

Ecological stochasticity and phage induction diversify bacterioplankton communities at the microscale

Rachel E. Szabo^{a,b}, Sammy Pontrelli^c, Jacopo Grilli^d, Julia A. Schwartzman^b, Shaul Pollak^b,
Uwe Sauer^c, Otto X. Cordero^{b,*}

^aMicrobiology Graduate Program, Massachusetts Institute of Technology, Cambridge, MA 02139, USA.

^bDepartment of Civil and Environmental Engineering, Massachusetts Institute of Technology, Cambridge, MA 02139, USA.

^cInstitute of Molecular Systems Biology, ETH Zürich, 8093 Zürich, Switzerland.

^dQuantitative Life Sciences, The Abdus Salam International Centre for Theoretical Physics, 34151 Trieste, Italy.

*To whom correspondence may be addressed. Email: ottox@mit.edu

Author Contributions: RE Szabo designed and performed experiments; analyzed and interpreted data; and wrote the manuscript. S Pontrelli performed metabolomics and the related analysis. J Grilli developed the population dynamics model and interpreted data. JA Schwartzman designed and performed the multi-particle incubation experiment and analyzed particle images. S Pollak developed the analysis pipeline for classifying MAGs into ecological roles. OX Cordero designed experiments, interpreted data, and wrote the manuscript. All authors edited and reviewed the manuscript.

Competing Interest Statement: The authors declare no competing interest.

This PDF file includes:

- Main Text
- Figures 1 to 5
- Main References
- Supplementary Text
- Figures S1 to S19
- Legends for Tables S1 to S7
- Supplementary References

Separate supplementary materials for this manuscript include:

- Tables S1 to S7

1 **Abstract**

2

3 In many natural environments, microorganisms self-assemble around heterogeneously distributed
4 resource patches. The growth and collapse of populations on resource patches can unfold within
5 spatial ranges of a few hundred micrometers or less, making such microscale ecosystems hotspots
6 of biological interactions and nutrient fluxes. Despite the potential importance of patch-level
7 dynamics for the large-scale evolution and function of microbial communities, we have not yet been
8 able to delineate the ecological processes that control natural populations at the microscale. Here,
9 we addressed this challenge in the context of microbially-mediated degradation of particulate
10 organic matter by characterizing the natural marine communities that assembled on over one
11 thousand individual microscale chitin particles. Through shotgun metagenomics, we found
12 significant variation in microscale community composition despite the similarity in initial species
13 pools across replicates. Strikingly, a subset of particles was highly populated by rare chitin-
14 degrading strains; we hypothesized that their conditional success reflected the impact of stochastic
15 colonization and growth on community assembly. In contrast to the conserved functional structures
16 that emerge in ecosystems at larger scales, this taxonomic variability translated to a wide range of
17 predicted chitinolytic abilities and growth returns at the level of individual particles. We found that
18 predation by temperate bacteriophages, especially of degrader strains, was a significant contributor
19 to the variability in the bacterial compositions and yields observed across communities. Our study
20 suggests that initial stochasticity in assembly states at the microscale, amplified through biotic
21 interactions, may have significant consequences for the diversity and functionality of microbial
22 communities at larger scales.

23

24

25 **Significance Statement**

26

27 The biogeochemical consequences of the degradation of particulate organic matter by
28 microorganisms represent the cumulative effect of microbial activity on individual microscale

29 resource patches. The ecological processes controlling community dynamics in these highly
30 localized microenvironments remain poorly understood. Here, we find that complex marine
31 communities growing on microscale resource particles diverge both taxonomically and functionally
32 despite assembling under identical abiotic conditions from a common species pool. We show that
33 this variability stems from bacteriophage predation and history-dependent factors in community
34 assembly, which create stochastic dynamics that are spatially structured at the microscale. This
35 microscale stochasticity may have significant consequences for the coexistence, evolution, and
36 function of diverse bacterial and viral populations in the global ocean.

37

38

39 **Main Text**

40

41 **Introduction**

42

43 A central challenge in microbial ecology is to connect the microscale world experienced by
44 microbial cells to observations of large-scale community functions (1, 2). In many environments –
45 ranging from soils (3) and sediments (4) to bioreactors (5) and hosts (6) – microbes live not in
46 homogeneous, well-mixed cultures, but rather in diverse, spatially-structured assemblages,
47 attached to surfaces and other cells in nutrient-dense patches on the order of 100 μ m in size.
48 Patches often exist in otherwise nutrient-limiting environments, creating hotspots of ecological
49 interactions and nutrient fluxes (7, 8).

50 A well-known example of micron-scale ecological hotspots is marine particulate organic
51 matter (POM), which is degraded by complex communities of bacteria, archaea, viruses, and
52 eukaryotes (9) with global biogeochemical consequences (10) (Fig. 1a). These interacting
53 community members can be broadly classified as primary degraders (that produce extracellular
54 enzymes to hydrolyze particle biopolymers), exploiters and scavengers (that are facilitated by
55 primary degraders) (11–14), and predators (such as bacteriophages (15) and grazers (16)).
56 Although these assemblages are often ephemeral, with organisms migrating through seawater

57 from patch to patch, cells can undergo multiple generations of growth while residing on a single
58 patch of nutrient-rich POM. This implies that the evolution and ecological functions of POM-
59 associated microbes are heavily influenced by their dynamics and interactions on microscale
60 particles. However, little is known about the processes governing community assembly at these
61 scales.

62 A major obstacle to understanding the factors that control populations at the patch level is
63 the difficulty of characterizing natural microscale communities with high replication. Microbial
64 communities are usually sampled at spatial scales orders of magnitude larger than those relevant
65 for microbial life (2), which homogenizes their inherent patchiness and results in inconsistent
66 inferences about ecological interactions (1, 17). Recent technological advances now permit the
67 sequencing of only thousands of cells (18, 19), presenting an opportunity to systematically
68 characterize microbial populations in units more closely approximating in scale the ecological
69 contexts experienced by microbes.

70 Here, we leveraged high-replicate sequencing of individual microscale communities to
71 evaluate the outcomes of assembly processes without the confounding effects of standard
72 sampling procedures. We employed a hybrid natural-laboratory approach that paired the
73 complexity of environmental microbial species pools with the controllability of synthetic
74 microparticles as discrete resource patches (11, 12). We immersed 1222 individual hydrogel
75 particles (85.0 ± 24.0 μm in diameter) made of chitin – a highly abundant biopolymer in marine POM
76 (9) – in samples of seawater containing microbes in their native states, which were then enriched
77 on particle surfaces. By incubating single particles separately under identical abiotic conditions,
78 each one became a microenvironment harboring a replicate community assembled from initially
79 similar species pools. We performed a comparative analysis across these microscale ecosystems
80 to investigate the natural variability in community composition and function among particles and to
81 identify biological processes that contribute to particle-level variability.

82

83

84 **Results**

85

86 **Bacterial community composition varies significantly across individual particles.**

87 To quantify the variation in community states across replicate microscale ecosystems, we
88 separately incubated single chitin particles in coastal seawater sampled from a common reservoir
89 (Fig. 1b; Methods). Assembly outcomes were assessed by removing particles from the seawater
90 at 13 time points over the course of 167 hours, a duration that aligns with previous measurements
91 (12, 20) of particle lifetimes (Methods). Shotgun metagenomic sequencing of individual particle-
92 attached communities was used to construct metagenome-assembled genomes (MAGs), which
93 were annotated to infer strains' potential ecological roles in a chitin-degrading community as
94 primary degraders, chitooligosaccharide exploiters, or metabolic byproduct scavengers (Fig. 1b;
95 Table S1; Methods). These MAGs served as the references for characterizing the community
96 composition of each particle.

97 We found a remarkable degree of compositional variability across individual particle
98 communities at the end of the time course ($n = 149$, after 154-167 hours of incubation). The
99 distributions of taxon (MAG) relative abundances across these late-stage particles spanned more
100 than three orders of magnitude (Fig. 2; Fig. S1a) and were approximately lognormal with a skew
101 towards high frequencies (Fig. S2). As a result, the community states observed at the single-particle
102 level diverged so significantly that the relative success of taxa across particles was poorly explained
103 by their average abundances (Fig. S1b; SI Methods). To assess whether non-ecological factors,
104 such as sampling bias in initial species pools, could have contributed to this compositional
105 divergence, we compared the variability in communities across late-stage particles to that across
106 unincubated aliquots of the seawater used as the inoculum (SI Methods). Inter-sample variation
107 was significantly higher across particles than across seawater samples (Mann-Whitney U test on
108 Aitchison distances: $p = 1.3 \times 10^{-13}$; Fig. S3), indicating that the observed variability stemmed more
109 from the community assembly process than from differences across inocula. Because other

110 technical sources of noise (Methods) also did not significantly impact the measured particle
111 compositions (Fig. S4) and all particles were chemically identical, we concluded that the variation
112 in taxon relative abundances across particles was due to biological and ecological factors that
113 amplified stochasticity in the initial assembly states of these communities.

114 The skew towards high frequencies in the relative abundance distributions implied that taxa
115 that were rare on average became dominant on a small number of particles (Fig. 2). As a result,
116 those particles harbored low-complexity communities (Fig. S5a-b) that diverged highly from the
117 average particle taxonomic composition (Fig. S5c). We termed the species that displayed this
118 phenomenon “jackpot taxa” for their simultaneous local success and global rarity (Methods). The
119 strains in this phylogenetically broad group of organisms included members of the
120 *Enterobacterales*, *Cytophagales*, *Pseudomonadales*, *Flavobacteriales*, *Rhodobacterales*,
121 *Fibrobacterales*, and *Chitinophagales* orders (Fig. S1, Table S1) and were mostly (87.9%)
122 classified as chitin degraders. Jackpot taxa were more prevalent across late-stage particles than
123 other taxa that were equally rare across inocula (Mann-Whitney U test: $p = 7.1 \times 10^{-3}$; Fig. S5d),
124 indicating that the probability of their success on particles, while influenced by their scarcity in
125 seawater, was also determined by ecological factors during community assembly. Notably, while
126 taxon-specific interactions did not explain the abundance patterns observed across particles (Fig.
127 S6; SI Methods), the most variable strains were likely to be degraders enriched in genes encoding
128 chitinases (Fig. 2; Fig. S7; coefficient of variation vs. chitinase copy number, Spearman's $\rho = 0.44$,
129 $p = 8.5 \times 10^{-7}$). These observations indicated that the conditional success of specialized degraders
130 from a diverse initial species pool contributed to the differentiation of the many rare community
131 states found at the single-particle level.

132

133 **Taxonomic variability translates to divergent community-level productivity.**

134 In contrast to the functionally similar gene content profiles predicted when microbial ecosystems
135 are characterized at the macroscale (21, 22), we found that the communities formed on particles
136 in individual microscale ecosystems were highly functionally divergent (Fig. 3a; Fig. S8). By the
137 end of the time course, most particles (63.8%) – and especially particles dominated by jackpot taxa

138 – harbored majority-degrader communities (Fig. 3a), highlighting the importance of degraders for
139 establishing and maintaining chitin-associated communities. However, the percentage of putative
140 degraders on each late-stage particle was as low as 13.1% and as high as 97.3%, indicating that
141 chitin degrading communities did not self-assemble to “optimized” or conserved ratios of ecological
142 roles after a fixed incubation period. Read mapping to chitinase protein sequences rather than
143 MAGs supported our interpretation that variability in the estimated proportion of degraders was not
144 due to the use of MAGs as reference genomes (Fig. S9; SI Methods). We hypothesized that this
145 extensive variability in community composition, primed by stochasticity in assembly processes,
146 could have had significant consequences for overall community function.

147 Consistent with this hypothesis, individual particles sustained highly variable particle-
148 attached biomass levels that were correlated with their community compositions (Fig. S10). The
149 number of bacterial cells in each late-stage community, estimated using qPCR of the 16S rRNA
150 region (Methods), ranged from approximately 1,000 to nearly 200,000 cells (Fig. 3b) and was
151 strongly correlated with the overall frequency of degraders (Spearman’s $\rho = 0.45$, $p = 1.6 \times 10^{-8}$).
152 Accordingly, particles that displayed the jackpot phenomenon had significantly higher cell counts
153 (Fig. 3b; Mann-Whitney U test: $p = 2.3 \times 10^{-7}$), revealing that jackpot taxa were dominant not only in
154 terms of relative abundances but also absolute abundances. The distribution of cells per particle
155 was approximately lognormal with a skew towards low cell numbers, indicating that some particles
156 were highly productive while others harbored small populations even by the end of the incubation,
157 as corroborated by visualizing particle-attached cells using a DNA stain (Fig. 3c; Methods).
158 Importantly, the initial colonization of single particles incubated together in the same volume of
159 seawater, rather than individually, resulted in particle-associated cell biomass that also spanned
160 several orders of magnitude (Fig. S11; SI Methods). This variability in initial particle colonization
161 was observed across a range of particle densities 15-140 times more concentrated than the
162 conditions of the individual particle incubations, indicating that phenomena such as jackpot
163 colonization are not specific to the environmental regime established in our separate microscale
164 ecosystems. Collectively, these results suggested that a strain’s growth was highly influenced by

165 its assembly context, raising the question of which biological or ecological factors could explain the
166 large variance in species compositions and consequent yields across replicate particles.

167

168 **Predation by bacteriophages contributes to variability in community composition and yield.**

169 Our observation that most (63.7%) of our MAGs contained sequences homologous to those of
170 bacteriophages led us to investigate whether these entities impacted the abundances of bacteria
171 on single particles. Bacteriophages (or phages, i.e. viruses that infect bacteria) are ubiquitous and
172 abundant in marine ecosystems, making predation by phages one of the primary forms of top-down
173 control of bacterial populations (23). High viral densities have been measured on marine particles
174 relative to ambient seawater (15), but it is unknown to what extent this represents passive
175 adsorption as opposed to active proliferation with impacts on bacterial growth in a natural, particle-
176 associated context. Therefore, we sought to identify populations of actively replicating phages
177 within the single particle communities to determine if heterogeneous phage predation could explain
178 the variability in community composition and yield.

179 To detect replicating phages, we first classified contigs in our metagenomic dataset as
180 phage-derived or bacteria-derived using tools (24, 25) that annotate phages from mixed
181 metagenomes (Methods). We reasoned that contigs classified as phage-derived, especially those
182 belonging to the genomes of temperate phages, were likely to be binned into the MAGs of their
183 bacterial hosts. Phage *k*-mer signatures tend to be more similar to those of their specific hosts than
184 to those of random bacteria (25, 26), and phages in a lysogenic cycle will have the same
185 sequencing read coverage patterns as their hosts across samples. Therefore, phages that were
186 lysogenic in most single-particle communities would tend to be binned with their hosts and have
187 similar coverage levels, reflected in an inferred virus-to-microbial cell ratio (VMR) close to 1 (Fig.
188 4a, top left). In contrast, phages in a productive cycle (lytic or chronic) would have higher coverage
189 than their hosts because of the multiple virion copies produced per bacterial cell (27, 28) (Fig. 4a,
190 top right). Therefore, we considered a phage-derived contig to be productive in a sample if it was
191 one of the most highly covered elements of its MAG (Methods).

192 Through this pipeline, we identified 256 phage contigs with coverage patterns consistent
193 with lysogenic infections in all samples and 263 phage contigs with coverage patterns consistent
194 with productive infections in a subset of samples (Table S2). Because our approach relied on
195 comparisons between co-binned phages and MAGs, phages that exclusively employ a lytic cycle
196 were unlikely to be detected. The VMRs of three representative examples of lysogenic and
197 productive phage contigs are shown for each particle in Figure 4b. Comparing the coverage
198 patterns of phage- and bacteria-derived contigs provided evidence that variable phage coverage
199 was not due to sequencing noise, lending confidence to our estimates of VMRs for specific phages
200 (Fig. S12; SI Methods). Using the VMRs of individual productive phages, we calculated the total
201 productive VMR per particle as a measure of overall phage replication in each community
202 (Methods).

203 The total productive VMRs of particle-associated communities sharply increased during
204 the early stages of particle incubation in seawater (Fig. 4c), consistent with the phenomenon of
205 rapid bacterial growth and high host densities driving the lysogeny-lysis switch in some phages
206 (29–33). The mean productive VMR was lowest for the initial seawater inocula and rose sharply
207 until the middle of the incubation period (59 hours), suggesting that phages became induced as
208 their particle-associated hosts began to grow. Concomitant with this increase in productive VMRs,
209 we observed the accumulation of metabolites in the seawater surrounding each particle until 59
210 hours of incubation, followed by a decrease in metabolite concentrations (Fig. 4d; Fig. S13; Table
211 S3; Methods). These observations could be explained by metabolite release upon the initiation of
212 bacterial growth (34) or lysis by phages (35) and by subsequent metabolite consumption by the
213 remaining viable bacteria (36). The coinciding timescales of metabolite liberation and rising VMRs
214 are consistent with our hypothesis that a particle-associated lifestyle among bacteria promoted
215 phage proliferation; therefore, we sought to assess the impact of variable phage induction on each
216 community's composition and consequent yield.

217 There was a striking negative relationship between cell counts and productive VMRs on
218 late-stage particles (Fig. 4e main, red data: Spearman's $\rho = -0.56$, $p = 3.3 \times 10^{-13}$), suggesting that
219 phage predation impacted bacterial growth success on particles upon induction. The degrader

220 populations contributed the most to this signal, indicating that strains among this trophic level may
221 have been especially prone to phage activation (Fig. 4e inset; Fig. S14). Importantly, jackpot
222 degrader taxa had lower productive VMRs than non-jackpot degraders (Fig. S15a; Mann-Whitney
223 U test: $p = 1.3 \times 10^{-49}$). This translated to jackpot particles having significantly lower productive VMRs
224 than non-jackpot particles (Mann-Whitney U test: $p = 4.3 \times 10^{-8}$), even controlling for differences in
225 cell counts between these groups of communities (ANCOVA: $F(1, 139) = 16.92$, $p = 4.1 \times 10^{-4}$, partial
226 $\eta^2 = 0.09$; Fig. S15b). Therefore, jackpot degraders may have been locally successful on a minority
227 of particles in part because they experienced less predation, supporting the hypothesis that top-
228 down population control by phages contributed to the large variability in the bacterial compositions
229 and thus yields observed across communities.

230 While cell counts were significantly correlated with both phage abundances and community
231 compositions, these features explained, respectively, 23% (Fig. 4e) and 34% (Fig. S10) of the
232 observed variation in yields, indicating that other factors also contributed to variable growth returns.
233 Therefore, we sought a more general framework in which to understand the key quantitative
234 features of the data – namely, the lognormal-like distributions of relative taxon abundances (with
235 right skews consistent with jackpot taxa) and of absolute cell abundances (with a left skew
236 corresponding to low-biomass communities). Incorporating (i) stochastic cell arrival on particles, (ii)
237 degraders as population founders, and (iii) noisy growth rates into a simple mathematical model of
238 community development on single particles was sufficient to reproduce these features (Figs. S16-
239 S19; SI Text). Taken together with our experimental data, this model indicates that the biological
240 processes which contribute to the stochasticity of particle colonization and growth rates – and
241 especially those processes that affect degraders – will result in variable growth returns for strains
242 across particles.

243

244

245 **Discussion**

246

247 While there is an abundance of evidence showing that the marine environment as
248 experienced by microbial cells is biologically, chemically, and physically heterogeneous (7),
249 characterizing the ecological processes controlling community assembly and development at these
250 scales remains a fundamental challenge, particularly *in situ*. Our study takes a step toward
251 addressing this problem using a hybrid natural-laboratory experiment that monitored the assembly
252 outcomes of complex marine communities across hundreds of individual chitin-based resource
253 particles. In accordance with prior work demonstrating small-scale heterogeneity on aquatic
254 resource particles (37), we found that bacterial compositions and absolute abundances varied to
255 such an extent across replicate particles that key community features – namely, species
256 composition and functional potential – were not conserved. Our results contrast with those of
257 previous studies (11, 13) that describe rapid ecological successions within particle systems that
258 are reproducible across batches. Despite this apparent reproducibility, biomass distributions in our
259 single-particle and our multi-particle incubations suggest that particle colonization is likely
260 heterogeneous in both systems. Thus, the reproducible dynamics previously observed in particle
261 systems could reflect the increasingly homogenizing effect of exchange between particles over
262 time or the mean of a process that is highly variable on the individual-particle scale. Future work is
263 required to determine the effect of dispersal and “cross-colonization” on the dynamics of particle
264 systems.

265 Stochastic factors are anticipated to strongly influence community assembly for
266 populations that are localized to small scales (38), such as in the microscale ecosystems on
267 resource particles (12, 20). The first step in community assembly – the arrival of cells to a particle
268 – is an intrinsically random process dependent on encounter probabilities. Our population dynamics
269 model demonstrated that historical contingencies (created by stochastic arrival times and the
270 growth dependency of non-degraders on degraders) magnified through noisy growth rates were
271 sufficient to reproduce the distributions of bacterial abundances observed across individual
272 particles. Because this chitin microparticle ecosystem is subject to conditions that have been shown
273 to promote strong priority effects (e.g. a large regional species pool, rapid local growth dynamics,
274 high resource overlap, and a dependence of late-arriving organisms on early-arriving ones) (38),

275 we hypothesize that biotic factors amplified this initial stochasticity in each assembly context and
276 influenced subsequent community development.

277 One key biological contributor to noisy growth returns may have been variable predation
278 by temperate bacteriophages. Phages became increasingly and differentially activated during
279 community development on particles, with elevated virus-to-microbial cell ratios (VMRs) in low-
280 biomass communities implicating phage-mediated lysis as one factor explaining the biomass
281 variability on late-stage particles. These results align with those of previous studies documenting
282 extensive variation in VMRs at small spatial scales (23) and an inverse relationship between VMRs
283 and cell densities in marine environments (39). Because phage induction was significantly less
284 associated with jackpot degrader strains, we hypothesize that the jackpot phenomenon –
285 characterized by globally rare yet locally productive degraders – was partially a reflection of lower
286 levels of phage-driven population collapse in those community contexts. Therefore, top-down
287 control by phages may link the highly variable community compositions and yields observed among
288 particles.

289 A synthesis of our mathematical model with our observations of bacterial and phage
290 abundances suggests a conceptual framework for key processes promoting variability in
291 microscale community composition and function (Fig. 5). We posit that stochastic arrival on
292 particles diversifies initial assembly states; that the timescale and magnitude of degrader
293 colonization determine the extent to which scavengers and exploiters are supported; and that
294 phage induction and subsequent host lysis, primarily among degraders, contribute to noisy growth
295 returns. Therefore, in this conceptual framework, the high-biomass jackpot particles are those in
296 which degraders arrive early and resist phage induction, leading to high relative and absolute
297 degrader abundances (Fig. 5, top). By contrast, low-biomass particles are those in which degraders
298 are not able to proliferate, either because phage induction leads to their population collapse (Fig.
299 5, middle) or because they become established on a particle relatively late (Fig. 5, bottom).

300 In addition to the growing body of evidence that marine aggregates can stimulate the
301 production of virulent phages (i.e. phages that exclusively engage in lytic cycles) (40), our study
302 suggests that resource particles may be replication hotspots also for temperate phages (i.e. those

303 that conditionally employ both lytic and lysogenic cycles). In marine environments, lysogeny is
304 promoted under conditions that limit bacterial growth while the lytic cycle is favored during periods
305 of high bacterial activity (29, 41, 42), indicating that rapid host growth and abundance can regulate
306 the lysogeny-lysis switch in some temperate phages (30–33). Therefore, in a patchy nutrient
307 landscape, temperate phages may employ lysogeny as a survival strategy when their bacterial
308 hosts are at low densities and are foraging for nutrients, hitchhiking with their hosts onto resource
309 particles. Robust bacterial growth on particles may induce prophages at a time when abundant
310 host resources can be co-opted and many susceptible cells are nearby, resulting in lytic
311 suppression of the bacterial population and the release of virions into the surrounding seawater.
312 Factors such as the variable presence of prophages in the flexible genomes of strains growing on
313 different particles (43), the co-occurrence of bacterial competitors that trigger induction (44, 45),
314 and phenotypic heterogeneity resulting in differential induction (46, 47) may all contribute to the
315 varying levels of phage activation observed on individual particles in our microscale ecosystems.
316 Further research is required on the mechanisms underlying prophage induction in complex
317 communities in order to understand how lysogeny and lysis on particle hotspots shape the
318 dynamics of marine microbial communities.

319 Our observations of wild marine communities, though made in a laboratory setting, may
320 provide insights on the ecosystem-level consequences of microscale stochastic assembly
321 dynamics. First, the stochasticity in bacterial growth, amplified through spatial structuring at the
322 microscale, may promote the maintenance of a diverse regional species pool. This is because the
323 variability in growth returns can effectively offset differences in relative fitness between competing
324 strains or species (48). Second, the variability in microscale community states could be reflected
325 in larger-scale biogeochemical patterns since the cumulative process of POM degradation can be
326 approximated as the sum of degradation events on individual particles. We found that late-stage
327 communities did not converge to a fixed proportion of chitin degraders or to a fixed amount of
328 biomass per particle; both measures are positively correlated with the rate of particle degradation
329 (12), suggesting that historical contingencies in community assembly promote functional
330 divergence (38, 49). These results contrast with those of previous studies on the replicability of

331 microbial community assembly at the functional level (21, 22) likely because of the homogenizing
332 effect of macroscale sampling. Intriguingly, the lognormal-like distribution of biomass on individual
333 particles aligns with observations and predictions of lognormally-distributed global marine organic
334 matter export and remineralization rates; these distributions may repeatedly emerge as a reflection
335 of the multiplicative effects of stochastic variables in ecological settings (50–52). Although our
336 experimental system significantly simplified the process of POM degradation in the ocean, our
337 approach provides a quantitative link between the microscale and larger-scale processes,
338 highlighting the importance of considering local variability when investigating the mechanisms
339 behind microbial community development in a spatially structured environment.

340

341

342 **Materials and Methods**

343

344 Abridged Methods are provided below; details and additional information are provided in SI
345 Methods.

346

347 **Seawater collection and individual chitin particle incubation.** Nearshore coastal seawater was
348 collected from Nahant, MA; filtered (63 μ m) to remove large particulate matter; gently concentrated
349 via centrifugation at 4000 \times *g* for 5 minutes; and aliquoted for incubations and sequencing. Chitin
350 magnetic particles (New England Biolabs, #E8036L) were washed in sterile artificial seawater
351 (Sigma-Aldrich, #S9883) and individually selected beneath a dissecting microscope in a laminar
352 flow hood. Single chitin particles (85.0 \pm 24.0 μ m in diameter) were transferred to sterile 96-well
353 plates (Thermo Fisher, #AB0600L), with one chitin particle per well. Plates were inoculated
354 consecutively with 175 μ L of filtered, centrifuged seawater per well; sealed (VWR, #89092-056);
355 and rotated end-over-end (7.5rpm) at room temperature. The particles in an entire plate were
356 harvested at each time point (after 12, 22.75, 34.5, 46, 59, 69, 82, 92, 103, 116.75, 113, 153.5, and
357 166.5 hours of incubation) by inspection and pipetting under a dissecting microscope in a laminar
358 flow hood. Each particle was transferred into sterile 96-well plates (Thermo Fisher, #AB0600L)

359 containing TE buffer and stored at -20°C. The seawater surrounding each harvested particle was
360 also saved in 96-well plates and stored at -20°C.

361

362 **Mock communities and negative controls.** To quantify the technical error associated with
363 creating metagenomic libraries from low DNA inputs, mock communities were simulated by
364 combining the DNA of two strains previously isolated from a chitin particle enrichment (11).
365 Libraries from three technical replicates of mock communities totaling 50pg or 5pg of DNA (SI
366 Methods), as well as from six negative controls (containing only nuclease-free water), were
367 prepared and analyzed with the same protocols used for individual chitin particle-attached
368 communities.

369

370 **DNA extraction and metagenomic sequencing.** DNA extractions were performed for twelve
371 175µL-volume aliquots of the initial, unincubated seawater and for particles harvested after 34.5,
372 59, 103, 116.75, 113, 153.5, and 166.5 hours of incubation (see Table S5 for sample metadata).
373 DNA was extracted from all samples with the Agencourt DNAAdvance Genomic DNA Isolation Kit
374 (Beckman Coulter; modifications noted in SI Methods). Metagenomic libraries were prepared with
375 the Nextera XT DNA Library Prep Kit and index primers (Illumina) using the protocol developed by
376 Rinke *et al.* (18) for low DNA inputs (SI Methods). Libraries were quantified on an Agilent 4200
377 TapeStation system with High Sensitivity D5000 ScreenTapes (Agilent Technologies) and pooled
378 by time point in equimolar amounts. Sequencing was performed on an Illumina HiSeq 2500
379 machine (250bp paired-end reads) at the Whitehead Institute for Biomedical Research (Cambridge,
380 MA).

381

382 **Metagenome-assembled genome (MAG) generation, taxonomic assignment, and role**
383 **classification.** Raw sequencing reads were quality trimmed with Trimmomatic v0.36 (53). Reads
384 mapping to the PhiX and human genomes were filtered out using BBDuk v38.16 (54) and BBDuk
385 v38.16, respectively (SI Methods). Trimmed, filtered reads that were error-corrected using
386 BayesHammer (55) were pooled within each time point and co-assembled using MEGAHIT v1.2.9

387 (56). Bins were generated with MaxBin v2.2.7 (57) and CONCOCT v1.1.0 (58); consolidated and
388 filtered using DAS Tool v1.1.1 (59); and evaluated for completeness and contamination using
389 CheckM v1.1.2 (60). The resulting 251 bins were used as reference MAGs ($\geq 50\%$ complete, $\leq 10\%$
390 contaminated; median completeness 93.7%, median contamination 3.9%; Table S1). Highly similar
391 MAGs obtained from separate co-assemblies were grouped into 132 clusters (SI Methods). MAG
392 taxonomic classifications were made using GTDB-Tk v1.1.1 (61). MAGs were functionally
393 annotated using a custom database of profile hidden Markov models (HMMs) of proteins involved
394 in growth on chitin (SI Methods; Table S6). Ecological roles for MAGs (as degraders,
395 chitooligosaccharide exploiters, or metabolic byproduct scavengers) were defined based on the
396 gene content patterns observed for sequenced and phenotyped (14) strains previously isolated
397 (11, 13) from particle enrichments (SI Methods).

398

399 **Read mapping to MAGs for relative abundance estimation.** Trimmed, filtered reads were
400 mapped competitively against the MAGs generated from sequencing particle-attached
401 communities, initial seawater samples, and negative controls. Read mapping was performed using
402 the approach described in Leventhal *et al.* (62) (SI Methods). Reads that best mapped to predicted
403 contaminant MAGs (SI Methods) were removed from consideration. MAG relative abundances
404 were calculated for each sample by (1) tallying the hits to all MAGs in each MAG cluster; (2)
405 normalizing the tally by the average genome length of all MAGs in each MAG cluster; and (3)
406 dividing the normalized tallies for each MAG cluster by their sum for each sample. Therefore, for
407 MAGs clustered together based on similarity, their relative abundances are represented in that of
408 the entire MAG cluster to which they belong; this calculation circumvents the artificial
409 underestimation of MAG relative abundances that would otherwise be obtained with a non-
410 dereplicated reference set. The relative abundances of organisms occupying the three ecological
411 roles (degrader, exploiter, scavenger) on each particle were calculated by summing the relative
412 abundances of MAGs classified into each role.

413

414 **Definitions of jackpot MAGs and jackpot particles.** A jackpot score was calculated for each
415 MAG cluster to quantitatively reflect the properties of rarity across most particles and dominance
416 on a few particles (SI Methods) such that MAGs with high scores strongly displayed the jackpot
417 phenomenon. Each particle's jackpot score was calculated as the weighted average of MAG
418 jackpot scores (i.e. the sum of the relative abundance of each MAG cluster multiplied by its jackpot
419 score). Particles with high jackpot scores and low Pielou's evenness were categorized as "jackpot
420 particles" (SI Methods).

421

422 **Bacteriophage analyses.** Binned contigs were classified as phage-derived or bacteria-derived
423 using VirSorter v1.0.3 with its RefSeqABVir database (24) and VirFinder v1.1 (25), two tools
424 designed to detect phage sequences among mixed metagenomes (SI Methods). We used a read
425 coverage-based approach to categorize phage-derived contigs as productive or lysogenic in
426 particle-attached communities (Table S2; see SI Methods for analysis controls). Based on read
427 mapping to MAGs, per-base coverage values for all binned contigs were computed with BEDTools
428 v2.27.0 (63) and were used to calculate contig-wide average coverage values. For each MAG and
429 for each sample, a phage-derived contig was considered to be productive if its coverage was
430 greater than the coverage of the 95th percentile bacteria-derived contig in the same MAG;
431 otherwise, it was considered to be lysogenic in that sample. The VMR of an individual phage contig
432 in one sample is defined as the phage contig coverage divided by average coverage of the MAG
433 with which it is binned (which was calculated using only the bacteria-derived contigs). Total VMRs
434 – i.e. the total number of phage copies relative to the total number of bacterial MAG copies in an
435 entire sample – were calculated separately for productive and lysogenic phage contigs. The total
436 productive VMR for a sample was defined as:

437

$$438 \quad \sum_i^n \left[\left(\frac{\text{average coverage of productive phage contigs in } MAG_i}{\text{average } MAG_i \text{ coverage}} \right) \times (MAG_i \text{ relative abundance}) \right]$$
$$439 \quad = \frac{\text{total \# phage copies (due to productive infections)}}{\text{total \# bacterial genome copies}}$$

440

441 where n is the number of MAGs found in a sample. This calculation is equivalent to

442

443
$$\frac{\sum_i^n (\text{average coverage of productive phage contigs in } MAG_i)}{\sum_i^n (\text{average } MAG_i \text{ coverage})}$$

444

445 where n is the number of MAGs found in a sample. Total lysogenic VMRs were calculated using
446 the same formula while considering only lysogenic-annotated contigs. VMRs for each ecological
447 role (i.e. for the subpopulation in a community that belongs to one of the three roles of degrader,
448 exploiter, or scavenger) were calculated using the same formula considering only the MAGs of
449 each role and their associated phages.

450

451 **Cell count estimation.** Bacterial DNA extracted from individual particle-attached communities was
452 quantified through qPCR of the 16S rRNA gene using the Femto Bacterial DNA Quantification Kit
453 (Zymo Research), which has a lower limit of detection of 20fg. Two sets of standards and negative
454 controls were included in each qPCR run. The number of bacterial cells for each particle was
455 estimated from the absolute DNA amounts based on measurements indicating a mean of 2.5fg
456 DNA per bacterial cell in seawater samples (64).

457

458 **Imaging of chitin particle colonization.** Subsets of chitin particles incubated individually in
459 seawater were stained at time points by adding the DNA stain SYTO9 (Invitrogen, #S34854) at a
460 final concentration of 500nM directly to the particle incubations. Particles were incubated in the
461 dark at room temperature for 15 minutes before being mounted separately on microscope slides
462 and imaged with a Zeiss epifluorescence microscope at 100X magnification.

463

464 **Metabolomics.** We performed untargeted metabolomics of the seawater that surrounded each
465 harvested chitin particle and of the initial, unincubated seawater (SI Methods). We used a binary
466 LC pump (Agilent Technologies) and an MPS2 Autosampler (Gerstel) coupled to an Agilent 6520

467 time-of-flight mass spectrometer (Agilent Technologies) operated in negative mode, at 2GHz,
468 extended dynamic range, with an m/z (mass/charge) range of 50-1000. Ions (Table S3) were
469 annotated against a curated library of metabolites present in marine microbes, based on the BioCyc
470 database (65). For metabolites that exceeded the limit of detection (SI Methods), the intensities of
471 each ion were normalized between 0 (the limit of detection) and 1 (the highest measured intensity
472 of a given ion). Weighted ion intensities for each timepoint were calculated by taking the sum of all
473 normalized intensities of ions in all samples for each timepoint.

474

475

476 **Data Sharing Plans**

477

478 All data will be made publicly available before publication. Sequencing data will be deposited to the
479 National Center for Biotechnology Information (NCBI) as a BioProject, with raw reads uploaded to
480 the Sequence Read Archive (SRA) and metagenome-assembled genomes (MAGs) uploaded to
481 the Whole Genome Shotgun (WGS) database. All mass spectra files from the metabolomics will
482 be accessible from MassIVE (<ftp://MSV000087936@massive.ucsd.edu>) before publication. MAG
483 relative abundances for each sample and metadata for samples, MAGs, phages, and detected
484 metabolites are provided as Supplementary Tables. All code and files used to generate figures will
485 be made available at personal GitHub pages before publication.

486

487

488 **Acknowledgments**

489

490 We extend our gratitude to all past and present members of the Cordero lab for their support and
491 critical feedback, as well as members of the Simons Collaboration on Principles of Microbial
492 Ecosystems for stimulating discussions. In particular, we would like to thank: Manoshi S. Datta, for
493 contributing to the genesis of this project; José T. Saavedra, for developing the DNA extraction and
494 metagenomic library preparation pipeline for single particle-attached communities; Gabriel E.

495 Leventhal and Jakob Russel for bioinformatic mentorship and assistance; Matti Gralka, for
496 feedback on analyses and this manuscript; Elise Ledieu, for quantifying chitin particle sizes;
497 Anthony Gaca, for advice on metagenomic library preparation; Fatima Aysha Hussain, for feedback
498 on the bacteriophage analysis; Emily Zakem, for insights on global particle remineralization rates;
499 Akshit Goyal, for comments on this manuscript; and Sara Szabo and William Mandella, for
500 assistance with seawater sampling. This material is based upon work supported by the National
501 Science Foundation Graduate Research Fellowship under Grant No. #174530. This project was
502 supported by the Simons Collaboration: Principles of Microbial Ecosystems (PriME) award number
503 542395. S Pontrelli was supported by a grant from the Simons Foundation (ID608247) as part of
504 PriME. S Pollak was supported by the EMBO ALTF Grant No. #800-2017.

505

506

507 **Main References**

508

- 509 1. O. X. Cordero, M. S. Datta, Microbial interactions and community assembly at microscales.
510 *Curr Opin Microbiol* 31, 227–234 (2016).
- 511 2. D. R. Nemergut, *et al.*, Patterns and Processes of Microbial Community Assembly.
512 *Microbiol Mol Biol R* 77, 342–356 (2013).
- 513 3. A. G. O'Donnell, I. M. Young, S. P. Rushton, M. D. Shirley, J. W. Crawford, Visualization,
514 modelling and prediction in soil microbiology. *Nat Rev Microbiol* 5, 689–699 (2007).
- 515 4. S. E. McGlynn, G. L. Chadwick, C. P. Kempes, V. J. Orphan, Single cell activity reveals
516 direct electron transfer in methanotrophic consortia. *Nature* 526, 531–535 (2015).
- 517 5. G. Gonzalez-Gil, C. Holliger, Aerobic granules: Microbial landscape and architecture,
518 stages, and practical implications. *Appl Environ Microb* 80, 3433–3441 (2014).
- 519 6. J. L. M. Welch, B. J. Rossetti, C. W. Rieken, F. E. Dewhirst, G. G. Borisy, Biogeography of
520 a human oral microbiome at the micron scale. *Proc National Acad Sci* 113, 791–800
521 (2016).
- 522 7. R. Stocker, Marine microbes see a sea of gradients. *Science* 338, 628–633 (2012).

- 523 8. L. M. Dann, *et al.*, Microbial micropatches within microbial hotspots. *Plos One* 13, 1–22
524 (2018).
- 525 9. H. Dang, C. R. Lovell, Microbial Surface Colonization and Biofilm Development in Marine
526 Environments. *Microbiol Mol Biol R* 80, 91–138 (2016).
- 527 10. N. Jiao, *et al.*, Microbial production of recalcitrant dissolved organic matter: Long-term
528 carbon storage in the global ocean. *Nat Rev Microbiol* 8, 593–599 (2010).
- 529 11. M. S. Datta, E. Sliwerska, J. Gore, M. F. Polz, O. X. Cordero, Microbial interactions lead to
530 rapid micro-scale successions on model marine particles. *Nat Commun* 7, 1–7 (2016).
- 531 12. T. N. Enke, G. E. Leventhal, M. Metzger, J. T. Saavedra, O. X. Cordero, Microscale ecology
532 regulates particulate organic matter turnover in model marine microbial communities. *Nat*
533 *Commun* 9, 2743 (2018).
- 534 13. T. N. Enke, *et al.*, Modular Assembly of Polysaccharide-Degrading Marine Microbial
535 Communities. *Curr Biol* 29, 1528-1535.e6 (2019).
- 536 14. S. Pontrelli, *et al.*, Hierarchical control of microbial community assembly by specialists.
537 bioRxiv [Preprint] (2021). <https://doi.org/10.1101/2021.06.22.449372> (accessed 24
538 September 2021).
- 539 15. M. G. Weinbauer, *et al.*, Viral ecology of organic and inorganic particles in aquatic systems:
540 Avenues for further research. *Aquat Microb Ecol* 57, 321–341 (2009).
- 541 16. M. Simon, H. P. Grossart, B. Schweitzer, H. Ploug, Microbial ecology of organic aggregates
542 in aquatic ecosystems. *Aquat Microb Ecol* 28, 175–211 (2002).
- 543 17. D. W. Armitage, S. E. Jones, How sample heterogeneity can obscure the signal of microbial
544 interactions. *Isme J* 13, 2639–2646 (2019).
- 545 18. C. Rinke, *et al.*, Validation of picogram- and femtogram-input DNA libraries for microscale
546 metagenomics. *Peerj* 2016, 1–28 (2016).
- 547 19. R. U. Sheth, *et al.*, Spatial metagenomic characterization of microbial biogeography in the
548 gut. *Nat Biotechnol* 37, 877–883 (2019).

- 549 20. M. H. Iversen, H. Ploug, Temperature effects on carbon-specific respiration rate and
550 sinking velocity of diatom aggregates – potential implications for deep ocean export
551 processes. *Biogeosciences* 10, 4073–4085 (2013).
- 552 21. S. Louca, *et al.*, High taxonomic variability despite stable functional structure across
553 microbial communities. *Nat Ecol Evol* 1, 1–12 (2016).
- 554 22. C. Huttenhower, *et al.*, Structure, function and diversity of the healthy human microbiome.
555 *Nature* 486, 207–214 (2012).
- 556 23. M. Breitbart, C. Bonnain, K. Malki, N. A. Sawaya, Phage puppet masters of the marine
557 microbial realm. *Nat Microbiol* 3, 754–766 (2018).
- 558 24. S. Roux, F. Enault, B. L. Hurwitz, M. B. Sullivan, VirSorter: Mining viral signal from microbial
559 genomic data. *PeerJ* 2015, 1–20 (2015).
- 560 25. J. Ren, N. A. Ahlgren, Y. Y. Lu, J. A. Fuhrman, F. Sun, VirFinder: a novel k-mer based tool
561 for identifying viral sequences from assembled metagenomic data. *Microbiome* 5, 69
562 (2017).
- 563 26. N. A. Ahlgren, J. Ren, Y. Y. Lu, J. A. Fuhrman, F. Sun, Alignment-free d2* oligonucleotide
564 frequency dissimilarity measure improves prediction of hosts from metagenomically-
565 derived viral sequences. *Nucleic Acids Res* 45, 39–53 (2017).
- 566 27. K. Kieft, K. Anantharaman, Deciphering active prophages from metagenomes. bioRxiv
567 [Preprint] (2021). <https://doi.org/10.1101/2021.01.29.428894> (accessed 24 September
568 2021).
- 569 28. R. F. von Boeselager, E. Pfeifer, J. Frunzke, Cytometry meets next-generation sequencing
570 – RNA-Seq of sorted subpopulations reveals regional replication and iron-triggered
571 prophage induction in *Corynebacterium glutamicum*. *Sci Rep-uk* 8, 1–13 (2018).
- 572 29. J. H. Paul, Prophages in marine bacteria: Dangerous molecular time bombs or the key to
573 survival in the seas? *Isme J* 2, 579–589 (2008).
- 574 30. J. R. Brum, B. L. Hurwitz, O. Schofield, H. W. Ducklow, M. B. Sullivan, Seasonal time
575 bombs: Dominant temperate viruses affect Southern Ocean microbial dynamics. *Isme J*
576 10, 437–449 (2016).

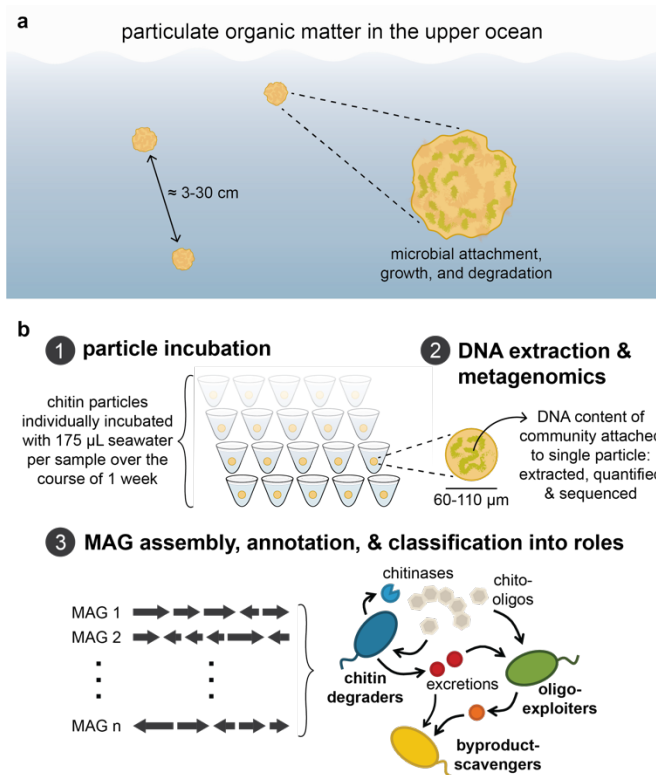
- 577 31. M. Touchon, A. Bernheim, E. P. C. Rocha, Genetic and life-history traits associated with
578 the distribution of prophages in bacteria. *Isme J* 10, 2744–2754 (2016).
- 579 32. J. E. Silpe, B. L. Bassler, A Host-Produced Quorum-Sensing Autoinducer Controls a Phage
580 Lysis-Lysogeny Decision. *Cell* 176, 268-280.e13 (2019).
- 581 33. L. Laganenka, *et al.*, Quorum sensing and metabolic state of the host control lysogeny-
582 lysis switch of bacteriophage T1. *Mbio* 10, 3–8 (2019).
- 583 34. B. E. Noriega-Ortega, *et al.*, Does the chemodiversity of bacterial exometabolomes sustain
584 the chemodiversity of marine dissolved organic matter? *Front Microbiol* 10, 1–13 (2019).
- 585 35. N. Y. D. Ankrah, *et al.*, Phage infection of an environmentally relevant marine bacterium
586 alters host metabolism and lysate composition. *Isme J* 8, 1089–1100 (2014).
- 587 36. S. Blasche, *et al.*, Metabolic cooperation and spatiotemporal niche partitioning in a kefir
588 microbial community. *Nat Microbiol* 6, 196–208 (2021).
- 589 37. M. Bizic-Ionescu, D. Ionescu, H.-P. Grossart, Organic particles: heterogeneous hubs for
590 microbial interactions in aquatic ecosystems. *Front Microbiol* Accepted, 1–15 (2018).
- 591 38. T. Fukami, Historical Contingency in Community Assembly: Integrating Niches, Species
592 Pools, and Priority Effects. *Annu Rev Ecol Evol Syst* 46, 1–23 (2015).
- 593 39. C. H. Wigington, *et al.*, Re-examination of the relationship between marine virus and
594 microbial cell abundances. *Nat Microbiol* 1, 4–11 (2016).
- 595 40. L. Riemann, H. P. Grossart, Elevated lytic phage production as a consequence of particle
596 colonization by a marine Flavobacterium (*Cellulophaga* sp.). *Microbial Ecol* 56, 505–512
597 (2008).
- 598 41. M. G. Weinbauer, I. Brettar, M. G. Höfle, Lysogeny and virus-induced mortality of
599 bacterioplankton in surface, deep, and anoxic marine waters. *Limnol Oceanogr* 48, 1457–
600 1465 (2003).
- 601 42. J. P. Payet, C. A. Suttle, To kill or not to kill: The balance between lytic and lysogenic viral
602 infection is driven by trophic status. *Limnol Oceanogr* 58, 465–474 (2013).
- 603 43. B. C. M. Ramisetty, P. A. Sudhakari, Bacterial “grounded” prophages: Hotspots for genetic
604 renovation and innovation. *Frontiers Genetics* 10, 1–17 (2019).

- 605 44. M. Jancheva, T. Böttcher, A Metabolite of Pseudomonas Triggers Prophage-Selective
606 Lysogenic to Lytic Conversion in Staphylococcus aureus. *J Am Chem Soc* 143, 8344–8351
607 (2021).
- 608 45. J. E. Silpe, J. W. H. Wong, S. V. Owen, M. Baym, E. P. Balskus, The gut bacterial natural
609 product colibactin triggers induction of latent viruses in diverse bacteria. bioRxiv [Preprint]
610 (2021). <https://doi.org/10.1101/2021.05.24.445430> (accessed 24 September 2021).
- 611 46. J. J. Dennehy, I. N. Wang, Factors influencing lysis time stochasticity in bacteriophage.
612 *BMC Microbiol* 11, 174 (2011).
- 613 47. L. Imamovic, E. Ballesté, A. Martínez-Castillo, C. García-Aljaro, M. Muniesa, Heterogeneity
614 in phage induction enables the survival of the lysogenic population. *Environ Microbiol* 18,
615 957–969 (2016).
- 616 48. A. Melbinger, M. Vergassola, The Impact of Environmental Fluctuations on Evolutionary
617 Fitness Functions. *Sci Rep* 5, 1–11 (2015).
- 618 49. L. S. Bittleston, M. Gralka, G. E. Leventhal, I. Mizrahi, O. X. Cordero, Context-dependent
619 dynamics lead to the assembly of functionally distinct microbial communities. *Nat Commun*
620 11, 1–10 (2020).
- 621 50. B. B. Cael, K. Bisson, C. L. Follett, Can Rates of Ocean Primary Production and Biological
622 Carbon Export Be Related Through Their Probability Distributions? *Global Biogeochem Cy*
623 32, 954–970 (2018).
- 624 51. E. J. Zakem, B. B. Cael, N. M. Levine, A unified theory for organic matter accumulation.
625 *Proc National Acad Sci* 118, e2016896118 (2021).
- 626 52. E. Limpert, W. A. Stahel, M. Abbt, Log-normal Distributions across the Sciences: Keys and
627 Clues. *BioScience* 51, 341–352 (2001).
- 628 53. A. M. Bolger, M. Lohse, B. Usadel, Trimmomatic: A flexible trimmer for Illumina sequence
629 data. *Bioinformatics* 30, 2114–2120 (2014).
- 630 54. B. Bushnell, BBMap: A Fast, Accurate, Splice-Aware Aligner (2014).
- 631 55. S. I. Nikolenko, A. I. Korobeynikov, M. A. Alekseyev, BayesHammer: Bayesian clustering
632 for error correction in single-cell sequencing. *Bmc Genomics* 14, S7 (2013).

- 633 56. D. Li, *et al.*, MEGAHIT v1.0: A fast and scalable metagenome assembler driven by
634 advanced methodologies and community practices. *Methods* 102, 3–11 (2016).
- 635 57. Y. W. Wu, B. A. Simmons, S. W. Singer, MaxBin 2.0: An automated binning algorithm to
636 recover genomes from multiple metagenomic datasets. *Bioinformatics* 32, 605–607 (2016).
- 637 58. J. Alneberg, *et al.*, Binning metagenomic contigs by coverage and composition. *Nat*
638 *Methods* 11, 1144–1146 (2014).
- 639 59. C. M. K. Sieber, *et al.*, Recovery of genomes from metagenomes via a dereplication,
640 aggregation and scoring strategy. *Nat Microbiol* 3, 836–843 (2018).
- 641 60. D. H. Parks, M. Imelfort, C. T. Skennerton, P. Hugenholtz, G. W. Tyson, CheckM:
642 Assessing the quality of microbial genomes recovered from isolates, single cells, and
643 metagenomes. *Genome Res* 25, 1043–1055 (2015).
- 644 61. P. A. Chaumeil, A. J. Mussig, P. Hugenholtz, D. H. Parks, GTDB-Tk: A toolkit to classify
645 genomes with the genome taxonomy database. *Bioinformatics* 36, 1925–1927 (2020).
- 646 62. G. E. Leventhal, *et al.*, Strain-level diversity drives alternative community types in
647 millimetre-scale granular biofilms. *Nat Microbiol* 3, 1295–1303 (2018).
- 648 63. A. R. Quinlan, I. M. Hall, BEDTools: A flexible suite of utilities for comparing genomic
649 features. *Bioinformatics* 26, 841–842 (2010).
- 650 64. D. K. Button, B. R. Robertson, Determination of DNA Content of Aquatic Bacteria by Flow
651 Cytometry. *Appl Environ Microb* 67, 1636–1645 (2001).
- 652 65. P. D. Karp, *et al.*, The BioCyc collection of microbial genomes and metabolic pathways.
653 *Brief Bioinform* 20, 1085–1093 (2017).

654 **Figures**

655



656

657 **Figure 1. Modeling particulate organic matter degradation with a laboratory system of**

658 **enriching of marine microbes on chitin particles. (a)** Microscale marine particles are spatially-

659 spatially-separated nutrient-rich habitats dynamically populated attached and degraded by complex communities of

660 heterotrophic bacteria. The interparticle distance range is estimated from data reported in Simon

661 *et al.* (16). **(b)** Schematic depicting experimental design and analysis. Microscale chitin particles

662 were individually incubated in seawater, and the DNA content of particle-attached communities

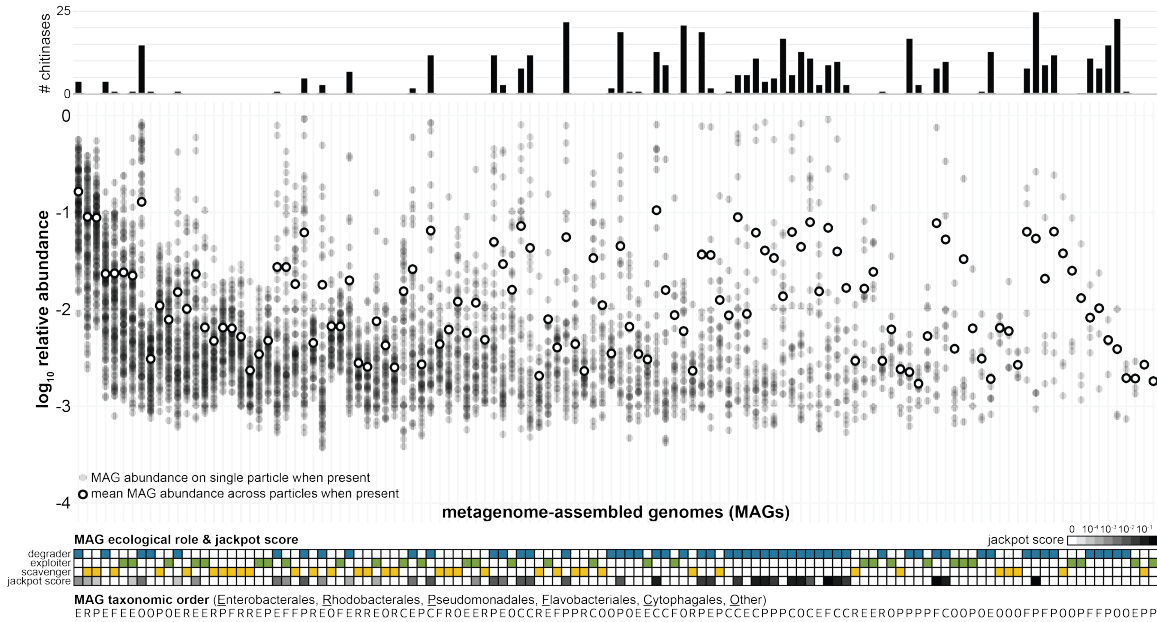
663 was quantified and submitted for shotgun metagenomic sequencing. Communities were

664 characterized using metagenome-assembled genomes (MAGs), which were classified into three

665 predicted ecological roles for this ecosystem: chitin degraders, chitooligosaccharide exploiters,

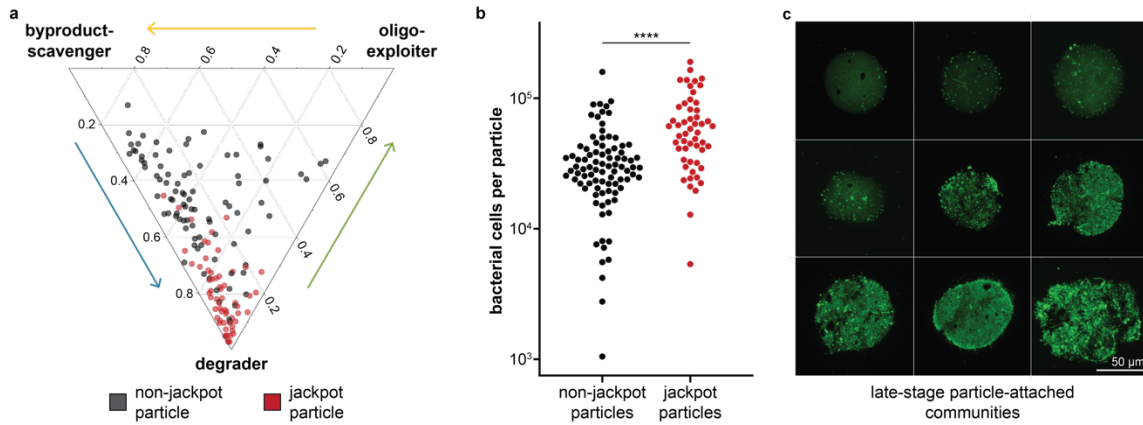
666 and metabolic byproduct scavengers.

667



668

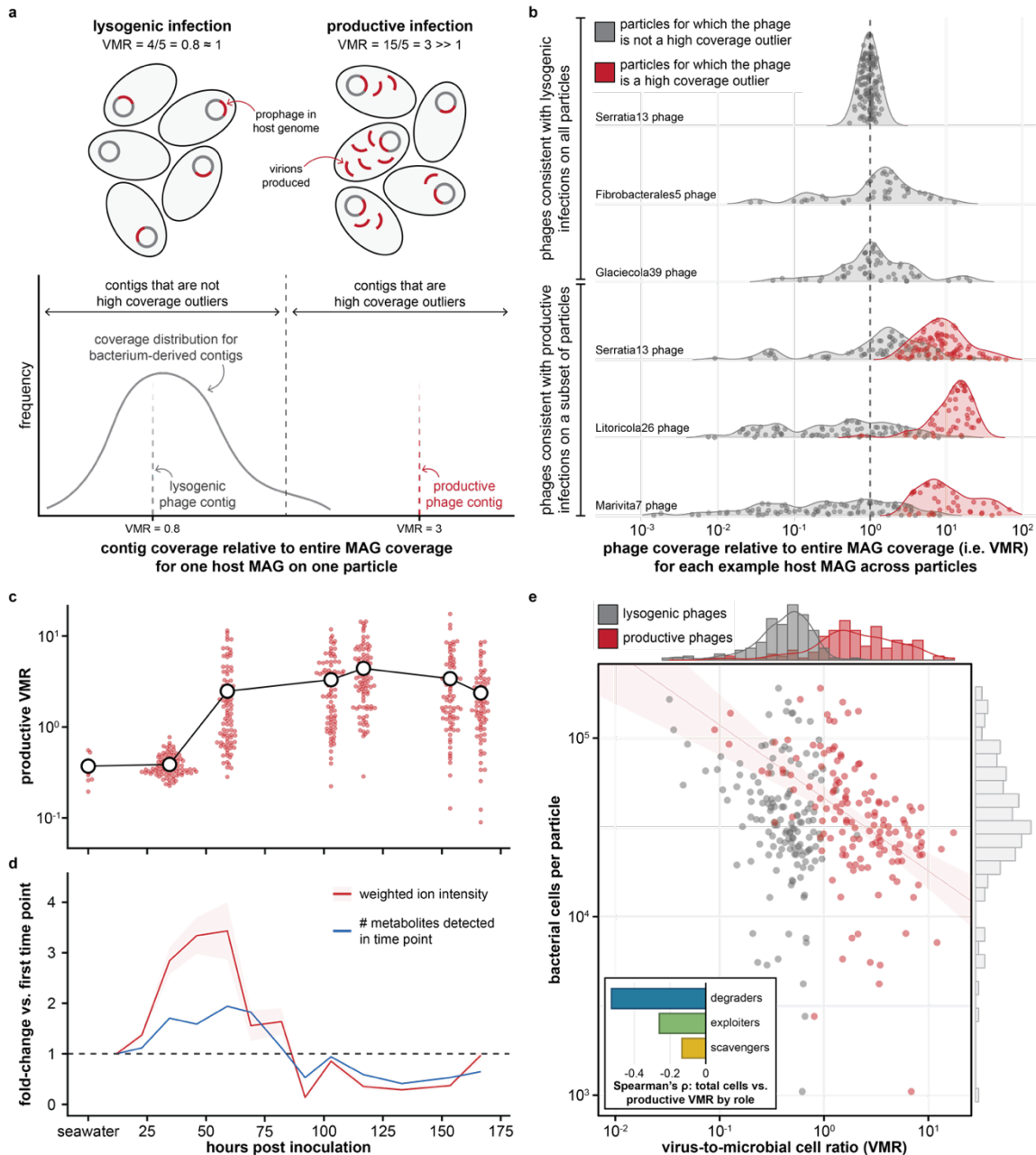
669 **Figure 2. High compositional variability across replicate late-stage particles is driven by**
 670 **conditionally rare degrader taxa.** Relative abundances of metagenome-assembled genomes
 671 (MAGs; $n = 120$) across late-stage particles. Smaller black dots indicate the relative abundance
 672 of each MAG per particle ($n = 149$). Larger white dots indicate the \log_{10} [mean relative abundance]
 673 across the particles on which the MAG was found. MAGs are sorted from left to right by their
 674 prevalence across particles (i.e. the number of particles on which they are detected). The bars
 675 above show the number of chitinases encoded in each MAG. The annotations below show each
 676 MAG's predicted ecological role (heatmap: blue = degrader, green = exploiter, yellow =
 677 scavenger); jackpot score (heatmap: white = low, black = high); and taxonomic order (E =
 678 Enterobacterales, R = Rhodobacterales, P = Pseudomonadales, F = Flavobacterales, C =
 679 Cytophagales, O = Other). See Fig. S1 for additional details.



680

681 **Figure 3. Late-stage particles diverge in community-level functional potential and biomass.**

682 **(a)** Ternary plot of the relative abundances of organisms occupying the three ecological roles
683 (degrader, exploiter, scavenger) on each late-stage particle ($n = 149$), calculated by summing the
684 relative abundances of MAGs classified into each role. Red dots represent jackpot particles, and
685 black ones represent non-jackpot particles. Jackpot particles harbored significantly higher degrader
686 populations than non-jackpot particles (79.8% vs. 47.4% on average; Mann-Whitney U test: $p <$
687 2.2×10^{-16}). **(b)** Estimates of absolute bacterial cell counts on late-stage particles through qPCR of
688 the 16S rRNA gene in DNA extracted from particle-attached communities. Jackpot particles (red
689 dots) harbored significantly higher numbers of cells (Mann-Whitney U test: $p = 2.3 \times 10^{-7}$) than non-
690 jackpot particles (black dots). **(c)** Representative images of late-stage particles that were harvested
691 after 167 hours of incubation in seawater and stained with the DNA-intercalating dye SYTO 9 (scale
692 bar, 50 μm). Particle-attached communities spanned a range of growth states, from sparsely to
693 densely populated.



694

695 **Figure 4. Bacteriophages become increasingly activated during community development**

696 **and contribute to variability in bacterial abundances on late-stage particles. (a)** Schematic

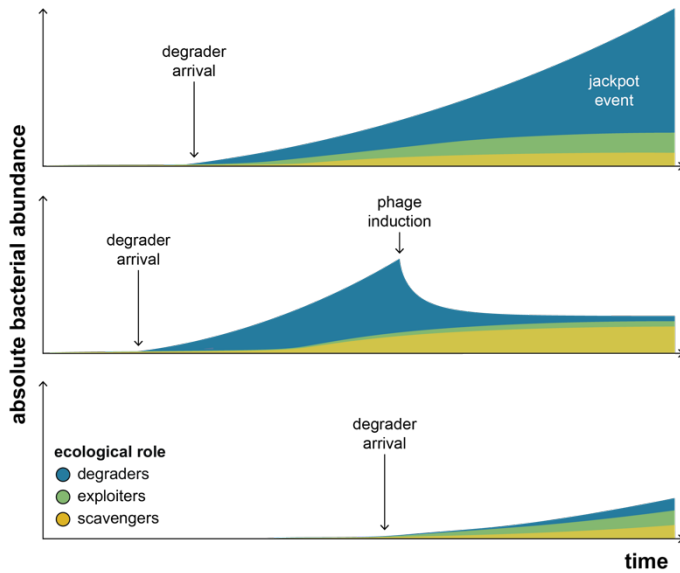
697 of approach to detect productive phage infections from metagenomic data. Left: during lysogenic

698 infections, prophages replicate with their bacterial hosts (virus-to-microbe ratio, VMR, ≈ 1 , top);

699 lysogenic phage contigs have read coverage values similar to those of most bacterial contigs of

700 their host MAG (bottom). Right: during productive infections, prophages replicate much more than

701 their hosts (VMR \gg 1, top); productive phage contigs have read coverage values much higher
702 than those of most bacterial contigs of their host MAG (bottom). **(b)** Representative examples of
703 phages with lysogenic coverage patterns on all late-stage particles (top three rows), and of
704 phages with productive coverage patterns on a subset of particles (bottom three rows). For each
705 phage contig, VMR is shown across late-stage particles on which each MAG is present. Gray
706 dots, particles on which the phage contig is not a coverage outlier; red dots, particles on which
707 the phage is a high coverage outlier. Dashed line: VMR = 1. **(c)** Total VMRs for productive
708 phages over time. The first time point shows productive VMRs of initial seawater samples;
709 subsequent time points show productive VMRs for chitin particle-attached communities incubated
710 in seawater. Smaller red dots, values for individual samples; larger white dots, mean VMR for
711 each time point. **(d)** Metabolomic profiles of the seawater surrounding chitin particles as a
712 function of incubation duration. Values are depicted in terms of fold-change at each time point
713 relative to the first time point (dashed line: no change). Red line (and shading): mean (\pm 1
714 standard deviation) weighted ion intensity (Methods). Blue line: number of unique metabolites. **(e)**
715 Main: Absolute bacterial cell counts on late-stage particles ($n = 142$), estimated through qPCR,
716 vs. each particle's total VMR for lysogenic phages (gray dots) and productive phages (red dots).
717 Cell counts were negatively correlated with productive VMRs (Spearman's $\rho = -0.56$, $p = 3.3 \times 10^{-13}$;
718 red line and shading: log-log linear regression and 95% confidence interval, $R^2 = 0.23$, $p =$
719 1.3×10^{-9}). Productive and lysogenic VMRs were decoupled (red vs. gray data: Spearman's $\rho =$
720 0.11 , $p = 0.18$). Marginal histograms: distributions of productive VMRs (red), lysogenic VMRs
721 (dark gray), and bacterial cell counts (light gray). Inset: Bar plot of values of Spearman's ρ
722 between cell counts and productive VMRs of bacterial populations by ecological role (blue =
723 degraders, green = exploiters, yellow = scavengers; see Fig. S14b for details).



724

725 **Figure 5. Conceptual model of key processes contributing to the diversification of**
726 **communities on microscale particles.** Schematics of community development over time are
727 shown for three example particles, with the absolute abundances depicted for bacterial
728 populations by ecological role (blue = degraders, green = exploiters, yellow = scavengers). Based
729 on our conceptual model (see Discussion), high-biomass jackpot particles are those on which
730 degraders arrive early and resist phage induction, leading to high relative and absolute degrader
731 abundances (top). By contrast, low-biomass particles are those on which degraders are not able
732 to proliferate, either because phage induction leads to their population collapse (middle) or
733 because they become established on a particle relatively late (bottom).

734 **Supplementary Information Text**

735

736 **Models of abundance fluctuations**

737

738 **Notation and context**

739

740 We considered a system with M MAGs and P particles. Let x_i be the abundance of MAG i on a
741 particle and $X = \sum_i x_i$ be the total abundance. The probability distribution $p(\underline{x})$ is the probability of
742 observing a given vector of abundance \underline{x} , while $p_i(x)$ is the probability that species i has abundance
743 x and $P(X)$ is the probability that the total abundance is X . We also define $y_i = \log x_i$ and $Y = \log X$
744 (\log means natural log everywhere).

745

746 Models #1-3 are reasonable models that nevertheless do not recapitulate the observed trends (i.e.
747 the right-skewed distributions of relative taxon abundances [Fig. S2] and the left-skewed
748 distribution of absolute cell abundances [Fig. 3b]), which model #4 (referenced in the main text)
749 does reproduce.

750

751 **Model #1: Stochastic arrival and exponential growth**

752

753 We assume that MAGs arrive stochastically to a particle and grow exponentially with a fixed MAG-
754 specific growth rate r_i . The log-abundance of MAG i at time t will therefore be $y_i = r_i(t - t_i^a)$,
755 where t_i^a is the arrival time of MAG i .

756

757 The only source of variation across particles is the intrinsic randomness in the arrival time, which
758 is exponentially distributed with (migration) rate λ_i . If we are considering only particles where i is
759 present, the probability should be normalized between 0 and the duration of the experiment t , which
760 leads to

761

762
$$\rho_i(t^a) = \frac{\lambda_i e^{-\lambda_i t^a}}{1 - e^{-\lambda_i t}}.$$

763

764 One can obtain the probability of observing a MAG with log-abundance y_i at time t simply by
765 inverting the relationship $y_i = r_i(t - t_i^a)$:

766

767
$$p_i(y) = \frac{\lambda_i \exp(-\lambda_i t + \lambda_i \frac{y}{r_i})}{r_i (1 - e^{-\lambda_i t})},$$

768

769 and, therefore, the probability of the abundance (conditioned on being present) reads

770

771
$$p_i(x) = \frac{\lambda_i \exp(-\lambda_i t)}{r_i (1 - e^{-\lambda_i t})} x^{\frac{\lambda_i}{r_i} - 1}.$$

772

773 Note that this distribution is normalized between 0 and $\tilde{x}_i = e^{r_i t}$. We can therefore rewrite this
774 expression as

775

776
$$p_i(x) = \frac{\lambda_i}{r_i} \frac{\lambda_i}{\tilde{x}_i^{r_i}} x^{\frac{\lambda_i}{r_i} - 1}.$$

777

778 Both the arrival rate λ_i and the growth rate m_i differ across MAGs. We set their values by drawing
779 them from two independent lognormal distributions. In particular, each λ_i for $i = 1, \dots, M$ was drawn
780 from a lognormal distribution with mean $\bar{\lambda}$ and log-variance s_λ^2 . Similarly, each r_i was drawn from a
781 lognormal with mean \bar{r} and variance s_r^2 .

782

783 Fig. S16 shows the distribution of collapsed MAG relative abundances and the distribution of total
784 abundances obtained with this model. Model #1 always predicts a relative log-abundance
785 distribution with negative skewness and a total log-abundance distribution with non-negative
786 skewness (contrarily to what observed in the data; see Fig. 3b).

787

788 **Model #2: Stochastic arrival and exponential growth with demographic stochasticity**

789

790 Model #2 assumes that MAGs arrive on particles with rate λ_i . The population growth that follows is
791 determined by a birth-death process with constant per-capita birth and death rates (b_i and d_i ,
792 respectively). The (average) growth rate r_i equals $b_i - d_i$.

793

794 Similar to the procedure of model #1, we assumed that the values of migration, growth, and death
795 rates of each MAG were initialized as lognormal random variables with means $\bar{\lambda}$, \bar{r} , and \bar{d} and log-
796 variances s_{λ}^2 , s_r^2 , and s_d^2 .

797

798 Fig. S17 shows that model #2 always predicts a total log-abundance distribution with positive
799 skewness, therefore failing in reproducing the empirical shape of the total abundance distribution.

800

801 **Model #3: Stochastic arrival and exponential growth with environmental stochasticity**

802

803 Model #3, similarly to model #1, assumes that MAGs arrive stochastically to a particle with arrival
804 rate λ_i and then grow exponentially. When a MAG arrives on a particle, it starts growing
805 exponentially. The growth rate of MAG i is not fixed, equal to r_i across all particles, but is itself a
806 random variable. In particular, the growth rates of MAG i across particles are normally distributed
807 with mean r_i and variance σ_i^2 proportional to the mean squared: $\sigma_i^2 = c_r^2 r_i^2$, where c_r is the
808 coefficient of variation.

809

810 As for model #1, the arrival rate λ_i and the average growth rates r_i are lognormally distributed with
811 means $\bar{\lambda}$ and \bar{r} and log-variances s_{λ}^2 and s_r^2 .

812

813 Fig. S18 shows that model #3 always predicts a total log-abundance distribution with positive
814 skewness, therefore failing in reproducing the empirical shape of the total abundance distribution.

815

816 **Model #4: Exponential growth with environmental stochasticity conditioned on degrader**
817 **presence**

818

819 In the previous models, the growth of all MAGs was only conditioned on arrival. This assumption
820 inevitably led to total log-abundance distributions with positive skewness, contrarily to the empirical
821 observation of negative skewness.

822

823 Model #4 assumes that, for a particle to become viable for growth, the presence of a degrader MAG
824 is required first. The arrival rate of a degrader is λ^d . All the cells that arrive to the particle after the
825 first arrival of the degrader are able to grow. The time at which the population of MAG i on a particle
826 will start to grow will be $t_g^i = t_a^i + t_d$, where t_d is the time of arrival of the degrader (an exponential
827 random variable with rate λ^d) and t_a^i is the time between arrival of the degrader and the arrival of
828 the MAG i (an exponential random variable with rate λ_i).

829

830 Starting at t_g^i , MAG i will start to grow exponentially with a random, normally distributed, growth
831 rate with mean r_i and coefficient of variation c_r . Similar to the previous models, the arrival rate λ_i
832 and the average growth rates r_i are lognormally distributed with means $\bar{\lambda}$ and \bar{r} and log-variances
833 s_λ^2 and s_r^2 .

834

835 Fig. S19 shows that the predictions of model #4 agrees with the empirical observations. The total
836 log-abundance distribution has a negative skewness, while the distribution of MAG relative
837 abundances has a positive skewness. The shape of the patterns is robust across different
838 parameters values. Only when the variation across MAGs is comparable to the growth rate
839 fluctuations across particles ($c_r \sim s_r \sim 1$) does the total log-abundance distribution display a
840 positive skewness.

841

842

843 **Extended methods**

844

845 **Sample collection and incubation with individual chitin particles**

846

847 **Seawater sampling and treatment.** Nearshore coastal ocean surface water samples were
848 collected on July 15, 2017 from Canoe Beach, Nahant, MA, USA (42°25'11.5" N, 70°54'26.0"
849 W). The seawater was immediately transported to Parsons Laboratory (MIT, Cambridge, MA,
850 USA) for processing. In order to decrease the degree of dissimilarity between seawater aliquots
851 used in incubations with chitin particles, large particulate matter was removed (using a 63µm
852 filter), and the flow-through was concentrated via gentle centrifugation in 1L batches at 4000 ×
853 g for 5 minutes. The lower 100mL of each 1L batch was saved and pooled; aliquots of this
854 water in 175µL volumes were either used for particle incubations or stored at -20°C for
855 downstream DNA extraction and metagenomic sequencing.

856

857 **Seawater incubation with individual chitin particles.** Artificial seawater (ASW), used for
858 washing and storing chitin particles, was prepared by dissolving 40g/L sea salts (Sigma-
859 Aldrich, #S9883) in Milli-Q deionized water and filtering the solution through a 0.22-µm filter.
860 Chitin magnetic particles (New England Biolabs, #E8036L) stored in 20% ethanol were washed
861 three times (2mL particles resuspended in 50mL ASW) using a magnet to pull down the
862 particles. Aliquots of washed chitin particles were further diluted in ASW in sterile petri dishes
863 and individually selected beneath a dissecting microscope in a laminar flow hood. Single chitin
864 particles were transferred in 3µL volumes of ASW into the wells of 96-well plates (Thermo
865 Fisher, #AB0600L; UV-sterilized; free from DNase, RNase, and human DNA), with one chitin
866 particle per well. The individual particles selected had a diameter of 85.0±24.0 µm, which was
867 quantified from a set of 60 particles on an ImageXpress Micro Confocal (Molecular Devices).
868 Plates containing individual particles were stored at 4°C until they were inoculated
869 consecutively with 175µL of filtered, centrifuged seawater per well. The plates were sealed

870 (VWR, #89092-056) and rotated end-over-end at 7.5 revolutions/minute at room temperature.
871 The particles in an entire plate were harvested at each time point (after 12, 22.75, 34.5, 46, 59,
872 69, 82, 92, 103, 116.75, 113, 153.5, and 166.5 hours of incubation) by pipetting the contents
873 of each well onto a sterile petri dish and inspecting the water under a dissecting microscope in
874 a laminar flow hood. Each particle was transferred in 1 μ L volumes into 96-well plates (Thermo
875 Fisher, #AB0600L) pre-filled with 100 μ L of TE buffer; plates with harvested particles were
876 stored at -20°C until downstream processing. The seawater surrounding each harvested
877 particle was also saved in 96-well plates (Thermo Fisher, #AB0600L) and stored at -20°C until
878 downstream processing.

879

880 **DNA extraction and metagenomic sequencing.** DNA extractions were performed for twelve
881 175 μ L-volume aliquots of the initial, unincubated seawater, as well as for particles harvested after
882 34.5, 59, 103, 116.75, 113, 153.5, and 166.5 hours of incubation. DNA was extracted from all
883 samples with the Agencourt DNAdvance Genomic DNA Isolation Kit (Beckman Coulter) using
884 reagent volumes 0.5X relative to those specified in the manufacturer's protocol, except for the
885 elution buffer, of which only 30 μ L was used for each sample to avoid over-diluting low DNA yields.
886 Metagenomic libraries were prepared with the Nextera XT DNA Library Prep Kit and index primers
887 (Illumina) using the protocol developed by Rinke *et al.* (1) for low DNA input samples. While the
888 results from the protocol in Rinke *et al.* were reproducible with as little as 100fg of input DNA, the
889 authors recommend using a minimum of 1pg as input. Based on our qPCR measurements of DNA
890 extracted from individual particle-attached communities (as described in the Methods section "Cell
891 count estimation"), only one of our libraries (with 0.44pg input) was created with less than 1pg DNA.
892 The modifications to the manufacturer's library preparation protocol included (i) diluting the
893 Amplicon Tagment Mix 1:10 in non-DEPC-treated nuclease-free water, and (ii) increasing the
894 number of PCR amplification cycles of the tagmented DNA from 12 to 20 cycles. Amplified libraries
895 were purified with 0.6X AMPure XP beads. Each library was quantified on an Agilent 4200
896 TapeStation system with High Sensitivity D5000 ScreenTapes (Agilent Technologies) following the
897 manufacturer's protocol, and successfully amplified libraries were pooled by time point in equimolar

898 amounts. Sequencing was performed on an Illumina HiSeq 2500 machine (250bp paired-end
899 reads) at the Genome Technology Core of the Whitehead Institute for Biomedical Research (MIT,
900 Cambridge, MA, USA). See Table S5 for all sample metadata.

901

902 **Metagenomic analyses**

903

904 **Read pre-processing.** Raw sequencing reads were clipped (to remove adapter sequences)
905 and trimmed for quality with Trimmomatic v0.36 (2) (parameters: LEADING:3, TRAILING:3,
906 SLIDINGWINDOW:10:20, MINLEN:36). Reads mapping to the PhiX genome were filtered out
907 with BBDuk v38.16 (3) (parameters: k=31, hdist=1) and those mapping to the human genome
908 (masked by Brian Bushnell at the Joint Genome Institute to prevent false positives) were
909 identified and removed using BBDuk v38.16 (parameters: minid=0.95 maxindel=3 bwr=0.16
910 bw=12 minhits=2 qtrim=rl trimq=10 untrim; reference genome:
911 hg19_main_mask_ribo_animal_allplant_allfungus.fa.gz).

912

913 **Metagenome assembly, binning, and MAG taxonomic assignment.** Default parameters
914 were used for all tools unless otherwise specified. Trimmed, filtered reads were error-corrected
915 using BayesHammer (4) (a component of the SPAdes v3.13.0 pipeline) in order to improve
916 contig assembly. Reads within each time point were pooled and co-assembled using MEGAHIT
917 v1.2.9 (5). Assembled contigs at least 1kb in length were binned using two complementary
918 tools – MaxBin v2.2.7 (6) and CONCOCT v1.1.0 (7). To provide CONCOCT with coverage
919 estimates, error-corrected reads were mapped to contigs using Bowtie 2 v2.3.4.1 (8) with the
920 parameters and approach described in Leventhal *et al.* (9). Bins generated with MaxBin and
921 CONCOCT were consolidated and filtered using DAS Tool v1.1.1 (10) and evaluated for
922 completeness and contamination with CheckM v1.1.2 (11). The resulting 251 bins that were at
923 least 50% complete and at most 10% contaminated were used as reference MAGs, with
924 median completeness and contamination values of 93.7% and 3.9%, respectively, across this
925 set of MAGs (Table S1). Taxonomic classifications from the Genome Taxonomy Database

926 (GTDB) (12) were assigned to MAGs using GTDB-Tk v1.1.1 (13). Highly similar MAGs obtained
927 from separate co-assemblies were identified and clustered through a pipeline developed by Dr.
928 Jakob Russel for performing whole-genome comparisons of each MAG against all others with
929 BLAT v36x2 (14). Briefly, a similarity score was calculated for each MAG relative to another by
930 dividing the combined length of its contigs at least 98% identical to those in the compared MAG
931 by the combined length of all its contigs. A threshold for distinguishing high similarity scores
932 from low ones was determined using Otsu's method (15) with code derived from the R (16)
933 package EBImage (17). 132 clusters of MAGs with mutually high similarity scores were
934 identified, and all MAGs in each cluster had consistent GTDB-based taxonomic assignments.
935 For one MAG cluster, one of the MAGs was classified as a different genus from the other
936 MAGs; this MAG was separated from the cluster. We chose to consider clustered MAGs as a
937 unit, rather than to dereplicate them, in order to retain potential strain-level microdiversity in our
938 reference set.

939

940 **MAG ecological role assignments.** For each MAG, protein-coding genes were predicted and
941 translated using Prodigal v2.6.3 (18). Predicted protein sequences were compared to a custom
942 database of profile hidden Markov models (HMMs) of proteins involved in growth on chitin using
943 the *hmmsearch* function of HMMER v3.3 with default parameters (19). Publicly-available
944 HMMs were downloaded from the Pfam v33.1 (20) or TIGRFAM v15.0 databases (21) (see
945 Table S6 for accession numbers). Custom HMMs were made by identifying experimentally-
946 verified proteins of interest (22, 23), finding their homologs in the UniProtKB/Swiss-Prot
947 v2020_06 database (24), creating a seed alignment using MAFFT v7 with default parameters
948 (25, 26), and building the profile HMMs using the *hmmbuild* function of HMMER with default
949 parameters (see Table S6 for details on each custom HMM). Protein-coding sequences were
950 annotated based on the *hmmsearch* results if the protein length was at least 100 amino acids,
951 the independent E-value was less than 1×10^{-9} , and the domain score was greater than 30. Only
952 the most significant annotation was used for each protein sequence. Gene copy numbers were
953 calculated for each MAG by tallying the number of annotations made for each protein group

954 (Table S1). Ecological roles (as degraders, chitooligosaccharide exploiters, or metabolic
955 byproduct scavengers) for MAGs were defined based on the gene content patterns observed
956 for strains previously isolated from particle enrichments (27, 28), fully sequenced, and
957 phenotyped according to their abilities to grow on colloidal chitin, chitobiose, and GlcNAc (29).
958 MAGs were classified as degrader genomes if they encoded at least 1 chitinase and at least 1
959 copy of any of the following genes: GlcNAc-specific methyl-accepting chemotaxis protein
960 (MCP), GlcNAc-specific phosphotransferase system IIBC component (PTS), GlcNAc-specific
961 TonB-dependent transporter (TBDT), N,N'-diacetylchitobiose phosphorylase, beta-N-
962 acetylhexosaminidase, or GlcNAc kinase. MAGs were classified as exploiter genomes if they
963 encoded 0 chitinases and had at least one of the following characteristics: more than 1 copy of
964 beta-N-acetylhexosaminidase or at least 1 copy of MCP, PTS, TBDT, or N,N'-
965 diacetylchitobiose phosphorylase. MAGs were classified as scavenger genomes if they
966 encoded 0 chitinase, MCP, PTS, TBDT, and N,N'-diacetylchitobiose phosphorylase copies,
967 and 1 or fewer copies of beta-N-acetylhexosaminidase. If MAGs clustered by similarity were
968 assigned different ecological roles by these heuristics, then either (i) the role assigned to all
969 MAGs defaulted to the role of the MAG with the lowest contamination and/or highest
970 completeness (which occurred for 4 MAG clusters), or (ii) the MAG cluster was split into two
971 subclusters (which occurred for 5 MAG clusters); these discrepancies are indicated in Table
972 S1. Following this MAG cluster curation, there were a total of 138 MAG clusters.

973

974 **Read mapping to MAGs for relative abundance estimation.** All trimmed, filtered reads were
975 mapped competitively against the MAGs created from sequencing particle-attached
976 communities; the initial, unincubated seawater; and the negative controls (see the Methods
977 section "Mock communities and negative controls"). Samples with fewer than 100,000 trimmed,
978 filtered reads were excluded from analyses. Read mapping was performed using Bowtie 2
979 v2.3.4.1 (8) with the parameters and approach described in Leventhal *et al.* (9) and post-
980 processed using SAMtools v1.7 (30). Reads that best mapped (based on alignment scores) to
981 MAGs obtained from the negative controls (which were contaminants from laboratory reagents)

982 and to MAGs obtained from particle sequences that were also likely environmental
983 contaminants (indicated in Table S1; determined through a literature search of each strain's
984 taxonomy in studies of the marine environment) were removed from consideration when
985 estimating community compositions. To avoid artifactually double-counting hits from paired
986 reads, only the best hit of the forward read was considered for read pairs that survived trimming
987 and quality filtering. Hits to completely bacteriophage-derived contigs (as opposed to
988 prophages integrated into bacterial genome contigs) were also excluded from estimates of
989 MAG relative abundances (see the Methods section "Bacteriophage analysis"). To minimize
990 spurious detection, MAGs were considered to be "present" in a sample if they recruited at least
991 0.05% of the reads in a sample; for MAGs that recruited reads below this threshold in a sample,
992 their abundance was set to 0 for that sample. MAG relative abundances for MAGs above this
993 threshold were calculated for each sample by (1) tallying the hits to all MAGs in each MAG
994 cluster; (2) normalizing the tally by the average genome length of all MAGs in each MAG
995 cluster; and (3) dividing the normalized tallies for each MAG cluster by their sum for each
996 sample. Therefore, for MAGs clustered together based on similarity (see the Methods section
997 "Metagenome assembly, binning, and MAG taxonomic assignment"), their relative abundances
998 are represented in that of the entire MAG cluster to which they belong; this calculation
999 circumvents the artificial underestimation of MAG relative abundances that would otherwise be
1000 obtained with a non-dereplicated reference set. The relative abundances of organisms
1001 occupying the three ecological roles (degrader, exploiter, scavenger) on each particle were
1002 calculated by summing the relative abundances of MAGs classified into each role. Based on
1003 information gathered from relative abundance estimation, particles harvested at 113 hours
1004 post-inoculation were excluded from analyses because of a clear batch effect at that time point
1005 characterized by high abundances of MAG *Serratia_liquefaciens93* (96.7% of particles on
1006 which *Serratia_liquefaciens93* was at least 10% abundant were from t=113h, which included
1007 98.9% of particles from that time point; this MAG was also the most abundant MAG on 81.1%
1008 of particles from t=113h and was not the most abundant MAG on any particles from other time
1009 points; see Table S7).

1010

1011 **Comparison of variability in seawater vs. particle-associated communities.** Inter-sample
1012 variability was estimated as the Aitchison distance between the community compositions of
1013 pairs of samples (i.e. the Euclidian distance between center log-ratio-transformed MAG relative
1014 abundance vectors). Aitchison distances were calculated between aliquots of the initial,
1015 unincubated seawater and between late-stage particle communities separately, and the
1016 distributions of distances between all pairs of samples were compared to each other.

1017

1018 **Definitions of jackpot MAGs and jackpot particles.** A jackpot score was calculated for each
1019 MAG cluster to quantitatively reflect the properties of rarity across most particles and
1020 dominance on a few particles. Based on relative abundances across late-stage particles, each
1021 MAG's jackpot score was defined as:

1022

1023
$$\frac{(\text{coefficient of variation of relative abundances}) * (\# \text{ particles on which MAG is the most abundant}) * (\text{highest relative abundance achieved})}{(\# \text{ particles on which MAG is present})^2}$$

1024

1025 Therefore, MAGs with high scores strongly display the jackpot phenomenon, whereas MAGs
1026 with low scores do not. The jackpot score for each particle was calculated as the weighted
1027 average of MAG jackpot scores (i.e. the sum of the relative abundance of each MAG cluster
1028 multiplied by its jackpot score). Each particle's jackpot score was compared to its species
1029 evenness (calculated as Pielou's evenness, i.e. the Shannon diversity index divided by the
1030 natural logarithm of species richness) with the expectation that particles that most strongly
1031 display the jackpot phenomenon have low species evenness. Particles were defined as
1032 "jackpot particles" if they have jackpot scores that exceed the threshold value above which log-
1033 transformed values of species evenness drop sharply (Fig. S5a); this value corresponds to the
1034 60th percentile of jackpot particle scores. For comparing the binary categories of "jackpot
1035 degraders" and "non-jackpot degraders" (Fig. S15a), "jackpot degraders" were those MAGs
1036 that had jackpot scores greater than zero and that were present on less than 75% of late-stage
1037 particles; this thresholding was done in order to exclude the MAG clusters Serratia13 and

1038 Fibrobacterales⁵ that had very low yet non-zero jackpot scores because of their high relative
1039 abundances on many particles (see Tables S1 and S7).

1040

1041 **Calculation of the percent variance explained in MAG abundances on individual**
1042 **particles by the MAG abundances theoretically obtained by sequencing particles in bulk.**

1043 To evaluate the extent to which community compositions at the single particle level diverged
1044 from that of a “bulk” measurement theoretically obtained by sequencing all particles together,
1045 we calculated “bulk” MAG abundances by (1) normalizing the mapped read counts to each
1046 MAG cluster by the total number of read counts for each sample; (2) summing the counts for
1047 each MAG cluster across samples; (3) normalizing the sum across samples by the average
1048 genome length of all MAGs in each MAG cluster; and (4) dividing the length-normalized counts
1049 for each MAG cluster by their sum. (These “bulk” MAG relative abundances are equivalent to
1050 the mean MAG relative abundances calculated across all particles, including those particles
1051 where MAGs are absent.) The percent variance in the abundance ranks of MAGs on single
1052 particles explained by the abundance ranks for the theoretical bulk measurement was
1053 calculated for each particle as the square of the Pearson correlation coefficient (between each
1054 individual vs. the bulk abundance rank), multiplied by 100.

1055

1056 **Multivariate analysis.** We inferred the number of conditional dependencies between MAGs
1057 from the estimated inverse covariance matrix of center log-ratio-transformed MAG relative
1058 abundances, repeating this process for 1000 randomizations of the data in which we permuted
1059 particle labels for each MAG but retained their abundance distributions. The inverse covariance
1060 matrices were estimated using a graphical lasso approach with the R package *glasso* (31) for
1061 several values of the regularization parameter ($\rho = 0.005$, $\rho = 0.001$, $\rho = 0.0005$, and $\rho =$
1062 0.0001).

1063

1064 **Read mapping to chitinases and calculating the chitinase-weighted means.** To evaluate
1065 whether the use of MAGs as reference genomes could have biased our estimate of the

1066 degrader population relative abundances in particle-attached communities, reads were also
1067 mapped to a reference set of chitinase genes (regardless of binning). All assembled contigs
1068 (binned and unbinned) were annotated for chitinase genes using the HMM-based approach
1069 described in the Methods section “MAG ecological role assignments.” A custom DIAMOND
1070 database of 3,370 translated chitinase genes was created using the *makedb* function of
1071 DIAMOND v0.9.10.111 (32) with default parameters. Because of the high sequence diversity
1072 of chitinase genes, we chose to make this custom database so that the chitinase sequences
1073 used as references would be representative of those found in this experiment. Trimmed,
1074 quality-filtered reads were mapped to this database using the *blastx* function of DIAMOND with
1075 default parameters. To avoid artifactually double-counting hits from paired reads, only the best
1076 hit of the forward read was considered for read pairs that survived trimming and quality filtering.
1077 Only the most significant hit was counted for each read and only if the E-value was less than
1078 or equal to 1×10^{-25} . The number of such hits was tallied for each sample and divided by the
1079 number of trimmed, quality-filtered reads used in the mapping step to yield the percent of reads
1080 in each sample mapping to chitinase genes (Table S5). If the degrader population relative
1081 abundance estimated by MAGs were a consistent approximation of the true degrader
1082 population abundance, then the wide range in the number of chitinases encoded in each
1083 degrader MAG (Table S1) would be reflected in the percent of reads in each community
1084 mapping to chitinase genes. Therefore, the community-weighted mean (CWM) for chitinases
1085 was calculated as another comparison to the percent of reads mapping to chitinases. The
1086 chitinase CWM was calculated by multiplying the relative abundance of each degrader MAG
1087 (or MAG cluster) by the number of chitinases encoded in it (or the mean number of chitinases
1088 for a MAG cluster), and finally by summing these values.

1089

1090 **Bacteriophage analysis.**

1091

1092 *Identifying phage-derived contigs.* Binned contigs assembled from our metagenomic
1093 dataset were first classified as phage-derived or bacteria-derived using tools designed to

1094 detect phage sequences among mixed metagenomes – namely, (i) VirSorter v1.0.3 (33)
1095 via the CyVerse platform (www.cyverse.org; National Science Foundation Awards DBI-
1096 0735191, DBI-1265383, DBI-1743442) using its RefSeqABVir database and default
1097 parameters; and (ii) VirFinder v1.1 (34) with default parameters. Contigs were classified as
1098 phage-derived if they met one of the following standards as employed in Gregory *et al.*
1099 (35): (i) they were classified by VirSorter as Category 1 or 2 (complete phage contig, higher
1100 confidence); (ii) they were classified by VirFinder with a score ≥ 0.9 and p-value < 0.05 ; or
1101 (iii) they were classified both by VirSorter as Category 3 (complete phage contig, lower
1102 confidence) and by VirFinder with a score ≥ 0.7 and p-value < 0.05 .

1103

1104 *Identifying productive vs. lysogenic phage-derived contigs.* We used a read coverage-
1105 based approach to categorize phage-derived contigs as productive or lysogenic in particle-
1106 attached communities. Phages in a productive cycle in a particular sample would have a
1107 higher coverage than the bacterial contigs of the MAG with which they were binned
1108 because of the multiple virion copies produced per bacterial cell. In contrast, phages in a
1109 lysogenic cycle would have coverage values comparable to those of the bacterial contigs
1110 of the MAG with which they were binned. We reasoned that contigs classified as phage-
1111 derived, especially those belonging to the genomes of temperate phages, were likely to be
1112 binned into the MAGs of their bacterial hosts because: i) phage *k*-mer signatures tend to
1113 be more similar to those of their specific hosts than to those of random bacteria (34, 36,
1114 37); ii) phages in a lysogenic cycle will have the same sequencing read coverage patterns
1115 as their hosts across samples; and iii) accordingly, both of the binning algorithms we
1116 employed clustered contigs based on their tetranucleotide frequencies and their coverage
1117 levels across multiple samples. Because our approach relied on comparisons between co-
1118 binned phages and MAGs, we considered in our analyses only phage-classified contigs at
1119 least 5kb in length, since the likelihood of mis-binning decreases with increasing contig
1120 length.

1121 Based on read mapping to MAGs (see the Methods section “Read mapping to
1122 MAGs for relative abundance estimation”), per-base coverage values for all binned contigs
1123 were computed with the *genomecov* function of BEDTools v2.27.0 (38) and were used to
1124 calculate contig-wide average coverage values. For each MAG and for each sample, a
1125 phage-derived contig was considered to be productive if its coverage was greater than the
1126 coverage of the 95th percentile bacteria-derived contig in the same MAG. A phage derived-
1127 contig was considered to be lysogenic in a sample if its coverage did not exceed the
1128 coverage of the 95th percentile bacteria-derived contig in the same MAG. Through this
1129 pipeline, we identified 263 phage contigs with coverage patterns consistent with productive
1130 infections in a subset of samples and 256 phage contigs with coverage patterns consistent
1131 with lysogenic infections in all samples (Table S2).

1132

1133 *Calculating virus-to-microbial cell ratios (VMRs).* The VMR of an individual phage contig in
1134 one sample is defined as the phage contig coverage divided by average coverage of the
1135 MAG with which it is binned (which was calculated using only the bacteria-derived contigs).
1136 Total VMRs – i.e. the total number of phage copies relative to the total number of bacterial
1137 MAG copies in an entire sample – were calculated separately for productive and lysogenic
1138 phage contigs. The total productive VMR for a sample was defined as:

1139

$$\begin{aligned} 1140 \quad & \sum_i^n \left[\left(\frac{\text{average coverage of productive phage contigs in } MAG_i}{\text{average } MAG_i \text{ coverage}} \right) \right. \\ 1141 \quad & \left. \times (MAG_i \text{ relative abundance}) \right] \\ 1142 \quad & = \frac{\text{total \# phage copies (due to productive infections)}}{\text{total \# bacterial genome copies}} \end{aligned}$$

1143

1144 where n is the number of MAGs found in a sample. This calculation is equivalent to

1145

1146
$$\frac{\sum_i^n (\text{average coverage of productive phage contigs in } MAG_i)}{\sum_i^n (\text{average } MAG_i \text{ coverage})}$$

1147

1148 where n is the number of MAGs found in a sample. Similarly, the total lysogenic VMR for
1149 a sample was defined as

1150

1151
$$\sum_i^n \left[\left(\frac{\text{average coverage of lysogenic phage contigs in } MAG_i}{\text{average } MAG_i \text{ coverage}} \right) \right.$$

1152
$$\left. \times (MAG_i \text{ relative abundance}) \right]$$

1153
$$= \frac{\text{total \# phage copies (due to lysogenic infections)}}{\text{total \# bacterial genome copies}}$$

1154

1155 where n is the number of MAGs found in a sample. VMRs for each ecological role (i.e. for
1156 the subpopulation in a community that belongs to one of the three roles of degrader,
1157 exploiter, or scavenger) were calculated using the same formulas as above while
1158 considering only the MAGs of each role and their associated phages. When calculating
1159 total VMRs, we used the average coverage value of all phage contigs in each MAG (rather
1160 than the sum of the coverage values for all phage contigs in each MAG) to obtain a more
1161 conservative estimate of phage copy number. For example, if two phage contigs belonged
1162 to the same phage genome but did not overlap in sequence, they would appear to be two
1163 separate phages; thus, using their sum would double the apparent phage copy number,
1164 while using their average would provide a more accurate representation of their
1165 abundance.

1166

1167 *Analysis controls.* Given that read coverage from metagenomic data is often noisy, it is
1168 conceivable that phage contigs identified as “productive” have high coverage relative to
1169 their associated bacterial MAGs simply due to sequencing noise. We performed two
1170 analyses to examine this possibility. Firstly, we considered that because productive phages

1171 are identified based on coverage, there is a chance that more productive phages would be
1172 found in samples with more reads. (Ensuring that this is not the case is one of the controls
1173 used in Kieft *et al.* (39), which also employs a coverage-based method for finding
1174 productive phages in mixed metagenomes.) Therefore, we calculated the Spearman's
1175 correlation coefficient between the number of reads in a sample and the number of phage-
1176 derived contigs with coverage values above the 95th percentile for their MAG (as described
1177 in the Methods section "Identifying productive vs. lysogenic phage-derived contigs"). We
1178 calculated these correlations for samples within time points to avoid spurious correlations
1179 created by systematic differences in the number of reads obtained across time points. The
1180 phage contigs from the MAGs that showed a significant correlation ($p < 0.05$) were
1181 excluded from analyses. Secondly, for each sample, we compared the average coverage
1182 of all phage-derived contigs (≥ 5 kb) with coverage values above the 95th percentile for their
1183 MAG with the average coverage of all bacteria-derived contigs (≥ 5 kb) with coverage
1184 values above the 95th percentile for their MAG. If the high coverage phage contigs have
1185 comparable average coverage to the bacterial contigs, that would indicate that the phage
1186 contigs had high coverage only due to sequencing noise. The average bacterial coverage
1187 is larger than the average phage coverage in only 2.0% (3/149) of late-stage particles (Fig.
1188 S12). Therefore, for samples with high total productive VMRs, phage contigs with high
1189 coverage values likely represent phages that were replicating more than their bacterial
1190 hosts, rather than representing contigs with randomly higher coverage values.

1191

1192 **Mock communities and negative controls.** In order to quantify the technical error associated with
1193 creating metagenomic libraries from low DNA inputs, mock communities were simulated by
1194 combining the DNA of two strains isolated from a previous chitin particle enrichment experiment
1195 using seawater from the same location sampled for this project (27). The total genomic DNA of
1196 *Vibrio splendidus* strain 1A01 (BioProject #PRJNA414740, Accession #PDUR00000000) and
1197 *Maribacter sp.* 6B07 (BioProject #PRJNA414740, Accession #PDUT00000000) was extracted
1198 using the MasterPure DNA Purification Kit (Epicentre), and double-stranded DNA content was

1199 quantified using the Quant-iT PicoGreen dsDNA Assay Kit (Invitrogen). The DNA of each strain
1200 was mixed in equimolar amounts and serially diluted to either 50pg or 5pg total (to include a range
1201 of expected DNA input concentrations from extractions of communities attached to single chitin
1202 particles). Metagenomic libraries from three technical replicates of each concentration, as well as
1203 from six negative controls (containing only nuclease-free water), were prepared using the same
1204 protocol used for individual chitin particle-attached communities (as described in the Methods
1205 section “DNA extraction and metagenomic sequencing”). The results from the mock community
1206 sequencing are shown in Fig. S4. Of the six negative control libraries, only three amplified; the
1207 MAGs recovered from these samples included *Delftia acidivorans* and a *Brevundimonas sp.*, which
1208 belong to taxonomic groups previously found as contaminants in laboratory reagents used in DNA
1209 extractions and sequencing (40, 41). These MAGs were included as references for the
1210 metagenomic analysis, and the reads best mapping to them (based on alignment scores) were
1211 removed from consideration when estimating community compositions.

1212

1213 **Cell count estimation.** Bacterial DNA extracted from individual particle-attached communities was
1214 quantified through qPCR of the 16S rRNA gene using the Femto Bacterial DNA Quantification Kit
1215 (Zymo Research), which has a lower limit of detection of 20fg. Two sets of standards and negative
1216 controls were included in each qPCR run. The number of bacterial cells for each particle was
1217 estimated from the absolute DNA amounts based on measurements indicating a mean of 2.5fg
1218 DNA per bacterial cell in seawater samples (42).

1219

1220 **Metabolomics experiments and analyses.** We performed untargeted metabolomics of the
1221 seawater that surrounded each chitin particle (after removing the chitin particles at each time point)
1222 and of the initial, unincubated seawater (t=0). All samples were first diluted 1:100 in nuclease-free
1223 water (in two serial 1:10 dilutions). We used a binary LC pump (Agilent Technologies) and an MPS2
1224 Autosampler (Gerstel) coupled to an Agilent 6520 time-of-flight mass spectrometer (Agilent
1225 Technologies) operated in negative mode, at 2GHz, extended dynamic range, with an *m/z*
1226 (mass/charge) range of 50-1000. The mobile phase consisted of isopropanol:water (60:40, v/v) pH

1227 9, with the addition of 5mM ammonium fluoride and a flow rate of 150 μ l/min. Raw data were
1228 processed and analyzed using preprocessing raw mass spectrometry data functions contained in
1229 the bioinformatics toolbox of MATLAB (43, 44). We detected 5714 ions, of which 121 were
1230 annotated against a curated library of metabolites that are present in marine microbes, based on
1231 the BioCyc database (45). Certain ions were matched with multiple isomeric or isobaric compounds
1232 (as noted in Table S3). Detectable metabolites were those with ion intensities that passed the
1233 detection threshold above the inoculum [sample ion intensity > (mean ion intensity at t=0) +
1234 (3*standard deviation of ion intensity at t=0)]. For metabolites that exceeded the limit of detection,
1235 the intensities of each ion were normalized between 0 and 1, where 0 is the limit of detection and
1236 1 is the highest intensity measured of a given ion. Weighted ion intensities for each timepoint were
1237 calculated by taking the sum of all normalized intensities of ions in all samples for each timepoint.

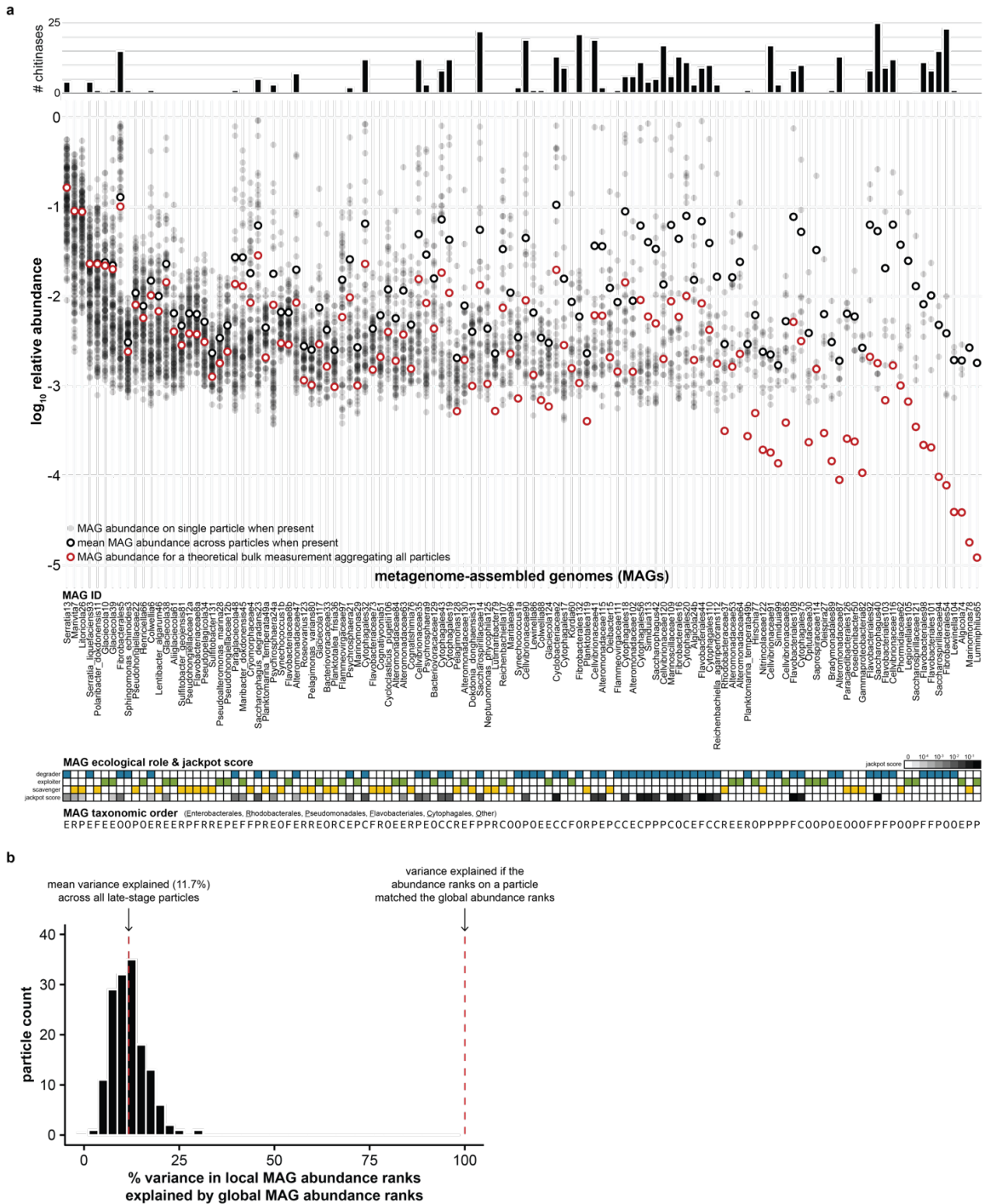
1238

1239 **Sample collection and incubation with many chitin particles.** Seawater was collected on the
1240 day of the experiment from Canoe Beach, Nahant, MA, USA (42°25'11.5" N, 70°54'26.0" W), the
1241 same source as seawater used elsewhere in this study. Chitin magnetic particles (New England
1242 Biolabs, #E8036S) were collected on a 40 μ m cell strainer then passed through a 100 μ m cell
1243 strainer to restrict the size range of the particles to 40-100 μ m (Corning). The size selected particles
1244 were then resuspended in 0.2 μ m-filtered natural seawater to create three suspensions: 807 \pm 99
1245 particles/mL (\pm indicates standard deviation, $n = 3$), 182 \pm 28 particles/mL (sd, $n = 3$), or 88 \pm 3
1246 particles/mL (sd, $n = 3$). Unfiltered natural seawater, containing microbes, was left undiluted, or
1247 diluted 1:10, or diluted 1:100 into 0.2 μ m-filtered natural seawater to create three different initial
1248 densities of bacterioplankton. All combinations of particles and cells were combined by adding 5mL
1249 particle mixture to 10mL cell mixture to create a matrix of 9 separate conditions. Particle/cell
1250 mixtures were incubated in 15mL polystyrene tubes (Falcon) with end-over end rotation at a rate
1251 of 8 revolutions/minute on a Stuart SB3 rotator at room temperature (21-25°C).

1252

1253 **Imaging and quantification of chitin particle colonization by natural seawater bacteria.** For
1254 the experiment incubating chitin particles individually in seawater (see the Methods section

1255 “Seawater incubation with individual chitin particles”), at each time point, the communities on a
1256 subset of particles (that were not sequenced) were stained with the DNA stain SYTO9 (Invitrogen,
1257 #S34854) at a final concentration of 500nM. SYTO9 was added directly to the wells containing the
1258 particles and seawater, which were subsequently incubated in the dark at room temperature for 15
1259 minutes before the individual chitin particles were harvested (as described in the Method section
1260 “Seawater incubation with individual chitin particles”) and mounted separately on microscope
1261 slides. Particles were imaged with a Zeiss epifluorescence microscope at 100X magnification. For
1262 the experiment incubating many particles together in seawater (see the Methods section “Sample
1263 collection and incubation with many chitin particles”), after 24 hours of incubation, 200µl samples
1264 of each condition were stained with SYTO9 at a final concentration of 5µM. The SYTO9-stained
1265 samples were transferred to a black-walled Greiner Bio-One µClear 96-well plate. Samples were
1266 imaged on an ImageXpress Micro Confocal (Molecular Devices) in widefield mode using a Nikon
1267 10x Plan Apo lambda objective (NA 0.45) and FITC filter (ex 482/35, em 536/40, dichroic 506 nm)
1268 with blue LED illumination from a Lumencore Light Engine. Nine fields of view capturing the entire
1269 well were acquired to quantify all particles present in each well. For chitin particles incubated both
1270 individually and in bulk, a custom analysis script was written in MATLAB vR2019a (The Mathworks)
1271 to quantify the area of each chitin particle colonized by cells. The code defines chitin particle area,
1272 and the area of each particle covered by cells using intensity-based thresholds. Code and original
1273 data will be publicly available before publication at the following GitHub page:
1274 https://github.com/jaschwartzman/seawater_colonize



1275

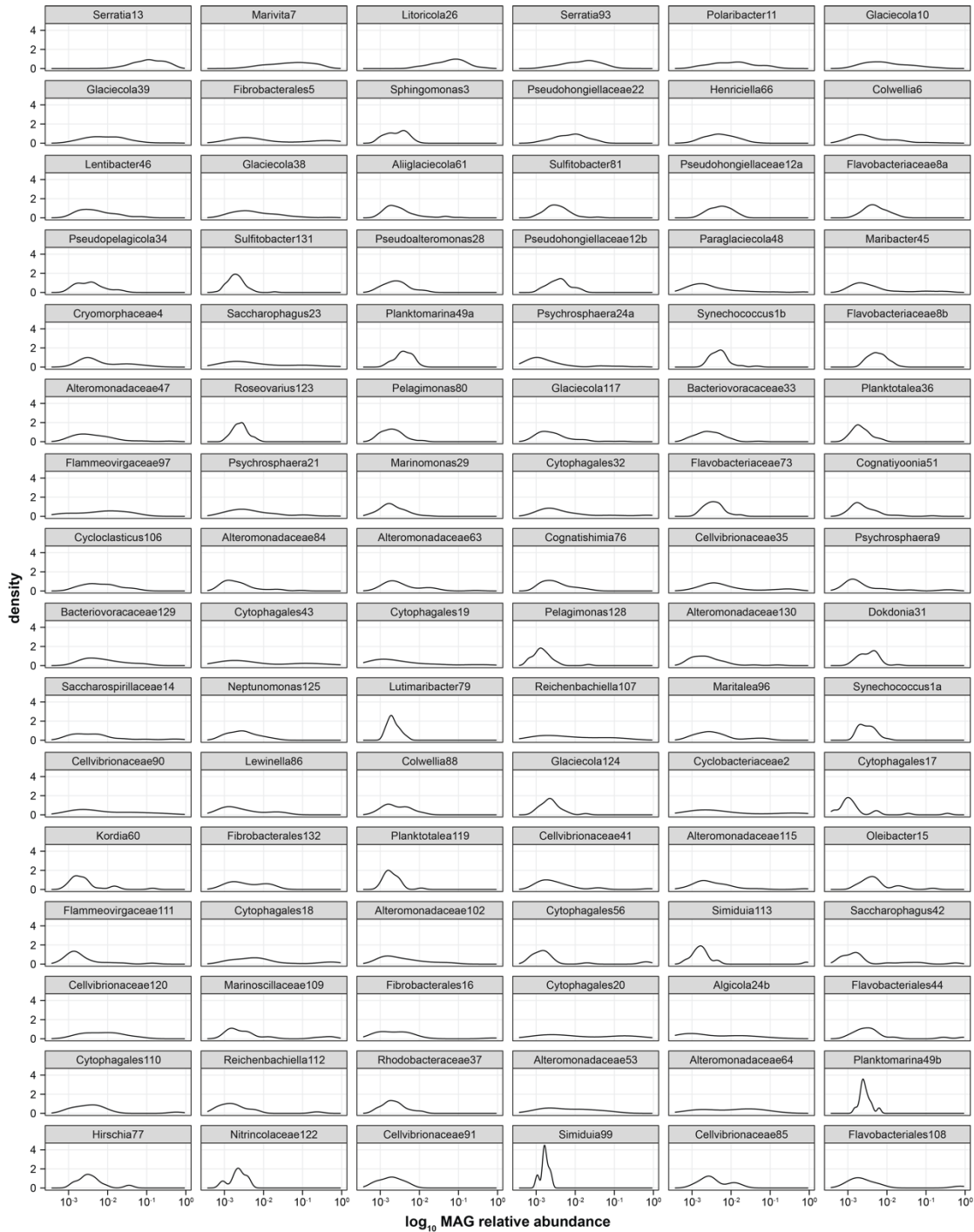
1276 **Fig. S1. Taxonomic abundances averaged across particles do not represent the**

1277 **compositions of communities on individual particles. (a)** Extended version of Fig. 2. Relative

1278 abundances of metagenome-assembled genomes (MAGs; $n = 120$) across late-stage particles.

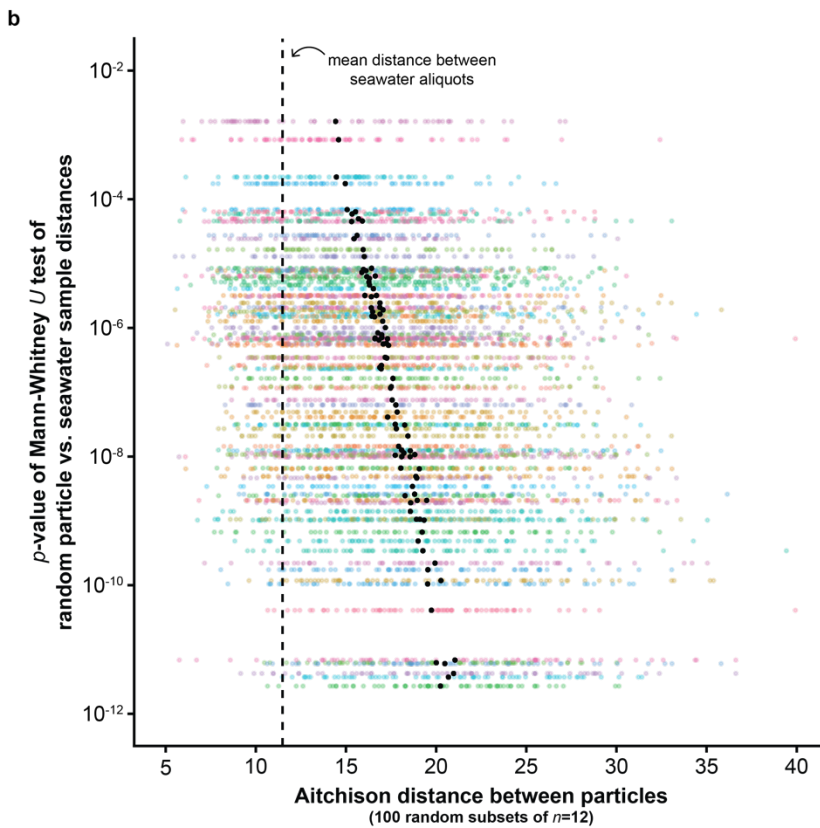
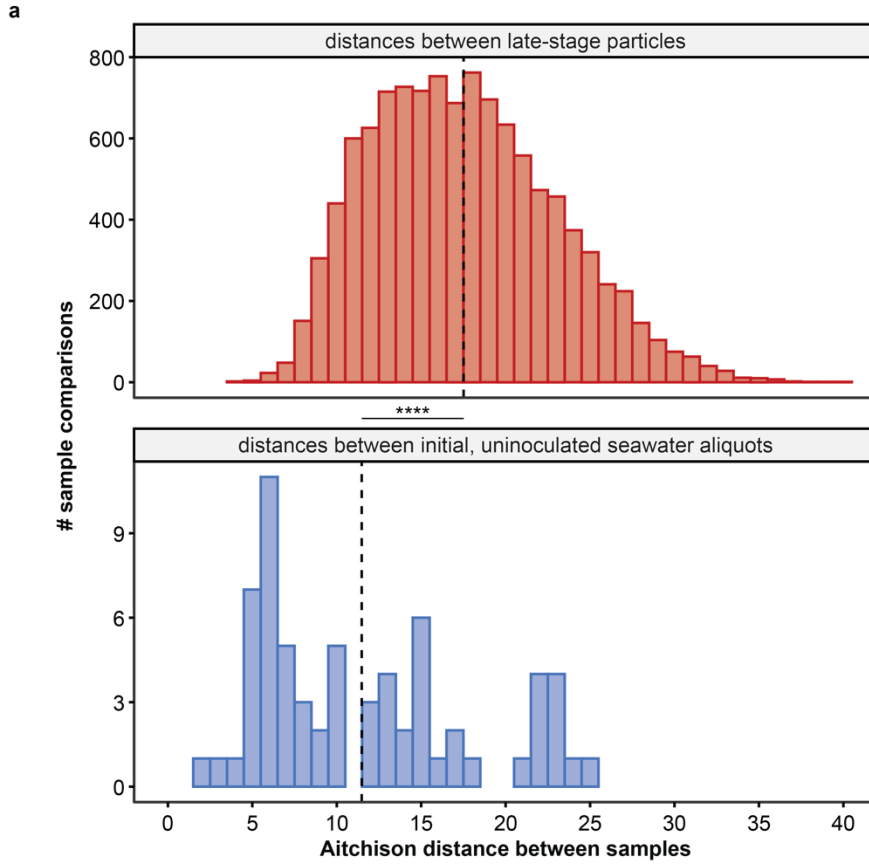
1279 Smaller black dots indicate the relative abundance of each MAG per particle ($n = 149$). Larger black

1280 circles indicate the \log_{10} [mean relative abundance] across the particles on which the MAG was
1281 found. Larger red circles indicate the \log_{10} [mean relative abundance] across all the particles (i.e.
1282 the MAG abundances for a theoretical bulk measurement aggregating all particles). MAGs are
1283 sorted from left to right by their prevalence across particles (i.e. the number of particles on which
1284 they are detected). The bars above show the average number of chitinases encoded in each cluster
1285 of highly similar MAGs (see Methods). The annotations below show each MAG's taxonomic ID
1286 (matching Table S1); predicted ecological role (heatmap: blue = degrader, green = exploiter, yellow
1287 = scavenger); jackpot score (heatmap: white = low, black = high); and taxonomic order (E =
1288 Enterobacterales, R = Rhodobacterales, P = Pseudomonadales, F = Flavobacteriales, C =
1289 Cytophagales, O = Other). **(b)** Histogram of the percent variance explained in the abundance ranks
1290 of MAGs on each late-stage particle by the abundance ranks for a theoretical bulk measurement
1291 aggregating all particles (which is equivalent to the average abundance across all particles). If the
1292 MAG abundance rank of a single particle's community matched that of the theoretical average, the
1293 percent variance explained would be 100% (right dashed line); however, the ensemble scale
1294 explained only an average of 11.7% (left dashed line) of the variance in abundance ranks at the
1295 single particle level.

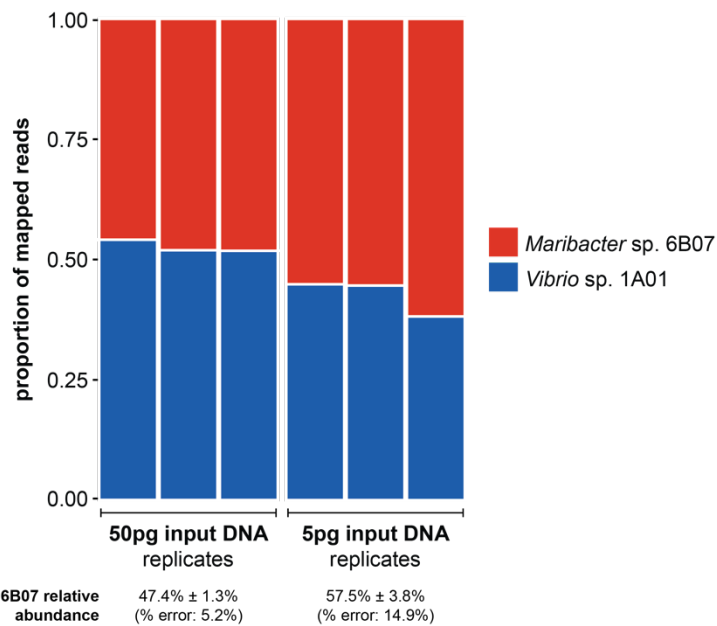


1296

1297 **Fig. S2. Distributions of MAG relative abundances on late-stage particles are approximately**
1298 **lognormal and right-skewed (i.e. towards high frequencies).** Distributions are shown as
1299 Gaussian kernel density estimates for MAGs present on at least 10 late-stage particles. The area
1300 under each curve equals one.

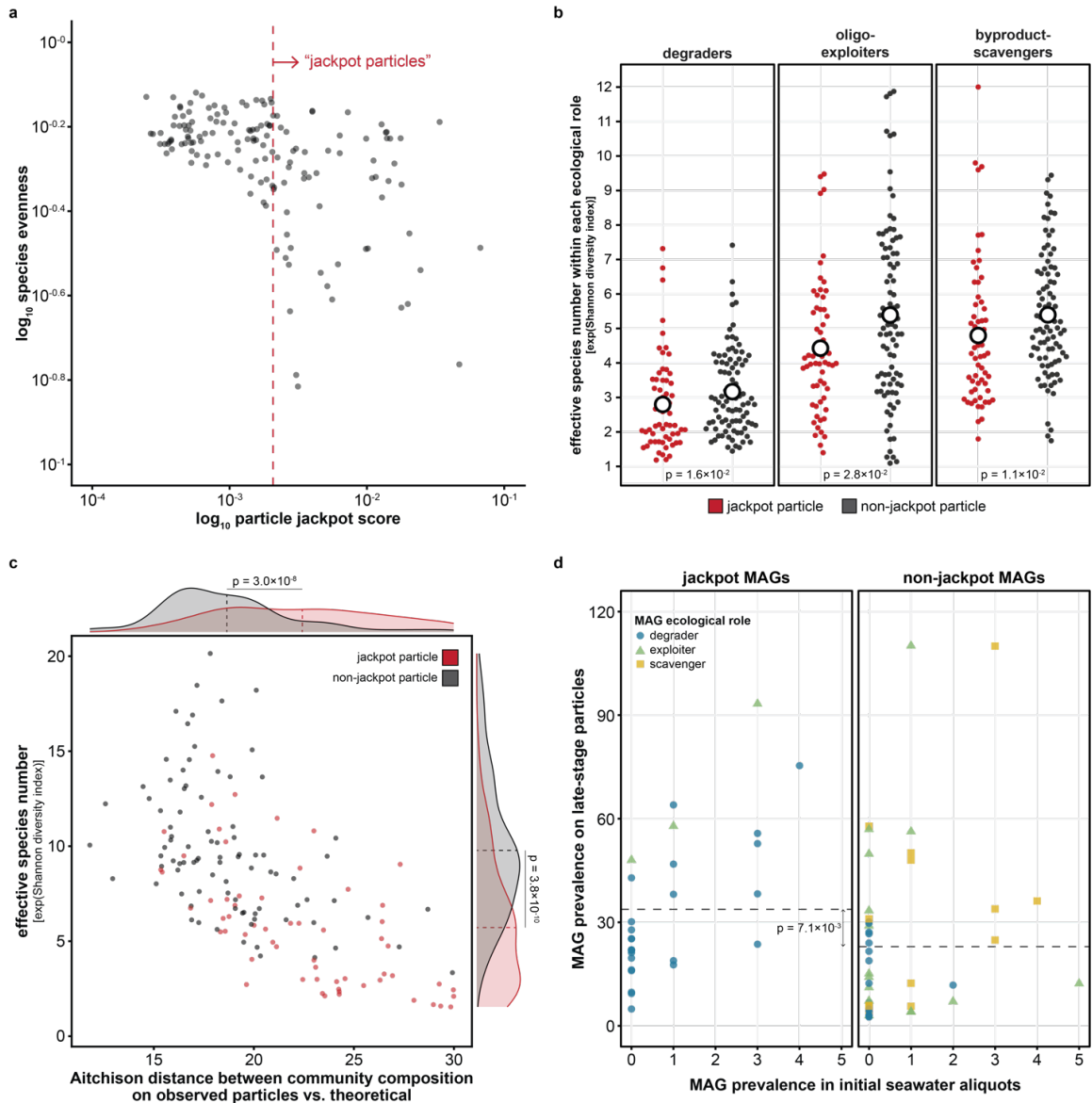


1302 **Fig. S3. Taxonomic variability in the initial seawater does not significantly account for**
1303 **variability observed across late-stage particles. (a)** Distributions of the Aitchison distances
1304 (Methods) calculated between all pairs of communities on late-stage particles ($n = 149$, red
1305 histogram) and between all pairs of aliquots of unincubated, initial seawater ($n = 12$, blue
1306 histogram). Dashed vertical lines represent the means of each distribution. Late-stage particles
1307 were significantly more dissimilar from one another than initial seawater samples (Mann-Whitney
1308 U test: $p = 1.3 \times 10^{-13}$). **(b)** The amount of inter-sample variability detected could depend on sample
1309 size, and many more pairs of particles than pairs of seawater samples were assessed in **(a)**.
1310 Therefore, we calculated the Aitchison distances between random subsets of 12 late-stage
1311 particles and compared those distributions to that of the seawater samples. Small points represent
1312 inter-particle Aitchison distances calculated for 100 random subsets (each with its own point color),
1313 and black dots indicate the mean value for each subset. The inter-particle distances for each subset
1314 are plotted against the p -value from a Mann-Whitney U test comparing the particle and seawater
1315 distributions. The dashed vertical line indicates the mean Aitchison distance between seawater
1316 samples (the same value as in the blue histogram in **(a)**). For all particle subsets, inter-particle
1317 distances were significantly higher than inter-seawater distances.



1318

1319 **Fig. S4. Mock communities sequenced with same protocols as particle-attached**
1320 **communities show relatively little deviation from expected strain abundances.** See Methods
1321 for details on the preparation of the mock communities, which contained equal proportions of
1322 *Marinobacter* sp. 6B07 genomic DNA and *Vibrio* sp. 1A01 genomic DNA. Relative abundances
1323 estimated from metagenomic libraries prepared using 50 pg of input DNA showed 5.2% error,
1324 whereas libraries prepared with 5 pg of input DNA showed 14.9% error.



1325

1326 **Fig. S5. Communities on jackpot particles are dominated by globally rare and locally**

1327 **abundant strains. (a)** Particles were defined as “jackpot particles” if they had high jackpot scores

1328 (indicating high relative abundances of jackpot taxa; see Methods for details). Each dot represents

1329 one late-stage particle ($n = 149$), and the red dashed line indicates the particle jackpot score

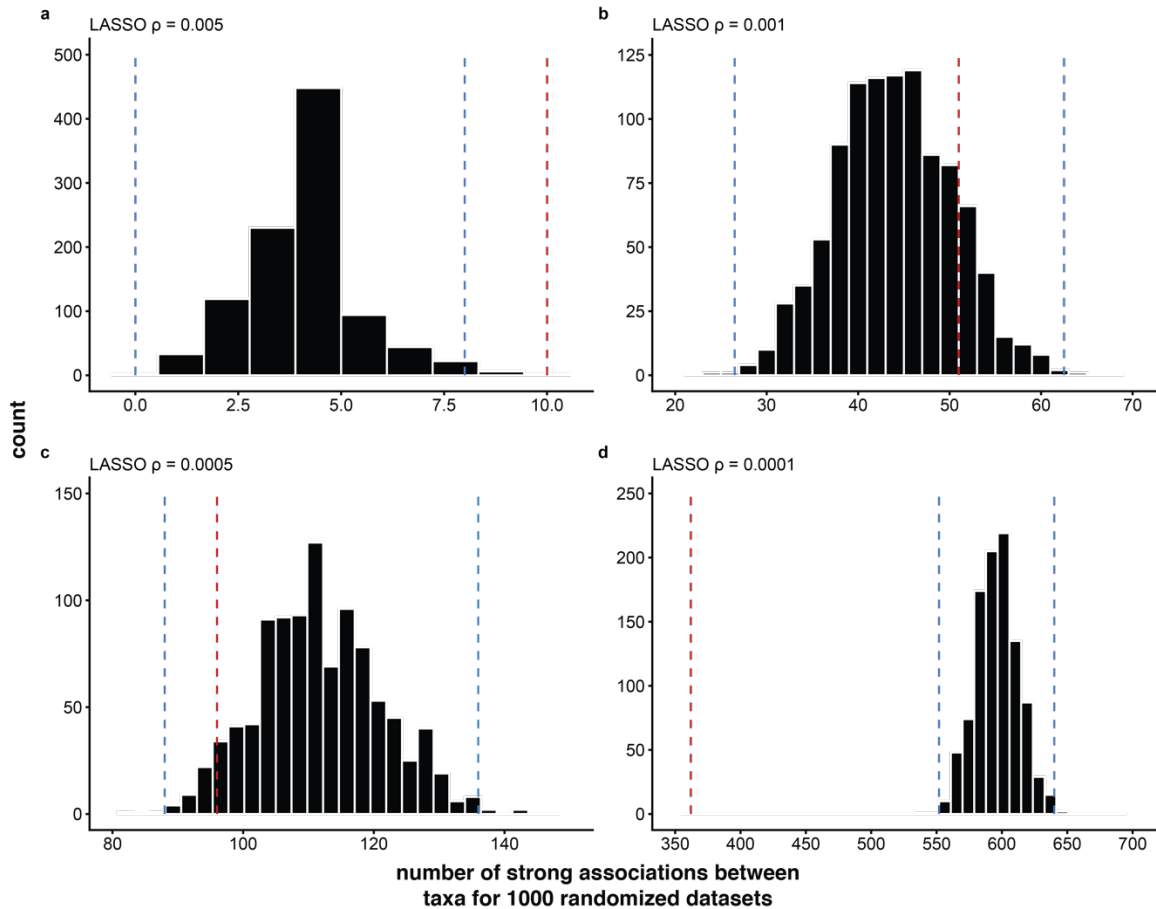
1330 threshold above which log-transformed values of Pielou’s species evenness drop sharply. **(b)** The

1331 effective species numbers (calculated from the Shannon diversity index) within each ecological

1332 role. Each smaller dot represents a late-stage particle ($n = 149$), and dot color indicates whether

1333 the particle was a jackpot particle (red) or a non-jackpot particle (black). Larger white dots represent

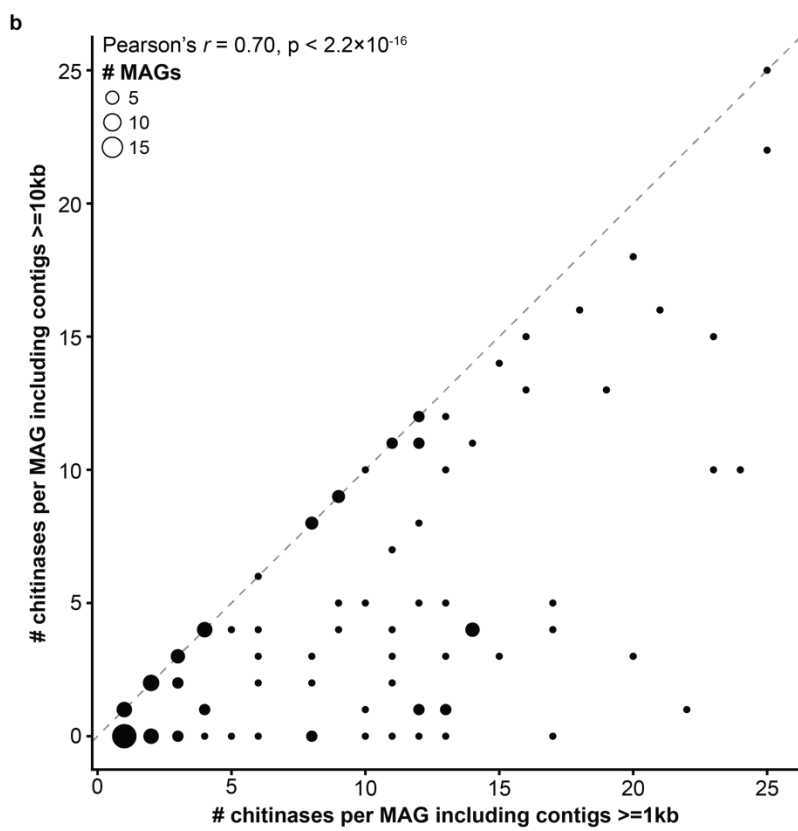
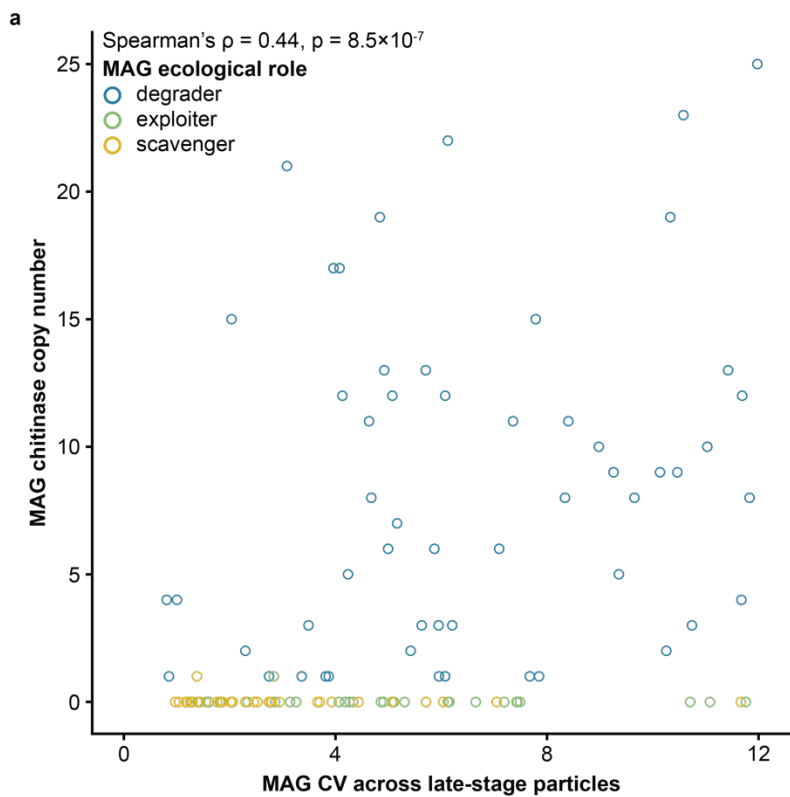
1334 the mean of each distribution. The diversity on jackpot particles was significantly lower than on non-
1335 jackpot particles for each of the roles (Mann-Whitney U test: degraders $p = 1.6 \times 10^{-2}$, exploiters $p =$
1336 2.8×10^{-2} , scavengers $p = 1.1 \times 10^{-2}$). **(c)** Community diversity (represented as effective species
1337 number, calculated from the Shannon diversity index) was inversely correlated (Spearman's $\rho = -$
1338 0.68 , $p < 2.2 \times 10^{-16}$) with the Aitchison distance between the community composition observed on
1339 each particle and the composition of the theoretical average particle (see larger red circles in Fig.
1340 S1a). Each dot represents a late-stage particle ($n = 149$), and dot color indicates whether the
1341 particle was a jackpot particle (red) or a non-jackpot particle (black; see Methods for definitions).
1342 Jackpot particle communities were significantly less diverse (Mann-Whitney U test: $p = 3.8 \times 10^{-10}$)
1343 and more divergent from the theoretical average particle (Mann-Whitney U test: $p = 3.0 \times 10^{-8}$) than
1344 non-jackpot particles. **(d)** Jackpot taxa (left panel) were significantly more prevalent across late-
1345 stage particles (Mann-Whitney U test: $p = 7.1 \times 10^{-3}$) than non-jackpot taxa (right panel) that were
1346 equally rare across aliquots of the initial, unincubated seawater (prevalence in seawater samples,
1347 Mann-Whitney U test: $p = 0.33$; mean abundance in seawater samples, Mann-Whitney U test: $p =$
1348 0.49). Each point represents a MAG that was detected on fewer than half of the seawater aliquots,
1349 with the point color and shape indicating its predicted ecological role (blue circle = degrader, green
1350 triangle = exploiter, yellow square = scavenger).



1351

1352 **Fig. S6. Late-stage particles exhibit little specific taxonomic structure.** For the observed data,
1353 as well as for 1000 randomizations of the data, we inferred the number of conditional dependencies
1354 between taxa from the estimated inverse covariance matrix of center log-ratio-transformed relative
1355 abundances across late-stage particles. The inverse covariance matrices were calculated using a
1356 graphical lasso approach with a regularization parameter of (a) $\rho = 0.005$, (b) $\rho = 0.001$, (c) $\rho =$
1357 0.0005 , or (d) $\rho = 0.0001$. Each plot shows the distribution of the number of conditional
1358 dependencies (with strengths ≥ 0.2 or ≤ -0.2) between MAGs inferred for the randomizations of the
1359 data. The red lines indicate the number of conditional dependencies (with strengths ≥ 0.2 or ≤ -0.2)
1360 inferred from the observed data. The blue lines indicate thresholds beyond which values are
1361 considered outliers relative to the distribution calculated for the randomized datasets (using the
1362 interquartile range [IQR] method – left lines indicate the value of $[Q_1 - 1.5 \times \text{IQR}]$, and right lines
1363 indicate the value of $[Q_3 + 1.5 \times \text{IQR}]$). The choice of the regularization parameter in the analyses

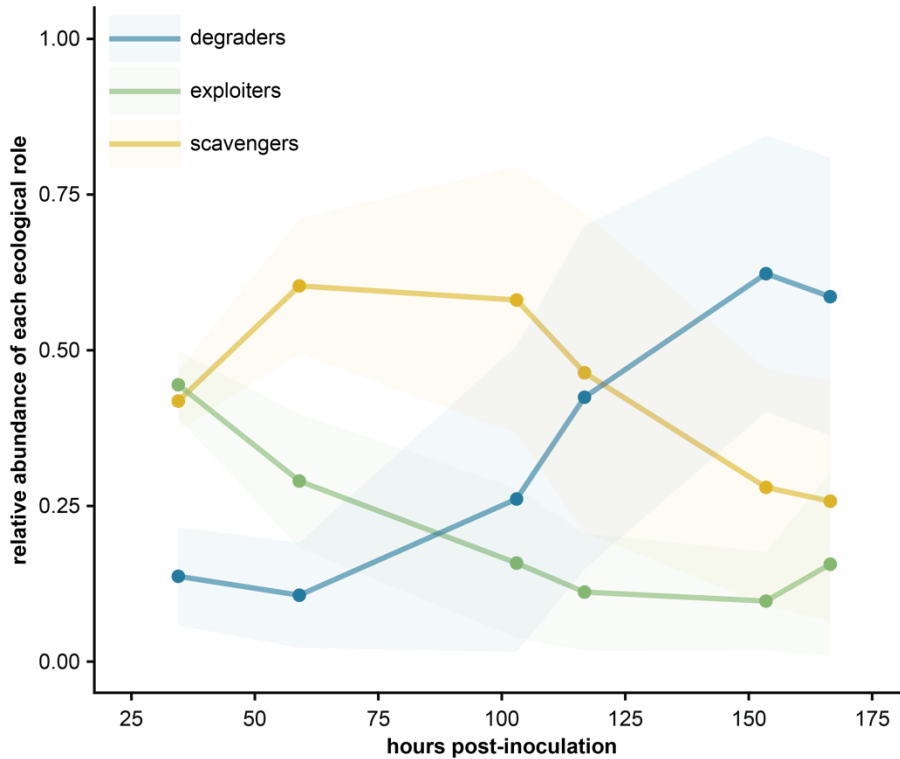
1364 used to estimate the number conditional dependencies between taxa resulted in more strong
1365 associations being inferred for the observed communities than the randomized ones only when a
1366 trivially small number of associations were inferred (panel **a**). Therefore, in terms of the number of
1367 strain-specific associations, the observed particles were either indistinguishable from, or less
1368 structured than, random communities.



1369

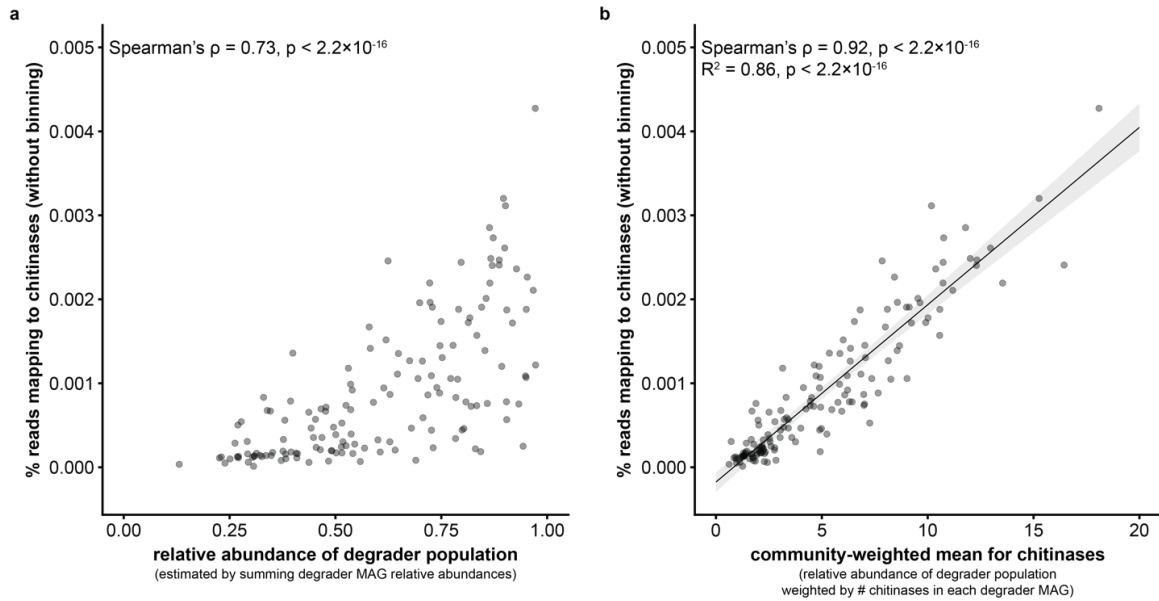
1370 **Fig. S7. Highly variable taxa are often degrader strains encoding many chitinase genes. (a)**

1371 There was a significant positive correlation between the coefficients of variation of MAG relative
1372 abundances across late-stage particles and the number of chitinase genes encoded in MAGs
1373 (Spearman's $\rho = 0.44$, $p = 8.5 \times 10^{-7}$; calculated for 120 MAGs across 149 particles). Each open dot
1374 represents a MAG, with the color indicating its predicted ecological role (blue = degrader, green =
1375 exploiter, yellow = scavenger). **(b)** Comparison of the number of chitinase genes encoded by each
1376 MAG when considering contigs ≥ 10 kb (which are binned more reliably than shorter contigs) vs.
1377 considering contigs ≥ 1 kb (the minimum length of binned contigs). There was a strong correlation
1378 between chitinase copy numbers when considering contigs ≥ 10 kb vs. contigs ≥ 1 kb (Pearson's r
1379 = 0.70, $p < 2.2 \times 10^{-16}$), lending confidence to estimates of high chitinase copy numbers in certain
1380 bins. Dot sizes indicate the number of MAGs at each coordinate.



1381

1382 **Fig. S8. The proportion of predicted degraders on particles increases and becomes more**
1383 **variable over time.** The relative abundances of the predicted ecological roles (degrader = blue,
1384 exploiter = green, scavenger = yellow) on particles harvested after varying incubation durations
1385 were calculated by summing the relative abundances of MAGs classified into each role. Points
1386 indicate mean values across particles at each time point, with shading representing ± 1 standard
1387 deviation. The number of particles considered at each time point were (in order from early to late)
1388 88, 88, 80, 90, 76, and 73.

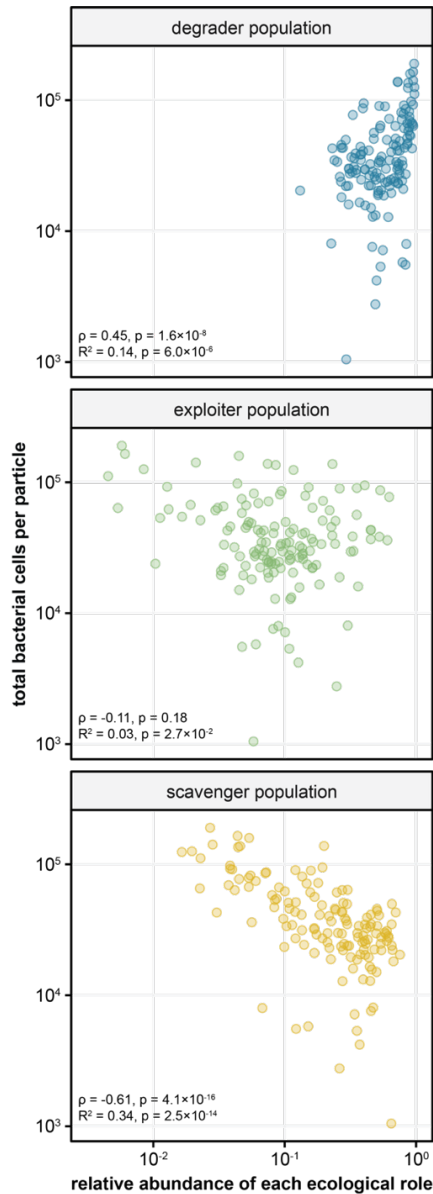


1389

1390 **Fig. S9. The wide range of degrader population relative abundances estimated for late-stage**
1391 **particles holds when genes are used as read mapping references rather than MAGs.**

1392 Conceivably, the use of MAGs as reference genomes could have biased our estimate of the
1393 degrader population abundance; therefore, we also mapped reads to a reference set of chitinase
1394 genes (Methods). **(a)** There was a strong correlation (Spearman's $\rho = 0.73$, $p < 2.2 \times 10^{-16}$) between
1395 the degrader population relative abundances estimated using MAGs and the percent of reads that
1396 mapped to chitinase genes. Each dot represents one late-stage particle ($n = 149$). **(b)** If the
1397 degrader population relative abundance estimated by MAGs were a consistent approximation of
1398 the true degrader population abundance, then the wide range in the number of chitinases encoded
1399 in each degrader MAG (Fig. S6, Table S1) would be reflected in the percent of reads in each
1400 community mapping to chitinase genes. Therefore, we weighted the degrader population relative
1401 abundances by the number of chitinases in each MAG to calculate the community-weighted mean
1402 for chitinases of each late-stage particle (Methods). There was an even stronger correlation
1403 (Spearman's $\rho = 0.92$, $p < 2.2 \times 10^{-16}$) between the chitinase community-weighted mean estimated
1404 using MAGs and the percent of reads that mapped to chitinase genes. Each dot represents one
1405 late-stage particle ($n = 149$). The black line represents the linear regression line ($R^2 = 0.86$, $p <$
1406 2.2×10^{-16} ; shading indicates the 99% confidence interval). Therefore, reference MAGs captured a

1407 representative subsample of the degraders within particle-attached communities, and predictions
1408 of chitinolytic potential were consistent between MAG- and gene-based approaches.



1409

1410 **Fig. S10: The overall yield of late-stage particles is correlated with community composition.**

1411 There was a strong negative correlation between the proportion of scavengers and the number of

1412 bacterial cells in late-stage communities, estimated through qPCR (yellow dots, $n = 142$;

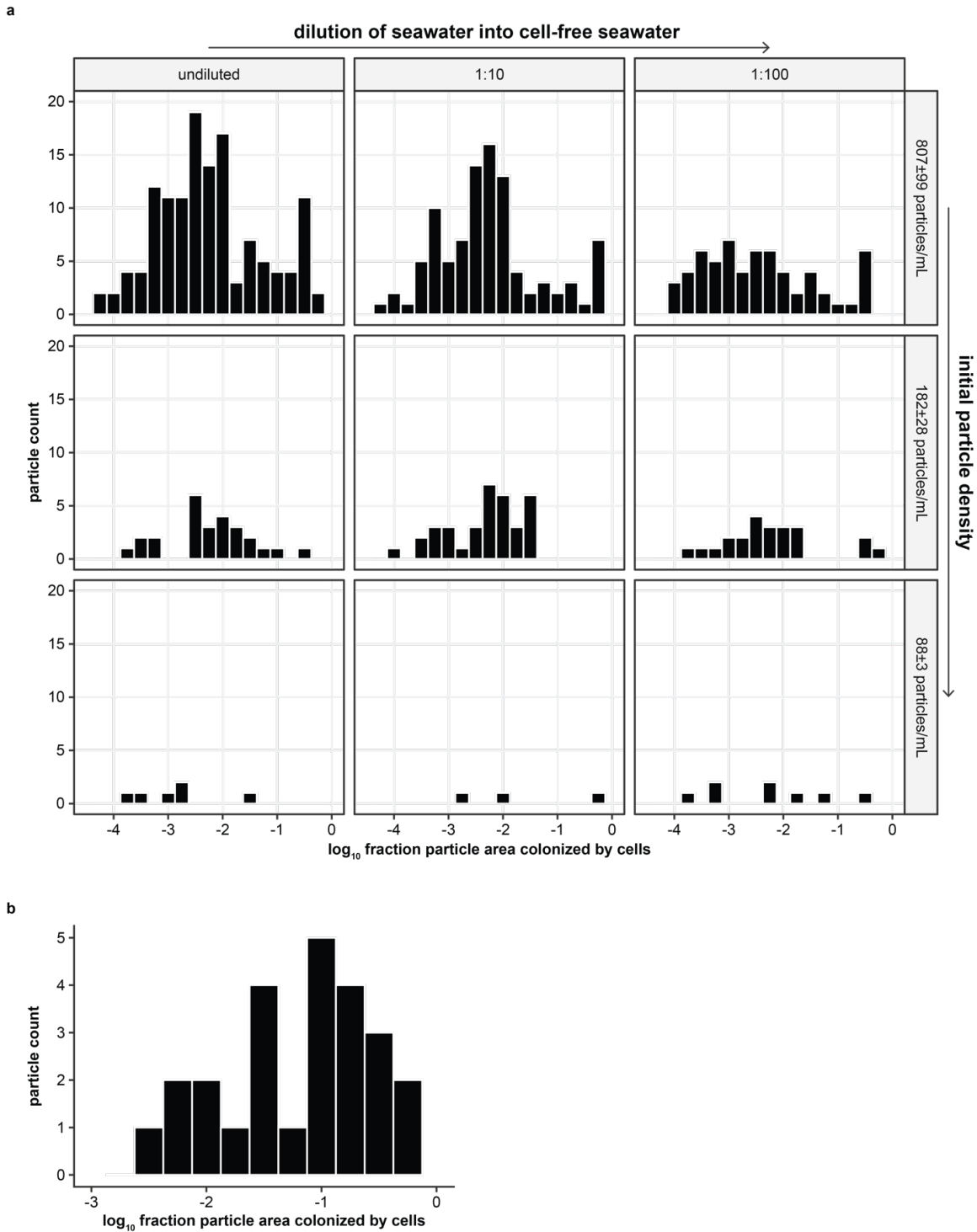
1413 Spearman's $\rho = -0.61, p = 4.1 \times 10^{-16}$; log-log linear regression: $R^2 = 0.34, p = 2.5 \times 10^{-14}$). There was

1414 a less strong, though still highly significant, positive correlation between biomass and the proportion

1415 of degraders (blue dots; Spearman's $\rho = 0.45, p = 1.6 \times 10^{-8}$; log-log linear regression: $R^2 = 0.14, p$

1416 $= 6.0 \times 10^{-6}$), and there was no correlation with the exploiter population (green dots; Spearman's ρ

1417 $= -0.11, p = 0.18$; log-log linear regression: $R^2 = 0.03, p = 2.7 \times 10^{-2}$).



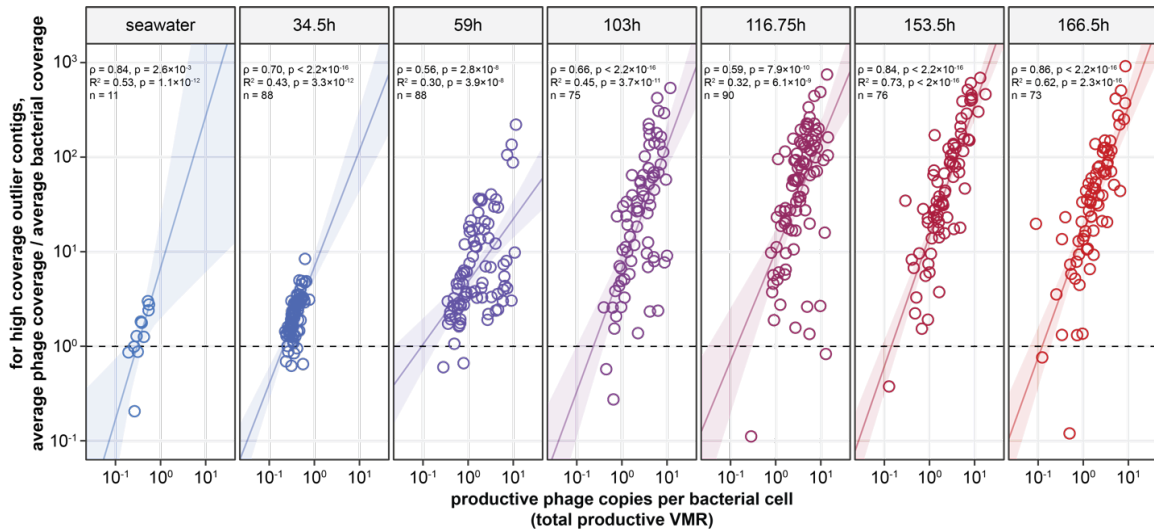
1418

1419 **Fig. S11. The cell counts on single particles incubated in the same volume of seawater span**

1420 **several orders of magnitude, matching the range estimated for single particles incubated**

1421 **separately. All plots show the distributions of the proportion of a particle's area occupied by cells**

1422 (transformed on a \log_{10} scale), estimated by visualizing particles stained with the DNA intercalating
1423 dye SYTO 9. **(a)** Cell count distributions for communities on single particles incubated in seawater
1424 together for 24 hours at various initial particle concentrations (top row: 807 ± 99 particles/mL; middle
1425 row: 182 ± 28 particles/mL; bottom row: 88 ± 3 particles/mL; , \pm indicates 1 standard deviation for $n =$
1426 3 replicates throughout) and at various initial cell concentrations (left column: undiluted natural
1427 seawater; middle column: seawater inoculum diluted 1:10 into $0.2\mu\text{m}$ -filtered natural seawater; right
1428 column: diluted 1:100). **(b)** Cell count distributions for communities on single particles incubated in
1429 seawater separately and harvested after 154-167 hours (i.e. late-stage communities).



1430

1431 **Fig. S12. Phage-derived contigs that are coverage outliers have much higher average read**

1432 **coverage than bacteria-derived contigs that are coverage outliers.** Given that read coverage

1433 values from metagenomic data are often noisy, it is conceivable that productive phage contigs had

1434 unusually high coverage simply due to sequencing noise. However, for contigs that were high

1435 coverage outliers, the ratio for each particle (open dots) of the average coverage of phage contigs

1436 to the average coverage of bacterial contigs was often much greater than 1 (the black horizontal

1437 dashed line). Notably, these coverage ratios were overall lowest in the initial seawater inocula and

1438 rose during the incubation period, coinciding with the timescale of increasing mean productive

1439 VMRs (Fig. 4c). This indicates that phage contigs with high coverage values represented phages

1440 that were replicating more than their bacterial hosts, rather than representing contigs with randomly

1441 higher coverage values. Furthermore, there were strong positive relationships between this

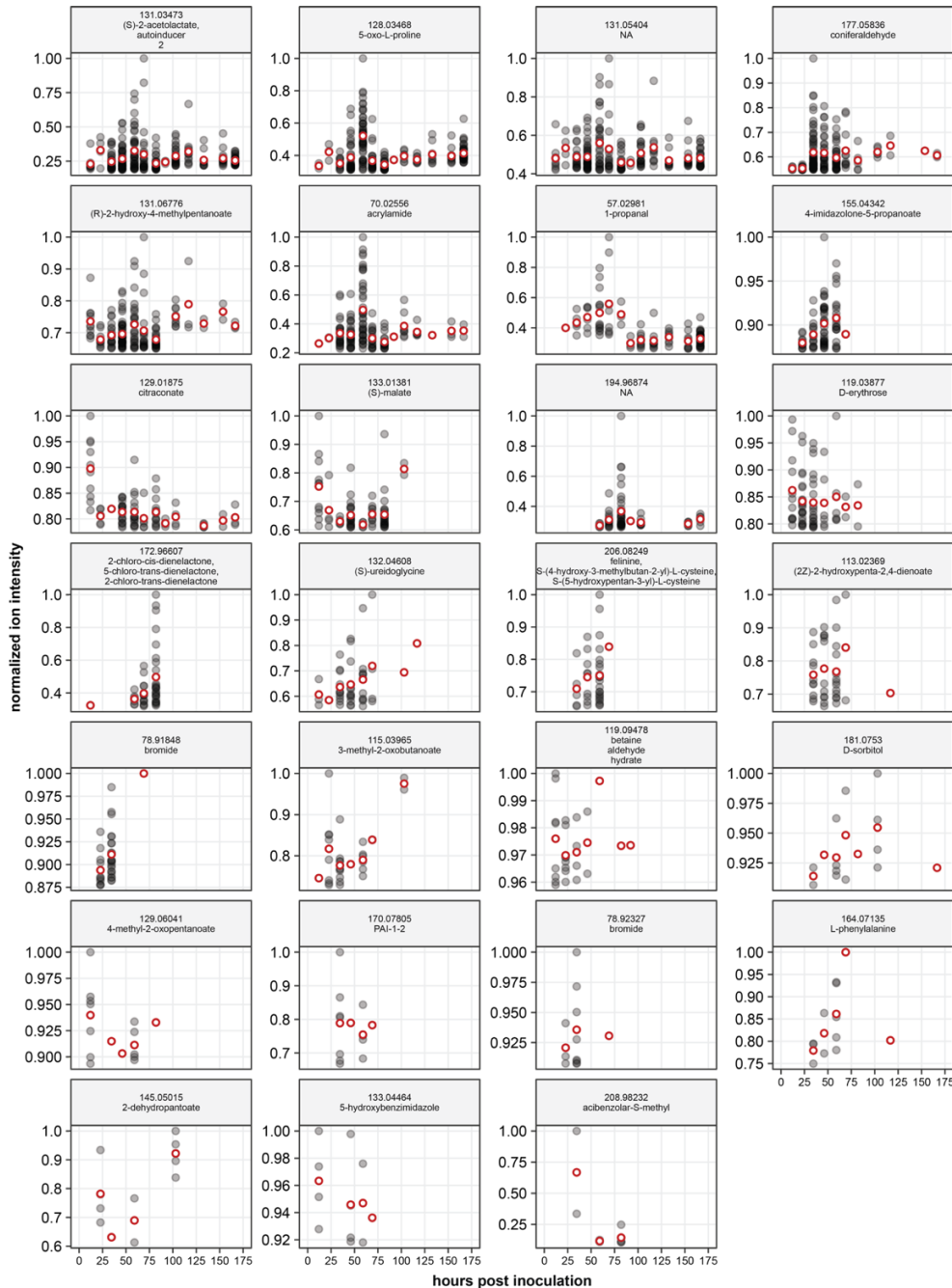
1442 coverage ratio and the total productive VMR at each time point (see each subplot for significance

1443 values; the solid lines represent the log-log linear regression lines, and shading indicates the 95%

1444 confidence intervals). Thus, as expected, the particles on which high outlier phage contig coverage

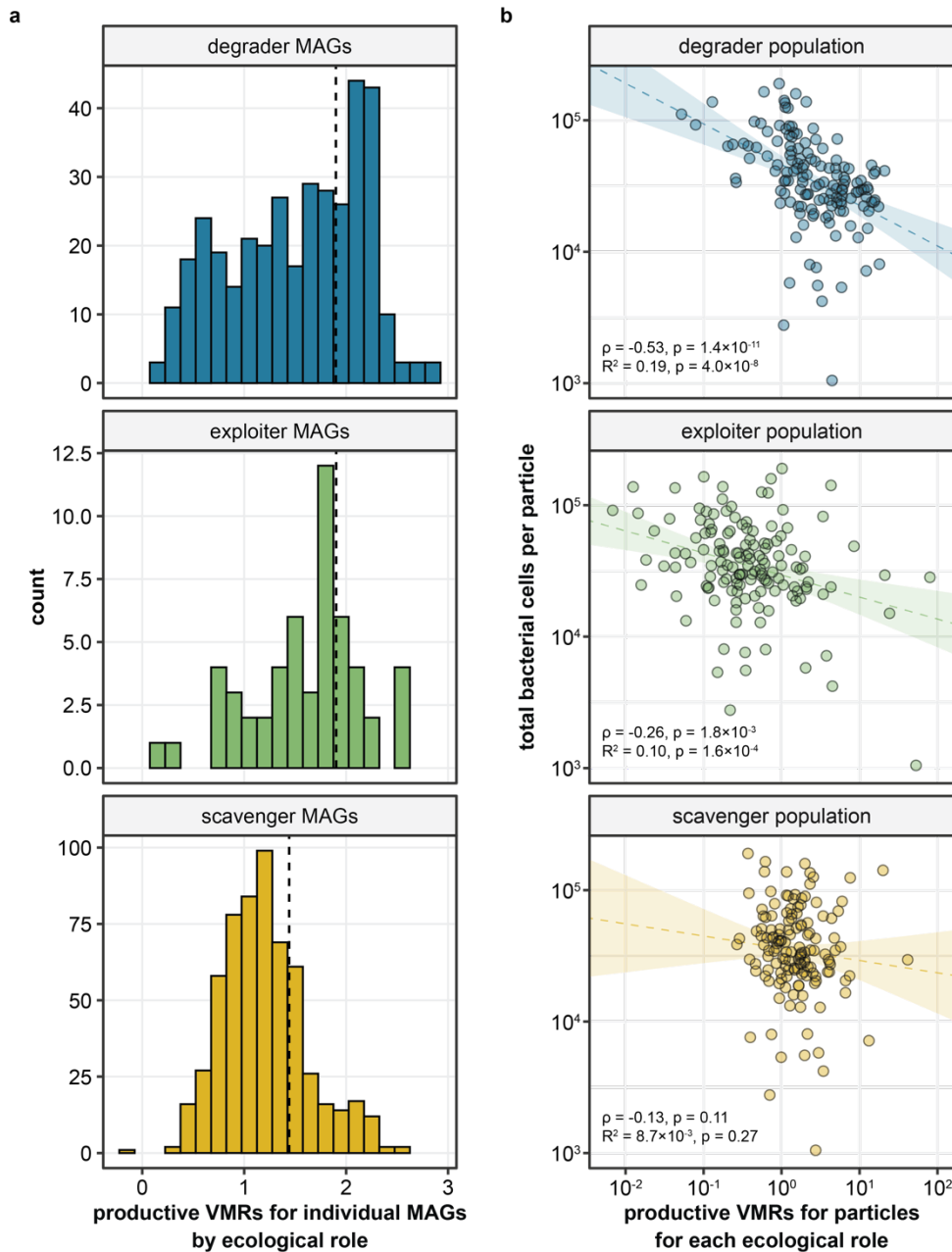
1445 was indistinguishable from high outlier bacterial contig coverage were mostly those with low

1446 productive VMRs.



1447

1448 **Fig. S13. Normalized intensities of individual ions over time.** The normalized intensities of ions
 1449 across particles harvested after varying incubation durations are shown for ions that were
 1450 significantly enriched (relative to the initial seawater) on at least 10 particles (see Methods). Gray
 1451 dots indicate measurements for individual particles, and red circles represent the mean normalized
 1452 intensities at each time point. Panel labels include the *m/z* ratio and predicted annotation for each
 1453 ion (see Methods).



1454

1455 **Fig. S14. The degrader population contributes significantly to particle-level productive**

1456 **VMRs. (a)** Degrader MAGs (top panel, blue distribution) and exploiter MAGs (middle panel, green

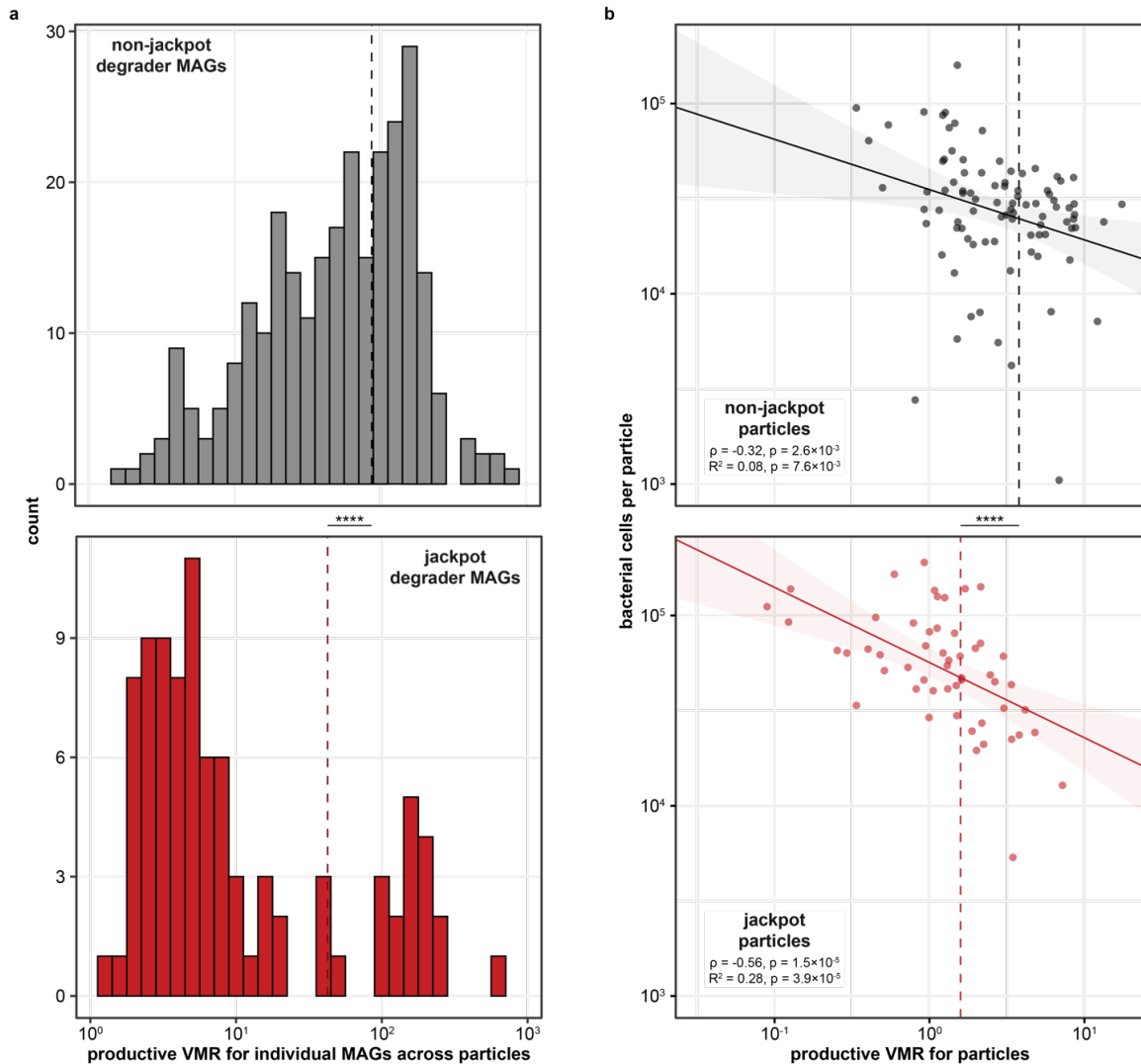
1457 distribution) had significantly higher productive VMRs across late-stage particles than scavenger

1458 MAGs (bottom panel, yellow distribution) when phages were productive. Distributions depict the

1459 non-zero VMRs for each group of MAGs (dashed lines represent the means of each distribution,

1460 with degraders having a mean VMR of 76.8, exploiters 78.6, and scavengers 27.3; one-way

1461 ANOVA: $F(998,2) = 47.4$, $p = 2.1 \times 10^{-20}$; Tukey's HSD test: degrader-exploiter $p = 0.40$; degrader-
1462 scavenger $p < 1.0 \times 10^{-7}$; exploiter-scavenger $p = 1.0 \times 10^{-7}$). When instances of VMRs equaling zero
1463 are included in the distributions, degraders had a mean VMR of 5.9, exploiters 1.3, and scavengers
1464 2.6 (one-way ANOVA: $F(14049,2) = 42.0$, $p = 6.3 \times 10^{-19}$; Tukey's HSD test: degrader-exploiter $p <$
1465 1.0×10^{-7} ; degrader-scavenger $p < 1.0 \times 10^{-7}$; exploiter-scavenger $p = 3.5 \times 10^{-2}$). This suggests that
1466 degraders overall experienced the most phage activation. **(b)** Absolute bacterial cell counts on late-
1467 stage particles ($n = 142$), estimated through qPCR, vs. each particle's productive VMR for the
1468 MAGs in each ecological role. Cell counts were negatively correlated with productive VMRs most
1469 strongly and significantly for degraders (top panel, blue dots; Spearman's $\rho = -0.53$, $p = 1.4 \times 10^{-11}$)
1470 and less so for exploiters (middle panel, green dots; Spearman's $\rho = -0.26$, $p = 1.8 \times 10^{-3}$), and there
1471 was no correlation between cell counts and productive VMRs for scavengers (bottom panel, yellow
1472 dots; Spearman's $\rho = -0.13$, $p = 0.11$). Dashed lines represent the log-log linear regression lines
1473 between cell counts and productive VMR (degraders: $R^2 = 0.19$, $p = 4.0 \times 10^{-8}$; exploiters: $R^2 = 0.10$,
1474 $p = 1.6 \times 10^{-4}$; scavengers: $R^2 = 8.7 \times 10^{-3}$, $p = 0.27$; shading indicates the 95% confidence intervals).
1475 The productive VMRs for each ecological role were also significantly different from each other, with
1476 degraders having the highest mean VMR (one-way ANOVA: $F(423,2) = 96.6$, $p = 2.9 \times 10^{-35}$; Tukey's
1477 HSD test: degrader-exploiter $p < 1.0 \times 10^{-7}$; degrader-scavenger $p = 1.7 \times 10^{-3}$; exploiter-scavenger
1478 $p = < 1.0 \times 10^{-7}$). This suggests that the effect of phage activation on particle yield was largely driven
1479 by the degrader trophic level.



1480

1481 **Fig. S15. The jackpot growth phenomenon is associated with less phage activation. (a)**

1482 Jackpot degrader MAGs (bottom panel, red distribution) had lower productive VMRs across late-

1483 stage particles than non-jackpot degraders (top panel, grey distribution). Distributions depict the

1484 non-zero VMRs for each group of MAGs (dashed lines represent the means of each distribution,

1485 with jackpot degraders having a mean VMR of 42.4 and non-jackpot degraders having a mean of

1486 88.0; Mann-Whitney U test: $p = 6.2 \times 10^{-15}$). When instances of VMRs equaling zero are included in

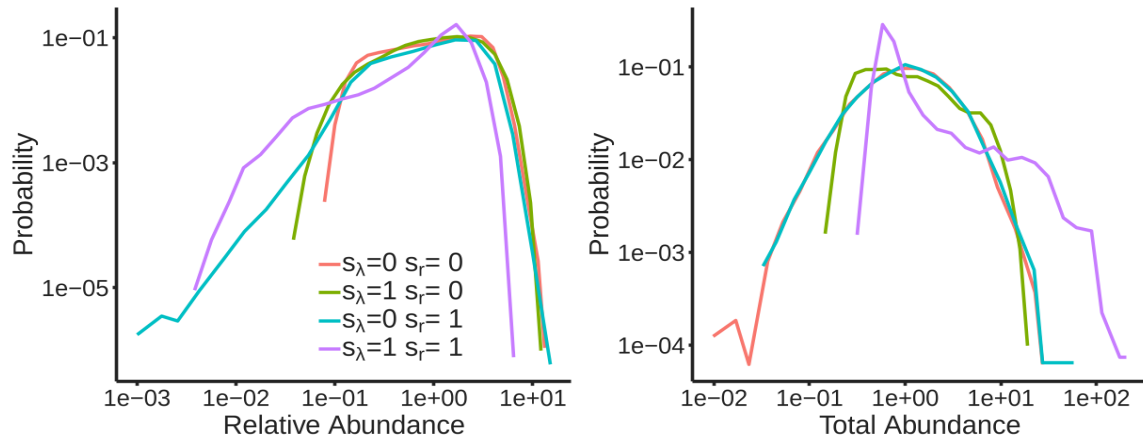
1487 the distributions, jackpot degraders still have a lower mean VMR (1.32 vs. 12.9; Mann-Whitney U

1488 test: $p = 1.3 \times 10^{-49}$). **(b)** Jackpot particles had lower productive VMRs than non-jackpot particles.

1489 Modified version of Fig. 4e in which jackpot particles (bottom panel, red dots) are shown separately

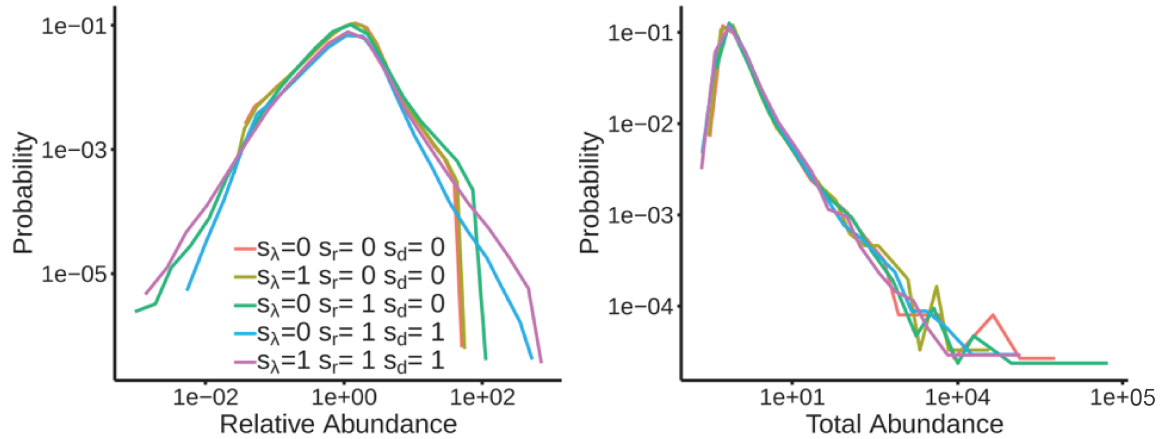
1490 from non-jackpot particles (top panel, dark grey dots). Jackpot particles had significantly lower

1491 productive VMRs than non-jackpot particles (dashed lines represent the means of each distribution;
1492 Mann-Whitney U test: $p = 4.3 \times 10^{-8}$), even controlling for differences in biomass between these
1493 groups of particles (ANCOVA: $F(1,139) = 16.92$, $p = 4.1 \times 10^{-4}$, partial $\eta^2 = 0.09$). Both groups of
1494 particles showed significant negative relationships between biomass (estimated through qPCR)
1495 and productive VMR (jackpot particles: Spearman's $\rho = -0.56$, $p = 1.5 \times 10^{-5}$; non-jackpot particles:
1496 Spearman's $\rho = -0.32$, $p = 2.6 \times 10^{-3}$). The solid lines represent the log-log linear regression lines
1497 between cell counts and productive VMRs (red line for jackpot particles: $R^2 = 0.28$, $p = 3.9 \times 10^{-5}$;
1498 black line for non-jackpot particles: $R^2 = 0.08$, $p = 7.6 \times 10^{-3}$; shading indicates the 95% confidence
1499 intervals).



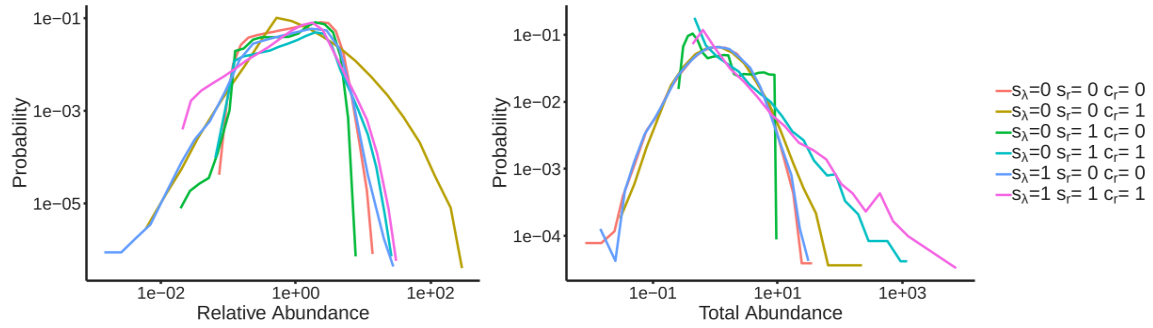
1500

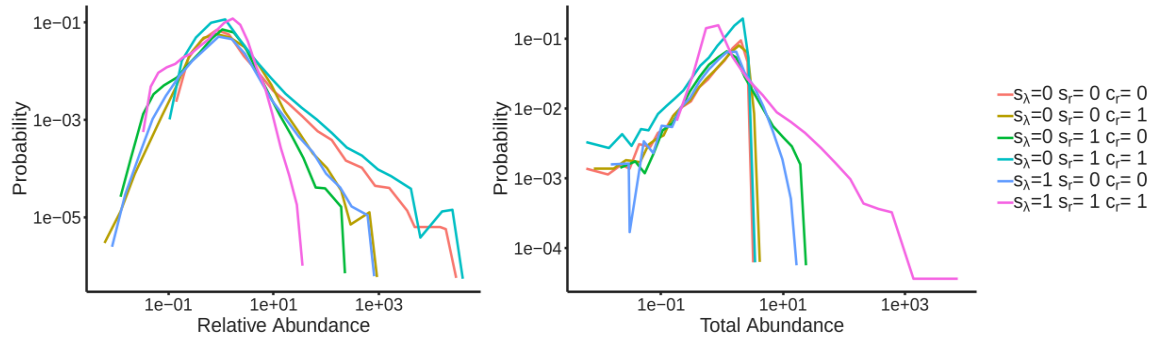
1501 **Fig. S16. Numerical simulations of mathematical model #1.** The left panel shows the distribution
1502 of rescaled relative abundances averaged over MAGs. For each MAG the logarithm of the relative
1503 abundances across particles was rescaled by mean and variance, so that it had mean zero and
1504 unit variance. Lines represent averages over MAGs. Colors refer to a particular parameterization.
1505 The right plot shows the distribution of the total biomass across particles. In all the simulations we
1506 set $\bar{\lambda} = 1$ and $\bar{r} = 1$ (where the total time of the experiment was also set to be equal to 1). Different
1507 values of s_r and s_λ correspond to different colors.



1508

1509 **Fig. S17. Numerical simulations of mathematical model #2.** The left panel shows the distribution
1510 of rescaled relative abundances averaged over MAGs. For each MAG the logarithm of the relative
1511 abundances across particles was rescaled by mean and variance, so that it had mean zero and
1512 unit variance. Lines represent averages over MAGs. Colors refer to a particular parameterization.
1513 The right plot shows the distribution of the total biomass across particles. In all the simulations we
1514 set $\bar{\lambda} = 1$ and $\bar{r} = 1$ (where the total time of the experiment was also set to be equal to 1). Different
1515 values of s_r and s_λ correspond to different colors.





1522 **Fig. S19. Numerical simulations of mathematical model #4.** The panels show the same
1523 distributions as in Fig. S16. The right plot shows the distribution of the total biomass across
1524 particles. In all the simulations we set $\lambda^d = 1$, $\bar{\lambda} = 1$, and $\bar{r} = 1$ (where the total time of the
1525 experiment was also set to be equal to 1). Different values of s_r , s_λ and c_r correspond to different
1526 colors.

1527 **Table S1 (separate file).** Metadata accompanying metagenome-assembled genomes (MAGs)
1528 from this study.

1529

1530 **Table S2 (separate file).** Metadata accompanying bacteriophage-annotated sequences from this
1531 study.

1532

1533 **Table S3 (separate file).** Metadata accompanying metabolomics performed in this study.

1534

1535 **Table S4 (separate file).** Statistics on the distributions of bacteriophage-annotated sequences in
1536 accompanying metagenome-assembled genomes (MAGs) from this study according to the
1537 predicted MAG ecological role.

1538

1539 **Table S5 (separate file).** Metadata accompanying metagenomic samples collected in this study.

1540

1541 **Table S6 (separate file).** Accession numbers and methods for the creation of custom profile hidden
1542 Markov models (HMMs) used to annotate chitin metabolism-related genes in metagenome-
1543 assembled genomes in this study.

1544

1545 **Table S7 (separate file).** Relative abundances of metagenome-assembled genomes (MAGs) in
1546 each sample collected in this study.

1547

1548

1549 **Supplementary References**

1550

- 1551 1. C. Rinke, *et al.*, Validation of picogram- and femtogram-input DNA libraries for microscale
1552 metagenomics. *PeerJ* 2016, 1–28 (2016).
- 1553 2. A. M. Bolger, M. Lohse, B. Usadel, Trimmomatic: A flexible trimmer for Illumina sequence
1554 data. *Bioinformatics* 30, 2114–2120 (2014).

- 1555 3. B. Bushnell, BMap: A Fast, Accurate, Splice-Aware Aligner (2014).
- 1556 4. S. I. Nikolenko, A. I. Korobeynikov, M. A. Alekseyev, BayesHammer: Bayesian clustering
1557 for error correction in single-cell sequencing. *Bmc Genomics* 14, S7 (2013).
- 1558 5. D. Li, *et al.*, MEGAHIT v1.0: A fast and scalable metagenome assembler driven by
1559 advanced methodologies and community practices. *Methods* 102, 3–11 (2016).
- 1560 6. Y. W. Wu, B. A. Simmons, S. W. Singer, MaxBin 2.0: An automated binning algorithm to
1561 recover genomes from multiple metagenomic datasets. *Bioinformatics* 32, 605–607 (2016).
- 1562 7. J. Alneberg, *et al.*, Binning metagenomic contigs by coverage and composition. *Nat*
1563 *Methods* 11, 1144–1146 (2014).
- 1564 8. B. Langmead, S. L. Salzberg, Fast gapped-read alignment with Bowtie 2. *Nat Methods* 9,
1565 357–359 (2012).
- 1566 9. G. E. Leventhal, *et al.*, Strain-level diversity drives alternative community types in
1567 millimetre-scale granular biofilms. *Nat Microbiol* 3, 1295–1303 (2018).
- 1568 10. C. M. K. Sieber, *et al.*, Recovery of genomes from metagenomes via a dereplication,
1569 aggregation and scoring strategy. *Nat Microbiol* 3, 836–843 (2018).
- 1570 11. D. H. Parks, M. Imelfort, C. T. Skennerton, P. Hugenholtz, G. W. Tyson, CheckM:
1571 Assessing the quality of microbial genomes recovered from isolates, single cells, and
1572 metagenomes. *Genome Res* 25, 1043–1055 (2015).
- 1573 12. D. H. Parks, *et al.*, A standardized bacterial taxonomy based on genome phylogeny
1574 substantially revises the tree of life. *Nat Biotechnol* 36, 996 (2018).
- 1575 13. P. A. Chaumeil, A. J. Mussig, P. Hugenholtz, D. H. Parks, GTDB-Tk: A toolkit to classify
1576 genomes with the genome taxonomy database. *Bioinformatics* 36, 1925–1927 (2020).
- 1577 14. W. J. Kent, BLAT---The BLAST-Like Alignment Tool. *Genome Res* 12, 656–664 (2002).
- 1578 15. N. Otsu, A threshold selection method from gray-level histograms. *Ieee Transactions Syst*
1579 *Man Cybern* 9, 62–66 (1979).
- 1580 16. R. C. Team, R: A Language and Environment for Statistical Computing (2021).
- 1581 17. G. Pau, F. Fuchs, O. Sklyar, M. Boutros, W. Huber, EBImage-an R package for image
1582 processing with applications to cellular phenotypes. *Bioinformatics* 26, 979–981 (2010).

- 1583 18. D. Hyatt, *et al.*, Prodigal: prokaryotic gene recognition and translation.
- 1584 19. S. R. Eddy, Accelerated profile HMM searches. *Plos Comput Biol* 7, e1002195 (2011).
- 1585 20. J. Mistry, *et al.*, Pfam: The protein families database in 2021. *Nucleic Acids Res* 49, 1–8
- 1586 (2020).
- 1587 21. D. H. Haft, *et al.*, TIGRFAMs: A protein family resource for the functional identification of
- 1588 proteins. *Nucleic Acids Res* 29, 41–43 (2001).
- 1589 22. K. L. Meiborn, *et al.*, The *Vibrio cholerae* chitin utilization program. *P Natl Acad Sci Usa*
- 1590 101, 2524–2529 (2004).
- 1591 23. S. Eisenbeis, S. Lohmiller, M. Valdebenito, S. Leicht, V. Braun, NagA-dependent uptake
- 1592 of N-acetyl-glucosamine and N-acetyl-chitin oligosaccharides across the outer membrane
- 1593 of *Caulobacter crescentus*. *J Bacteriol* 190, 5230–5238 (2008).
- 1594 24. A. Bateman, UniProt: A worldwide hub of protein knowledge. *Nucleic Acids Res* 47, D506–
- 1595 D515 (2019).
- 1596 25. K. Katoh, D. M. Standley, MAFFT multiple sequence alignment software version 7:
- 1597 Improvements in performance and usability. *Mol Biol Evol* 30, 772–780 (2013).
- 1598 26. F. Madeira, *et al.*, The EMBL-EBI search and sequence analysis tools APIs in 2019.
- 1599 *Nucleic Acids Res* 47, W636–W641 (2019).
- 1600 27. M. S. Datta, E. Sliwerska, J. Gore, M. F. Polz, O. X. Cordero, Microbial interactions lead to
- 1601 rapid micro-scale successions on model marine particles. *Nat Commun* 7, 1–7 (2016).
- 1602 28. T. N. Enke, *et al.*, Modular Assembly of Polysaccharide-Degrading Marine Microbial
- 1603 Communities. *Curr Biol* 29, 1528-1535.e6 (2019).
- 1604 29. S. Pontrelli, *et al.*, Hierarchical control of microbial community assembly by specialists.
- 1605 *Biorxiv*, 1–18 (2021).
- 1606 30. H. Li, *et al.*, The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25, 2078–
- 1607 2079 (2009).
- 1608 31. J. Friedman, T. Hastie, R. Tibshirani, glasso: Graphical Lasso: Estimation of Gaussian
- 1609 Graphical Models (2019).

- 1610 32. B. Buchfink, C. Xie, D. H. Huson, Fast and sensitive protein alignment using DIAMOND.
1611 *Nat Methods* 12, 59–60 (2015).
- 1612 33. S. Roux, F. Enault, B. L. Hurwitz, M. B. Sullivan, VirSorter: Mining viral signal from microbial
1613 genomic data. *Peerj* 2015, 1–20 (2015).
- 1614 34. J. Ren, N. A. Ahlgren, Y. Y. Lu, J. A. Fuhrman, F. Sun, VirFinder: a novel k-mer based tool
1615 for identifying viral sequences from assembled metagenomic data. *Microbiome* 5, 69
1616 (2017).
- 1617 35. A. C. Gregory, *et al.*, Marine DNA Viral Macro- and Microdiversity from Pole to Pole. *Cell*
1618 177, 1109-1123.e14 (2019).
- 1619 36. N. A. Ahlgren, J. Ren, Y. Y. Lu, J. A. Fuhrman, F. Sun, Alignment-free d2* oligonucleotide
1620 frequency dissimilarity measure improves prediction of hosts from metagenomically-
1621 derived viral sequences. *Nucleic Acids Res* 45, 39–53 (2017).
- 1622 37. R. A. Edwards, K. McNair, K. Faust, J. Raes, B. E. Dutilh, Computational approaches to
1623 predict bacteriophage-host relationships. *Fems Microbiol Rev* 40, 258–272 (2016).
- 1624 38. A. R. Quinlan, I. M. Hall, BEDTools: A flexible suite of utilities for comparing genomic
1625 features. *Bioinformatics* 26, 841–842 (2010).
- 1626 39. K. Kieft, K. Anantharaman, Deciphering active prophages from metagenomes. *bioRxiv*,
1627 2021.01.29.428894 (2021).
- 1628 40. M. R. Olm, *et al.*, The source and evolutionary history of a microbial contaminant identified
1629 through soil metagenomic analysis. *Mbio* 8, 1–12 (2017).
- 1630 41. N. Dumont-Leblond, M. Veillette, C. Racine, P. Joubert, C. Duchaine, Development of a
1631 robust protocol for the characterization of the pulmonary microbiota. *Commun Biology* 4,
1632 1–9 (2021).
- 1633 42. D. K. Button, B. R. Robertson, Determination of DNA Content of Aquatic Bacteria by Flow
1634 Cytometry. *Appl Environ Microb* 67, 1636–1645 (2001).
- 1635 43. MATLAB, 9.7.0.1190202 (R2019b) (The MathWorks Inc.).

- 1636 44. T. Fuhrer, D. Heer, B. Begemann, N. Zamboni, High-throughput, accurate mass
1637 metabolome profiling of cellular extracts by flow injection-time-of-flight mass spectrometry.
1638 *Anal Chem* 83, 7074–7080 (2011).
- 1639 45. P. D. Karp, *et al.*, The BioCyc collection of microbial genomes and metabolic pathways.
1640 *Brief Bioinform* 20, 1085–1093 (2017).