

1 Discovery of an Antarctic ascidian-associated uncultivated
2 Verrucomicrobia that encodes antimelanoma palmerolide biosynthetic
3 capacity

4

5 Alison E. Murray^{a*}, Chien-Chi Lo^b, Hajnalka E. Daligault^b, Nicole E. Avalon^c, Robert W. Read^a,
6 Karen W. Davenport^b, Mary L. Higham^a, Yuliya Kunde^b, Armand E.K. Dichosa^b, Bill J. Baker^{c*},
7 Patrick S.G. Chain^{b*}

8 ^a Division of Earth and Ecosystem Science, Desert Research Institute, Reno, Nevada, USA.

9 ^b Bioscience Division, Los Alamos National Laboratory, Los Alamos, New Mexico, USA.

10 ^c Department of Chemistry, University of South Florida, Tampa, Florida, USA.

11 Corresponding authors:

12 Alison E. Murray, Division of Earth and Ecosystem Sciences, Desert Research Institute, Reno, NV
13 89512, USA. +1 775-673-7361. alison.murray@dri.edu.

14 Bill J. Baker, Department of Chemistry, University of South Florida, Tampa, FL 33620, USA. +1-
15 813-974-1967, bjbaker@usf.edu

16 Patrick S.G. Chain, Bioscience Division, Los Alamos National Laboratory, Los Alamos, NM 87545,
17 USA. +1-505-665-4019. pchain@lanl.gov

18 **Author Contributions:** The team that came together to conduct this research included A.E.M.
19 B.J.B. and P.S.G.C. who designed this research. A.E.M., C-C.L., H.E.D., N.E.A., R.W.R., K.W.D.,
20 M.L.H., Y.K., A.E.K.D. performed the research. P.S.G.C., C-C.L. and Y.K. contributed analytic
21 tools and reagents. A.E.M., C-C.L., H.E.D., N.E.A., and K.W.D., analyzed the data and A.E.M., C-
22 C.L., H.E.D., N.E.A., A.E.K.D., B.J.B. and P.S.G.C. wrote the paper.

23 **Competing Interest Statement:** The authors have no competing interests

24 **Classification:** Biological sciences-environmental sciences

25 **Keywords:** Antarctic, palmerolide, anticancer, microbiome metagenome, *Synoicum adareanum*

26 **This PDF file includes:**

27 Main Text
28 Figures 1 to 4
29 Tables 1
30
31

32 Abstract

33 The Antarctic marine ecosystem harbors a wealth of biological and chemical innovation that has
34 risen in concert over millennia since the isolation of the continent and formation of the Antarctic
35 circumpolar current. Scientific inquiry into the novelty of marine natural products produced by
36 Antarctic benthic invertebrates led to the discovery of a bioactive macrolide, palmerolide A, that
37 has specific activity against melanoma and holds considerable promise as an anticancer
38 therapeutic. While this compound was isolated from the Antarctic ascidian *Synoicum adareanum*,
39 its biosynthesis has since been hypothesized to be microbially mediated, given structural
40 similarities to microbially-produced hybrid non-ribosomal peptide-polyketide macrolides. Here, we
41 describe a metagenome-enabled investigation aimed at identifying the biosynthetic gene cluster
42 (BGC) and palmerolide A-producing organism. A 74kb candidate BGC encoding the multimodular
43 enzymatic machinery (hybrid Type I-*trans*-AT polyketide synthase-non-ribosomal peptide
44 synthetase and tailoring functional domains) was identified and found to harbor key features
45 predicted as necessary for palmerolide A biosynthesis. Surveys of ascidian microbiome samples
46 targeting the candidate BGC revealed a high correlation between palmerolide-gene targets and a
47 single 16S rRNA gene variant ($R=0.83 - 0.99$). Through repeated rounds of metagenome
48 sequencing followed by binning contigs into metagenome-assembled genomes, we were able to
49 retrieve a near-complete genome (10 contigs) of the BGC organism, a novel verrucomicrobium
50 within the *Opitutaceae* family that we propose here as *Candidatus* *Synoicohabitans*
51 *palmerolidicus*. The refined genome assembly harbors five highly similar BGC copies, along with
52 structural and functional features that shed light on the host-associated nature of this unique
53 bacterium.

54 Significance Statement

55 Palmerolide A has potential as a chemotherapeutic agent to target melanoma. We interrogated
56 the microbiome of the Antarctic ascidian, *Synoicum adarenum*, using a cultivation-independent
57 high-throughput sequencing and bioinformatic strategy. The metagenome-encoded biosynthetic
58 machinery predicted to produce palmerolide A was found to be associated with the genome of a
59 member of the *S. adareanum* core microbiome. Phylogenomic analysis suggests the organism
60 represents a new deeply-branching genus, *Candidatus* *Synoicohabitans* *palmerolidicus*, in the
61 *Opitutaceae* family of the Verrucomicrobia phylum. The *Ca. S. palmerolidicus* 4.29 Mb genome
62 encodes a repertoire of carbohydrate-utilizing and transport pathways enabling its ascidian-
63 associated lifestyle. The palmerolide-producer's genome also contains five distinct copies of the
64 large palmerolide biosynthetic gene cluster that may provide structural complexity of palmerolide
65 variants.

66

67

68 Main Text

69

70 Introduction

71

72 Across the world's oceans, marine benthic invertebrates harbor a rich source of natural products
73 that serve metabolic and ecological roles in situ. These compounds provide a multitude of
74 medicinal and biotechnological applications to science, health and industry. The organisms
75 responsible for their biosynthesis are often not clear (1, 2). Increasingly, the products, especially
76 in the polyketide class (*trans*-AT in particular), are found to be produced by microbial counterparts
77 associated with the invertebrate host (3-5). Invertebrates including sponges, corals, and ascidians
78 for example, are increasingly being recognized to harbor a wealth of diverse microbes, few of
79 which have been cultivated (e.g., 6-8). Genomic tools, in particular, are revealing biochemical
80 pathways potentially critical in the host-microbe associations (9). Microbes that form persistent
81 mutualistic (symbiotic) associations provide key roles in host ecology, such as provision of

82 metabolic requirements, production of adaptive features such as photoprotective pigments,
83 bioluminescence, or antifoulants, and biosynthesis of chemical defense agents.

84 Antarctic marine ecosystems harbor species-rich macrobenthic communities (10-12),
85 which have been the subject of natural products investigations over the past 30 years resulting in
86 the identification of > 600 metabolites (13). Initially, it was not known whether the same selective
87 pressures (namely predation and competition, e.g. (14)) that operate in mid and low latitudes
88 would drive benthic organisms at the poles to create novel chemistry (15). However, this does
89 appear to be the case, and novel natural products have been discovered across algae, sponges,
90 corals, nudibranchs, echinoderms, bryozoans, ascidians, and increasingly amongst
91 microorganisms (16) for which the ecological roles have been deduced in a number of cases (13,
92 17). Studies of Antarctic benthic invertebrate-microbe associations however, pale in comparison
93 to studies at lower latitudes, yet the few studies that have been reported suggest these
94 associations (i) harbor an untapped reservoir of biological diversity (18-21) including fungi (22),
95 (ii) are host species-specific (23, 24), (iii) provide the host with sources of nitrogen and fixed
96 carbon (25) and (iv) have biosynthetic functional potential (26, 27).

97 This study was specifically motivated by our desire to understand the biosynthetic origins
98 of a natural product, palmerolide A, given its potent anticancer activity (28), that is found to be
99 associated with the polyclinid Antarctic ascidian, *Synoicum adareanum* (Fig. 1 a and b). Ascidians
100 are known to be rich sources of bioactive natural products (9). They have been found to harbor
101 polyketide, terpenoid, peptide, alkaloid and a few other classes of natural products, of which the
102 majority have cytotoxic and/or antimicrobial activities. In addition to palmerolide A, a few other
103 natural products derived from Antarctic ascidians have been reported (29-31). Ascidian-
104 associated microbes responsible for natural product biosynthesis have been shown to be
105 affiliated with bacterial phyla including Actinobacteria (which dominates the recognized diversity),
106 Cyanobacteria, Firmicutes, Proteobacteria (both Alphaproteobacteria and Gammaproteobacteria)
107 and Verrucomicrobia in addition to many fungi (32, 33). Metagenome-enabled studies have been
108 key in linking natural products to the organisms producing them in a number of cases, e.g.
109 patellamide A and C to Cyanobacteria-affiliated *Prochloron* spp. (4), the tetrahydroisoquinoline
110 alkaloid ET-743 to Gammaproteobacteria-affiliated *Candidatus* Endoecteinascidia frumentensis
111 (34), patellazoles to Alphaproteobacteria-affiliated *Candidatus* Endolissoclinum faulkneri (35), and
112 mandelalides to Verrucomicrobia-affiliated *Candidatus* Didemnitutus mandela (36). However, this
113 is most certainly an under-representation of the diversity of ascidian-associated microorganisms
114 with capabilities for synthesizing bioactive compounds, given the breadth of ascidian biodiversity
115 (37). These linkages have been yet to be investigated for Antarctic ascidians.

116 Palmerolide A has anticancer properties with selective activity against melanoma when
117 tested in the National Cancer Institute 60 cell-line panel (28). This result is of particular interest,
118 as there are few natural product therapeutics for this devastating form of cancer. Palmerolide A
119 inhibits vacuolar ATPases, which are highly expressed in metastatic melanoma. Given the current
120 level of understanding that macrolides often have microbial biosynthetic origins, that the holobiont
121 metagenome has biosynthetic potential (26) and a diverse, yet persistent core microbiome is
122 found in palmerolide A-containing *S. adareanum* (27) we have hypothesized that a microbe
123 associated with *S. adareanum* is responsible for the biosynthesis of palmerolide A.

124 The core microbiome of the palmerolide A-producing ascidian *S. adareanum* in samples
125 collected across the Anvers Island archipelago (n=63 samples; (27)) is comprised of five bacterial
126 phyla including Proteobacteria (dominating the microbiome), Bacteroidetes, Nitrospirae,
127 Actinobacteria and Verrucomicrobia. A few candidate taxa in particular, were suggested to be
128 likely palmerolide A producers based on relative abundance and biosynthetic potential
129 determined by analysis of lineage-targeted biosynthetic capability (genera: *Microbulbifer*,
130 *Pseudovibrio*, *Hoeflea*, and the family *Opiritaceae* (27). This motivated interrogation of the *S.*
131 *adareanum* microbiome metagenome, with the goals of determining the metagenome-encoded

132 biosynthetic potential, identifying candidate palmerolide A biosynthetic gene cluster (BGC)(s) and
133 establishing the identity of the palmerolide A producing organism.

134

135

136 Results and Discussion

137

138 **Identification of a putative palmerolide A biosynthetic gene cluster.** Microbial-enriched
139 fractions of *S. adareanum* metagenomic DNA sequence from 454 and Ion Proton next generation
140 sequencing (NGS) libraries (almost 18 billion bases in all) were assembled independently, then
141 merged, resulting in ~ 145 MB of assembled bases distributed over 86,387 contigs (referred to as
142 CoAssembly 1; *SI Appendix*, Table S1). As the metagenome sequencing effort was focused on
143 identifying potential BGCs encoding the machinery to synthesize palmerolide A, the initial steps of
144 analysis specifically targeted those contigs in the assembly that were > 40 kb, as the size of the
145 macrolide ring with 24 carbons would require a large number of polyketide modules to be encoded.
146 This large fragment subset of CoAssembly 1 was submitted to antiSMASH v.3 (38) and more
147 recently to v.5 (39). The results indicated a heterogeneous suite of BGCs, including a bacteriocin,
148 two non-ribosomal peptide synthetases (NRPS), two hybrid NRPS -Type I PKS, two terpenes, and
149 three hybrid *trans*-AT-PKS hybrid NRPS clusters (*SI Appendix*, Table S2).

150 We predicted several functional characteristics of the BGC that would be required for
151 palmerolide A biosynthesis which aided our analysis (see (40) for details). This included evidence
152 of a hybrid nonribosomal peptide-polyketide pathway and enzymatic domains leading to placement
153 of two distinct structural features of the polyketide backbone, a carbamoyl transferase that appends
154 a carbamate group at C-11 of the macrolide ring, and an HMGC_oA synthetase that inserts a methyl
155 group on an acetate C-1 position of the macrolide structure (C-25). The antiSMASH results
156 indicated that two of the three predicted hybrid NRPS *trans*-AT-Type I PKS contained the predicted
157 markers. Manual alignment of these two contigs suggested near-identical overlapping sequence
158 (36,638 bases) and, when joined, the merged contig resulted in a 74,672 Kb BGC (Fig. 1c). The
159 cluster size was in the range of other large *trans*-AT PKS encoding BGCs including pederin (54 Kb;
160 41), leinamycin (135.6 Kb; 42) as well as a *cis*-acting AT-PKS, jamaicamide (64.9 Kb; 43). The
161 combined contigs encompassed what appeared to be a complete BGC that was flanked at the start
162 with a transposase and otherwise unlinked in the assembly to other contiguous DNA. The cluster
163 lacked phylogenetically informative marker genes from which putative taxonomic assignment could
164 be attributed.

165 The antiSMASH results suggested that the BGC appears to be novel with the highest
166 degree of relatedness to pyxipyrrolone A and B (encoded in the *Pyxidicoccus* sp. MCy9557 genome
167 (44), to which only 14% of the genes have a significant BLAST hit to genes in the metagenome-
168 encoded cluster. The ketosynthase (KS) sequences (13 in all) fell into three different sequence
169 groups (40). One was nearly identical (99% amino acid identity) to a previously reported sequence
170 from a targeted KS study of *S. adareanum* microbiome metagenomic DNA (26). The other two were
171 most homologous to KS sequences from *Allochromatium humboldtianum*, and *Dickeya dianthicola*
172 in addition to a number of hypothetical proteins from environmental sequence data sets.

173

174 **Taxonomic inference of palmerolide A BGC.** Taxonomic attribution of the BGC was inferred
175 using a real time PCR strategy targeting three coding regions of the putative palmerolide A BGC
176 spanning the length of the cluster (acyltransferase, AT1; hydroxymethylglutaryl Co-A synthase,
177 HCS, and the condensation domain of the non-ribosomal peptide synthase NRPS, Fig. 1c) to assay
178 a *Synoicum* microbiome collection of 63 samples that have been taxonomically classified using
179 Illumina SSU rRNA gene tag sequencing (27). The three gene targets were present in all samples
180 ranging within and between sites at levels from $\sim 7 \times 10^1 - 8 \times 10^5$ copies per gram of host tissue
181 (Fig. 2a). The three BGC gene targets co-varied across all samples ($r^2 > 0.7$ for all pairs), with the
182 NRPS gene copy levels slightly lower overall (mean: 0.66 and 0.59 copies per ng host tissue for
183 NRPS:AT1 and NRPS:HCS respectively, n=63). We investigated the relationship between BGC
184 gene copies per ng host tissue for each sample and palmerolide A levels determined for the same
185 samples using mass spectrometry, however no correlation was found ($R < 0.03$, n=63; (27)). We

186 then assessed the semi-quantitative relationship between the occurrence of SSU rRNA amplicon
187 sequence variants (ASV, n=461; (27)) and the abundance of the three palmerolide A BGC gene
188 targets. Here, we found a robust correlation ($R=0.83 - 0.99$) between all 3 gene targets and a single
189 amplicon sequence variant (ASV15) in the core microbiome (*SI Appendix*, Fig. S1). This ASV is
190 affiliated with the *Opitutaceae* family of the Verrucomicrobium phylum. The *Opitutaceae* family ASV
191 (SaM_ASV15) was a member of the core microbiome as it was detected in 59 of the 63 samples
192 surveyed at varying levels of relative abundance and displayed strong correlations with the
193 abundances of BGC gene targets (Fig. 2b, $r^2 = 0.68$ with AT1, 0.97 with HCS, and 0.69 with NRPS,
194 n=63 for all). The only other correlations $R > 0.5$ were ASVs associated within the “variable” fraction
195 of the microbiome, e.g., one low abundance ASV was present in 24 of 63 samples (*SI Appendix*,
196 Fig. S1).

197 This result supports the finding of Murray et al. (27) in which gene abundance and natural
198 product chemistry do not reflect a 1:1 ratio in this host-associated system. Neither the semi-
199 quantitative measure of ASV copies nor the real time PCR abundance estimates of the three
200 biosynthetic gene targets correlated with the mass-normalized levels of palmerolide A present in
201 the same samples. As discussed (27), this is likely a result of bioaccumulation in the ascidian
202 tissues. This result provided strong support that the genetic capacity for palmerolide A production
203 was associated with a novel member of the *Opitutaceae*, a taxonomic family with representatives
204 found across diverse host-associated and free-living ecosystems. Although the biosynthetic
205 capacity of this family is not well known (45), recent evidence (36) suggests this family may be a
206 fruitful target for cultivation efforts and natural product surveys.

207
208 **Assembly of the palmerolide BGC-associated *Opitutaceae*-related metagenome assembled**
209 **genome (MAG).** With metagenomes some genomes come together easily – while others present
210 compelling puzzles to solve. Assembly of the *pal* BGC-containing *Opitutaceae* genome was the
211 result of a dedicated effort of binning contigs, gene searches, additional sequencing of samples
212 with high BGC titer, and manual, targeted assembly. Binning efforts with CoAssembly 1 did not
213 result in association of the *pal* BGC with an associated metagenome assembled genome (*SI*
214 *Appendix* text). Therefore, a further round of metagenome sequencing using long read technology
215 (Pacific Biosciences Sequel Systems technology; PacBio) ensued.

216 The 16S rRNA gene ASV occurrence (27) and real time PCR data were used to guide *S.*
217 *adarenum* sample selection for sequencing. Two ascidian samples (Bon-1C-2011 and Del-2B-
218 2011) with high *Opitutaceae* ASV occurrences (ASV_015; > 1000 sequences each – relative
219 abundance of ~ 13.3-15.3 % compared to an overall average of $1.3 \pm 2.77\%$ across the 63 samples
220 respectively, (27)) and high BGC gene target levels ($> 6.9 \times 10^5$ and $> 2.0 \times 10^5$ for the NRPS
221 respectively; *SI Appendix*, Table S3) were selected for PacBio sequencing. This effort generated
222 28 GB of data that was used to create a new hybrid CoAssembly 2 which combined all three
223 sequencing technologies. Similar to the assembly with the *Mycale hentscheli*-associated polyketide
224 producers (46), the long-read data set improved the assembly metrics, and subsequent binning
225 resulted in a highly resolved *Opitutaceae*-classified bin (*SI Appendix*, Fig. S2, Table S4).
226 Interestingly however, the palmerolide BGC contigs still did not cluster with this bin, which we later
227 attributed to binning reliance on sequence depth.

228 We used PacBio circular consensus sequence (CCS) reads to generate and manually edit
229 the assembly for our *Opitutaceae* genome of interest. The resulting 4.3 MB genome (Fig. 3a) had
230 a GC content of 58.7% and was resolved into a total of 10 contigs. Five of the contigs were unique
231 and the other five contigs represented highly similar repeated units of the *pal* BGC (labeled *pal*
232 BGC 1, 2, 3, 4, and 5) with broken ends resulting in linkage gaps. Nucmer alignment of contigs to
233 the longest palmerolide-containing BGC revealed a long (36,198 kb) repeated region that was
234 shared between all 5 contigs with some substantial differences at the beginning of the cluster and
235 only minor differences at the end, indicating 3 full length, and 2 shorter palmerolide BGC-containing
236 contigs (Fig. 1 and 3). This was consistent with coverage estimates based on read-mapping that
237 suggested lower depth at the beginning of the cluster (Fig. 3b). BGC 1 and 3 are nearly identical
238 (over 86,135 bases) with only 2 single nucleotide polymorphisms (SNPs) and an additional 1,468
239 bases in BGC1 (237 bases at the 5' end and 1231 bases at the 3' end). BGC 4 is 13,470 bases

240 shorter than BGC1 at the 5' end, and 5 bases longer than BGC 1 at the 3' end. Alignment of the
241 real time PCR gene targets to the 5 *pal* BGCs provided independent support for the different lengths
242 of the 5 BGCs, as the region targeted by the NRP primers was missing in two of the *pal* BGCs,
243 thus explaining lower NRP: AT or NRP:HCS gene dosages reported above.

244 Interestingly, precedent for naturally occurring multi-copy BGCs to our knowledge, has only
245 been found in another ascidian (*Lissoclinium* sp.)-associated *Opiritaceae*, *Candidatus*
246 *Didemnitutus mandela* which have been linked to cytotoxic mandelalides (36). Likewise, we can
247 invoke a rationale similar to (36) that multiple gene clusters may be linked to biosynthesis of
248 different palmerolide variants, see Avalon et al. (40) for retro-biosynthetic predictions of these
249 clusters. Gene duplication, loss, and rearrangement processes over evolutionary time, likely
250 explain the source of the multiple copies. At present we do not yet understand the regulatory
251 controls, whether all five are actively transcribed, if there is a producing-organism function and how
252 this may vary amongst host microbiomes.

253

254 **Phylogenomic characterization of the *Opiritaceae*-related MAG.** The taxonomic relationship
255 of the *Opiritaceae* MAG to other Verrucomicrobiota was assessed using distance-based analyses
256 with 16S rRNA and average amino acid identity (AAI). Then it was classified using the GTDB-Tk
257 tool (47), and a phylogenomic analysis based on concatenated ribosomal protein markers.
258 Comparison of 16S rRNA gene sequences amongst other Verrucomicrobia with available genome
259 sequences (that also have 16S rRNA genes; *SI Appendix*, Fig. S3) suggests that the nearest
260 relatives are *Cephalotococcus primus* CAG34 (similarity of 0.9138), *Opiritus terrae* PB90-1
261 (similarity of 0.9132) and *Geminisphaera coliterminum* TAV2 (similarity of 0.9108). The
262 *Opiritaceae*-affiliated MAG sequence is identical to a sequence (uncultured bacterium clone Tun-
263 3b A3) reported from the same host (*S. adareanum*) in a 2008 study (26); bootstrapping supported
264 a deep branching position in the *Opiritaceae* family.

265 When characterizing the MAG using AAI metrics (average nucleotide identity, ANI, found
266 no closely related genomes) the closest genomes were environmental metagenome assemblies
267 from the South Atlantic TOBG_SAT_155 (53.08 % AAI) and WB6_3A_236 (52.71 % AAI); and the
268 two closest isolate type genomes were *Nibricoccus aquaticus* str. NZ CP023344 (52.82 % AAI) and
269 *Opiritus terrae* str. PB90 (52.75 % AAI). The Microbial Genome Atlas (MiGA) support for the MAG
270 belonging in the *Opiritaceae* family was weak (p-values of 0.5). Attempts to classify this MAG using
271 GTDB-Tk (47) were hampered by the fact we have no real representative in the genome databases,
272 resulting in low confidence predictions at the species or genus levels (see the *SI Appendix* text for
273 details).

274 Verrucomicrobia exhibit free-living and host-associated lifestyles in a multitude of terrestrial
275 and marine habitats on Earth. We performed a meta-analysis of Verrucomicrobia genomes, with
276 an emphasis on marine and host-associated *Opiritaceae*, to establish more confidence in the
277 phylogenetic position of the *Opiritaceae* MAG. The analysis was based on 24 conserved proteins
278 – 21 ribosomal proteins and three additional conserved proteins (InfB, lepA, pheS). The diversity
279 of the *Opiritaceae* family, and of Verrucomicrobia in general, is largely known from uncultivated
280 organisms in which there are 20 genera in GTDB (release 05-RS95), 2 additional genera in the
281 NCBI taxonomy database, and numerous unclassified single amplified genomes (SAGs); in all,
282 only eight genera have cultivated representatives. Given the uneven representations of the 24
283 proteins across all (115) genomes assessed (MAGs and SAGs are often incomplete), we selected
284 a balance of 16 proteins across 48 genomes to assess phylogenomic relatedness across the
285 *Opiritaceae* (Fig 4). Here too, as seen with the 16S rRNA gene phylogenetic tree, the *S. adareanum*-
286 *Opiritaceae* MAG held a basal position compared to the other *Opiritaceae* genomes in the
287 analysis.

288

289 ***Opiritaceae*-related MAG relative abundance estimates and ecological inference.** The
290 relative abundance of *Opiritaceae* bin 8 was estimated in the shotgun metagenomic samples by
291 mapping the NGS reads back to the assembled MAG across the four *S. adareanum* samples
292 collected. This indicated varying levels of genome coverage in the natural samples, with the two
293 samples selected based on real time PCR-quantified high BGC copy number being clearly enriched

294 in this strain (44.70 % of reads mapped to Bon-1C-2011 and 36.78 % to Del-2B-2011, Table 1).
295 These levels are higher than estimates of relative abundance derived from the 16S rRNA gene
296 amplicon surveys (estimated at 13.33 and 15.34 % respectively) for the same samples. This is
297 likely a result of the single-copy nature of the ribosomal operon in *Opiritaceae* bin 8 vs. other taxa
298 with multiple rRNA operon copies that could thus be over-represented in the core microbiome
299 library (e.g., *Pseudovibrio* sp. str. PSC04-5.14 has 9 and *Microbulbifer* sp. is estimated at 4.1 ± 0.8
300 based on 9 finished *Microbulbifer* genomes available at the Integrated Microbial Genomes
301 Database). All host *S. adareanum* lobes surveyed (n=63) in the Anvers Island regional survey
302 contained high levels (0.49 – 4.06 mg palmerolide A x g⁻¹ host dry weight) of palmerolide A (27),
303 and variable, yet highly concordant levels of the *pal* BGCs and 16S rRNA ASV levels (Fig. 2).
304 Despite the natural population structure sampled here (four single host lobes), the bin-level
305 sequence variation was low (ranging from 72-243 SNPs) when the PacBio reads were mapped
306 back to the *Opiritaceae* bin 8 (Table 1). This suggests maintenance of a relatively invariant
307 population at the spatial and temporal scales of this coastal Antarctic region while highlighting our
308 limited understanding of the biogeographical extent of the *S. adareanum*-symbiont-palmerolide
309 relationship across a larger region of the Southern Ocean.

310 Several questions remain with regard to the in situ function of palmerolide A (a eukaryotic
311 V-ATPase inhibitor in human cell line assays (28)) in this cryohabitat: how and why is it
312 bioaccumulated by the host? Overall, the study of natural products in high latitude marine
313 ecosystems is in its infancy. This palmerolide producing, ascidian-associated, *Opiritaceae*
314 provides the first Antarctic example in which a well-characterized natural product has been linked
315 to the genetic information responsible for its biosynthesis. Gaining an understanding of
316 environmental and biosynthetic regulatory controls, establishing integrated transcriptomic,
317 proteomic, and secondary metabolome expression in the environment will also reveal whether the
318 different clusters are expressed in situ. In addition to ecological pursuits, the path to clinical studies
319 of palmerolide will require genetic or cultivation efforts. At present, we hypothesize that cultivation
320 of *Opiritaceae* bin 8 may be possible, given the lack of genome reduction or of other direct evidence
321 for host-associated dependencies.

322

323 ***Candidatus Synoicohabitans palmerolidicus* genome attributes.** The Antarctic ascidian,
324 *Synoicum adareanum*, harbors a dense community of bacteria that has a conserved core set of
325 taxa (27). The near complete ~4.30 Mbp *Opiritaceae* bin 8 metagenome assembled genome (Fig.
326 3) represents one of the core members. This MAG is remarkable in that it encodes for five 36-74
327 kb copies of the candidate BGCs that are implicated in biosynthesis of palmerolide A and possibly
328 other palmerolide compounds. Intriguingly, this genome does not seem to show evidence of
329 genome reduction as found in *Candidatus Didemnitutus mandela* (36); the other ascidian-
330 associated *Opiritaceae* genome currently known to encode multiple BGC gene copies. This is the
331 first *Opiritaceae* genome characterized from a permanently cold, ~ -1.8 - 2 °C, often ice-covered
332 ocean ecosystem. This genome encodes one rRNA operon, 45 tRNA genes, and an estimated
333 5058 coding sequences. Based on the low (< 92%) SSU rRNA gene identity and low (< 54% AAI)
334 values to other genera in the *Opiritaceae*, along with the phylogenomic position of the *Opiritaceae*
335 bin 8, the provisional name “*Candidatus Synoicohabitans palmerolidicus*” (*Ca. S. palmerolidicus*) is
336 proposed for this novel verrucomicrobium. The genus name *Synoicohabitans* (*Syn.o.i.ci.ha'bitans*.
337 N.L. neut. N. *Synoicum* a genus of ascidians; L. pres. part *habitans* inhabiting; N.L. masc. n.)
338 references this organism as an inhabitant of the ascidian genus *Synoicum*. The species name
339 *palmerolidicus* (*pal.me.ro.li'di.cus*. N.L. neut. n. *palmerolidum* palmerolide; N.L. masc.
340 adj.) designates the species as pertaining to palmerolide.

341 The GC content of 58.7% is rather high compared to other marine *Opiritaceae* genomes
342 (ave. 51.49 s.d. 0.02, n=12), yet is ~ average for the family overall (61.58 s.d. 0.06, n=69; *S/*
343 *Appendix*, Table S5). MetaERG includes metagenome assembled genomes available in the GTDB
344 as a resource for its custom GenomeDB that new genomes are annotated against. This was a clear
345 advantage in annotating the *Ca. S. palmerolidicus* genome as Verrucomicrobia genomes are widely
346 represented by uncultivated taxa. Likewise, antiSMASH was an invaluable tool for *pal* BGC

347 identification and domain structure annotation. This formed the basis to derive a predicted step-
348 wise mechanism of *pal* biosynthesis (40).

349

350 **Ca. S. palmerolidicus genome structure, function and host-associated features.** Beyond the
351 *pal* BGCs, the *Ca. S. palmerolidicus* genome encodes a variety of additional interesting structural
352 and functional features that provide a window into its lifestyle. Here we will only provide a brief
353 synopsis. In addition to the repeated BGCs, three additional repeats with two nearly identical copies
354 each (15.3 Mb, 17.0 Mb, 27.4 Mb) were identified during the assembly process (Fig. 3a, *SI*
355 *Appendix*, Table S6). These coded for 20, 25 and 41 CDs respectively, were in some cases flanked
356 by transposase/integrases (both internal and proximal) and had widespread homology with
357 Verrucomicrobia orthologs. The contents of the three repetitive elements were unique.

358 Annotations were assigned to a little more than half of the CDSs in the 15.3 Mb repeat in
359 which support for xylose transport, two sulfatase copies, two endonuclease copies and a MacB-
360 like (potential macrolide export) periplasmic core domain were encoded. Xylan might be sourced
361 from seaweeds (48) or the even ascidian as it is a minor component of the tunic cellulose (49).
362 Related to this, an endo-1,4-beta-xylanase which has exoenzyme activity in some microorganisms
363 (50) was identified elsewhere in the genome. Altogether, eight sulfatase copies were identified in
364 this host-associated organism (four in the 15.5 Mb repeat elements). These may be involved in
365 catabolic activities of sulfonated polysaccharides, and possibly as *trans*-acting elements in
366 palmerolide biosynthesis (40). In addition to the MacB-like CDSs found in this repeat, 13 different
367 MacB-homologs were present in the genome – none of which were associated with the *pal* BGCs
368 (*SI Appendix*, Fig. S4). MacB is a primary component of the macrolide tripartite efflux pump that
369 operates as a mechanotransmission system which is involved both in antibiotic resistance and
370 antibiotic export depending on the size of the macrolide molecule (51). However, two additional
371 elements required for this pump to be functional, an intramembrane MacA and an outer membrane
372 protein TolC, were not co-located in the genome. MacA may be missing, as hits to two other
373 verrucomicrobia-associated MacA CDS were not identified using BLAST (Peat Soil MAG SbV1
374 SBV1_730043 and *Ca. Udaeobacter copiosis* KAF5408997.1; (52)). At least nine MacB CDS were
375 flanked by a FstX-like permease family protein; the genomic structure of which were quite complex
376 including several with multiple repeated domains. Detailed transporter modeling is beyond the
377 scope of this work, but it is likely that these proteins are involved in signaling of cell division
378 machinery rather than macrolide transport (53).

379 Predicted CDSs in the 17.0 Mb repeat included sugar binding and transport domains, as
380 well as domains encoding rhamnosidase, arabinofuranosidase, and other carbohydrate catabolism
381 functions. About half proteins encoded in the 27.4 Mb repeat were unknown in function, and those
382 characterized suggested diverse potential functional capacities. For example, a zinc
383 carboxypeptidase (1 of 3 in the genome), multidrug and toxic compound transporter (MatE/NorM),
384 and an exodeoxyribonuclease were identified.

385 The *Ca. S. palmerolidicus* MAG has a number of features that suggest it is adapted to a
386 host-associated lifestyle, several of these features were reported recently for two related sponge-
387 associated Opitutales metagenome bins (*Petrosia ficiformis*-associated bins 0 and 01, Fig. 4; (54)).
388 These include identification of a bacterial microcompartment (BMC) ‘super locus’. Such loci were
389 recently reported to be enriched in host-associated Opitutales genomes when compared to free-
390 living relatives. The structural proteins for the BMC were present as were other conserved
391 Planctomyces-Verrucomicrobia BMC genes (55). As in the sponge *Pectoria ficiformis* metagenome
392 bins, enzymes for carbohydrate (rhamnose) catabolism and modification were found adjacent to
393 the BMC locus (*SI Appendix*, Fig. S5), in addition to the two that were found in the 27.4 Mb repeat.
394 The genome did not appear to encode the full complement of enzymes required for fucose
395 metabolism, though a few alpha-L-fucosidases were identified. Further evidence for carbohydrate
396 metabolism was supported through classification of the genome using the CAZY database (56),
397 including 7 carbohydrate binding modules, a carbohydrate esterase, 14 glycoside hydrolases, 6
398 glycosyl transferases and a polysaccharide lyase. In addition, three bacterial cellulases (PF00150,
399 cellulase family A; glycosyl hydrolase family 5) were identified, all with a canonical conserved
400 glutamic acid residue. These appear to have different evolutionary histories in which each variant

401 has nearest neighbors in different bacterial phyla (*SI Appendix*, Fig. S6) matching between 68%
402 identity for Protein J6386_03765 to *Lacunisphaera limnophila*, 57.5% identity for Protein J6386
403 22340 with a cellulase from a shipworm symbiont *Alteromonadaceae* (*Terridinibacter* sp.), and
404 37.5% sequence identity to a Bacteroidetes bacterium. This suggests the potential for cellulose
405 degradation – which is consistent with ascidians being the only animals known to produce cellulose
406 where it acts as a skeletal structure (49). In addition to the BGCs, the enzymatic resources in this
407 genome (e.g., xylan and cellulose hydrolysis) are a treasure trove rich with biotechnological
408 potential.

409 Other indicators of host-association in the *Opiritales* include T-A domains, which were
410 prevalent in the *Petrosia ficiformis*-associated bins 0 and 01 (54). The *Ca. S. palmerolidicus*
411 genome encoded at least 22 TA-related genes including multiple MazG and AbiEii toxin type IV TA
412 systems, AbiEii-Phd_YefM type II toxin-antitoxin systems, along with genes coding for PIN
413 domains, Zeta toxin, RelB, HipA, MazE and MraZ. This analysis also resulted in identifying a
414 putative AbiEii toxin (PF13304) with homology to SyrD, a cyclic peptide ABC type transporter that
415 was present in all 5 BGCs (Fig. 1c; BLAST percent identity 52.7% to a *Desulfamplus* sp. homolog
416 over the full length of the protein, and a variety of other bacteria including an *Opiritaceae*-related
417 strain at similar levels of identity). These genes are encoded downstream of the BGC following the
418 acyl transferase domains and precede the predicted *trans*-acting domains at the 3' end of the BGC.
419 Given the proximity adjacent to the primary biosynthetic gene clusters, this protein is a candidate
420 for palmerolide transport. Further research is needed however to discern the details of *Ca. S.*
421 *palmerolidicus*' cellular biology, localization of palmerolide production, transport, and resistance
422 mechanisms to the potent vacuolar ATPase as well as products made by others in the *S.*
423 *adareanum* microbiome. Along these lines, in addition to the MatE (found in the 27.4 Mb repeat)
424 two other multidrug export systems with homology to MexB and MdtB were identified.

425 Unlike in *Ca. D. mandela* (36), there does not appear to be ongoing genome reduction,
426 which may suggest that the *S. adareanum*-*Ca. S. palmerolidicus* relationship is more recent, and/or
427 that the relationship is commensal rather than interdependent. Likewise, we suspect that the
428 pseudogene content may be high as several CDS appear to be truncated, in which redundant CDS
429 of varying lengths were found in several cases (including the MacB). There is evidence of lateral
430 gene transfer acquisitions of cellulase and numerous other enzymes that may confer ecological
431 advantages through the evolution of this genome. Likewise, the origin of the *pal* BGCs and how
432 recombination events play out in the success of this Antarctic host-associated system in terms of
433 adaptive evolution (57), not to mention the ecology of *S. adareanum* is a curiosity. This phylum
434 promises to be an interesting target for further culture-based and cultivation-free studies –
435 particularly in the marine environment.

436 Together, it appears that the genome of *Ca. S. palmerolidicus* is equipped for life in this
437 host-associated interactive ecosystem that stands to be one of the first high latitude marine
438 invertebrate-associated microbiomes with a genome-level understanding – and one that produces
439 a highly potent natural product, palmerolide A. This system holds promise for future research now
440 that we have identified the producing organism and *pal* BGC. We still have much to learn about the
441 ecological role of palmerolide A – if it is involved predation avoidance, antifouling, antimicrobial
442 defense or some other yet to be recognized aspect of life in the frigid, often ice-covered and
443 seasonally light-limited waters of the Southern Ocean.

444
445

446 **Materials and Methods**

447

448 **Sample Collection.** *S. adareanum* lobes were collected in the coastal waters off Anvers Island,
449 Antarctica and stored at -80 °C until processing (*SI Appendix*, Table S1). See the *SI Appendix*
450 text and (27) for details of sample collection, microbial cell preparation and DNA extraction.

451

452 **Metagenome sequencing.** Three rounds of metagenome sequencing were conducted, the
453 details of which are in the *SI Appendix* text. This included an initial 454 pyrosequencing effort with
454 a bacterial-enriched metagenomic DNA preparation from *S. adareanum* lobe (Nor2c-2007). Next,

455 an Ion Proton System was used to sequence a metagenomic DNA sample prepared from *S.*
456 *adareanum* lobe Nor2a-2007. Then two additional *S. adareanum* metagenome DNA samples
457 (Bon-1C-2011 and Del-2b-2011) selected based on high copy numbers of the palmerolide A BGC
458 (see real time PCR Methods, *SI Appendix* text) and sequenced using Pacific Biosciences Sequel
459 Systems technology.

460
461 **Metagenome assembly, annotation and binning.** Raw 454 metagenomic reads (1,570,137
462 single end reads, 904,455,285 bases) were assembled by Newbler (58) v2.9 (Life Technologies,
463 Carlsbad, CA, flags: -large -rip -mi 98 -ml 80), while Ion Proton metagenomic reads (89,330,870
464 reads, 17,053,251,055 bases) were assembled using SPAdes (59) v3.5 (flags: --iontorrent). Both
465 assembled datasets were merged with MeGAMerge (60) v1.2 and produced 86,387 contigs with
466 a maximum contig size 153,680 and total contig size 144,953,904 bases (CoAssembly 1). To
467 achieve more complete metagenome coverage and facilitate metagenome assembled genome
468 assembly, a Circular Consensus Sequence (CCS) protocol (PacBio) was used to obtain high
469 quality long reads on two samples Bon-1C-2011 and Del-2b-2011. The 5,514,426 PacBio reads
470 were assembled with aforementioned assembled contigs (CoAssembly 1) on EDGE
471 Bioinformatics using wtdbg2 (61), a fast and accurate long-read assembler. The contigs were
472 polished with three rounds of polishing by Racon (62) into a second Coassembly (CoAssembly 2)
473 which has 4,215 contigs with a maximum contig size 2,235,039 and total size 97,970,181 bases.
474 Lastly, A manual approach was implemented to arrive at assembly of the MAG of interest the
475 details of which are described in the *SI Appendix* text.

476 The contigs from both co-assemblies 1 and 2 were submitted initially to the EDGE
477 bioinformatics platform (63) for sequence annotation using Prokka (64) v1.13 and taxonomy
478 classification using BWA (65) mapping to NCBI RefSeq (Version: NCBI 2017 Oct 3).
479 Bioinformatic predictions of natural product potential was performed using the antibiotics and
480 secondary metabolite analysis shell (antiSMASH, bacterial versions 3.0, 4.0 and 5.0 (45, 70)).
481 This tool executed contig identification, annotation and analysis of secondary metabolite
482 biosynthesis gene clusters on both CoAssemblies 1 and 2 (> 1 kb and > 40 kb data sets). As
483 most of our attention was focused on analysis the *Ca. S. palmerolidicus* assembled metagenome,
484 we also used MetaERG (66) as the primary pipeline for metagenome annotation of the ten final
485 contigs in addition to NCBI's PGAP pipeline. There were 5186 coding sequences predicted in the
486 MetaERG annotation and 5186 in NCBI's PGAP annotation.

487 MaxBin (67) and MaxBin2 (68) were used to form metagenome bins for both
488 CoAssembly 1 and 2. CheckM v1.1.11 (69) and v.1.1.12, and GTDB-Tk v.1.0.2 (47) were used to
489 verify bin quality and taxonomic classification. See *SI Appendix* text for details. In order to assess
490 the representation of assembled *Opitutaceae* genome across the 4 environmental samples used
491 for metagenome sequencing (resulting from MaxBin2 binning of CoAssembly 2), we used BWA to
492 map the CCS reads to each metagenome data set.

493
494 **Real time PCR.** Gene targets (non-ribosomal peptide synthase, acyltransferase, and 3-
495 hydroxymethylglutaryl coenzyme A synthase) were selected at different positions along the length
496 of the candidate BGC. *SI Appendix*, Table S7 lists the primer and the GBlocks synthetic positive
497 control sequence. Metagenomic DNA extracts from a large *S. adareanum* sample set (n=63 *S.*
498 *adareanum* lobes from 21 colonies), all containing high levels of palmerolide A (27), were
499 screened with the real time PCR assays on a Quant Studio 3 (Thermo Fisher Scientific, Inc.; see
500 *SI Appendix* text for details of controls and analysis).

501
502 **Phylogenomic analyses.** A phylogenomic analysis of the assembled *Opitutaceae* MAG was
503 conducted based on shared ribosomal RNA and ribosomal proteins amongst 46 and 48 reference
504 genomes respectively, out of 115 genomes in total, mined from various databases (NCBI, GTDB
505 and IMG) for uncultivated and cultivated microorganisms identified in the Verrucomicrobia phylum
506 (*SI Appendix*, Table S8). The details of these analyses are described in the *SI Appendix* text. In
507 addition, we used MiGA (NCBI Prokaryotic taxonomy and the environmental TARA Oceans

508 (Tully) databases; accessed August 2020) and GTDB-Tk (ver. 1.3.0) tools for MAG taxonomic
509 classification.

510 Phylogenetic analysis of the MacB CDS sequences were retrieved from MetatERG
511 annotated *Ca. S. palmerolidicus* contigs and homologs were retrieved from the NCBI based on
512 BLAST results. Maximum likelihood analysis was conducted on 994 aligned (MUSCLE) positions
513 using RAxML v.8.2.12 using the PROTGAMMALG model and 550 bootstrap replicates. For the
514 phylogenetic analysis of the cellulase CDS, homologs were retrieved from the NCBI based on
515 BLAST results resulting in 19 sequences and 496 aligned positions (ClustalOmega) was also
516 conducted using RAxML v.8.2.12 under the PROTGAMMALG model of evolution with 1000
517 bootstraps.

518

519 Acknowledgments

520

521 Support for this research was provided in part by the National Institute of Health award
522 (CA205932) to A.E.M., B.J.B., and P.S.G.C., with additional support from National Science
523 Foundation awards (OPP-0442857, ANT-0838776, and PLR-1341339) to B.J.B., and Desert
524 Research Institute (Institute Project Assignment) to A.E.M..

525 The assistance of several collaborators and students are acknowledged including C.
526 Amsler, M. Amsler, L. Bishop, J. Cuce, B. Dent, N. Ernster, C. Gleasner, A. Maschek, A. Shilling,
527 L. Siao, S. Thomas, and the Palmer Station science support staff. We especially would like to
528 acknowledge the help of A. Oren in assisting with bacterial nomenclature proposed here.
529 Likewise, we thank J.T. Hollibaugh and A.L. Reysenbach for comments on previous drafts of the
530 manuscript.

531

532

533 References

534

- 535 1. T. L. Simmons *et al.*, Biosynthetic origin of natural products isolated from marine
536 microorganism-invertebrate assemblages. *Proc. Natl. Acad. Sci. U.S.A.* **105**, 4587-4594
537 (2008).
- 538 2. M. Morita, E. W. Schmidt, Parallel lives of symbionts and hosts: chemical mutualism in
539 marine animals. *Nat. Prod. Rep.* **35**, 357-378 (2018).
- 540 3. J. Piel, Metabolites from symbiotic bacteria. *Nat. Prod. Rep.* **21**, 519-538 (2004).
- 541 4. E. W. Schmidt *et al.*, Patellamide A and C biosynthesis by a microcin-like pathway in
542 *Prochloron didemni*, the cyanobacterial symbiont of *Lissoclinum patella*. *Proc. Natl.*
543 *Acad. Sci. U.S.A.* **102**, 7315-7320 (2005).
- 544 5. J. Piel, Metabolites from symbiotic bacteria. *Nat. Prod. Rep.* **26**, 338-362 (2009).
- 545 6. S. Sunagawa, C. M. Woodley, M. Medina, Threatened corals provide underexplored
546 microbial habitats. *Plos One* **5**, e9554 (2010).
- 547 7. L. Behrendt *et al.*, Microbial diversity of biofilm communities in microniches associated
548 with the didemnid ascidian *Lissoclinum patella*. *ISME J.* **6**, 1222-1237 (2012).
- 549 8. N. S. Webster, M. W. Taylor, Marine sponges and their microbial symbionts: love and
550 other relationships. *Environ. Microbiol.* **14**, 335-346 (2012).
- 551 9. M. McFall-Ngai *et al.*, Animals in a bacterial world, a new imperative for the life
552 sciences. *Proc. Natl. Acad. Sci. U.S.A.* **110**, 3229-3236 (2013).
- 553 10. J. Gray, Antarctic marine benthic biodiversity in a world-wide latitudinal context. *Polar*
554 *Biol.* **24**, 633-641 (2001).
- 555 11. A. Clarke, Antarctic marine benthic diversity: patterns and processes. *J. Exp. Mar. Biol.*
556 *Ecol.* **366**, 48-55 (2008).

- 557 12. D. Piepenburg *et al.*, Towards a pan-Arctic inventory of the species diversity of the
558 macro- and megabenthic fauna of the Arctic shelf seas. *Mar. Biodivers.* **41**, 51-70 (2011).
- 559 13. J. von Salm, K. Schoenrock, J. McClintock, C. Amsler, B. Baker, "The status of marine
560 chemical ecology in Antarctica: form and function of unique high-latitude chemistry" in
561 Chemical Ecology: The ecological impacts of marine natural products. (CRC Press, Boca
562 Raton, 2018), 10.1201/9780429453465, pp. 69.
- 563 14. J. B. McClintock, C. D. Amsler, B. J. Baker, Introduction to the symposium: Antarctic
564 marine biology. *Am. Zool.* **41**, 1-2 (2001).
- 565 15. G. J. Bakus, G. Green, Toxicity in sponges and holothurians - geographic pattern. *Science*
566 **185**, 951-953 (1974).
- 567 16. S. Soldatou, B. J. Baker, Cold-water marine natural products, 2006 to 2016. *Nat. Prod.*
568 *Rep.* **34**, 585-626 (2017).
- 569 17. C. Avila, S. Taboada, L. Nunez-Pons, Antarctic marine chemical ecology: what is next?
570 *Mar. Ecol.-Evol. Persp.* **29**, 1-71 (2008).
- 571 18. N. S. Webster, A. P. Negri, M. Munro, C. N. Battershill, Diverse microbial communities
572 inhabit Antarctic sponges. *Environ. Microbiol.* **6**, 288-300 (2004).
- 573 19. N. S. Webster, D. Bourne, Bacterial community structure associated with the Antarctic
574 soft coral, *Alcyonium antarcticum*. *FEMS Microbiol. Ecol.* **59**, 81-94 (2007).
- 575 20. A. E. Murray *et al.*, Microbiome composition and diversity of the ice-dwelling sea
576 anemone, *Edwardsiella andrillae*. *Integr. Comp. Biol.* **56**, 542-555 (2016).
- 577 21. S. Rodriguez-Marconi *et al.*, Characterization of bacterial, archaeal and eukaryote
578 symbionts from Antarctic sponges reveals a high diversity at a three-domain level and a
579 particular signature for this ecosystem. *Plos One* **10**, e0138837 (2015).
- 580 22. V. M. Godinho *et al.*, Diversity and distribution of hidden cultivable fungi associated with
581 marine animals of Antarctica. *Fungal Biol.-UK* **123**, 507-516 (2019).
- 582 23. G. Steinert *et al.*, Prokaryotic diversity and community patterns in Antarctic continental
583 shelf sponges. *Front. Mar. Sci.* **6**, 297 (2019).
- 584 24. O. Sacristan-Soriano, N. P. Criado, C. Avila, Host species determines symbiotic
585 community composition in Antarctic sponges (Porifera: *Demospongiae*). *Front. Mar. Sci.*
586 **7**, 474 (2020).
- 587 25. M. Moreno-Pino, A. Cristi, J. F. Gillooly, N. Trefault, Characterizing the microbiomes of
588 Antarctic sponges: a functional metagenomic approach. *Sci. Rep.* **10**, 645 (2020).
- 589 26. C. S. Riesenfeld, A. E. Murray, B. J. Baker, Characterization of the microbial community
590 and polyketide biosynthetic potential in the palmerolide-producing tunicate, *Synoicum*
591 *adareanum*. *J. Nat. Prod.* **71**, 1812-1818 (2008).
- 592 27. A. E. Murray *et al.*, Uncovering the core microbiome and distributions of palmerolide in
593 *Synoicum adareanum* across the Anvers Island archipelago, Antarctica. *Mar. Drugs* **18**,
594 298 (2020).
- 595 28. T. Diyabalanage, C. D. Amsler, J. B. McClintock, B. J. Baker, Palmerolide A, a cytotoxic
596 macrolide from the Antarctic tunicate *Synoicum adareanum*. *J. Am. Chem. Soc.* **128**,
597 5630-5631 (2006).
- 598 29. L. H. Franco *et al.*, Indole alkaloids from the tunicate *Aplidium meridianum*. *J. Nat. Prod.*
599 **61**, 1130-1132 (1998).
- 600 30. Y. Miyata *et al.*, Ecdysteroids from the antarctic tunicate *Synoicum adareanum*. *J. Nat.*
601 *Prod.* **70**, 1859-1864 (2007).

- 602 31. A. M. Seldes, M. F. Rodriguez Brasco, L. H. Franco, J. A. Palermo, Identification of two
603 meridianins from the crude extract of the tunicate *Aplidium meridianum* by tandem
604 mass spectrometry. *Nat. Prod. Rep.* **21**, 555–563 (2007).
- 605 32. L. Chen, J. S. Hu, J. L. Xu, C. L. Shao, G. Y. Wang, Biological and chemical diversity of
606 ascidian-associated microorganisms. *Mar. Drugs* **16**, 362 (2018).
- 607 33. X. Dou, B. Dong, Origins and Bioactivities of Natural Compounds Derived from Marine
608 Ascidiaceans and Their Symbionts. *Mar. Drugs* **17**, 670 (2019).
- 609 34. C. M. Rath *et al.*, Meta-omic characterization of the marine invertebrate microbial
610 consortium that produces the chemotherapeutic natural product ET-743. *ACS Chem.*
611 *Biol.* **6**, 1244-1256 (2011).
- 612 35. J. C. Kwan *et al.*, Genome streamlining and chemical defense in a coral reef symbiosis.
613 *Proc. Natl. Acad. Sci. U.S.A.* **109**, 20655-20660 (2012).
- 614 36. J. Lopera, I. J. Miller, K. L. McPhail, J. C. Kwan, Increased biosynthetic gene dosage in a
615 genome-reduced defensive bacterial symbiont. *mSystems* **2**, e00096-00017 (2017).
- 616 37. N. Shenkar, B. J. Swalla, Global diversity of Ascidiacea. *PLoS One* **6**, e20657 (2011).
- 617 38. T. Weber *et al.*, antiSMASH 3.0-a comprehensive resource for the genome mining of
618 biosynthetic gene clusters. *Nucleic Acids Res.* **43**, W237-W243 (2015).
- 619 39. K. Blin *et al.*, antiSMASH 5.0: updates to the secondary metabolite genome mining
620 pipeline. *Nucleic Acids Res.* **47**, W81-W87 (2019).
- 621 40. N. Avalon *et al.*, Palmerolide PKS-NRPS gene clusters characterized from the microbiome
622 of an Antarctic ascidian. *bioRxiv*, doi.org/10.1101/2021.1104.1105.438531.
- 623 41. J. Piel, A polyketide synthase-peptide synthetase gene cluster from an uncultured
624 bacterial symbiont of *Paederus* beetles. *Proc. Natl. Acad. Sci. U.S.A.* **99**, 14002-14007
625 (2002).
- 626 42. Y. Q. Cheng, G. L. Tang, B. Shen, Identification and localization of the gene cluster
627 encoding biosynthesis of the antitumor macrolactam leinamycin in *Streptomyces*
628 *atroolivaceus* S-140. *J. Bacteriol.* **184**, 7013-7024 (2002).
- 629 43. D. J. Edwards *et al.*, Structure and biosynthesis of the jamaicamides, new mixed
630 polyketide-peptide neurotoxins from the marine cyanobacterium *Lyngbya majuscula*.
631 *Chem. Biol.* **11**, 817-833 (2004).
- 632 44. L. Kjaerulff *et al.*, Pyxipyrrolones: structure elucidation and biosynthesis of cytotoxic
633 myxobacterial metabolites. *Angew. Chem. Int. Ed. Engl.* **56**, 9614-9618 (2017).
- 634 45. K. Blin *et al.*, The antiSMASH database version 2: a comprehensive resource on
635 secondary metabolite biosynthetic gene clusters. *Nucleic Acids Res.* **47**, D625-D630
636 (2019).
- 637 46. M. A. Storey *et al.*, Metagenomic exploration of the marine sponge *Mycale hentscheli*
638 uncovers multiple polyketide-producing bacterial symbionts. *Mbio* **11**, e02997-02919
639 (2020).
- 640 47. P.-A. Chaumeil, A. Mussig, P. Hugenholtz, D. Parks, GTDB-Tk: a toolkit to classify
641 genomes with the Genome Taxonomy Database. *Bioinformatics* **36**, 1925–1927 (2020).
- 642 48. F. I. Qeshmi *et al.*, Xylanases from marine microorganisms: A brief overview on scope,
643 sources, features and potential applications. *BBA-Proteins Proteom.* **1868**, 140312
644 (2020).
- 645 49. Y. Zhao, J. Li, Excellent chemical and material cellulose from tunicates: diversity in
646 cellulose production yield and chemical and morphological structures from different
647 tunicate species. *Cellul. Chem. Technol.* **21**, 3427-3441 (2014).

- 648 50. V. Juturu, J. C. Wu, Microbial exo-xylanases: a mini review. *Appl. Biochem. Biotechnol.*
649 **174**, 81-92 (2014).
- 650 51. N. P. Greene, E. Kaplan, A. Crow, V. Koronakis, Antibiotic resistance mediated by the
651 MacB ABC transporter family: a structural and functional perspective. *Front. Microbiol.*
652 **9**, 950 (2018).
- 653 52. S. F. Altschul, W. Gish, W. Miller, E. W. Myers, D. J. Lipman, Basic local alignment search
654 tool. *J. Mol. Biol.* **215**, 403-410 (1990).
- 655 53. A. Crow, N. P. Greene, E. Kaplan, V. Koronakis, Structure and mechanotransmission
656 mechanism of the MacB ABC transporter superfamily. *Proc. Natl. Acad. Sci. U.S.A.* **114**,
657 12572-12577 (2017).
- 658 54. S. Sizikov *et al.*, Characterization of sponge-associated Verrucomicrobia:
659 microcompartment-based sugar utilization and enhanced toxin-antitoxin modules as
660 features of host-associated Opitutales. *Environ. Microbiol.* **22**, 4669-4688 (2020).
- 661 55. O. Erbilgin, K. L. McDonald, C. A. Kerfeld, Characterization of a planctomycetal organelle:
662 a novel bacterial microcompartment for the aerobic degradation of plant saccharides.
663 *Appl. Environ. Microbiol.* **80**, 2193-2205 (2014).
- 664 56. Lombard V, H. G. Ramulu, E. Drula, P. Coutinho, B. Henrissat, The carbohydrate-active
665 enzymes database (CAZy) in 2013. *Nucleic Acids Res* **42**, D490–D495 (2014).
- 666 57. M. G. Chevrette *et al.*, Evolutionary dynamics of natural product biosynthesis in
667 bacteria. *Nat. Prod. Rep.* **37**, 566-599 (2020).
- 668 58. M. J. Chaisson, P. A. Pevzner, Short read fragment assembly of bacterial genomes.
669 *Genome Res.* **18**, 324-330 (2008).
- 670 59. S. Nurk *et al.*, Assembling single-cell genomes and mini-metagenomes from chimeric
671 MDA products. *J. Comput. Biol.* **20**, 714-737 (2013).
- 672 60. M. Scholz, C. C. Lo, P. S. G. Chain, Improved assemblies using a source-agnostic pipeline
673 for metagenomic assembly by merging (MeGAMerge) of contigs. *Sci. Rep-UK* **4**, 6480
674 (2014).
- 675 61. J. Ruan, L. Li, Fast and accurate long-read assembly with wtdbg2. *Nat. Methods* **17**, 155-
676 158 (2019).
- 677 62. R. Vaser, I. Sovic, N. D. Nagarajan, M. Sikic, Fast and accurate de novo genome assembly
678 from long uncorrected reads. *Genome Res.* **27**, 737-746 (2017).
- 679 63. P.-E. Li *et al.*, Enabling the democratization of the genomics revolution with a fully
680 integrated web-based bioinformatics platform. *Nucleic Acids Res.* **45**, 67-80 (2017).
- 681 64. T. Seemann, Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068-
682 2069 (2014).
- 683 65. H. Li, R. Durbin, Fast and accurate short read alignment with Burrows-Wheeler
684 transform. *Bioinformatics* **25**, 1754-1760 (2009).
- 685 66. X. L. Dong, M. Strous, An integrated pipeline for annotation and visualization of
686 metagenomic contigs. *Front. Genetics* **10**, 999 (2019).
- 687 67. Y. W. Wu, Y. H. Tang, S. G. Tringe, B. A. Simmons, S. W. Singer, MaxBin: an automated
688 binning method to recover individual genomes from metagenomes using an
689 expectation-maximization algorithm. *Microbiome* **2**, 26 (2014).
- 690 68. Y. W. Wu, B. A. Simmons, S. W. Singer, MaxBin 2.0: an automated binning algorithm to
691 recover genomes from multiple metagenomic datasets. *Bioinformatics* **32**, 605-607
692 (2016).

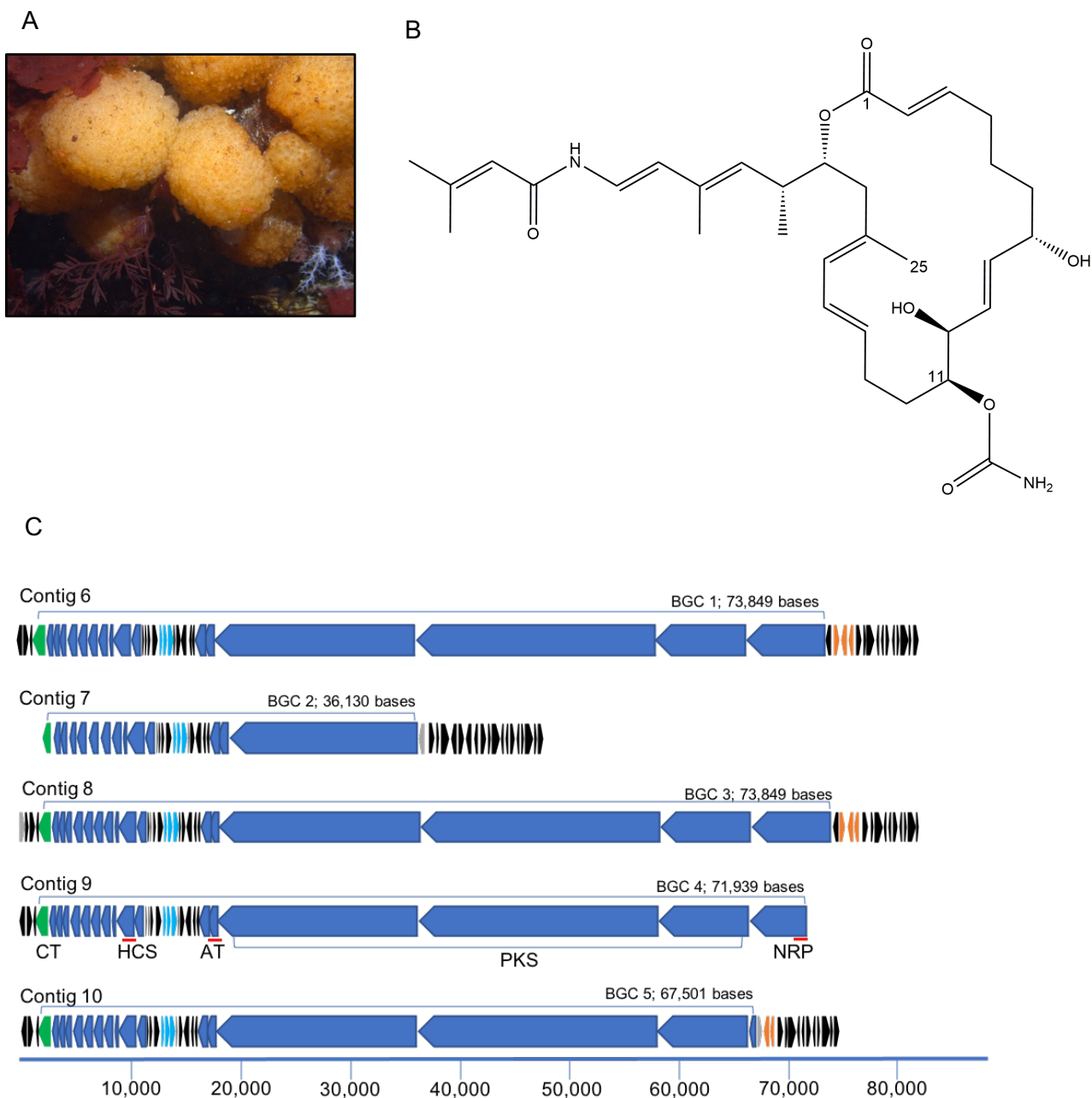
- 693 69. D. H. Parks, M. Imelfort, C. T. Skennerton, P. Hugenholtz, G. W. Tyson, CheckM:
694 assessing the quality of microbial genomes recovered from isolates, single cells, and
695 metagenomes. *Genome Res.* **25**, 1043-1055 (2015).
696 70. T. Weber *et al.*, antiSMASH 3.0 - a comprehensive resource for the genome mining of
697 biosynthetic gene clusters. *Nucleic Acids Res.* **43**, W237-W243 (2015).
698

699 **Figures and Tables**

700

701

702

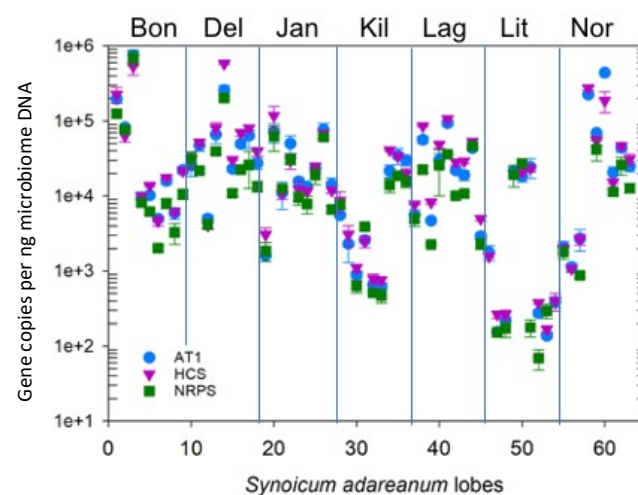


708

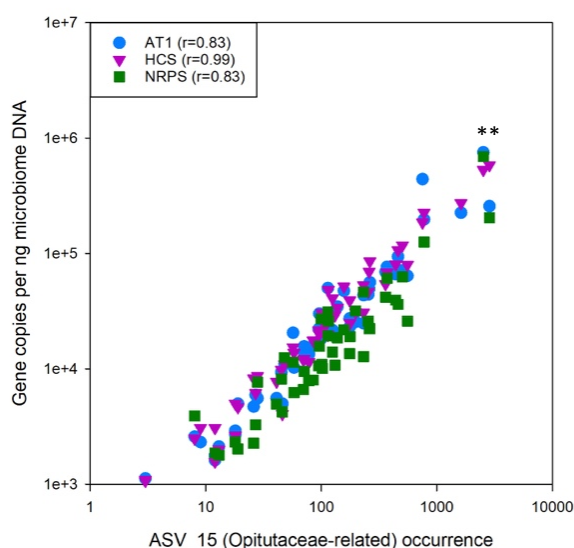
709 **Figure 1.** Palmerolide A, a cytotoxic, macrolide with anti-melanoma activity is found in the tissues of
710 of *Synoicum adareanum* in which a candidate biosynthetic gene cluster has been identified. (A) *S.*
711 *adareanum* occurs on rocky coastal seafloor habitats in the Antarctic; this study focused on the
712 region off-shore of Anvers Island in the Antarctic Peninsula. (B) Palmerolide A, is the product of a
713 hybrid PKS-NRPS system in which biosynthesis begins with a PKS starter unit followed by
714 incorporation of a glycine subunit by an NRPS module. Subsequent elongation, cyclization and
715 termination steps follow. Two additional features of the molecule include a methyl group on C-25

716 and a carbamate group on C-11. (C) Five repeats encompassing candidate palmerolide
717 biosynthetic gene clusters were identified. The BGC (in blue) is defined as starting with the NRP
718 unit and ending at the carbamoyltransferase (green). Candidate palmerolide A biosynthetic gene
719 cluster BGC4 was identified from initial metagenome library assemblies. The other four clusters
720 were identified following a third round of sequencing, assembly and manual finishing. Primary
721 BGC coding sequences (CDS) and a conserved tailoring cassette are in blue. Light blue CDS are
722 an ATP transporter with homology to an antibiotic transporter, SyrD. All black CDS are repeated
723 amongst the BGCs. Orange CDS are transposase/integrase domains. Gray CDS are unique,
724 non-repeated; and in BGC2 and 5, the unique CDS encode transposases, distinct from the
725 predicted amino acid sequences of those in orange. The red lines associated with Contig 9
726 indicate targeted quantitative PCR regions.
727

728 A
 729
 730
 731
 732
 733
 734
 735
 736
 737
 738
 739
 740
 741 B

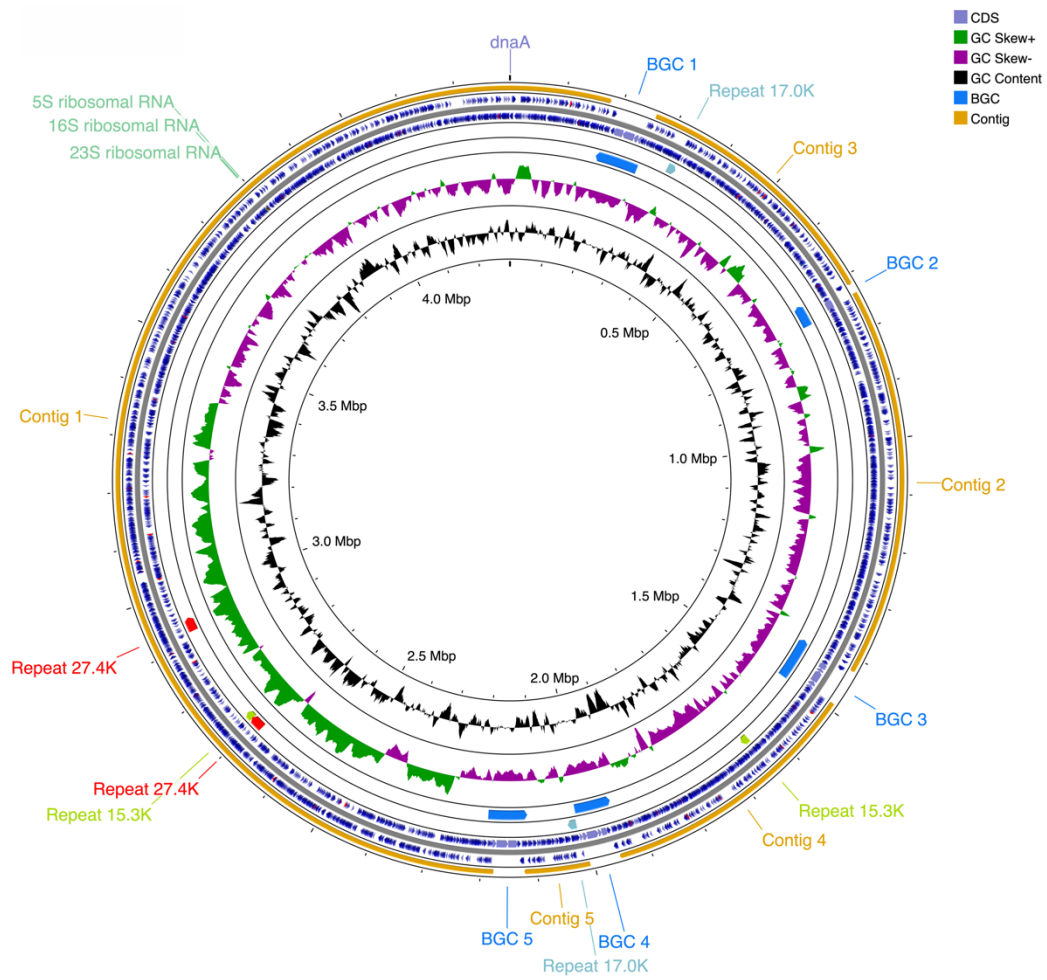


742
 743
 744
 745
 746
 747
 748
 749
 750
 751
 752
 753
 754
 755
 756
 757
 758
 759
 760
 761
 762
 763



764 **Figure 2.** Abundances of real time PCR-targeted coding regions in the candidate pal biosynthetic
 765 gene cluster in Antarctic ascidian samples. (A) Gene copies estimated for three targeted coding
 766 regions (Acyltransferase, AT1; 3-hydroxy-methyl-glutaryl coenzyme A synthase, HCS; and the
 767 condensation domain of a non-ribosomal peptide synthase; NRPS) in the candidate *palA*
 768 biosynthetic gene cluster surveyed over 63 DNA extracts derived from microbial cell preparations
 769 enriched from the Antarctic ascidian *Synoicum adareanum*. Nine samples were collected at each
 770 of seven sites: Bon, Bonaparte Point; Del, Delaca Island; Jan, Janus Island, Kil, Killer Whale
 771 Rocks; Lag, Laggard Island, Lit, Litchfield Island; Nor, Norsel Point (27). (B) Relationship between
 772 gene copy number for the three gene targets and the 16S rRNA gene ASV occurrences of
 773 Opiritaceae-related ASV_15 across a 63 *S. adareanum* microbial DNA sample set. * indicates
 774 samples Bon-1C-2011 and Del-2b-2011 that were selected for PacBio sequencing.
 775

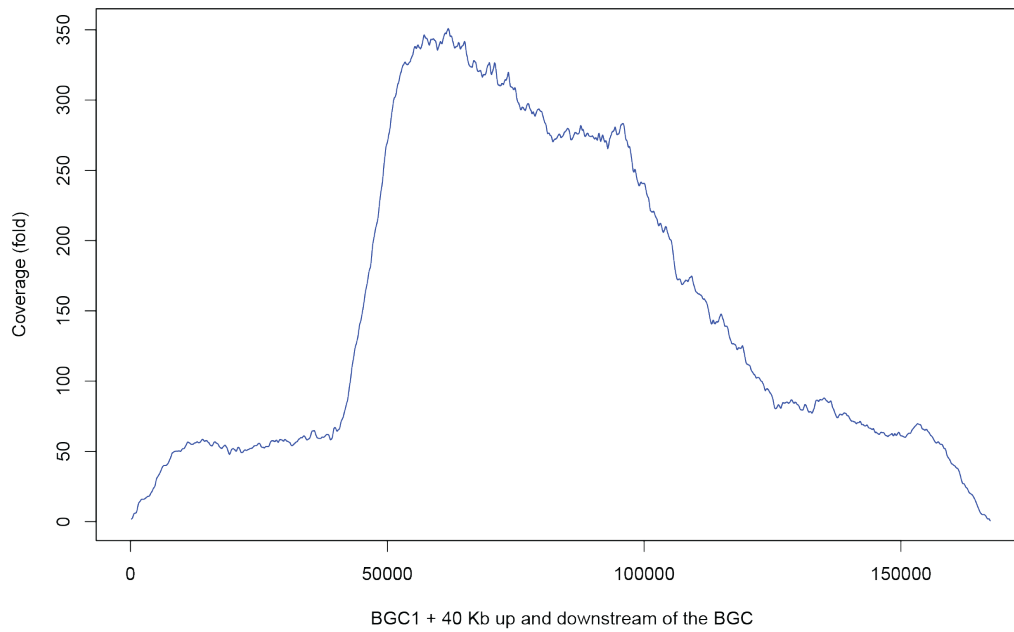
A



776

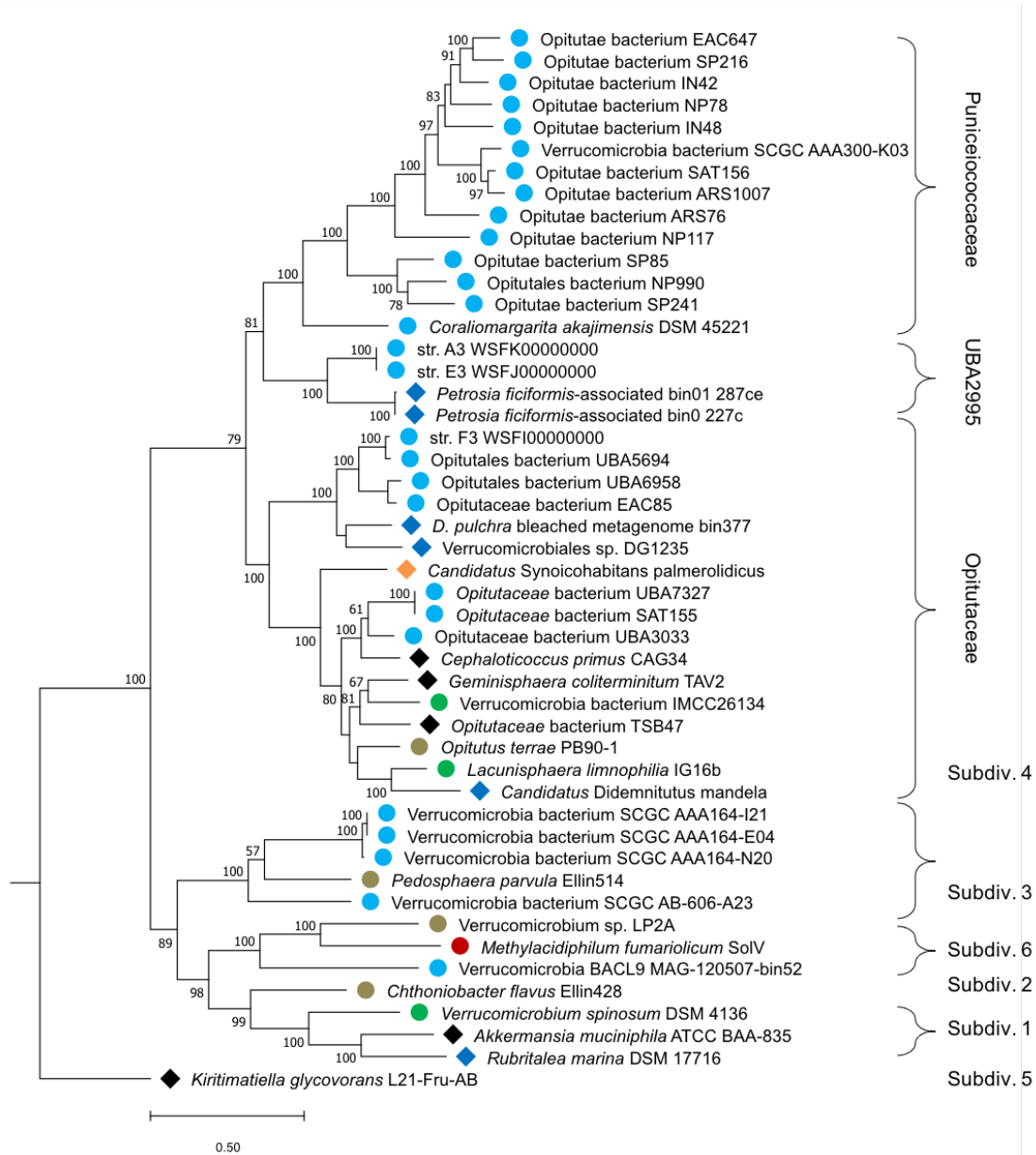
777

B



778

779 **Figure 3.** Genome maps of assembled MAG, *Candidatus* Synoicohabitans palmerolidicus and
780 evidence of multi-copy biosynthetic gene clusters. (a) The 4,297,084 bp gene map is oriented to
781 dnaA at the origin. One possible assembly scenario of the *Ca. S. palmerolidicus* genome is
782 shown as the order of the contigs and palmerolide BGC's are not currently known. In addition to
783 the 5 BGCs, three other internally repetitive regions were identified (15.3, 17.0 and 27.4 kb). The
784 genes and orientation are shown with blue and tRNA indicated in red. (b) To demonstrate depth
785 of coverage outside and inside the BGC regions, CCS reads from sample Bon-1C-2011 and Del-
786 2b-2011 were mapped to a 167.6 Kb region. The profile extends 40 kb into the genome on either
787 side of the BGC where depth of coverage averages 60 fold, while in the BGC depth of coverage
788 varies across the BGC given differences in cover across the BGC, the highest cover is 5X, or ~
789 300 fold, supporting the finding of 5 repeats encoding the BGC.
790



791

792 **Figure 4.** Maximum likelihood phylogenomic tree showing 48 verrucomicrobia genomes
 793 Phylogenomic relationship of *Candidatus Synoicohabitans palmerolidicus* (*Opitutaceae* bin 8)
 794 with respect to other mostly marine, and host-associated verrucomicrobia subdivision 4 and
 795 other genomes. The tree is based on 16 concatenated ribosomal proteins (5325 amino acids)
 796 common across 48 Verrucomicrobia genomes. Distance was estimated with RAxML with 300
 797 bootstrap replicates. Symbols designate environmental origins of the organisms: free-living are
 798 represented by circles: light blue - marine, green – freshwater, red – hydrothermal mud, brown –
 799 soils. Host-associated taxa from marine systems with blue diamonds, and from terrestrial systems
 800 with black diamonds.

801 **Table 1.** Metagenomic reads from 4 different samples were mapped back to the *Ca. S.*
 802 *palmerolidicus* MAG.
 803

	<i>Synocicum adareanum</i> samples			
	Bon-1c-2011	Del-2b-2011	Nor-2c-2007	Nor-2a-2007
Technology	PacBio CCS reads (1 cell)	PacBio CCS reads (3 cells)	454	Ion Torrent Proton
Number of reads	48,298	9,576	1,570,126	89,330,870
Mapped reads	21,591	3,522	23,993	15,979,084
Mapped reads (%)	44.70	36.78	1.53	17.89
Base coverage	99.98	99.89	90.15	99.98
Average fold	58.38	8.43	2.79	708.79
Gaps	2	3	1,734	8
Gap bases	644	4,618	422,870	774
SNPs	72	196	168	243
Indels	126	64	17	68

804