1 **Title:** Uncovering Lasonolide A biosynthesis using genome-resolved metagenomics

- 2 **Running title (54 characters)**: Elucidation of Lasonolide A biosynthesis
- 3 Authors: Siddharth Uppal^a, Jackie L. Metz^b, René K.M. Xavier^b, Keshav Nepal^b, Dongbo Xu^b,
- 4 Guojun Wang^{b*}, Jason C. Kwan^{a*}
- 5 Affiliation: ^aDivision of Pharmaceutical Sciences, School of Pharmacy, University of
 6 Wisconsin—Madison, Madison, Wisconsin, USA
- ⁷ ^bHarbor Branch Oceanographic Institute, Florida Atlantic University, Florida, USA
- 8 *Address correspondence to Guojun Wang, <u>guojunwang@fau.edu</u>; Jason C. Kwan,
 9 jason.kwan@wisc.edu

10 Abstract

11 Invertebrates, in particular sponges, have been a dominant source of new marine natural 12 products. For example, lasonolide A (LSA) is a potential anti-cancer molecule isolated from the 13 marine sponge *Forcepia* sp., with nanomolar growth inhibitory activity and a unique cytotoxicity 14 profile against the National Cancer Institute 60 cell line screen. Here, we identified the putative 15 biosynthetic pathway for LSA. Genomic binning of the Forcepia sponge metagenome revealed a 16 gram-negative bacterium belonging to the phylum Verrucomicrobia as the candidate producer of 17 LSA. Phylogenetic analysis showed this bacterium, herein named *Candidatus* Thermopylae 18 lasonolidus, only has 88.78% 16S rRNA identity with the closest relative Pedosphaera parvula 19 Ellin514, indicating it represents a new genus. The lasonolide A (*las*) biosynthetic gene cluster 20 (BGC) was identified as a *trans*-AT polyketide synthase (PKS) pathway. When compared with 21 its host genome, the las BGC exhibits a significantly different GC content and penta-nucleotide

22 frequency, suggesting a potential horizontal acquisition of the gene cluster. Furthermore, three 23 copies of the putative *las* pathway were identified in the candidate producer genome. Differences 24 between the three *las* repeats were observed including the presence of three insertions, two 25 single-nucleotide polymorphisms and the absence of a stand-alone acyl carrier protein in one of 26 the repeats. Even though the Verrucomicrobial producer shows signs of genome-reduction, its 27 genome size is still fairly large (about 5Mbp) and when compared to its closest free-living 28 relative contains most of the primary metabolic pathways, suggesting that it is in the early stages 29 of reduction.

30 Importance

31 While sponges are valuable sources of bioactive natural products, a majority of these compounds 32 are produced in small amounts by uncultured symbionts, hampering the study and clinical 33 development of these unique compounds. Lasonolide A (LSA), isolated from marine sponge 34 Forcepia sp., is a cytotoxic molecule active at nanomolar concentrations and causes premature 35 chromosome condensation, blebbing, cell contraction and loss of cell adhesion, indicating a 36 novel mechanism of action and making it a potential anti-cancer drug lead. However, its limited 37 supply hampers progression to clinical trials. We investigated the microbiome of *Forcepia* sp. 38 using culture-independent DNA sequencing to uncover how an uncultured bacterium produces 39 LSA. This provides future opportunities for heterologous expression and cultivation efforts that 40 may minimize LSA's supply problem.

41 Keywords: Lasonolide A, horizontal gene transfer, multiple repeats, Verrucomicrobia, *trans*-AT
42 PKS, genome reduction

43 Introduction

44 Lasonolide A (LSA) is a cytotoxic polyketide derived from the marine sponge *Forcepia* sp. (Fig. 45 **1A and 1B**) (1). Out of its analogs (B–G) (**Fig. 1C**), LSA is the most potent (2) and exhibits IC_{50} 46 values in the nanomolar range against certain cell lines in the National Cancer Institute 60 cell 47 line screen (3). Furthermore, its unique mechanism of action - induction of premature 48 chromosome condensation, loss of cell adhesion, activation of the RAF1 kinase in Ras pathway, 49 along with cell blebbing and contraction (3–5) - makes it a promising candidate as a scaffold for 50 future pharmaceutical development. However, a major challenge to its clinical development is 51 the lack of availability. Scarcity and limited accessibility of the sponge prevent it from being a 52 sustainable source of lasonolide A. Furthermore, the chemical synthesis of LSA is tedious and 53 has poor yields, limiting its scalability (6–8).

It is well known that bacteria living in a symbiotic relationship with higher animals are valuable 54 55 sources of novel bioactive secondary metabolites (9). In many instances, these molecules serve a 56 protective function for the host but the identity of the microbial producer remains unknown (9– 57 12). Furthermore, attempts to isolate these associated microbes are hampered by low cultivation 58 success; it is estimated less than 1% of bacteria are currently culturable from the environment 59 (13–15). These drawbacks have created the need to genetically engineer surrogate hosts for the 60 sustainable and sufficient production of the desired natural products in the laboratory. The first 61 step in engineering microbes for production of bioactive compounds is to identify the genes 62 responsible for natural product synthesis, which can be elucidated through metagenomic analysis 63 and cloning (16, 17). Based on its potent antitumor activity, it is likely that LSA also acts as a 64 chemical defense within its host sponge. The structure of LSA very likely arises from an 65 assembly-line type polyketide synthase (PKS), rather than the iterative PKSs that predominate in

fungi and other eukaryotes, and therefore the source is likely bacterial (18–20). Identifying the
bacterium responsible for synthesizing LSA and elucidating its biosynthetic pathway will allow
us to explore routes for LSA's heterologous expression, and potentially facilitate the synthesis of
analogs.

70 Here, we describe a *trans*-AT PKS pathway (*las*) that is likely responsible for the biosynthesis of 71 LSA. Furthermore, the entire las BGC has been captured on five overlapping fosmid clones and 72 reassembled for the purpose of heterologous expression. We propose that the *las* BGC is present 73 in a vet uncultivated bacterium belonging to a novel genus under the phylum Verrucomicrobia. 74 Additionally, evidence suggests the las BGC is repeated thrice within the Verrucomicrobia 75 symbiont with minor sequence variations between them. We also suggest that the las BGC has 76 been horizontally acquired and has a codon adaptation index comparable to that of highly expressed genes. Finally, we show that the Verrucomicrobia symbiont is in very early stages of 77 78 genome reduction and is likely to further reduce its size.

79 Results and discussion

80 Identification and capture of *las* BGC

In our initial studies, we constructed a high-capacity metagenomic DNA library consisting of ~600,000 *cfu* from *Forcepia* sp. sponges collected from the Gulf of Mexico (**Fig. 2A**) to search for potential *las* biosynthetic genes. The structure of LSA contains two tetrahydropyran rings and two β -methylations (21, 22) at C-13 and C-35 (**Fig. 2B**). These structural features have been identified in a variety of *trans*-AT PKS pathways but are rarely found in *cis*-AT PKS systems (23, 24) thus hinting that LSA is produced by a *trans*-AT PKS pathway (24). Therefore, clade-

87 guided degenerate primers targeted to conserved *trans*-AT PKS genes such as 3-hydroxy-3-88 methyglutaryl-CoA (HMG-CoA) synthase, free-standing ketosynthase (KS), acyl carrier protein 89 (ACP), and two enoyl-CoA hydratases (ECH) were utilized for initial screening of the Forcepia 90 fosmid library (Table S1A). From the metagenomic library, five fosmids were identified using 91 these primers (fosmids 5-16, 6-71, 3-46, 1-80, and 4-77) resulting in the capture of 92 approximately 48kb of the putative *las* BGC at its 3' end (Fig. S1A). However, minimal progress 93 was made toward capturing the remaining half of the BGC as primer walking failed to produce 94 new hits in the region upstream of fosmid 5-16. Therefore, we sequenced the metagenome of 95 Forcepia sp. and searched for trans-AT PKS BGCs. DNA was extracted from two different 96 regions (referred to as Forcepia_v1 and Forcepia_v2) of the same sponge, and subjected to whole 97 genome shotgun metagenomic sequencing. The reads were trimmed, assembled and then binned 98 into metagenome assembled genomes (MAGs). The metagenomes were found to be abundant in 99 Acidobacteria, Proteobacteria and Chloroflexota (Fig. 2C and S1B), with 56 and 55 MAGs 100 recovered from the two metagenomes. Based on MiMAG (25) standards for completeness and 101 contamination, 11 and 6 MAGs were high quality and 21 and 19 MAGs were medium quality, 102 for Forcepia_v1 and Forcepia_v2, respectively (Table S2).

A tBLASTN (26) search of KS domains from publicly available *trans*-AT PKS pathways against our assembled metagenome was performed. In the case of Forcepia_v1, the top hits were all to a contig of length 98kbp labeled gnl|UoN|bin5_1_edit_8, thus strongly suggesting that this contig contains *trans*-AT PKS genes and may possess the potential LSA pathway. Contig gnl|UoN|bin5_1_edit_8 was manually inspected and corrected for sequence gaps (**Text S1**). With the exception of a 1.1kbp contig annotated as containing a *trans*-AT PKS pathway with a truncated condensation domain (in bin3674_131), analysis of the metagenome using

110 AntiSMASH (27) (Fig. 2D) did not reveal any other BGC with plausible size and genes for the 111 synthesis of LSA. Contig gnl|UoN|bin5 1 edit 128 (3.6 kbp) which was connected to the 5 \Box 112 end of gnl|UoN|bin5 1 edit 8 was found to contain a stand-alone ACP domain and about 47 113 amino acid residues which completed the terminal KS domain of gnl|UoN|bin5 1 edit 8 (see 114 multiple repeats of the las BGC). Both the contigs were assembled together and annotation of 115 genes and biosynthetic domains within this assembly re-affirmed that they are likely involved in 116 LSA synthesis. We termed the gene cluster deemed relevant to LSA biosynthesis as las BGC v1. 117 Furthermore, the sequence of *las* BGC v1 was also in alignment with fosmids identified from 118 the metagenomic library. A screening strategy was then developed for isolating the previously 119 missing 5' end of the pathway from the metagenomic library (fosmids 5-41, 2-18 and 2-13) (Fig. 120 S1A and Table S1B).

121 Inspection of the MAGs revealed that gnl|UoN|bin5_1_edit_8 binned with genome bin75_1. 122 However, to our surprise, visual inspection of the assembly graph (Fig. S1C) in BANDAGE (29) 123 indicated that gnl|UoN|bin5 1 edit 8 is present between contigs belonging to bin5 1 (phylum 124 Verrucomicrobia). Furthermore, mapping the paired-end reads on the genomic bin (Fig. S1D) 125 showed that multiple read pairs aligned across the contig junction. The terminal connections 126 between contig gnl|UoN|bin5_1_edit_8 and several contigs in bin5_1 were verified via PCR 127 (Table S1C) and Sanger sequencing of the amplicons using metagenomic DNA as the template. 128 Based on this evidence gnl|UoN|bin5 1 edit 8 was manually placed with bin5 1, as well as 129 additional contigs (Text S1).

In the case of Forcepia_v2, tBLASTN of KS domains hit to eight different contigs which could
be assembled together (<u>https://www.geneious.com</u>) (Fig. S1E). Except for contig
gnl|UoN|bin4_1_edit_10 the other seven contigs assembled into a single large contig of 102kbp

6

(*las* BGC_v2). Similar to *las* BGC_v1, inspection of the assembly graph (Fig. S1F) and mapping
of paired-end reads (Fig. S1G) revealed that contigs forming *las* BGC_v2 have been binned
incorrectly and should be part of the bin4_1 (phylum Verrucomicrobia). As a result, the contigs
comprising *las* BGC_v2, as well as additional contigs (Text S1) were manually placed with
bin4_1. No other contig containing a *trans*-AT PKS pathway was identified in the metagenome
(Fig. S1H).

Alignment of *las* BGC from both Forcepia_v1 and Forcepia_v2 using clinker (28) revealed that these pathways are highly similar (**Fig. 2E**). The amino acid identity is 100% for most of the genes except for *lasJLO* where the amino acid identity is 98.37%, 99.84% and 99.83% respectively. The slightly lower identity of *lasJLO* is due to the insertion sequence present in *las* BGC_v2 but absent in *las* BGC_v1. These insertion variants were later identified to be present in some repeats of *las* BGC_v1 as well.

145 The putative symbiont genome carrying the las BGC (Forcepia v1 bin5 1 and Forcepia v2 146 bin4_1) was identified to belong to phylum Verrucomicrobiota, order Pedosphaerales, and genus 147 UBA2970 by GTDB-TK v1.5.0 (database r202) (30). Excluding the las genes, the ANI of 148 Forcepia_v1 bin5_1 and Forcepia_v2 bin4_1 is 99.9%, suggesting little strain heterogeneity 149 between the sites in the sponge, beyond a small amount perhaps attributable to sequencing errors. 150 To our knowledge, this is the first time a *trans*-AT PKS BGC has been reported in an organism 151 belonging to order Pedosphaerales. A phylogenetic tree of 51 different Verrucomicrobia 152 genomes (Fig. S1I) placed the LSA producer in subdivision 3 (NCBI taxonomy). The closest 153 relative of the symbiont with a publicly available genome is *Pedosphaera parvula* Ellin514 154 (GCA 000172555.1), with 88.78% identity to the 16S ribosomal RNA sequence. As per the 16S 155 taxonomic cutoffs proposed by Yarza et al. (31), this represents a new genus within the family

156 AAA164-E04 (as classified by GTDB-Tk (30)). We named the bacterium "Candidatus 157 Thermopylae lasonolidus": Thermopylae is a tribute to the 300 Spartan hoplites and other Greek 158 soldiers that fought at the battle of Thermopylae. The Spartans fought to protect Greece from 159 Persians and the LSA-producing bacterium with its three copies of the las BGC (see below) is 160 proposed to be protecting the host sponge from predators. Lasonolidus suggests the bacterium is 161 associated with lasonolide A and also rhymes with the Spartan king of the 300 hoplites, 162 Leonidas. Despite being the putative producer of LSA, "Ca T. lasonolidus" is not highly 163 abundant in the metagenome, having a relative abundance of just over 2.65% in Forcepia v1 and 164 1.78% in Forcepia_v2 (Fig. 2C and Fig. S1B).

With the aid of metagenome sequence, additional fosmids covering the $5\square$ -end were acquired (Fig. S1), which enabled us to capture the *las* BGC minimally on 5 fosmids (Fig. 3) and to subsequently reassemble the BGC into a plasmid for heterologous expression. Network analysis with BiG-SCAPE (32) revealed no shared families with MIBiG reference BGCs indicating the novelty of the *las* BGC.

170 Model for lasonolide biosynthesis by *las* BGC

The proposed biosynthetic scheme for the synthesis of LSA based on the *las* BGC is shown in Fig. 3. The complete *las* BGC consists of six *trans*-AT PKS proteins (*lasHJLMNO*), ten accessory genes (*lasCDEFIKPQRS*) and five genes with no or unknown role in LSA synthesis (*lasABGTU*). Phylogenetic analysis of 944 different KS domains (Fig. S2) was used to predict KS substrate specificity (33), which was found to be similar to the proposed biosynthetic model. The pathway is predicted to be collinear with the first KS domain of *lasH* clustering into the same clade as other starter KS domains in the KS phylogenetic tree. Moreover, the last *trans*-AT

PKS protein (*lasO*) contains a condensation domain, similar to those found in nonribosomal
peptide synthetase pathways, as its terminal domain, which is proposed to be responsible for
cyclizing and cleaving the final PKS product (24).

181 An acylhydrolase (AH) domain is often used in *trans*-AT PKS systems for proofreading by 182 cleaving the acyl units from stalled sites (34, 35). Alls are closely related to acyltransferase (AT) 183 domains, which are involved in the addition of malonyl-S-coenzyme A extender units on the 184 phosphopantetheine arms on ACP domains (24, 36). LasE (AH) and LasF (AT) were correctly 185 identified as AH and AT domains respectively, based on the presence of active site residues (Fig. 186 S3A) and phylogeny (35) (Fig. S3B). The accessory proteins LasCDRS include enzymes known to be involved in β -branch formation at module 1 and 10 (21). The ACPs at module 1 and 10 187 188 contain a conserved tryptophan which is involved in interacting with β -branching enzymes (37, 189 38). LasR was identified to be responsible for dehydration (ECH1) while LasH ECHb and LasO 190 ECHb to be responsible for decarboxylation (ECH2) during β -formation (39) (Fig. S3C). Due to 191 their truncated size and lack of homology to the conserved sites needed for oxyanion hole 192 formation LasH ECHa and LasO ECHa are proposed to be inactive (40, 41) (Fig. S3D and S3E). 193 An endo- β -methyl (\Box , β -unsaturated β -methyl) is predicted to form on module 10. The presence 194 of a truncated ECH domain just upstream of ECH2 domain has been commonly observed with 195 the formation of exo- β -methylene (β , γ -unsaturated β -methylene), but to our knowledge this is 196 the first time such an architecture has been reported to form a endo- β -methyl (38). Based on the 197 collinearity of the pathway suggested above and the split module architecture (KS-DH MOX 198 ACP-KS) associated with different Baeyer-Villiger (BV) monooxygenases as seen in oocydin 199 and sesbanimide biosynthesis (42-44) we propose LasI to be involved in BV oxidation and LasK 200 in the addition of a hydroxyl group. Based on the recent reports that the most common

transformation by cytochrome P450 enzymes in PKS biosynthesis is C-H hydroxylation (45) we
suggest LasP to be oxidizing C-31. Another accessory protein, the enoylreductase (ER) domain
LasQ (46) is proposed to be acting in *trans* as observed in other pathways including lagriamide
(47), patellazoles (48) and bacillaene (24, 49).

205 Due to the disruption of the catalytically active residues (CHH, Fig. S3F), we predict certain KS domains to be inactive (LasL KS1, LasL KS4, LasM KS5 and LasO KS7). We propose that the 206 207 ACP domain of LasL directly takes the molecule from the first ACP of LasJ and thus we predict 208 the KS domain in LasJ to be catalytically inactive despite the presence of catalytic residues, as 209 observed in lagriamide, lankacidin, and etnangien pathways (24, 36). Likewise, the alignment of 210 ketoreductase (KR) domains (Fig. S3G) allowed us to identify the ones lacking the KSY 211 catalytic triad and thus spot the inactive KR domain in module 2 (LasL KR1). Additionally, it 212 was found that the predicted stereoconfiguration of KR products (50, 51) in the las BGC, 213 matched the configuration of the equivalent moieties within the LSA structure produced by total 214 synthesis (8). The absence of a KR domain required in module 14 is proposed to be compensated 215 by a *trans*-acting KR likely from the following module as proposed in the patellazole (48) 216 pathway.

We were able to identify two pyran synthase (PS) domains (in module 7 and module 13) based on their phylogeny (**Fig. S4A**) and alignment (52, 53) (**Fig. S4B**). These PS domains are at the correct position in the *las* BGC to insert the pyran rings required to synthesize LSA. Even though module 13 lacks a DH domain required for pyran ring formation, we predict this role to be played by a *trans*-acting DH domain as commonly seen in *trans*-AT PKS pathways (24). Similarly, we were able to identify double bond-shifting DH domains in module 4 (LasM DH1) and 8 (LasN DH1) by the absence of both proline at the HxxxGxxxxP motif and of Glutamine/

10

224 Histamine at the DxxQ/H (Fig. S4C) (54) motif. Moreover, alignment of the DH domains 225 allowed us to identify the presence of inactive DH domains in module 2 (LasH DH1) and 6 226 (LasM DH2) by the absence of catalytic histidine at the HxxxGxxxxP motif and catalytic 227 aspartic acid at the DxxxQ/H motif (Fig. S4D). LasL DH3 has a serine in place of proline in its 228 HxxxGxxxxP motif. Alignment with different DH domains with serine in the HxxxGxxxxP 229 motif revealed a mixture of domains annotated as active and inactive (Fig. S4E). The majority of 230 times, when the DH domain had the conserved histidine in the HxxxGxxxxP motif it was 231 annotated as active. Based on this we propose LasL DH3 to be active. Specific primers were 232 designed based on the las BGC sequence and used to identify additional fosmids so that the 233 whole pathway could be assembled from five overlapping clones for future heterologous 234 expression.

235 For the biosynthesis of other LSA analogs, we propose that all of them except for lasonolide D 236 are modified post PKS (Fig. 4). The cytochrome P450 LasP is predicted to oxidize LSA at C-37 237 and C-36 leading to the synthesis of lasonolide B and C respectively. However, in the complete 238 biosynthesis of lasonolide B it is unclear how the methyl group is transferred from C-38 in LSA 239 to C-36 in lasonolide B. Recently it was shown that serine hydrolase activity of lipid droplet-240 associated hydrolase is responsible for cleaving the ester bond in LSA and yielding the active 241 form of the molecule, i.e. lasonolide F (55). Due to its hydrophobicity LSA is able to easily 242 diffuse into the plasma membrane and into lipid droplets, where it is converted into lasonolide F, 243 which is more hydrophilic and therefore able to diffuse out of the lipid droplet and into the 244 cytoplasm to exhibit its cytotoxic effect (55). Lasonolide C seems to undergo an esterification 245 reaction with a long-chain fatty acid ($CH_3(CH_2)_{10}COOH$) to produce lasonolide G. We suggest 246 that lasonolide E is also biosynthesized by a trans-esterification reaction, by reacting with an

ethanol molecule. As with the production of LSA, we suggest that for the biosynthesis of lasonolide D the molecule passes through the entire *las* BGC, however, the starter molecule in this case is an acetate instead of a malonate that gets loaded on the ACP of LasJ, with LasH and LasI being inactive.

251 Multiple repeats of the *las* BGC

252 The k-mer coverage of the las BGC (400.165× for las BGC v1 and 159.02× for las BGC v2) is 253 roughly three times that of "Ca. T. lasonolidus" (135.16x in Forcepia v1 and 48.24x in 254 Forcepia v2). The $3 \times$ coverage suggests three repeats of the putative *las* BGC. Visual inspection 255 of the assembly graph as well as mapping of the paired-end reads onto "Ca. T. lasonolidus" 256 allowed us to identify three connections on the 3' end of las BGC but only two connections on 257 the 5' end of the pathway (contig 7 and 8) (Fig. 5 and Table S3). Another contig (contig 5) was 258 observed to be connected to *las* BGC about 3 kbp (3.6 kbp for *las* BGC_v1 and 3.7 kbp for *las* 259 BGC_v2) from the $5\square$ end of *las* BGC. This suggests that the majority of *las* BGC (about 260 98kbp) is repeated thrice with a 3 kbp segment of the pathway (contig 6) is repeated twice (Fig. 261 5). The two repeats of contig 6 were further verified by more than twice paired end reads 262 mapping to it as compared to contig 5 (56) as well as its $2 \times$ coverage when compared to "Ca. T. 263 lasonolidus". All the connections between the las BGC and the bacterial genome were verified 264 using PCR and Sanger sequencing of the amplicons. We believe that the three repeats of the *las* 265 BGC might be involved in contributing to the increased expression of LSA through increased 266 gene dosage (57).

267 On comparing the three repeats it was observed that the *las* BGC repeat connected to contig 5 268 lacks *lasC* (ACP domain, highlighted area in **Fig. 5**), which is predicted to play an important role

269 in β -branch formation. Furthermore, the same repeat which lacks *lasC* also shows the presence 270 of an incomplete *lasD* (decarboxylating KS domain used in β -branching). Although this KS 271 domain has the catalytic active site residues SHH, characteristic of decarboxylating KSs (38), it 272 lacks about 47 amino acids that are present in the KS domain of the other two repeats connected 273 to contig 6. On further investigation with GATK HaplotypeCaller (58, 59) we were able to detect 274 three insertions and two single nucleotide polymorphisms (SNPs) between the three repeats of 275 las BGC v1 (Fig. 6 and Table 1). This was further supported by the allelic depth (AD) -276 informative reads supporting each allele - and phred-scaled likelihoods (PL) of the possible 277 genotypes. The genotype quality (GQ) which represents the confidence in the PL values was 99 278 for all five variants, which is the maximum value GATK reports for GQ. Furthermore, alignment 279 of las BGC_v1 with las BGC_v2 revealed that las BGC_v2 contains all the variants that were 280 called by GATK, thus further supporting their presence. All the three insertions are multiples of 281 three base pairs (60bp, 24bp and 54bp), and are thus not causing any frame-shift mutations. 282 Moreover, all the three insertions lie between *trans*-AT PKS domains, suggesting they do not 283 contribute to functional differences. Change in one base from G to A at 93,995 bp does not result 284 in a change in the amino acid sequence as both codons (TAT and TAC) encode for tyrosine. 285 Finally, the change in base from A to G at 95,154 bp lies just outside *lasS*, i.e. in the non-coding 286 region. The above-mentioned differences in the three repeats of *las* BGC indicate that the repeats 287 have been present long enough to allow divergence. However, the differences between the three 288 repeats are not predicted to affect the function of the las BGC.

289 Evidence for horizontal gene transfer

290 During the binning process by Autometa (60), Barnes-Hut Stochastic Neighbor Embedding (BH-291 tSNE) was used to reduce 5-mer frequencies to two dimensions. Generally, contigs belonging to 292 the same genome would have similar 5-mer frequency and would be expected to cluster close to 293 each other (61, 62). Visualization of the dimension-reduced data (Fig. 7A-B and S5A-B), 294 revealed that the las BGC contigs significantly differ in their 5-mer frequency from "Ca. T. 295 lasonolidus", suggesting that the las BGC could have been recently horizontally acquired. 296 Furthermore, the GC% of the *las* BGC is significantly different (p < 0.05, ANOVA followed by 297 Tukey HSD) from annotated, hypothetical and pseudogenes (Fig. 7C and S5C) providing further 298 evidence for horizontal transfer of the las BGC.

Codon adaptation index (CAI) compares the synonymous codon usage of a gene and that of a reference set along with measuring the synonymous codon usage bias (63). The CAI for the *las* BGC was significantly different (p < 0.05, ANOVA followed by Tukey HSD) from the annotated, hypothetical and pseudogenes, but matched that of highly expressed genes (i.e. ribosomal proteins) (**Fig. S5D-E**). Thus, despite its horizontal acquisition, the BGC's codon usage has been adapted to be efficiently translated even though the 5-mer composition is still different when compared to the rest of the "*Ca*. T. lasonolidus" genome.

306 The genome of the putative lasonolide producing symbiont

307 "*Ca.* T. lasonolidus", with multiple *las* BGC repeats, represents an important addition to the 308 growing collection of symbiotic Verrucomicrobia ("*Candidatus* Didemnitutus mandela" and 309 "*Candidatus* Synoicihabitans palmerolidicus") being identified with repeated *trans*-AT PKS

310 BGCs (57, 64, 65). Recently, two simultaneous studies have also identified a trans-AT PKS 311 BGC for pateamine in a bacterium ("Candidatus Patea custodiens") belonging to phylum 312 Kiritimatiellaeota (66, 67), a recently proposed phylum which was previously classified within 313 Verrucomicrobia (68). This highlights the importance of this understudied phylum as an 314 important producer of natural products. "Ca. T. lasonolidus" is little over 5 Mbp long and has 315 GC percentage of about 53%. It is estimated to be 99% complete, is 1.35% contaminated (69), 316 has tRNAs for all amino acids and complete 5S, 16S and 23S rRNA genes. Based on MIMAG 317 standards (25) the bin is classified as a high-quality MAG. Detailed statistics of the putative LSA 318 producer are provided in **Table 2**.

Eukaryotic-like proteins (ELPs) are known to be present in genomes of sponge symbionts and have been found to play an important role in regulating their interaction with the host sponge (70–73). It is hypothesized that interaction with ELPs allow the symbiotic bacteria to evade phagocytosis by the sponge, thus allowing discrimination between food and symbiont bacteria (72, 74). A number of ELPs were identified in "*Ca*. T. lasonolidus" (**Table 3, and Table S4A**), thus suggesting a symbiotic relationship of the bacterium with the *Forcepia* sp.

325 Bacterial microcompartments (BMCs) are organelles that enclose enzymes within a selectively 326 permeable proteinaceous shell (75), and they are rare among bacteria. Members of the phyla 327 Planctomycetes and Verrucomicrobia have a unique BMC gene cluster called the 328 Planctomycetes-Verrucomicrobia bacterial microcompartment (PV BMC) which is responsible 329 for production of microcompartment shell proteins BMC-P and BMC-H, as well as degradation 330 of L-rhamnose, L-fucose and fucoidans (71, 76, 77). Genes encoding the PV BMC cluster were 331 identified in "Ca. T. lasonolidus" (Table S4B), and the respective gene clusters in Forcepia v1 332 "Ca. T. lasonolidus" and Forcepia_v2 "Ca. T. lasonolidus" were compared using clinker (28)

333 and were found to be 100% identical to each other. One interesting finding was that the 334 identified PV BMC clusters had a DNA-methyltransferase and PVUII endonuclease gene 335 between the first and the second BMC-H genes. This is different from the usual arrangement of 336 the PV BMC gene cluster where both the BMC-H genes lie next to each other and the cluster 337 lacks DNA-methyltransferase and PVUII endonuclease genes (Fig. 8). The presence of PV BMC 338 genes in the "Ca. T. lasonolidus" genome suggests that it possesses bacterial microcompartments 339 and that they might be involved in L-fucose and L-rhamnose degradation. Despite repeated 340 attempts, we only found rhamnulokinase and fumarylacetoacetate hydrolase family proteins in 341 the "Ca. T. lasonolidus" genomes and failed to identify other complementary enzymes involved 342 in the degradation of L-fucose and L-rhamnose. However, other enzymes involved in 343 carbohydrate metabolism including glycoside hydrolases, carbohydrate binding module, 344 polysaccharide lyase, carbohydrate esterases and glycoside transferase were detected (Table 345 **S4C**) indicating that "*Ca*. T. lasonolidus" is capable of polysaccharide degradation, something 346 that is observed in a number of marine Verrucomicrobia (78–80).

347 A characteristic of obligate host-symbiont relationships is the loss of symbiont genes which are 348 required for independent survival. The early stages of genome reduction are characterized by 349 reduced coding density, and a high number of pseudogenes (81–83). We compared "Ca T. 350 with its closest free-living relative - Pedosphaera parvula Ellin514 lasonolidus" 351 (GCA_000172555.1). The draft genome of P. paruva Ellin514 is 7.41Mbp long, about 2.2 Mbp 352 longer than "Ca. T. lasonolidus". Furthermore, in P. paruva Ellin514 only 0.5% of total ORFs 353 were found to be pseudogenes (57, 84, 85) as opposed to about 16% in "Ca. T. lasonolidus" (Fig. 354 **9A-B**). Another indication of ongoing genome reduction comes from the fact that a much smaller 355 percentage of genes with annotated function were identified in "Ca. T. lasonolidus" as compared

to P. paruva Ellin514 (Fig. 9C), perhaps indicating sequence degradation and divergence from 356 357 functionally-annotated genes. Moreover, when compared with P. paruva Ellin514, "Ca. T. 358 lasonolidus" lacks genes involved in DNA repair, DNA replication, chemotaxis and nucleotide 359 metabolism (Fig. 9D), a trend which is commonly observed in symbionts undergoing genome 360 reduction (81). However, "Ca. T. lasonolidus" contains most of the primary metabolic pathways 361 (Fig. 9E) when compared to *P. paruva* Ellin514, and has a fairly large genome to be classified as 362 reduced. Based on the above evidence we suggest that "Ca. T. lasonolidus" is in early stages of 363 genome reduction. This hypothesis is also supported by its low coding density of \sim 72% (without 364 pseudogenes), relative to the average coding density of 85-90% for free-living bacteria (81) 365 which suggests a recent transitional event, such as host restriction (81).

Due to its potency and unique mechanism of action, LSA is considered a potential anti-cancer drug lead; however, its limited supply has hampered its transition to clinical trials. The evidence provided here suggests that LSA is synthesized by a yet uncultured Verrucomicrobial symbiont, which harbors three copies of the putative *las* BGC. The detailed analysis of the biosynthetic scheme, genome characteristics of the putative producer as well as the capture of the *las* BGC on a plasmid will aid future cultivation and heterologous expression efforts.

372 Acknowledgments

The authors acknowledge Amy Wright and Shirley Pomponi for providing sponge specimens, and Amy Wright, Shirley Pomponi and Peter McCarthy for valuable discussions during the project. The authors also thank Samantha C. Waterworth for fruitful discussions, John Barkei for discussions on heterologous expression strategies, and Chase Clark for providing DH sequences for sequence comparisons.

378 Funding

Support was provided by NCI (R21 CA209189) and a Start-up Fund from Harbor Branch Oceanographic Institute Foundation. The sample used in the study was collected with funds from a grant from the State of Florida Board of Education awarded to Florida Atlantic University for the Center of Excellence in Biomedical and Marine Biotechnology. This material is based upon work supported by the National Science Foundation under Grant No. DBI 1845890.

384 Data availability

The data associated with this study is deposited under BioProject RJNA833117. The WGS reads have been deposited in the sequence read archive (SRA) with accessions SRR18966768 (Forcepia_v1) and SRR18966767 (Forcepia_v2).

388 Methods

389 For full details see **Text S1**.

390 Sponge collection

Forcepia sp. (class, Demospongiae; order, Poecilosclerida; family, Coelosphaeridae) was
collected in August of 2005 using the HBOI Johnson Sea Link submersible. Samples were
collected at a depth of 70m from the Gulf of Mexico (26.256573N, 83.702772W) on the Pulley
Ridge (https://shinyapps.fau.edu/app/bmr). The sponge samples were immediately frozen at
-80°C. The sample ID was 12-VIII-05-1-006 200508121006 2005-08-12 JSL I-4837 (HBOI) *Forcepia* sp. 131921.

397 DNA purification and sequencing

398 The sponge hologenome was extracted using a modified cetyl trimethylammonium bromide 399 (CTAB) DNA extraction method (87) and then size-fractionated by low melting point gel 400 electrophoresis. DNA fragments greater than 40 kb were recovered from the gel and used for 401 fosmid library preparation (Text S1) as well as metagenomic sequencing. Two rounds of 402 sequencing were performed for different DNA extracts from the Forcepia sp. sponge. For the 403 first round (referred to as Forcepia_v1) Illumina TruSeq DNA libraries were prepared and 404 sequenced by RTL Genomics using an Illumina MiSeq sequencer giving us 108 million paired-405 end reads with length of 151bp. For the second round of sequencing (referred to as Forcepia v^{2}) 406 Illumina Nextera libraries were prepared and sequenced using a NovaSeq 6000 sequencer giving 407 us 303 million paired-end reads with length of 150 bp. Fosmids were sequenced by RTL 408 Genomics and Genewiz.

409 Identification and annotation of *las* BGC

Identification of the *las* BGC was done using tBLASTN (26), where KS domains from different *trans*-AT PKS pathways were used as a query against the metagenomic assembly (assembled using MetaSpades (88), see **Text S1**). Genes for each bin were called and annotated using Prokka v1 (89, 90). MetaSpades contig headers were replaced by their respective Prokka headers to maintain consistency with the annotation file submitted to NCBI. Genes on contigs making up the *las* BGC were not called correctly by Prokka (89, 90) and were thus annotated manually in Artemis (91) with the help of AntiSMASH (27), CDD (92) and SMART (93, 94).

417 Functional analysis of the "*Ca*. T. lasonolidus" genome

418 Genes called using Prokka v1 were used for all functional analysis (89, 90). PV-BMC clusters 419 were identified in "Ca. T. lasonolidus" using Interproscan v5.52-86.0 (95) and CDD (92). Initial 420 identification of ELPs was done using Diamond BLASTP against the diamond-formatted nr 421 database (using -k 1 --max-hsps 1 options) (96) and Interproscan v5.52-86.0 (95). This was 422 followed by verification of non-pseudogenes using CDD (92). Enzymes involved in 423 carbohydrate metabolism were detected using dbCAN2 (97) where genes annotated by ≥ 2 tools 424 (out of HMMER, Diamond and Hotpep) were kept. COG categories were identified using the 425 eggNOG mapper online server (98, 99).

The genome of *P. paruva* Ellin514 was downloaded from Genbank (GCA_000172555.1) and genes were called and annotated using Prokka v1 (89, 90). Primary metabolic pathways were identified for non-pseudogenes with kofamscan using the --mapper flag (100) and annotated against the KEGG database (101–103). The matrix with presence/absence of different enzymes was constructed in RStudio (104). Completeness of metabolic pathways was identified using KEGG-Decoder (86).

432 References

- Horton PA, Koehn FE, Longley RE, McConnell OJ. 1994. Lasonolide A, a new cytotoxic
 macrolide from the marine sponge *Forcepia* sp. J Am Chem Soc 116:6015–6016.
- 435 2. Wright AE, Chen Y, Winder PL, Pitts TP, Pomponi SA, Longley RE. 2004. Lasonolides C436 G, five new lasonolide compounds from the sponge *Forcepia* sp. J Nat Prod 67:1351–1355.
- 437 3. Isbrucker RA, Guzmán EA, Pitts TP, Wright AE. 2009. Early effects of lasonolide A on

438 pancreatic cancer cells. J Pharmacol Exp Th	her 331:733–739.
---	------------------

	439	4.	Jossé R.	Zhang	Y-W.	Giroux V	7. C	Shosh AK.	Luo J.	Pommier	Y.	2015.	Activation	of l	RA	١	71
--	-----	----	----------	-------	------	----------	------	-----------	--------	---------	----	-------	------------	------	----	---	----

- 440 (c-RAF) by the marine alkaloid lasonolide A induces rapid premature chromosome
- 441 condensation. Mar Drugs 13:3625–3639.
- 442 5. Zhang Y-W, Ghosh AK, Pommier Y. 2012. Lasonolide A, a potent and reversible inducer
 443 of chromosome condensation. Cell Cycle 11:4424–4435.
- 444 6. Yang L, Lin Z, Shao S, Zhao Q, Hong R. 2018. An enantioconvergent and concise synthesis

of lasonolide A. Angew Chem Int Ed Engl 57:16200–16204.

- Yang L, Lin Z, Shao S, Zhao Q, Hong R. 2019. Corrigendum: An enantioconvergent and
 concise synthesis of lasonolide A. Angew Chem Int Ed Engl 58:4431.
- 8. Trost BM, Stivala CE, Fandrick DR, Hull KL, Huang A, Poock C, Kalkofen R. 2016. Total
 synthesis of (-)-lasonolide A. J Am Chem Soc 138:11690–11701.

450 9. Piel J. 2009. Metabolites from symbiotic bacteria. Nat Prod Rep 26:338–362.

- 451 10. Lopanik NB. 2014. Chemical defensive symbioses in the marine environment. Funct Ecol
 452 28:328–340.
- 453 11. Flórez LV, Biedermann PHW, Engl T, Kaltenpoth M. 2015. Defensive symbioses of
 454 animals with prokaryotic and eukaryotic microorganisms. Nat Prod Rep 32:904–936.
- 455 12. Oliver KM, Smith AH, Russell JA. 2014. Defensive symbiosis in the real world advancing

456 ecological studies of heritable, protective bacteria in aphids and beyond. Funct Ecol

457 28:341–355.

- 458 13. Bodor A, Bounedjoum N, Vincze GE, Erdeiné Kis Á, Laczi K, Bende G, Szilágyi Á,
- Kovács T, Perei K, Rákhely G. 2020. Challenges of unculturable bacteria: Environmental
 perspectives. Rev Environ Sci Technol 19:1–22.
- 461 14. Hofer U. 2018. The majority is uncultured. Nat Rev Microbiol 16:716–717.
- 462 15. Vartoukian SR, Palmer RM, Wade WG. 2010. Strategies for culture of "unculturable"
 463 bacteria. FEMS Microbiol Lett 309:1–7.
- 16. Stevens DC, Hari TPA, Boddy CN. 2013. The role of transcription in heterologous

465 expression of polyketides in bacterial hosts. Nat Prod Rep 30:1391–1411.

- 466 17. Trindade M, van Zyl LJ, Navarro-Fernández J, Abd Elrazak A. 2015. Targeted
- 467 metagenomics as a tool to tap into marine natural product diversity for the discovery and
- 468 production of drug candidates. Front Microbiol 6:890.
- 18. Nivina A, Yuet KP, Hsu J, Khosla C. 2019. Evolution and diversity of assembly-line
 polyketide synthases. Chem Rev 119:12524–12547.
- 471 19. Cox RJ. 2007. Polyketides, proteins and genes in fungi: Programmed nano-machines begin
 472 to reveal their secrets. Org Biomol Chem 5:2010–2026.
- 473 20. Jenke-Kodama H, Sandmann A, Müller R, Dittmann E. 2005. Evolutionary implications of
 474 bacterial polyketide synthases. Mol Biol Evol 22:2027–2039.
- 475 21. Calderone CT, Kowtoniuk WE, Kelleher NL, Walsh CT, Dorrestein PC. 2006.
- 476 Convergence of isoprene and polyketide biosynthetic machinery: Isoprenyl-S-carrier
- 477 proteins in the *pksX* pathway of *Bacillus subtilis*. Proc Natl Acad Sci U S A 103:8977–

478 8982.

- 479 22. Calderone CT. 2008. Isoprenoid-like alkylations in polyketide biosynthesis. Nat Prod Rep
 480 25:845–853.
- 481 23. Gu L, Wang B, Kulkarni A, Geders TW, Grindberg RV, Gerwick L, Håkansson K, Wipf P,
 482 Smith JL, Gerwick WH, Sherman DH. 2009. Metamorphic enzyme assembly in polyketide
 483 diversification. Nature 459:731–735.
- 484 24. Helfrich EJN, Piel J. 2016. Biosynthesis of polyketides by *trans*-AT polyketide synthases.
 485 Nat Prod Rep 33:231–316.
- 486 25. Bowers RM, Kyrpides NC, Stepanauskas R, Harmon-Smith M, Doud D, Reddy TBK,
- 487 Schulz F, Jarett J, Rivers AR, Eloe-Fadrosh EA, Tringe SG, Ivanova NN, Copeland A,
- 488 Clum A, Becraft ED, Malmstrom RR, Birren B, Podar M, Bork P, Weinstock GM, Garrity
- 489 GM, Dodsworth JA, Yooseph S, Sutton G, Glöckner FO, Gilbert JA, Nelson WC, Hallam
- 490 SJ, Jungbluth SP, Ettema TJG, Tighe S, Konstantinidis KT, Liu W-T, Baker BJ, Rattei T,
- 491 Eisen JA, Hedlund B, McMahon KD, Fierer N, Knight R, Finn R, Cochrane G, Karsch-
- 492 Mizrachi I, Tyson GW, Rinke C, Lapidus A, Meyer F, Yilmaz P, Parks DH, Murat Eren A,
- 493 Schriml L, Banfield JF, Hugenholtz P, Woyke T. 2017. Minimum information about a
- 494 single amplified genome (MISAG) and a metagenome-assembled genome (MIMAG) of
- 495 bacteria and archaea. Nat Biotechnol 35:725–731.
- 26. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009.
 BLAST+: Architecture and applications. BMC Bioinform 10:1–9.
- 498 27. Blin K, Shaw S, Steinke K, Villebro R, Ziemert N, Lee SY, Medema MH, Weber T. 2019.

499	antiSMASH 5.0: Updates to the secondary metabolite genome mining pipeline. Nucl	eic
500	Acids Res 47:W81–W87.	

- 501 28. Gilchrist CLM, Chooi Y-H. 2021. Clinker & clustermap.js: Automatic generation of gene
 502 cluster comparison figures. Bioinformatics 37:2473–2475.
- 503 29. Wick RR, Schultz MB, Zobel J, Holt KE. 2015. Bandage: Interactive visualization of *de*504 *novo* genome assemblies. Bioinformatics 31:3350–3352.
- 505 30. Chaumeil P-A, Mussig AJ, Hugenholtz P, Parks DH. 2020. GTDB-Tk: A toolkit to classify

genomes with the Genome Taxonomy Database. Bioinformatics 36:1925–1927.

507 31. Yarza P, Yilmaz P, Pruesse E, Glöckner FO, Ludwig W, Schleifer K-H, Whitman WB,

- Euzéby J, Amann R, Rosselló-Móra R. 2014. Uniting the classification of cultured and
 uncultured bacteria and archaea using 16S rRNA gene sequences. Nat Rev Microbiol
 12:635–645.
- - 511 32. Navarro-Muñoz JC, Selem-Mojica N, Mullowney MW, Kautsar SA, Tryon JH, Parkinson
 - 512 EI, De Los Santos ELC, Yeong M, Cruz-Morales P, Abubucker S, Roeters A, Lokhorst W,
 - 513 Fernandez-Guerra A, Cappelini LTD, Goering AW, Thomson RJ, Metcalf WW, Kelleher
 - 514 NL, Barona-Gomez F, Medema MH. 2020. A computational framework to explore large-
 - scale biosynthetic diversity. Nat Chem Biol 16:60–68.
 - 516 33. Nguyen T, Ishida K, Jenke-Kodama H, Dittmann E, Gurgui C, Hochmuth T, Taudien S,
 - 517 Platzer M, Hertweck C, Piel J. 2008. Exploiting the mosaic structure of *trans*-
 - 518 acyltransferase polyketide synthases for natural product discovery and pathway dissection.
 - 519 Nat Biotechnol 26:225–233.

520	34.	Jensen K, Niederkrüger H, Zimmermann K, Vagstad AL, Moldenhauer J, Brendel N, Frank
521		S, Pöplau P, Kohlhaas C, Townsend CA, Oldiges M, Hertweck C, Piel J. 2012. Polyketide
522		proofreading by an acyltransferase-like enzyme. Chem Biol 19:329–339.
523	35.	Jenner M, Afonso JP, Kohlhaas C, Karbaum P, Frank S, Piel J, Oldham NJ. 2016. Acyl
524		hydrolases from <i>trans</i> -AT polyketide synthases target acetyl units on acyl carrier proteins.
525		Chem Commun 52:5262–5265.
526	36.	Piel J. 2010. Biosynthesis of polyketides by <i>trans</i> -AT polyketide synthases. Nat Prod Rep
527		27:996–1047.
528	37.	Haines AS, Dong X, Song Z, Farmer R, Williams C, Hothersall J, Płoskoń E, Wattana-
529		Amorn P, Stephens ER, Yamada E, Gurney R, Takebayashi Y, Masschelein J, Cox RJ,
530		Lavigne R, Willis CL, Simpson TJ, Crosby J, Winn PJ, Thomas CM, Crump MP. 2013. A
531		conserved motif flags acyl carrier proteins for β -branching in polyketide synthesis. Nat
532		Chem Biol 9:685–692.
533	38.	Walker PD, Weir ANM, Willis CL, Crump MP. 2021. Polyketide β-branching: Diversity,
534		mechanism and selectivity. Nat Prod Rep 38:723–756.
535	39.	Slocum ST, Lowell AN, Tripathi A, Shende VV, Smith JL, Sherman DH. 2018.
536		Chemoenzymatic dissection of polyketide β -branching in the bryostatin pathway, p. 207–
537		236. In Moore, BS (ed.), Marine Enzymes and Specialized Metabolism - Part A. Academic
538		Press.
539	40.	Gu L, Jia J, Liu H, Håkansson K, Gerwick WH, Sherman DH. 2006. Metabolic coupling of
540		dehydration and decarboxylation in the curacin A pathway: Functional identification of a

541 mechanistically diverse enzyme pair. J Am Chem Soc 128:9014	-9015.
---	--------

542	41.	Matilla MA, Stöckmann H, Leeper FJ, Salmond GPC. 2012. Bacterial biosynthetic gene
543		clusters encoding the anti-cancer haterumalide class of molecules: Biogenesis of the broad
544		spectrum antifungal and anti-oomycete compound, oocydin A. J Biol Chem 287:39125-
545		39138.
546	42.	Kačar D, Cañedo LM, Rodríguez P, González EG, Galán B, Schleissner C, Leopold-Messer
547		S, Piel J, Cuevas C, de la Calle F, García JL. 2021. Identification of <i>trans</i> -AT polyketide
548		clusters in two marine bacteria reveals cryptic similarities between distinct symbiosis
549		factors. Environ Microbiol 23:2509–2521.
550	43.	Meoded RA, Ueoka R, Helfrich EJN, Jensen K, Magnus N, Piechulla B, Piel J. 2018. A
551		polyketide synthase component for oxygen insertion into polyketide backbones. Angew
552		Chem Int Ed Engl 57:11644–11648.
553	44.	Hemmerling F, Meoded RA, Fraley AE, Minas HA, Dieterich CL, Rust M, Ueoka R,
554		Jensen K, Helfrich EJN, Bergande C, Biedermann M, Magnus N, Piechulla B, Piel J. 2022.
555		Modular halogenation, α -hydroxylation, and acylation by a remarkably versatile polyketide
556		synthase. Angew Chem Int Ed Engl 61:e202116614.
557	45.	Greule A, Stok JE, De Voss JJ, Cryle MJ. 2018. Unrivalled diversity: The many roles and
558		reactions of bacterial cytochromes P450 in secondary metabolism. Nat Prod Rep 35:757-
559		791.

560 46. Bumpus SB, Magarvey NA, Kelleher NL, Walsh CT, Calderone CT. 2008. Polyunsaturated
561 fatty-acid-like *trans*-enoyl reductases utilized in polyketide biosynthesis. J Am Chem Soc

26

562 130:11614–11616.

563	47.	Flórez LV, Scherlach K, Miller IJ, Rodrigues A, Kwan JC, Hertweck C, Kaltenpoth M.
564		2018. An antifungal polyketide associated with horizontally acquired genes supports
565		symbiont-mediated defense in Lagria villosa beetles. Nat Commun 9:2478.
566	48.	Kwan JC, Donia MS, Han AW, Hirose E, Haygood MG, Schmidt EW. 2012. Genome
567		streamlining and chemical defense in a coral reef symbiosis. Proc Natl Acad Sci USA
568		109:20655–20660.
569	49.	Chen X-H, Vater J, Piel J, Franke P, Scholz R, Schneider K, Koumoutsi A, Hitzeroth G,
570		Grammel N, Strittmatter AW, Gottschalk G, Süssmuth RD, Borriss R. 2006. Structural and
571		functional characterization of three polyketide synthase gene clusters in Bacillus
572		amyloliquefaciens FZB 42. J Bacteriol 188:4024-4036.
573	50.	Keatinge-Clay AT. 2007. A tylosin ketoreductase reveals how chirality is determined in
574		polyketides. Chem Biol 14:898–908.
575	51.	Caffrey P. 2003. Conserved amino acid residues correlating with ketoreductase
576		stereospecificity in modular polyketide synthases. ChemBioChem 4:654–657.
577	52.	Pöplau P, Frank S, Morinaka BI, Piel J. 2013. An enzymatic domain for the formation of
578		cyclic ethers in complex polyketides. Angew Chem Int Ed Engl 52:13215–13218.
579	53.	Wagner DT, Zhang Z, Meoded RA, Cepeda AJ, Piel J, Keatinge-Clay AT. 2018. Structural
580		and functional studies of a pyran synthase domain from a trans-acyltransferase assembly
581		line. ACS Chem Biol 13:975–983.

582	54.	Gay DC, Spear PJ, Keatinge-Clay AT. 2014. A double-hotdog with a new trick: Structure
583		and mechanism of the <i>trans</i> -acyltransferase polyketide synthase enoyl-isomerase. ACS
584		Chem Biol 9:2374–2381.

55. Dubey R, Stivala CE, Nguyen HQ, Goo Y-H, Paul A, Carette JE, Trost BM, Rohatgi R.
2020. Lipid droplets can promote drug accumulation and activation. Nat Chem Biol
16:206–213.

588 56. Albertsen M, Hugenholtz P, Skarshewski A, Nielsen KL, Tyson GW, Nielsen PH. 2013.

Genome sequences of rare, uncultured bacteria obtained by differential coverage binning of
multiple metagenomes. Nat Biotechnol 31:533–538.

- 57. Lopera J, Miller IJ, McPhail KL, Kwan JC. 2017. Increased biosynthetic gene dosage in a
 genome-reduced defensive bacterial symbiont. mSystems 2:e00096–17.
- 593 58. Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, Del Angel G, Levy-Moonshine A,
- Jordan T, Shakir K, Roazen D, Thibault J, Banks E, Garimella KV, Altshuler D, Gabriel S,
- 595 DePristo MA. 2013. From FastQ data to high-confidence variant calls: The genome analysis

toolkit best practices pipeline. Curr Protoc Bioinform 43:11.10.1–11.10.30.

597 59. DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, Philippakis AA, del

598 Angel G, Rivas MA, Hanna M, McKenna A, Fennell TJ, Kernytsky AM, Sivachenko AY,

- 599 Cibulskis K, Gabriel SB, Altshuler D, Daly MJ. 2011. A framework for variation discovery
- and genotyping using next-generation DNA sequencing data. Nat Genet 43:491–498.
- 60. Miller IJ, Rees ER, Ross J, Miller I, Baxa J, Lopera J, Kerby RL, Rey FE, Kwan JC. 2019.
- 602 Autometa: Automated extraction of microbial genomes from individual shotgun

603 metagenomes. Nucleic Acids Res 47:e57.

604	61.	Dick GJ, Andersson AF, Baker BJ, Simmons SL, Thomas BC, Yelton AP, Banfield JF.
605		2009. Community-wide analysis of microbial genome sequence signatures. Genome Biol
606		10:R85.
607	62.	Laczny CC, Pinel N, Vlassis N, Wilmes P. 2014. Alignment-free visualization of
608		metagenomic data by nonlinear dimension reduction. Sci Rep 4:4516.
609	63.	Sharp PM, Li W-H. 1987. The codon adaptation index-a measure of directional
610		synonymous codon usage bias, and its potential applications. Nucleic Acids Res 15:1281-
611		1295.
612	64.	Avalon NE, Murray AE, Daligault HE, Lo C-C, Davenport KW, Dichosa AEK, Chain PSG,
613		Baker BJ. 2021. Bioinformatic and mechanistic analysis of the palmerolide PKS-NRPS
614		biosynthetic pathway from the microbiome of an Antarctic ascidian. Front Chem 9:802574.
615	65.	Murray AE, Lo C-C, Daligault HE, Avalon NE, Read RW, Davenport KW, Higham ML,
616		Kunde Y, Dichosa AEK, Baker BJ, Chain PSG. 2021. Discovery of an Antarctic ascidian-
617		associated uncultivated Verrucomicrobia with antimelanoma palmerolide biosynthetic
618		potential. mSphere 6:e00759–21.
619	66.	Michael R, N. HEJ, F. FM, Pakjira N, M. FC, Christian R, Tomas K, J. PM, L. WV, Jörn K,
620		Shinichi S, Jörn P. 2020. A multiproducer microbiome generates chemical diversity in the
621		marine sponge Mycale hentscheli. Proc Natl Acad Sci U S A 117:9508–9518.

622 67. Storey MA, Andreassend SK, Bracegirdle J, Brown A, Keyzers RA, Ackerley DF,

623		Northcote PT, Owen JG. 2020. Metagenomic exploration of the marine sponge Mycale
624		hentscheli uncovers multiple polyketide-producing bacterial symbionts. mBio 11:e02997-
625		19.
626	68.	Spring S, Bunk B, Spröer C, Schumann P, Rohde M, Tindall BJ, Klenk H-P. 2016.
627		Characterization of the first cultured representative of Verrucomicrobia subdivision 5
628		indicates the proposal of a novel phylum. ISME J 10:2801–2816.
629	69.	Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. 2015. CheckM:
630		Assessing the quality of microbial genomes recovered from isolates, single cells, and
631		metagenomes. Genome Res 25:1043–1055.
632	70.	Burgsdorf I, Handley KM, Bar-Shalom R, Erwin PM, Steindler L. 2019. Life at home and
633		on the roam: Genomic adaptions reflect the dual lifestyle of an intracellular, facultative
634		symbiont. mSystems 4:e00057–19.
635	71.	Sizikov S, Burgsdorf I, Handley KM, Lahyani M, Haber M, Steindler L. 2020.
636		Characterization of sponge-associated Verrucomicrobia: Microcompartment-based sugar
637		utilization and enhanced toxin-antitoxin modules as features of host-associated Opitutales.
638		Environ Microbiol 22:4669–4688.
639	72.	Frank AC. 2019. Molecular host mimicry and manipulation in bacterial symbionts. FEMS
640		Microbiol Lett 366:fnz038.
641	73.	Díez-Vives C, Moitinho-Silva L, Nielsen S, Reynolds D, Thomas T. 2017. Expression of
642		eukaryotic-like protein in the microbiome of sponges. Mol Ecol 26:1432–1451.

30

643	74. Nguyen MTHD, Liu M, Thomas T. 2014. Ankyrin-repeat proteins from sponge symbion
644	modulate amoebal phagocytosis. Mol Ecol 23:1635–1645.

- 645 75. Kerfeld CA, Aussignargues C, Zarzycki J, Cai F, Sutter M. 2018. Bacterial
- 646 microcompartments. Nat Rev Microbiol 16:277–290.
- 647 76. Erbilgin O, McDonald KL, Kerfeld CA. 2014. Characterization of a planctomycetal
- 648 organelle: A novel bacterial microcompartment for the aerobic degradation of plant
- saccharides. Appl Environ Microbiol 80:2193–2205.
- 650 77. Sichert A, Corzett CH, Schechter MS, Unfried F, Markert S, Becher D, Fernandez-Guerra
- A, Liebeke M, Schweder T, Polz MF, Hehemann J-H. 2020. Verrucomicrobia use hundreds
- of enzymes to digest the algal polysaccharide fucoidan. Nat Microbiol 5:1026–1039.
- 653 78. Cardman Z, Arnosti C, Durbin A, Ziervogel K, Cox C, Steen AD, Teske A. 2014.
- 654 Verrucomicrobia are candidates for polysaccharide-degrading bacterioplankton in an arctic
- 655 fjord of Svalbard. Appl Environ Microbiol 80:3749–3756.
- 656 79. Martinez-Garcia M, Brazel DM, Swan BK, Arnosti C, Chain PSG, Reitenga KG, Xie G,
- 657 Poulton NJ, Lluesma Gomez M, Masland DED, Thompson B, Bellows WK, Ziervogel K,
- Lo C-C, Ahmed S, Gleasner CD, Detter CJ, Stepanauskas R. 2012. Capturing single cell
- 659 genomes of active polysaccharide degraders: An unexpected contribution of
- 660 Verrucomicrobia. PLoS One 7:e35314.
- 661 80. Herlemann DPR, Lundin D, Labrenz M, Jürgens K, Zheng Z, Aspeborg H, Andersson AF.
- 662 2013. Metagenomic *de novo* assembly of an aquatic representative of the verrucomicrobial
- class *Spartobacteria*. mBio 4:e00569–12.

664	81.	McCutcheon JP, Moran NA. 2011. Extreme genome reduction in symbiotic bacteria. Nat
665		Rev Microbiol 10:13–26.

- 666 82. Lo W-S, Huang Y-Y, Kuo C-H. 2016. Winding paths to simplicity: Genome evolution in
 667 facultative insect symbionts. FEMS Microbiol Rev 40:855–874.
- 668 83. Dietel A-K, Merker H, Kaltenpoth M, Kost C. 2019. Selective advantages favour high
- genomic AT-contents in intracellular elements. PLoS Genet 15:e1007778.
- 670 84. Lerat E, Ochman H. 2005. Recognizing the pseudogenes in bacterial genomes. Nucleic
- 671 Acids Res 33:3125–3132.
- 672 85. Waterworth SC, Flórez LV, Rees ER, Hertweck C, Kaltenpoth M, Kwan JC. 2020.
- 673 Horizontal gene transfer to a defensive symbiont with a reduced genome in a multipartite
- beetle microbiome. mBio 11:e02430–19.
- 675 86. Graham ED, Heidelberg JF, Tully BJ. 2018. Potential for primary productivity in a
 676 globally-distributed bacterial phototroph. ISME J 12:1861–1866.
- 677 87. Piel J, Hui D, Wen G, Butzke D, Platzer M, Fusetani N, Matsunaga S. 2004. Antitumor
- 678 polyketide biosynthesis by an uncultivated bacterial symbiont of the marine sponge
- 679 *Theonella swinhoei*. Proc Natl Acad Sci U S A 101:16222–16227.
- 88. Nurk S, Meleshko D, Korobeynikov A, Pevzner PA. 2017. metaSPAdes: A new versatile
 metagenomic assembler. Genome Res 27:824–834.
- 89. Seemann T. 2014. Prokka: Rapid prokaryotic genome annotation. Bioinformatics 30:2068–
 2069.

- 684 90. Hyatt D, Chen G-L, Locascio PF, Land ML, Larimer FW, Hauser LJ. 2010. Prodigal:
- 685 Prokaryotic gene recognition and translation initiation site identification. BMC Bioinform686 11:119.
- 687 91. Carver T, Harris SR, Berriman M, Parkhill J, McQuillan JA. 2012. Artemis: An integrated
 688 platform for visualization and analysis of high-throughput sequence-based experimental
 689 data. Bioinformatics 28:464–469.
- 690 92. Lu S, Wang J, Chitsaz F, Derbyshire MK, Geer RC, Gonzales NR, Gwadz M, Hurwitz DI,
- 691 Marchler GH, Song JS, Thanki N, Yamashita RA, Yang M, Zhang D, Zheng C, Lanczycki
- 692 CJ, Marchler-Bauer A. 2020. CDD/SPARCLE: The conserved domain database in 2020.
- 693 Nucleic Acids Res 48:D265–D268.
- 694 93. Letunic I, Khedkar S, Bork P. 2021. SMART: Recent updates, new developments and status
 695 in 2020. Nucleic Acids Res 49:D458–D460.
- 696 94. Letunic I, Bork P. 2018. 20 years of the SMART protein domain annotation resource.
 697 Nucleic Acids Res 46:D493–D496.
- 698 95. Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J,
- 699 Mitchell A, Nuka G, Pesseat S, Quinn AF, Sangrador-Vegas A, Scheremetjew M, Yong S-
- Y, Lopez R, Hunter S. 2014. InterProScan 5: Genome-scale protein function classification.
 Bioinformatics 30:1236–1240.
- 96. Buchfink B, Reuter K, Drost H-G. 2021. Sensitive protein alignments at tree-of-life scale
 using DIAMOND. Nat Methods 18:366–368.

704	97.	Zhang H, Yohe T, Huang L, Entwistle S, Wu P, Yang Z, Busk PK, Xu Y, Yin Y. 2018.
705		dbCAN2: A meta server for automated carbohydrate-active enzyme annotation. Nucleic
706		Acids Res 46:W95–W101.
707	98.	Cantalapiedra CP, Hernández-Plaza A, Letunic I, Bork P, Huerta-Cepas J. 2021. eggNOG-
708		mapper v2: Functional annotation, orthology assignments, and domain prediction at the
709		metagenomic scale. Mol Biol Evol 38:5825–5829.
710	99.	Huerta-Cepas J, Szklarczyk D, Heller D, Hernández-Plaza A, Forslund SK, Cook H, Mende
711		DR, Letunic I, Rattei T, Jensen LJ, von Mering C, Bork P. 2019. eggNOG 5.0: A
712		hierarchical, functionally and phylogenetically annotated orthology resource based on 5090
713		organisms and 2502 viruses. Nucleic Acids Res 47:D309–D314.
714	100.	Aramaki T, Blanc-Mathieu R, Endo H, Ohkubo K, Kanehisa M, Goto S, Ogata H. 2020.
715		KofamKOALA: KEGG ortholog assignment based on profile HMM and adaptive score
716		threshold. Bioinformatics 37:2251–2252.
717	101.	Kanehisa M, Sato Y. 2020. KEGG Mapper for inferring cellular functions from protein
718		sequences. Protein Sci 29:28–35.
719	102.	Kanehisa M, Sato Y, Kawashima M. 2021. KEGG mapping tools for uncovering hidden
720		features in biological data. Protein Sci 31:47–52.
721	103.	Kanehisa M, Goto S. 2000. KEGG: Kyoto encyclopedia of genes and genomes. Nucleic
722		Acids Res 28:27–30.
723	104.	Team R. 2020. RStudio: Integrated Development for R. RStudio, PBC, Boston, MA.

725 Figure Legends

Fig. 1. (A) Sponge *Forcepia* sp. as seen in the field. (B) The *Forcepia* sp. specimen used for
DNA extraction (sample ID: 12-VIII-05-1-6). Photo credit: HBOI Marine Biomedical and
Biotechnology Program. (C) The chemical structures of lasonolide A (LSA) and its analogs.

729 Fig. 2. (A) Collection site of Forcepia sp. sponge (dark red diamond, 26.256573N, 730 83.702772W). (B) Features in lasonolide A (LSA) characteristic of biosynthesis by a trans-AT 731 PKS pathway. (C) Relative abundance of different phyla in the sequenced Forcepia v1 732 metagenome. Each block shows the relative abundance of each metagenome-assembled genome 733 (MAG), with colors representing the phylum they belong to. The *las* biosynthetic gene cluster 734 (BGC)-carrying bin is marked with a star. (**D**) BGC distribution in *Forcepia* v1 sp. metagenome. 735 AntiSMASH (27) annotations of bacterial contigs greater than 500 bp are shown. Each bar 736 indicates a MAG. Bars have been grouped by phylum. The star represents the MAG possessing 737 las BGC. BGC annotations have been simplified into polyketide synthase (PKS), Type 1 PKS, 738 Type 3 PKS, trans-AT PKS, nonribosomal peptide synthetase (NRPS), ribosomally synthesized, 739 post-translationally modified peptide (RiPP), hgIE-KS, hgIE-KS-T1PKS, terpenes, RiPP-terpene 740 and others. (E) Comparison of las BGC_v1 and las BGC_v2 using clinker (28). V1 refers to las 741 BGC v1 while V2 refers to *las* BGC v2. Numbers in the boxes indicate amino acid identity as a 742 fraction of 1.

Fig. 3. Proposed LSA biosynthetic scheme. Colored lines above the *las* BGC represent an alignment of individual fosmids to the pathway, the fosmids were subsequently assembled together in a plasmid. A cross indicates a domain predicted to be catalytically inactive. Open reading frames colored in gray represent genes with unknown or no role in LSA synthesis.

36

747 Numbers below domains indicate the module number and 'A' and 'B' denote the predicted 748 stereoconfiguration of the KR product, as previously described (50, 51). Predicted substrate 749 specificity of KS domains, obtained through phylogeny (Fig. S3) (33), are shown above each 750 respective KS domain. Carbon 31 is highlighted in blue to represent the site where P450 LasP is 751 predicted to act. Abbreviations: ACP, acyl carrier protein, also denoted by a filled black circle; 752 AH, acylhydrolase; AT, acyltransferase; C, condensation; DH, dehydratase; ECH, enoyl-CoA 753 reductase; ER, enoylreductase; HMGS, 3-hydroxy-3-methylglutaryl-CoA synthase; KR, 754 ketoreductase; KS, ketosynthase; MOX, monooxygenase; PS, pyransynthase; P450, cytochrome 755 P450; THP, Tetrahydropyran.

Fig. 4. Proposed biosynthesis of different analogs from lasonolide A. We could not identify an
enzyme which would explain the migration of the methyl from C-38 in LSA to C-36 in
lasonolide B.

Fig. 5. Model for three repeats of the *las* BGC. The 5' end of *las* BGC is highlighted to demonstrate the location where one of the *las* BGC repeats lacks *lasC*. Contig(s) making up the 98 kbp segment of *las* BGC (one in *las* BGC_v1 and six in *las* BGC_v2) have been collectively referred to as contig 1. Contigs except for the *las* BGC are not shown to scale.

Fig. 6. Variants identified between the three repeats of *las* BGC_v1.

Fig. 7. (**A**) 2D visualization of Automata binning of Forcepia_v1. The "*Ca.* T. lasonolidus" genome is circled in black and the *las* BGC contigs are marked with an arrow. Axes represent dimension-reduced Barnes-Hut Stochastic Neighbor Embedding (BH-tSNE) values (BH-tSNE x and BH-tSNE y). (**B**) 3D visualization of contigs present in "*Ca.* T. lasonolidus". The *las* BGC is colored red. Axes represent BH-tSNE values (BH-tSNE x and BH-tSNE y) along with k-mer

coverage. (C) GC percentage of different sets of genes in Forcepia_v1 "*Ca*. T. lasonolidus". The *las* BGC genes are colored in red.

Fig. 8. Comparison of PV BMC gene cluster in Forcepia_v1 "*Ca.* T. lasonolidus" with the PV
BMC cluster from other Verrucomicrobia. "*Ca.* T. lasonolidus" has DNA-methyltransferase and
PVUII endonuclease genes (in gray, labeled 1 and 2) between the first and the second BMC-H
genes. This kind of arrangement was not observed in other PV BMC clusters.

775 Fig. 9. (A and B) Comparison of the gene length in (A) Forcepia v1 "Ca. T. lasonolidus" and 776 (B) Forcepia_v2 "Ca. T. lasonolidus", respectively with their closest homologs in the nr 777 database. Genes with length less than 80% of the closest homolog (below the lower black line) 778 are classified as putative pseudogenes (57, 84, 85). The graphs have been truncated for clarity, as 779 some genes are many kbp long. (C) Comparison of functional COG categories in Forcepia v1 780 "Ca. T. lasonolidus", Forcepia v2 "Ca. T. lasonolidus" and P. paruva Ellin514 for non-781 pseudogenes. A gene is considered to have a functional annotation when it belongs to a COG 782 category, except for category S which represents unknown function. (D) Comparison of genes in 783 different metabolic pathways for Forcepia_v1 "Ca. T. lasonolidus", Forcepia_v2 "Ca. T. 784 lasonolidus" and P. paruva Ellin514, including only non-pseudogenes. Colored squares represent 785 presence of a gene while white squares represent absence of gene. *K00940 is involved in both 786 purine and pyrimidine metabolism. Genes absent in all three genomes have been removed. (E) 787 Comparison of completeness of different metabolic pathways in Forcepia_v1 "Ca. T. 788 lasonolidus", Forcepia_v2 "Ca. T. lasonolidus" and P. paruva Ellin514 (including only non-789 pseudogenes) as determined by KEGG decoder (86). Pathways have been grouped into 790 categories wherever possible. Pathways absent in all three genomes have been removed. V1 and 791 V2 refer to Forcepia_v1 "Ca. T. lasonolidus" and Forcepia_v2 "Ca. T. lasonolidus" respectively.

Table 1. Description of the variants identified between the three repeats of *las* BGC_v1. Both

AD and PL values are represented in the manner "reference, variant". A lower PL value

represents a higher likelihood of the sample being that genotype.

ID	Location	Change	Length	Allelic depth	PL
	(bp)		(bp)		
Insertion 1	Between	+GGAGGATGGGGTGGAGGA	60	7, 15	1129, 0
	15,246 and	TGGGGTGGAGGATGGGGTG			
	15,247	GAGGATGGGGTGGAGGATG			
		GGGT			
Insertion 2	Between	+GGGGTCGGATGGGGGGGTC	24	8, 57	3092, 0
	24,776 and	GGATGG			
	24,777				
Insertion 3	Between	+GCGGCGGTTGAGGCGGAG	54	23, 164	5987, 0
	67,976 and	GCGGCGGTTGAGGCGGAGG			
	67,977	CGGCGGTTGAGGCGGAG			
SNP 1	93,995	$G \rightarrow A$	1	363, 457	2973, 0
SNP 2	95,154	$A \rightarrow G$	1	175, 621	16101, 0

Table 2. Genome statistics for "*Ca*. T. lasonolidus". *Coding density is weighted by length
taking into account the 97.11% coding density of *las* BGC repeats.

Characteristic	Forcepia_v1 "Ca. T.	Forcepia_v2 "Ca. T.			
	lasonolidus"	lasonolidus"			
Size	4.85 Mbp	4.93 Mbp			
Size after adding the three <i>las</i>	5.05 Mbp	5.13 Mbp			
repeats					
checkM completeness	99.24%	99.32%			
checkM contamination	1.35%	1.35%			
No. of contigs	144	92			
Longest contig	204,102 bp	649,894 bp			
N50	52,980 bp	96,223			
Average GC percentage	53.81%	53.88%			
Percentage of pseudogenes	16.31% of total ORFs	16.62% of total ORFs			
Transposase genes	6	15			
Coding density*	79.45%	79.41%			
Coding density without	72.58%	72.38%			
pseudogenes*					
Eukaryotic like proteins					

Ankyrin repeats	3	3
Tetratricopeptide repeat	43 (9 were Sel-1 repeats)	42 (9 were Sel-1 repeats)
Pyrrolo-quinoline quinone	21	21
Leucine-rich repeat	16	16
WD40	4	5

797

798

799 Legend for supplementary material

800 Fig S1. (A) Alignment of fosmids to the las BGC. Fosmids are depicted as arrows above the las 801 BGC. Fosmids captured before WGS are colored orange (3-46, 5-16, 6-17, 4-77 and 1-80), 802 whereas fosmids captured after WGS are colored blue (5-41, 2-18, and 2-13). (B) Relative 803 abundance of different phyla in the sequenced Forcepia v2 metagenome. Each block shows the 804 relative abundance of each metagenome-assembled genome (MAG), with different colors 805 representing the phylum they belong to. The las BGC-carrying bin is highlighted is marked with 806 a star. (C) Assembly graph of *las* BGC_v1 visualized in BANDAGE (Wick RR, Schultz MB, 807 J, KE, Zobel Holt **Bioinformatics** 31:3350-3352, 2015, 808 https://doi.org/10.1093/bioinformatics/btv383). (D) Mapping of paired-end reads to contigs 809 making up *las* BGC v1. Contigs in green boxes represent the *las* BGC, red boxes represent the 5' 810 end of las BGC and blue boxes represent the 3' end of las BGC. (E) Assembly of the seven 811 contigs making up las BGC_v2. (F) Assembly graph of las BGC_v2 visualized in BANDAGE 812 (Wick RR, Schultz MB, Zobel J, Holt KE, Bioinformatics 31:3350-3352, 2015,

https://doi.org/10.1093/bioinformatics/btv383). (G) Mapping of paired-end reads to las BGC_v2. 813 814 Contigs in green boxes represent the *las* BGC, red boxes represent the $5\Box$ end of *las* BGC and 815 blue boxes represent the $3\Box$ end of *las* BGC. Panels D, E, G, and H were edited for clarity by 816 removing contigs which had either very few paired-end read connections, were mapping to 817 themselves or were very small. (H) BGC distribution in Forcepia_v2 sp. Metagenome. 818 AntiSMASH (Blin K, Shaw S, Steinke K, Villebro R, Ziemert N, Lee SY, Medema MH, Weber 819 T, Nucleic Acid Res 47:W81–W87, 2019, https://doi.org/10.1093/nar/gkz310) annotations of 820 bacterial contigs greater than 3000bp are shown. Each bar indicates a MAG, grouped by phylum. 821 The star represents the MAG containing the las BGC. BGC annotations have been simplified 822 into polyketide synthase (PKS), Type 1 PKS, Type 3 PKS, trans-AT PKS, nonribosomal peptide 823 synthetase (NRPS), ribosomally synthesized, post-translationally modified peptide (RiPP), hgIE-824 KS, hgIE-KS-T1PKS, terpenes, and others. (I) Phylogenetic tree of 51 different Verrucomicrobia 825 genomes. Bootstrap values were calculated using RaxML with 1000 bootstrap replicates.

Fig S2. Clades from a phylogenetic tree of 944 KS domains from *trans*-AT PKS and the erythromycin BGC as an outgroup, containing KS domains from the *las* BGC. Color within the individual clades corresponds to the chemical structure shown on its right. 'i' and 'c' in *lasD* KS1 indicate the incomplete and complete KS domain in *lasD* respectively.

Fig S3. (A) Alignment of *las* AT and AH domains with AT and AH domains from different *trans*-AT PKS pathways. Active sites as well as sites distinguishing AT and AH domains (Jenner
M, Afonso JP, Kohlhaas C, Karbaum P, Frank S, Piel J, Oldham NJ, Chem Commun 52:5262–
5265, 2016, <u>https://doi.org/10.1039/C6CC01453D</u>) have been marked. (B) Phylogenetic tree of
AT and AH domains. The different types of domain separate into different clades (Jenner M,
Afonso JP, Kohlhaas C, Karbaum P, Frank S, Piel J, Oldham NJ, Chem Commun 52:5262–5265,

836 2016, https://doi.org/10.1039/C6CC01453D). (C) Phylogenetic tree of ECH1 and ECH2 837 domains. Both the domains separate into different clades (Slocum ST, Lowell AN, Tripathi A, 838 VV. Smith JL. Sherman DH. Methods Enzvmol Shende 604:207-236. 2018. 839 https://doi.org/10.1016/bs.mie.2018.01.034). (D-E) Alignment of las (D) ECH1 and (E) ECH2 840 domains with respective ECH domains from other PKS pathways. Sequence that is required for 841 the formation of the oxyanion hole which stabilizes the enolate anions is marked. LasH_a and 842 LasO a are proposed to be inactive as they are truncated and show poor homology to the rest of 843 the ECH domains (Gu L, Jia J, Liu H, Håkansson K, Gerwick WH, Sherman DH, J Am Chem 844 Soc 128:9014–9015, 2006, http://doi.org/10.1021/ja0626382, Matilla MA, Stöckmann H, Leeper FJ. Salmond J 845 GPC. Biol Chem 287:39125-39138, 2012. 846 https://doi.org/10.1074/jbc.M112.401026). (F) Alignment of las KS domains with active site 847 (CHH) marked. 'i' and 'c' in the LasD KS indicate the incomplete and complete KS domain in 848 different repeats of LasD respectively. LasD is a decarboxylating KS which are known to lack 849 the active site cysteine (Walker PD, Weir ANM, Willis CL, Crump MP, Nat Prod Rep 38:723-850 756, 2021, https://doi.org/10.1039/D0NP00045K). (G) Alignment of las KR domains with two 851 from the erythromycin BGC to allow comparison. Active site residues and conserved motifs are 852 marked. The presence or absence of the second aspartate in the LDD motif is supposed to predict 853 the stereochemistry of the hydroxyl group (Keatinge-Clay AT, Chem Biol 14:898–908, 2007, 854 https://doi.org/10.1016/j.chembiol.2007.07.009, Caffrey P, ChemBioChem 4:654-657, 2003, 855 https://doi.org/10.1002/cbic.200300581). Figures have been truncated for clarity and to show 856 only the relevant sites. In phylogenetic trees, *las* BGC domains are highlighted in white.

Fig S4. (A) Phylogenetic tree of DH and PS domains, which separate into different clades
(Wagner DT, Zhang Z, Meoded RA, Cepeda AJ, Piel J, Keatinge-Clay AT, ACS Chem Biol

43

859 13:975–983, 2018, http://doi.org/10.1021/acschembio.8b00049). Las BGC DH/PS domains are 860 highlighted in white. (B) Alignment of PS domains identified in the las BGC with PS domains from other trans-AT PKS pathways. The DH domain from the erythromycin BGC is used for 861 862 comparison. LasO DH2 and LasM DH4 are annotated as putative PS domains. Generally, PS 863 domains have a Hx₄P motif instead of a Hx₈P and they lack the catalytic aspartate at the 864 DxxxQ/H motif (Wagner DT, Zhang Z, Meoded RA, Cepeda AJ, Piel J, Keatinge-Clay AT, ACS 865 Chem Biol 13:975–983, 2018, http://doi.org/10.1021/acschembio.8b00049, Pöplau P, Frank S, 866 Morinaka Int Ed BI. Piel J. Angew Chem Engl 52:13215-13218. 2013. https://doi.org/10.1002/anie.201307406). This was found to be true only for LasM DH4 and not 867 868 LasO DH2. However, identical variations from a traditional PS domain architecture are also seen 869 in PS domains found in the mandelalide pathway (MndC DH3 and MndD DH3) (Lopera J, 870 Miller IJ, **McPhail** KL, Kwan JC, mSystems 2:e00096-17, 2017. 871 https://doi.org/10.1128/mSystems.00096-17). (C) Alignment of double bond-shifting DH 872 domains identified in *las* BGC with similar domains found in other *trans*-AT PKS pathways. The 873 DH domain from the erythromycin BGC is used for comparison. LasM DH1 and LasN DH1 are 874 annotated as putative double bond-shifting DH domains. Generally, in DH shifting domains the 875 conserved proline (P) in Hx_8P motif is often replaced by either value (V) or leucine (L). In the 876 case of LasM DH1, a methionine (M) instead of V or L appears in the place of P, which is in line 877 with what is observed in difficidin biosynthesis as well (Chen X-H, Vater J, Piel J, Franke P, 878 Scholz R, Schneider K, Koumoutsi A, Hitzeroth G, Grammel N, Strittmatter AW, Gottschalk G, 879 Süssmuth RD, Borriss R, J Bacteriol 188:4024–4036, 2006, https://doi.org/10.1128/JB.00052-880 06). Furthermore, DH shifting domains are sometimes characterized by the replacement of the 881 conserved aspartic acid (D) with asparagine (N) and substitution of glutamine (Q) or histidine

882 (H) with V or L in the DxxxO/H motif. Even though LasN DH1 has an N in place of D in the 883 DxxxQ/H motif, it substitutes Q/H with a serine (S). This is unusual and not found in any other 884 double bond shifting DH. (D) Alignment of DH domains present in the las BGC with the DH 885 domain from the erythromycin BGC. Putative PS and double bond-shifting DH domains have 886 been excluded. LasH DH1 and LasM DH2 are annotated as inactive domains due to disrupted 887 catalytic motifs Hx_8P and DxxxQ/H. Even though in LasL DH1 the catalytic aspartic acid is 888 replaced by glutamic acid (DxxxQ/H motif), we propose it is active, as a similar mutation is 889 observed in the palmerolide BGC (Avalon NE, Murray AE, Daligault HE, Lo C-C, Davenport 890 KW, Dichosa AEK, Chain PSG, Baker BJ. Front Chem 9. 2021, 891 https://doi.org/10.3389/fchem.2021.802574). € Alignment of LasL DH3 with other DH domains 892 having a serine in place of proline in Hx₈P motif. The DH domain from the erythromycin BGC is 893 used for comparison. Sequence headers in blue represent DH domains annotated as active while 894 in red represent the ones annotated as inactive.

895 Fig S5. (A) 2D visualization of the initial Autometa binning of Forcepia v2. The "Ca. T. 896 lasonolidus" genome is circled in black and contigs making up the las BGC are marked with 897 arrows. Axes represent dimension-reduced Barnes-Hut Stochastic Neighbor Embedding (BH-898 tSNE) values (BH-tSNE x and BH-tSNE y). (B) 3D visualization of contigs present in the "Ca. 899 T. lasonolidus" genome. Contigs making up the las BGC are colored red. (C) GC percentage of 900 different sets of genes in Forcepia_v2 "Ca. T. lasonolidus". Las BGC genes are colored in red. 901 (D) and (E) Codon adaptation index (CAI) of different categories of genes present in 902 Forcepia_v1 "Ca. T. lasonolidus" and Forcepia_v2 "Ca. T. lasonolidus", respectively. Las BGC 903 genes are colored in red. P values for pairwise comparison between different categories of genes 904 are shown in the matrix below their respective plots. Values with p < 0.05 are considered

significant. Other non-significant p values are colored red. Annotated and hypothetical genes 905 906 represent the genes annotated with a function and genes annotated as hypothetical respectively 907 T, by Prokka (Seemann **Bioinformatics** 30:2068-2069, 2014, 908 https://doi.org/10.1093/bioinformatics/btu153).

909 Table S1. List of oligonucleotide primers used for different purposes. (A) Primers used for 910 screening the Forcepia sp. fosmid library before WGS. (B) Primers used for screening the 911 Forcepia sp. fosmid library after WGS. (C) Primers used for confirming the presence of terminal 912 connections with the las BGC.

913 Table S2. Metadata and taxonomic classification of all the MAGs

914 Table S3. Contigs making up the three repeats of the las BGC in "Ca. T. lasonolidus" in 915 Forcepia_v1 and Forcepia_v2. Contig IDs represent the labels in Fig. 5 of the main text.

916 Table S4. Gene annotation in Forcepia_v1 and Forcepia_v2 "Ca. T. lasonolidus". (A) Non-

917 pseudogenes annotated as Eukaryotic-like proteins. (B) Genes forming the PV BMC cluster. (C)

918 Non-pseudogenes annotated by dbCAN2

919 Text S1. Supplementary methods

920





24 R₂O 32 ОН Lasonopyran Skeleton

R,

Н

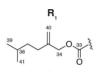
н

н

н

н

н



C

ОН Rí

Analog	Name

Lasonolide A

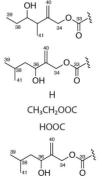
Lasonolide B

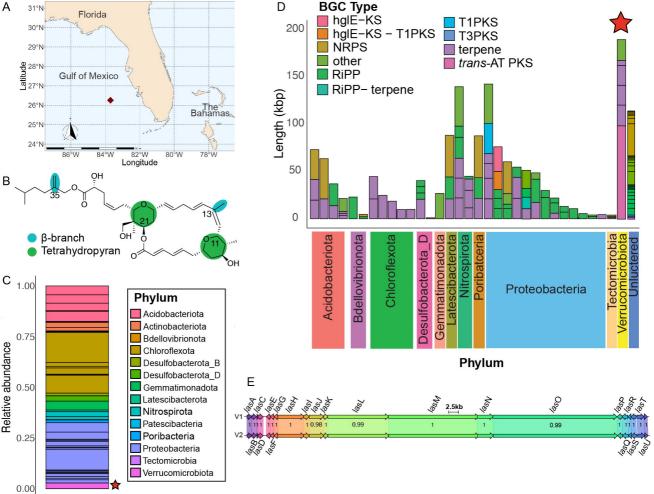
- Lasonolide C
 - Lasonolide D
 - Lasonolide E
 - Lasonolide F

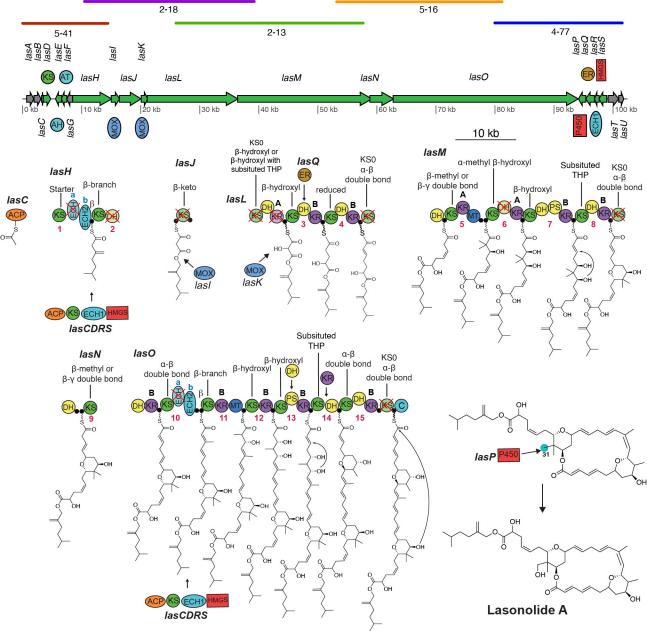
CH₃(CH₂)₁₀CO Lasonolide G

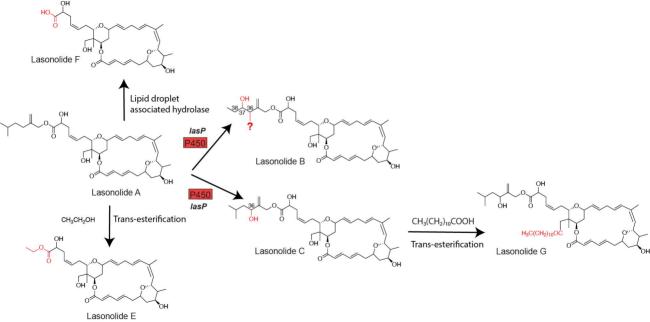
B

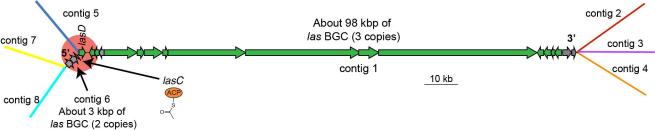


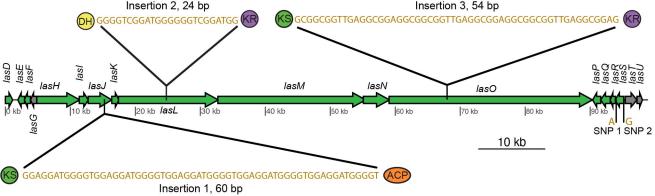


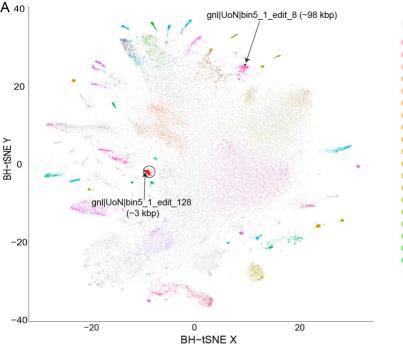




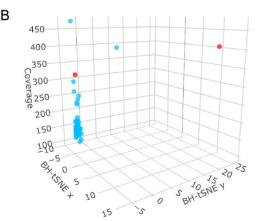


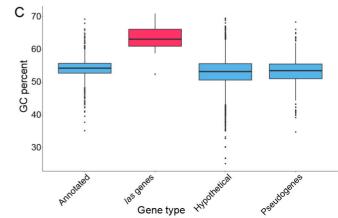


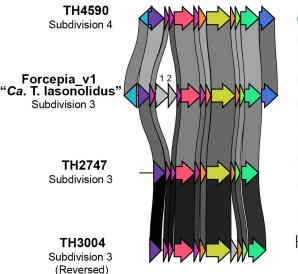




Bins			
obin1_1 📢	bin3_11	•	bin3678_217
bin120_0	bin3_12		bin4_1
bin132_0	bin3_13	•	bin4_2
o bin168_0 🛛	bin3_14	•	bin5_1
o bin178_16 🔵	bin3_2		bin5_15
o bin186_0 🛛	bin3_23		bin5_49
obin186_2	bin3_27	•	bin6_0
o bin187_35 🔵	bin3_3	۲	bin6_17
o bin2_0 🛛 🔵	bin3_5	٠	bin6_18
o bin2_1 🛛 🔵	bin3_6	•	bin6_4
o bin2_11 🛛	_		bin72_0
o bin2_3 🛛 🏼	bin3_8	•	bin75_1
o bin2_4 🛛 🏼 🍋	bin3_9		bin86_1
o bin2_5 🛛 🔹	bin3243_72	0	bin92_54
_	bin3673_21		_
o bin2_7 🛛 🏼	bin3674_127	•	refinement_1
bin216_1			refinement_2
bin217_2			refinement_3
bin3_10	bin3675_31		unclustered







- DeoR Transcriptional regulator
- Phosphate Propanoyltransferase
- BMC-H
- Acetate kinase
- BMC-P
- Aldehyde dehydrogenase
- BMC-P
- Aldolase class II
- L-lactate dehydrogenase
- 1: DNA-methyltransferase
- 2: PVUII Endonuclease, subunit A



