

## 1 **Supplemental Text S1**

### 2 **MATERIAL AND METHODS**

#### 3 **Metagenome datasets**

4 Pre-assembled metagenomes (n = 596) were downloaded between 2014 and 2017 from  
5 MG-RAST (<http://metagenomics.anl.gov/>), IMG/M (<https://img.jgi.doe.gov/>), EBI/ENA  
6 (<https://www.ebi.ac.uk/metagenomics/>), HMP (<http://hmpdacc.org/>), and the NCBI shotgun  
7 metagenomes database (<http://www.ncbi.nlm.nih.gov/Traces/wgs/>). Metagenome assembly  
8 quality was evaluated using MetaQUAST (Mikheenko *et al.*, 2016) (v 3.0). Protein coding  
9 genes were predicted using MetaProdigal (Hyatt *et al.*, 2012) (v 2.6.2). The number of  
10 genome equivalents within metagenomes was calculated using MicrobeCensus (Nayfach  
11 and Pollard, 2015) (v. 1.0.7) with read trimming to 100 bp. Metagenome sources were  
12 classified following accepted metagenome classification system (Ivanova *et al.*, 2010).  
13 Environmental metagenomes were classified based on ecosystem category (aquatic,  
14 terrestrial, air) and type (e.g. freshwater, marine) followed by other subcategories (e.g. depth,  
15 soil type). Host-associated communities were classified by host taxonomy and body site.  
16 Engineered communities were classified based on their function (e.g. wastewater treatment).

#### 17 **Generation of Hidden Markov models (HMM) and HMM search**

18 Twenty-three HMMs representing ten steroid-degradation protein families involved in  
19 degradation of steroid rings A and B, rings C and D (HIP), and HPDOA (see **Fig. 1**) were  
20 used for this analysis. HMMs used in this study were generated as previously described  
21 (Bergstrand *et al.*, 2016). Briefly, steroid-degradation proteins from the reference organisms  
22 *R. jostii* RHA1, *M. tuberculosis* H37Rv, *C. testosteronei* CNB-2 and *Pseudomonas* sp. strain  
23 Chol1 and from 256 newly identified putative steroid-degraders were clustered with CD-hit  
24 (Fu *et al.*, 2012) (v4.6.1) using a minimum sequence identity of 45%, a word size value of  
25 two, and all other parameters left at default. Homologous sequences from the resulting  
26 clusters were aligned using Mega (Tamura *et al.*, 2011) (v5.2.2) and manually trimmed. The  
27 aligned proteins were used to generate HMMs using HMMER (v3.1b1, <http://hmmer.org>).  
28 HMMs are available online  
29 ([https://github.com/MohnLab/Steroid\\_Degradation\\_Metagenomes\\_HMMs\\_2017](https://github.com/MohnLab/Steroid_Degradation_Metagenomes_HMMs_2017)). Proteins  
30 predicted from metagenomes were compared against our HMMs using HMMER with a  
31 maximum E-value of  $10^{-25}$  and a minimum coverage of 30%. These values were determined  
32 in our previous study to best identify known steroid catabolism genes in model organism  
33 genomes while providing maximum stringency against false positives (Bergstrand *et al.*,  
34 2016).

#### 35 **Taxonomic classification**

36 HMM hit protein sequences from the 105 selected metagenomes were aligned against the  
37 bacterial and archaeal non-redundant Ref-Seq protein databases (release 80 from 09

38 January 2017) using Diamond (Buchfink *et al.*, 2015) (v 0.7.11) with a maximum E-value of  
39  $10^{-5}$ . The taxonomies of the ten best hits for each protein were used for classification using  
40 the lowest common ancestor analysis with MEGAN (Huson *et al.*, 2016) (v6) against the  
41 NCBI taxonomy database (protein-accession numbers to taxon-ID reference file  
42 prot\_acc2tax-Nov2016.abin). Taxonomic assignments were curated manually where  
43 necessary. For better visualization of the taxonomic assignments we created KRONA charts  
44 for all analyzed metagenomes, which are available online  
45 ([https://github.com/MohnLab/Steroid\\_Degradation\\_Metagenomes\\_KRONA\\_charts\\_2017](https://github.com/MohnLab/Steroid_Degradation_Metagenomes_KRONA_charts_2017)).

#### 46 **Metagenome binning**

47 The 105 selected metagenomes were subjected to genome binning based on genomic  
48 signatures and marker genes using MyCC (Lin and Liao, 2016) with a minimum contig length  
49 of 2500 for clustering and the metagenome gene prediction mode of prodigal. The quality of  
50 recovered metagenome-assembled genomes (MAGs) was assessed using CheckM (Parks  
51 *et al.*, 2015) (v1.0.3) using lineage specific marker genes. Only MAGs with more than 25%  
52 genome completeness and less than 10% genome contamination and with HMM hits for at  
53 least five out of ten steroid-degradation protein families including at least one hit for KshA or  
54 HsaC were used for further analysis. The taxonomy of MAG contigs was assessed using the  
55 contig annotation tool CAT (Cambuy *et al.*, 2016). The taxonomy of whole MAGs was  
56 assessed by determining the weighted majority of taxonomic lineages based on contig length  
57 using a custom-made python script  
58 ([https://github.com/Holert/GitScripts/blob/master/annotate\\_cat\\_contigs.py](https://github.com/Holert/GitScripts/blob/master/annotate_cat_contigs.py)). Protein coding  
59 genes in MAGs were predicted using MetaProdigal (Hyatt *et al.*, 2012) (v 2.6.2). Best  
60 reciprocal BLAST hit analysis using BackBLAST (Bergstrand *et al.*, 2016) was used to  
61 compare the steroid-degradation proteome of MAGs to known and hypothetical steroid-  
62 degradation proteins from characterized steroid-degrading model organisms. A minimum  
63 identity filter of 25% and a maximum E-value of  $10^{-5}$  were used for analysis. MAGs annotated  
64 as Actinobacteria were compared to *R. jostii* RHA1 (Accessions NC\_008268.1,  
65 NC\_008269.1, NC\_008270.1, NC\_008271.1) and *M. tuberculosis* H37Rv (Accession  
66 NC\_000962.3). Because the phylogeny of steroid-degradation proteins in Proteobacteria  
67 does not follow the phylogeny of the respective 16S rRNA genes (Bergstrand *et al.*, 2016),  
68 MAGs annotated as Proteobacteria were compared to multiple model organisms, namely *C.*  
69 *testosteroni* CNB-2 (Accession NC\_013446.2), *Pseudomonas* sp. strain Chol1 (Accession  
70 NZ\_AMSL00000000.1), and *Pseudoalteromonas haloplanktis* TAC125 (Accessions  
71 NC\_007481.1, NC\_007482.1). Steroid degradation genes in the latter organism were  
72 recently identified by HMM analysis and best reciprocal BLAST analysis (Bergstrand *et al.*,  
73 2016). MAGs classified only to the bacterial domain were compared to all known steroid-  
74 degraders.

## 75 **Novelty estimation and phylogenetic reconstruction**

76 Protein novelty was assessed by analyzing sequence similarities of predicted KshA and  
77 HsaC homologs from predicted steroid-degradation MAGs to their best Diamond BLAST hit  
78 against the non-redundant Ref-Seq protein database (see above). Due to sequence identity  
79 of the corresponding proteins encoded in the two *Rhodococcus* MAGs from Antarctic dry  
80 valley metagenomes, we only analyzed sequences of one representative MAG (SOI\_12.1).  
81 For phylogenetic analysis, all predicted KshA and HsaC sequences from predicted steroid-  
82 degrader MAGs and all KshA and HsaC sequences from known and previously predicted  
83 steroid-degraders (Bergstrand et al., 2016) were aligned using the Muscle algorithm within  
84 MEGA7 (Kumar et al., 2016). KshA sequence alignment showed that the highly-conserved  
85 Rieske-dioxygenase and non-heme binding motifs typical for class IA terminal oxygenases  
86 (van der Geize et al., 2002) are conserved in all KshA sequences (not shown). Similarly, the  
87 highly-conserved iron-binding and active site motifs typical for meta-cleavage dioxygenases  
88 (Horinouchi et al., 2001) are conserved in all HsaC sequences. Phylogenetic maximum  
89 likelihood trees were computed using MEGACC (Kumar et al., 2012) with 1000 bootstrap  
90 repetitions and complete deletion. Phylogenetic trees were visualized with iTol (Letunic and  
91 Bork, 2016) (v3). For both analyses, only KshA and HsaC HMM hit proteins from MAGs were  
92 used, which had a protein sequence length of more than 70% of the median length of all  
93 KshA and HsaC proteins from analyzed MAGs.

## 94 **Isolation of steroid-degraders from marine sponges**

95 We isolated steroid-degrading bacteria from six sponge species (NCBI biosample numbers  
96 SAMN02192786, SAMN02192789, SAMN02192792, SAMN02192793, SAMN02192796, and  
97 SAMN02192803) collected off the coast of Santa Barbara, US, and British Columbia,  
98 Canada. Frozen sponge material was thawed and cut into pieces of approximately 1 cm<sup>3</sup>,  
99 which were washed in sterile artificial calcium- and magnesium-free sea water medium  
100 (doi:10.1101/pdb.rec12053) and minced into smaller pieces before sponge tissue was  
101 dissociated in 20 ml of the same medium for 20 min. Samples were centrifuged (3000 g for 5  
102 min) and 1 ml of supernatant was used to inoculate 9 ml of artificial seawater medium (36 g l<sup>-1</sup>  
103 sea salts (Instant Ocean Sea Salt, Aquarium Systems Inc., Blacksburg, VA, USA), 2.4 g l<sup>-1</sup>  
104 HEPES, 1 g l<sup>-1</sup> NH<sub>4</sub>Cl, 0.2 g l<sup>-1</sup> KH<sub>2</sub>PO<sub>4</sub>, trace elements (Bauchop and Elsdon, 1960) and  
105 vitamins (biotin 2 µg l<sup>-1</sup>, nicotinic acid 20 µg l<sup>-1</sup>, p-Aminobenzoate 10 µg l<sup>-1</sup>, D(+)-Pantothenate  
106 5 µg l<sup>-1</sup>, Pyridoxal 50 µg l<sup>-1</sup>, vitamin B12 (Cyanocobalamin) µg l<sup>-1</sup> mg), pH 7.2) containing 1  
107 mM cholesterol as the sole organic substrate solubilized with 0.5% (w/v) methyl-β-  
108 cyclodextrin. Cultures were incubated at 20°C at 180 rpm. Once the cultures turned turbid  
109 (between 7 and 14 days), 100 µl were transferred to 10 ml fresh medium. Enrichment  
110 transfers were repeated ten times for all samples. For selected transfers, substrate  
111 consumption was measured by organic extraction of culture supernatants and GC-MS

112 analysis, and biomass was determined using a bicinchoninic acid protein assay. Control  
113 experiments with artificial seawater medium containing 0.5% (w/v) methyl- $\beta$ -cyclodextrin but  
114 no cholesterol were carried out regularly for all enrichments but never turned turbid. Finally,  
115 cultures were serially diluted in artificial seawater medium and 100  $\mu$ l of selected dilutions  
116 were plated on artificial sea water medium agar (1.5%, (w/v)) supplemented with cholesterol  
117 and cyclodextrin. Representatives of morphologically different colony types were picked and  
118 further isolated on cholesterol or marine broth medium plates. 16S rRNA genes of purified  
119 strains were Sanger sequenced after PCR amplification using the primers 27f and 1492r.  
120 GenBank accession numbers for 16S rRNA sequences are MF770252 - MF770257. Isolates  
121 were taxonomically classified by aligning their 16S rRNA gene sequences against the SILVA  
122 SSU database using the SILVA Incremental Aligner (SINA) with default settings. Growth of  
123 pure strains was tested in the aforementioned medium containing 1 mM cholesterol and  
124 cyclodextrin. Substrate consumption was measured by organic extraction of culture  
125 supernatants and GC-MS analysis of outgrown cultures as described earlier (Casabon *et al.*,  
126 2013). The applied GC-MS analysis method detects cholesterol as well as typical side chain,  
127 and A- and B-ring degradation intermediates. Biomass was determined using a bicinchoninic  
128 acid protein assay. Control experiments without inoculum were included in all growth  
129 experiments.

## 130 **References**

- 131 Bauchop T, Elsdon SR. (1960). The growth of micro-organisms in relation to their energy supply. *J Gen*  
132 *Microbiol* **23**: 457–469.
- 133 Bergstrand LH, Cardenas E, Holert J, Van Hamme JD, Mohn WW. (2016). Delineation of Steroid-  
134 Degrading Microorganisms through Comparative Genomic Analysis. *MBio* **7**: e00166–16.
- 135 Buchfink B, Xie C, Huson DH. (2015). Fast and sensitive protein alignment using DIAMOND. *Nature*  
136 *Methods* **12**: 59–60.
- 137 Cambuy DD, Coutinho FH, Dutilh BE. (2016). Contig annotation tool CAT robustly classifies assembled  
138 metagenomic contigs and long sequences. *BioRxiv*: doi: <https://doi.org/10.1101/072868>
- 139 Casabon I, Crowe AM, Liu J, Eltis LD. (2013). FadD3 is an acyl-CoA synthetase that initiates catabolism  
140 of cholesterol rings C and D in actinobacteria. *Mol Microbiol* **87**: 269–283.
- 141 Fu L, Niu B, Zhu Z, Wu S, Li W. (2012). CD-HIT: accelerated for clustering the next-generation  
142 sequencing data. *Bioinformatics* **28**: 3150–3152.
- 143 Horinouchi M, Yamamoto T, Taguchi K, Arai H, Kudo T. (2001). Meta-cleavage enzyme gene *tesB* is  
144 necessary for testosterone degradation in *Comamonas testosteroni* TA441. **147**: 3367–3375.
- 145 Huson DH, Beier S, Flade I, Górska A, El-Hadidi M, Mitra S, *et al.* (2016). MEGAN Community Edition  
146 - Interactive Exploration and Analysis of Large-Scale Microbiome Sequencing Data. *PLoS Comput Biol*  
147 **12**: e1004957.
- 148 Hyatt D, LoCascio PF, Hauser LJ, Uberbacher EC. (2012). Gene and translation initiation site prediction

149 in metagenomic sequences. *Bioinformatics* **28**: 2223–2230.

150 Ivanova N, Tringe SG, Liolios K, Liu W-T, Morrison N, Hugenholtz P, *et al.* (2010). A call for standardized  
151 classification of metagenome projects. *Environ Microbiol* **12**: 1803–1805.

152 Kumar S, Stecher G, Peterson D, Tamura K. (2012). MEGA-CC: computing core of molecular  
153 evolutionary genetics analysis program for automated and iterative data analysis. *Bioinformatics* **28**:  
154 2685–2686.

155 Kumar S, Stecher G, Tamura K. (2016). MEGA7: Molecular Evolutionary Genetics Analysis Version 7.0  
156 for Bigger Datasets. *Mol Biol Evol* **33**: 1870–1874.

157 Letunic I, Bork P. (2016). Interactive tree of life (iTOL) v3: an online tool for the display and annotation  
158 of phylogenetic and other trees. *Nucleic Acids Res* **44**: W242–5.

159 Lin H-H, Liao Y-C. (2016). Accurate binning of metagenomic contigs via automated clustering  
160 sequences using information of genomic signatures and marker genes. *Sci Rep* **6**: 24175.

161 Mikheenko A, Saveliev V, Gurevich A. (2016). MetaQUAST: evaluation of metagenome assemblies.  
162 *Bioinformatics* **32**: 1088–1090.

163 Nayfach S, Pollard KS. (2015). Average genome size estimation improves comparative metagenomics  
164 and sheds light on the functional ecology of the human microbiome. *Genome Biol* **16**: 59–18.

165 Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. (2015). CheckM: assessing the quality  
166 of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res* **25**: 1043–  
167 1055.

168 Tamura K, Peterson D, Peterson N, Stecher G, Nei M, Kumar S. (2011). MEGA5: molecular evolutionary  
169 genetics analysis using maximum likelihood, evolutionary distance, and maximum parsimony methods.  
170 *Mol Biol Evol* **28**: 2731–2739.

171 van der Geize R, Hessels GI, van Gerwen R, van der Meijden P, Dijkhuizen L. (2002). Molecular and  
172 functional characterization of *kshA* and *kshB*, encoding two components of 3-ketosteroid 9 $\alpha$ -  
173 hydroxylase, a class IA monooxygenase, in *Rhodococcus erythropolis* strain SQ1. *Mol Microbiol* **45**:  
174 1007–1018.

175