1 **From sequence to molecules: Feature sequence-based genome mining**

2 **uncovers the hidden diversity of bacterial siderophore pathways**

3 Shaohua Gu[1,2#], Yuanzhe Shao[2#], Karoline Rehm[3], Laurent Bigler[3], Di Zhang[1], Ruolin He[1],

4 Jiqi Shao[1], Alexandre Jousset[4], Ville-Petri Friman[5], Zhong Wei[4*], Rolf Kümmerli[6*], Zhiyuan

5 Li[1,2*]

6 [1] Center for Quantitative Biology, Academy for Advanced Interdisciplinary Studies, Peking

7 University, Beijing, 100871, China

8 [2] Peking-Tsinghua Center for Life Sciences, Academy for Advanced Interdisciplinary Studies,

9 Peking University, Beijing, 100871, China

10 [3] University of Zurich, Department of Chemistry, Winterthurerstr. 190, 8057 Zurich,

11 Switzerland

12 [4] Jiangsu Provincial Key Lab for Organic Solid Waste Utilization, Key lab of organic-based

13 fertilizers of China, Nanjing Agricultural University, Nanjing, P R China

14 [5] University of Helsinki, Department of Microbiology, 00014, Helsinki, Finland

15 [6] University of Zurich, Department of Quantitative Biomedicine, Winterthurerstr. 190, 8057
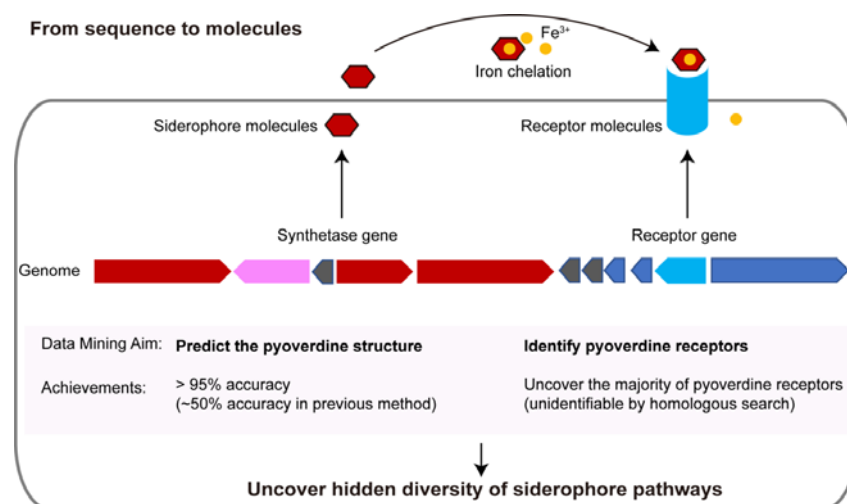
16 Zurich, Switzerland

17

18 [#] These authors contributed equally to this article.

19 [*] Corresponding authors (email: weizhong@njau.edu.cn; rolf.kuemmerli@uzh.ch;

20 zhiyuanli@pku.edu.cn )

21

22 **Abstract**

23 Microbial secondary metabolites have long been recognized as a rich source for

24 pharmaceutical compound discovery and to have crucial ecological functions. However, the

25 sequence-to-function mapping in microbial secondary metabolism pathways remains

26 challenging because neither protein function nor substrate specificity can accurately be

27 predicted from genome data. Here we focus on the iron-scavenging pyoverdines,

28 siderophores of *Pseudomonas* bacteria, as model system to develop a knowledge-guided

29 bioinformatic pipeline that extracts functional information of both the pyoverdine synthesis

30 machinery and uptake receptors from 1928 draft genomes. For pyoverdine synthesis, our

31 approach predicts the chemical structure of 188 different pyoverdines with nearly 100%

32 accuracy. For pyoverdine uptake, our pipeline uncovers 94 different pyoverdine receptor

33 groups. Our results demonstrate that combining feature sequence and phylogenetic

34 approaches is a powerful way to reconstruct bacterial secondary metabolism pathways based

35 on sequence data, unveiling an enormous yet overlooked diversity of siderophores and their

36 receptors.

37

**Introduction**

Rapid advancements in sequencing technologies have revolutionized our view on microbial communities. While amplicon sequencing provides information on community composition and diversity, shotgun and whole genome sequencing allow us to reliably anticipate evolutionary and ecological relationships between microbes and to obtain functional information on communities. Computational models assessing the metabolic capacity of individual members, or an entire consortium, have become very popular and powerful[1-3]. The major focus of such modelling approaches is typically on the primary metabolism of bacteria, as genes involved in core metabolic pathways are highly conserved and can be identified with relative ease[2,4]. Conversely, analysis of the secondary metabolites has attracted less attention, even though they include compounds such as antibiotics, toxins, siderophores, biosurfactants, all known to have important implications for community assembly[5,6] and to be important sources for pharmaceutical discoveries[7,8 9,10].

There are multiple challenges that currently prevent a detailed unravelling of secondary metabolism of bacteria based on genome data[5,11]. First, most secondary metabolites are produced by pathways comprised of modular enzymes such as non-ribosomal peptide synthetases (NRPSs) or polyketide synthases (PKS)[12,13]. Locating complete synthesis clusters and identifying all enzyme-encoding genes is challenging from highly fragmented metagenomic sequences or draft genomes with a high number of contigs. Second, functional predictions for coding regions within a cluster rely on homologous comparisons with experimentally characterized genes. Such information is often restricted to a limited number of model organisms, meaning that only a small portion of the existing secondary metabolism

60    pathways is covered by current data bases. Finally, given the complex multi-modular

61    synthesis machineries, it is challenging to precisely predict the secondary metabolites

62    produced even with accurately annotated NRPS or PKS clusters. The main challenge is that a

63    large pool of non-proteinogenic amino acids is used as substrates and the specificity of an

64    enzyme's A domain, connecting these unusual amino acids, is often poorly understood[14]. As a

65    result, new computational methods are needed to accurately reconstruct bacterial secondary

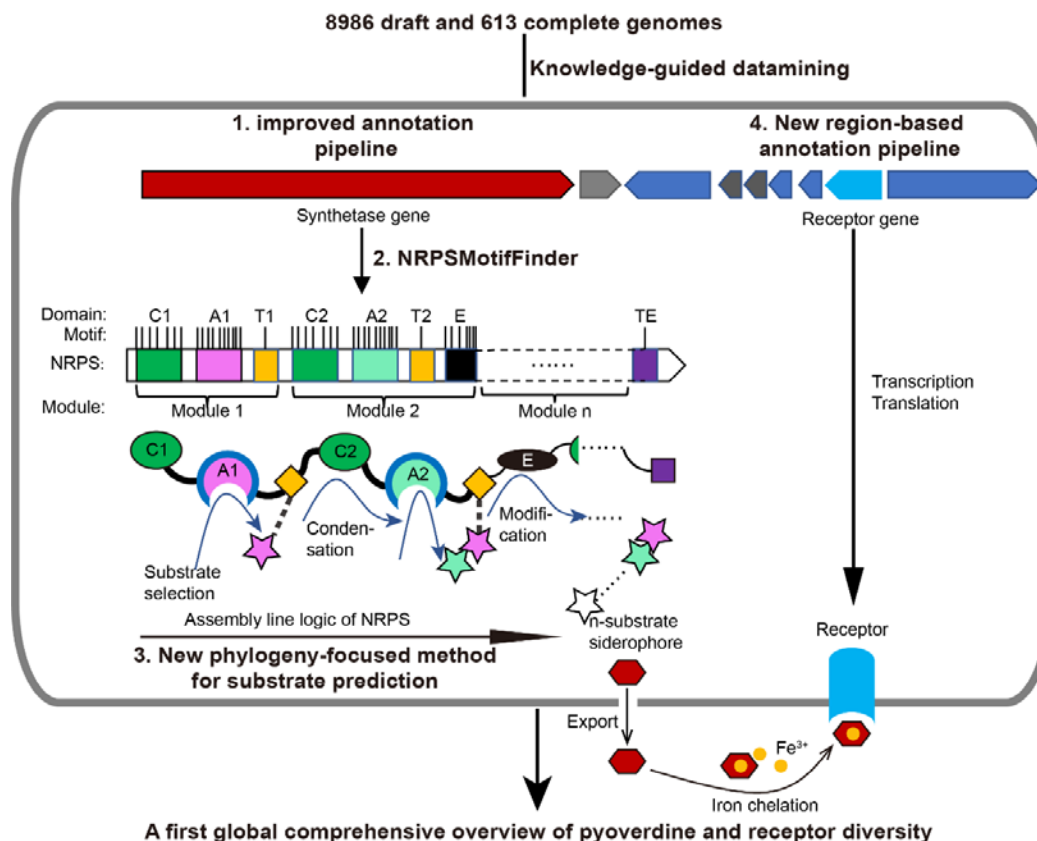66    metabolism from sequence data.

67        Here, we present a new bioinformatic pipeline that overcomes these challenges. We

68    specifically focus on a particular class of secondary metabolites (iron-scavenging

69    siderophores) as a case study to develop a bioinformatic workflow that predicts the chemical

70    structure of the produced metabolites with near 100% accuracy. Our pipeline is based on

71    improved gene annotation combined with a phylogeny- and feature sequence-based

72    substrate prediction techniques (Figure 1). In comparison with the currently available

73    databases and bioinformatic tools[14-16], the main advancement of our method is the more

74    accurate prediction of synthesized products based on NRPS clusters identified in genome

75    data.

76        Among siderophores, we focus on NPRS machineries that are responsible for the

77    synthesis of pyoverdines, a class of chemically diverse siderophores with high iron affinity,

78    produced by *Pseudomonas* bacteria[17,18]. While each *Pseudomonas* strain produces a single

79    type of pyoverdine, an enormous structural diversity has been described across strains and

80    species[19-22]. Pyoverdine types differ in their peptide backbone, meaning that the diversity

81    should be mirrored in NPRS enzyme diversity and their selectivity for the different amino acid

82   substrates[18]. Based on this knowledge, our pipeline entails the following steps (Figure 1): (i)

83   identification of the complete sequences of pyoverdine synthetase genes from fragmented

84   draft genomes, (ii) building the pyoverdine synthesis machinery *in silico* by extracting the

85   feature sequences for substrate specificity from motif-standardized NRPSs, and (iii) predicting

86   the precise chemical structure of pyoverdines followed by empirical verification.

87   An additional element of iron metabolism is that when siderophores are secreted and

88   bound to iron, bacteria rely on a specific receptor for their uptake into the cell. Pyoverdine

89   receptors are annotated as FpvA and it is known that receptor diversity matches pyoverdine

90   diversity[22,23]. Moreover, FpvA belongs to the family of TonB-dependent receptors and a single

91   *Pseudomonas* species often has many gene copies encoding these receptors. This poses an

92   additional bioinformatic challenge: how to find the gene encoding the specific pyoverdine

93   receptor among several potential receptor genes? To overcome this, we develop an algorithm

94   that focuses on sequence regions involved in pyoverdine recognition and translocation across

95   the outer membrane with supervised learning methods that locate the *fpvA* genes in the

96   fragmented genomes based on these regions (Figure 1). Altogether, our bioinformatic pipeline

97   uses knowledge-guided insights empowered by supervised learning to construct a first

98   systematic sequence-to-function mapping of a family of secondary metabolites (pyoverdine)

99   and their corresponding receptors. Our analysis unveils a yet unrecognized extraordinary

100   diversity of iron-scavenging machineries in pseudomonads.

101

102

**Figure 1 Scheme depicting our new genome mining pipeline to precisely predict the**

**synthesis, the molecular structure and the uptake machinery of pyoverdines, a family**

**of iron-scavenging siderophores produced by members of the *Pseudomonas* genus.**

The grey rounded outer rectangle represents a bacterial cell. The red and blue arrow-shaped boxes

stand for the synthetase and receptor genes for pyoverdines, respectively. Synthetase genes are

transcribed and translated to form the *n*-modular NRPS enzymes. These enzymes synthesize the

peptide backbone of pyoverdine through an assembly line using their repeating module units, with the A

domain being responsible for substrate selection and the E domain for chirality. The *n*-substrate

siderophores are then exported to the extracellular space for iron chelation. Membrane-embedded

TonB-dependent receptors recognize the ferri-siderophore complex and import it into the cell. Bold black

text and black arrows describe our multi-step computational methods developed to reconstruct the entire

114 process from genome sequence data. First, the annotation pipeline was improved (from antiSMASH) to

115 extract the complete sequence of pyoverdine synthetase genes from draft genomes. Second,

116 NRPSMotifFinder was used to define A- and E-domains and to determine the exact motif-intermotif

117 structure of the pyoverdine assembly line. Third, intermotif regions most indicative of substrate

118 specificity were used to develop a phylogeny-focused method for precise product prediction. Fourth, a

119 sequence-region-based annotation method was combined with genome architecture features to identify

120 the FpvA, receptors responsible for ferri-pyoverdine import.
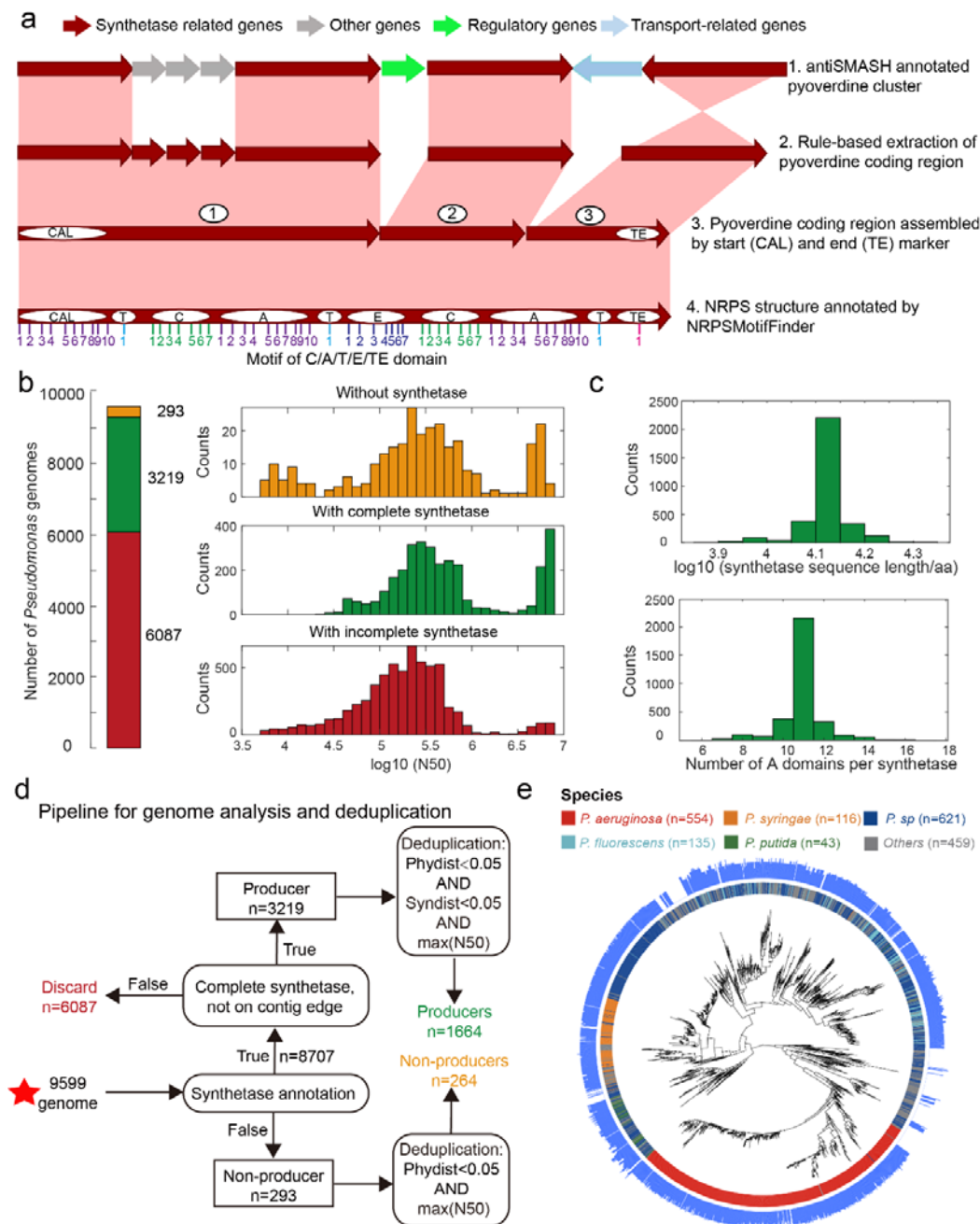
121

122 **Results**

123

124 **Section 1: Improved annotation pipeline reveals a vast reservoir of pyoverdine**

125 **synthetase genes**

126 The first step of our bioinformatic pipeline was to improve the annotation of pyoverdine

127 synthetase genes. The pyoverdine molecules is composed of a conserved fluorescent

128 chromophore (Flu) and a peptide chain (Pep), which are both synthesized by NRPS

129 enzymes[24]. There are already existing tools, such as antiSMASH, that can find and annotate

130 NRPS clusters in microbial genomes[25]. However, antiSMASH (and other popular annotation

131 platforms[15,16]) rely on accurate gene predictions, which are typically problematic for

132 fragmented genomes. Consequently, while antiSMASH can recognize and annotate certain

133 genes of an NRPS cluster, the precise reconstruction of a complete NRPS assembly line

134 often fails. This is particularly problematic because most available genomes are drafts and

135 any analysis suffers from the unavoidable issue of incomplete or misannotation of gene

136    fragments.

137



**Figure 2 Improved annotation pipeline reveals a vast diversity of pyoverdine**

**synthetase genes. a.** Improved annotation pipeline based on the raw annotation from antiSMASH. **b.**

The annotation pipeline was applied to 9599 *Pseudomonas* genomes (94% draft genomes). Genomes

could be separated into three categories. Yellow: genomes without pyoverdine cluster. Green: genomes

143     with a complete pyoverdine cluster. Red: genomes with incomplete pyoverdine synthetase cluster. The

144     red category involved genomes with truly incomplete clusters (lacking Flu or Pep synthetic genes) or

145     genomes with likely truncated synthetic genes at the edge of contigs. **c.** Distributions of the sequence

146     length (upper panel) and the number of A domains (lower panel) across all the genomes with a complete

147     synthetase cluster. **d.** Workflow applied to separate the 9599 *Pseudomonas* genomes into the three

148     categories described in b and removing of redundant genomes with high phylogenic similarity and

149     showing high similarity in pyoverdine synthetases. Red star indicates the start of the workflow. **e.**

150     Phylogenetic tree depicting the relationship among the 1928 non-redundant *Pseudomonas* strains (1664

151     producers and 264 non-producers) based on the concatenated alignment of 400 single-copy conserved

152     genes in their genomes. The inner ring depicts the taxonomical classification including the four most

153     prevalent species. The outer ring shows the number of A domains present in the pyoverdine synthetase

154     assembly line in each strain.

155        To overcome these issues, we developed an improved four-step annotation pipeline

156     starting with the raw annotation of the pyoverdine cluster obtained from antiSMASH (Figure

157     2a). First, we implemented a NRPS Hidden Markov Model (HMM) to re-annotate and extract

158     the entire nucleotide sequence of the pyoverdine synthetase cluster[26], including the genes

159     missed by antiSMASH. For this step, the nucleotide sequences were converted into amino

160     acid sequences to avoid erroneous gene predictions typically associated with antiSMASH.

161     Second, we assembled the entire re-annotated pyoverdine coding region into a single

162     sequence with a defined start (CAL) and/or end (TE) markers. Third, we used

163     NRPSMotifFinder to identify the C, A, T, E and TE motifs that are characteristic for the NRPS

164     structure of pyoverdine[14]. Finally, we applied a safety measure to ensure that the recovered

165   NRPS assembly line is complete (contains both Flu and Pep) and is not truncated, which can

166   occur when a synthetase coding region is at the edge of a contig. Consequently, we

167   dismissed all pyoverdine synthesis clusters located within 100 bp proximity to contigs' edges

168   and either lacked Flu or Pep synthetic genes.

169   Next, we applied our improved pyoverdine synthesis annotation pipeline to 9599

170   *Pseudomonas* genomes (including 613 complete and 8986 draft genomes) retrieved from the

171   Pseudomonas Genome Database[27]. We found the pyoverdine synthesis machinery in 97% of

172   the genomes (Figure 2b), indicating that the machinery is ubiquitous in *Pseudomonas*.

173   However, since 94% of the analyzed genomes were in draft form, the pyoverdine synthesis

174   machinery was likely truncated (i.e., on the edge of the contig) in 63.4% (6087) of the

175   genomes. These genomes were excluded from further analysis. Around 3.1% of retained

176   genomes (293) with high assembly completeness were missing pyoverdine synthetic genes,

177   indicating that these *Pseudomonas* strains were not able to produce pyoverdine ('non-

178   producers'). The rest of the genomes (33.5%; 3219 genomes) were classified as 'producers'

179   with complete pyoverdine NRPS assembly lines that meet all our quality controls. For these

180   3219 genomes, we used NRPSMotifFinder to find boundaries between the various synthesis

181   domains and to determine amino acid length and the number of A domains. The lengths of

182   pyoverdine synthetic genes ranged between 7690 and 21333 amino acids, and the number of

183   A domains per synthetase ranged between 6 and 17, with a total of 35,281 A domains being

184   present across all strains (Figure 2c and Figure S1). Overall, our analysis pipeline unveiled a

185   vast diversity of pyoverdine synthetase that goes far beyond of what has previously been

186   described in the literature.

187        Finally, we conducted a phylogenetic analysis based on 400 conserved genes with the

188        293 non-producers and the 3219 producers. We first removed redundant non-producers by

189        retaining the most integrative genome among strains with high phylogenic similarity. Then, we

190        removed redundant producers by retaining the most integrative genome among strains with

191        high phylogenic and pyoverdine synthetase similarity (Figure 2d). This data cleaning yielded a

192        total of 1928 *Pseudomonas* strains (403 complete and 1525 incomplete genomes),

193        segregating into 1664 pyoverdine producers and 264 non-producers. The phylogenetic tree

194        revealed that all major *Pseudomonas* species clades were present in our data set (Figure 2e).

195        Moreover, the number of A domains varied widely among species and even between strains

196        within species. For example, the number of *Pseudomonas aeruginosa* A domains ranges

197        between 7 and 14. In summary, by improving the synthetase annotation method, we

198        successfully obtained 1664 highly reliable pyoverdine synthetases (with a total of 18,292 A

199        domains) and 264 non-producers.

200

201        **Section 2: Phylogeny-focused substrate prediction for pyoverdine A domains**

202        Our next goal was to precisely predict the molecular structure of the pyoverdines produced by

203        the 1664 strains with complete synthetase gene clusters. The first essential step towards this

204        goal was to reliably predict the substrate selectivity of all A domains in the NRPS assembly

205        line. The A domain of each module selects for a single substrate among 22 proteinogenic and

206        hundreds of non-proteinogenic amino acids[28,29]. Moreover, whenever an E domain exists

207        downstream of an A domain, the chirality of the amino acid incorporated into the peptide chain

208        gets modified from L to D. Thus, the modularity combined with the selectivity of A domains

209    can promote an enormous diversity of pyoverdine molecule structures. To date, 73 pyoverdine

210    structures have been reported (Supplementary_table1) out of which 13 have their synthetase

211    genes sequenced (Supplementary_table2). In order to make reliable predictions, two

212    challenges must be addressed: (i) the extraction of relevant information from A domain

213    sequences for which the substrate is known, and (ii) the effective application of this

214    information to predict specificity of A domain sequences for which the substrate is unknown.

215        To address the first challenge, we built our analysis on the NRPS assembly lines of the

216    known 13 pyoverdines to extract relevant information from the A domain sequences. From

217    this dataset, we could identify 101 A domains that could be experimentally linked to 13 amino

218    acid substrates (Supplementary_table3). We next performed multisequence alignment of the

219    101 A domains to determine the "feature sequence distance", which is the most informative

220    for the substrate selectivity. To this end, we tested three different A domain regions, three

221    different sequence similarity measurements and seven different clustering methods for their

222    predictive power (Figure 3a). We found that the full A domain sequence is not informative for

223    substrate prediction (Figure 3b, left panel). Instead, our analysis indicated that information-

224    rich positions start with motif A4 and end before motif A5, consistent with the known role of

225    the A domain pocket in substrate selectivity[14]. Overall, the sequence region from motifs A4 to

226    A5 (termed "Amotif4-5"), in conjunction with Jukes-Cantor distance and Ward linkage

227    clustering, performed best in accurately distinguishing between different substrates and

228    maintaining homogeneity for identical substrates (Figure 3b).

229        To address the second challenge, we developed a "phylogeny-focused method" to apply

230    the feature sequence distance derived in the preceding paragraph to the 18,292 discovered A

231   domains. We realized that a direct construction of a phylogenetic tree including all 18,292

232   query A domains and the 101 reference A domains would be computationally too demanding

233   and impossible to scale up. Furthermore, such an approach would result in phylogeny-

234   interference issues, where domains would cluster not only based on their substrate

235   similarities but also based on overall species relatedness[14]. To minimize the effect of

236   phylogeny and speed up calculation, we took each of the 18,292 query A domains and

237   identified the two most similar A domain clusters within the 101 reference A domain set. We

238   then compared the feature distance between each query A domain and the two most similar

239   reference A domains in different clusters and assigned the query A domain to a substrate

240   specificity using the following rules (Figure 3c). (1) If the feature distance is below the 0.7

241   threshold (corresponding to 50% identity) for only one of the two reference A domains, then

242   the substrate of the query A domain is matched to the substrate of the more similar (lower

243   distance) reference A domain. (2a) If the feature distance is below the 0.7 threshold for both

244   reference A domains, then we considered the relative difference of the query A domain

245   towards the two reference A domains. If the relative difference is larger than 0.2, the query A

246   domain is matched to the substrate of the more similar reference A domain. (2b) If the relative

247   difference is smaller than 0.2, the substrate of the query A domain cannot unambiguously be

248   determined and is thus matched with both reference substrates. (3) If the feature distance is

249   above the 0.7 threshold (below 50% identity) for both reference A domains, then the substrate

250   of the query A domain is marked as "unknown". For most query A domains, rule (1) could be

251   applied (17880 cases), whereas rules (2) and (3) had to be used rarely (133 and 279 cases,

252   respectively). We applied our methodology termed "phylogeny-focused method" to all

253    following substrate and pyoverdine structure predictions.

254

255    **Section 3: Experimental validation of the annotation and prediction pipeline**

256    We tested whether our bioinformatic pipeline can reliably predict the structure of a set of yet

257    uncharacterized pyoverdines. To achieve this objective, we selected 20 *Pseudomonas* strains,

258    all known to produce pyoverdines, from a natural strain collection that was previously isolated

259    from soil and water[30]. We sequenced their genomes and subsequently applied our annotation

260    and prediction pipeline to generate predicted pyoverdine structures for all 20 strains harboring

261    a total of 237 A domains. Then, we elucidated the chemical structure of the 20 pyoverdines

262    using culture-based methods combined with UHPLC-HR-MS/MS[31]. We found a near-perfect

263    match (96.2%) between the predicted and the observed pyoverdine chemical structures and

264    were able to accurately assign amino acids in 228 out of 237 cases (Figure 3d). Our method

265    demonstrated a substantial improvement comparing to the prediction accuracy of AntiSMASH

266    in pyoverdines (46.0%), which could accurately assign correct amino acids only in 109 out of

267    237 cases (Supplementary_table4). The nine non-matching cases of our method segregated

268    into three groups. In three cases (1.3%), our algorithm could not distinguish between the

269    substrates Lysine and Ornithine, as these two amino acids are highly similar both in terms of

270    their chemical structures and corresponding A domain sequences. This is the only sensitivity

271    issue that is associated with our approach. In four cases (1.7%), our technique assigned an

272    "unknown" substrate to amino acids that turned out to be valine, citrulline and histidine.

273    Indeed, these three amino acids have not been reported in pyoverdines before and are

274    therefore not yet present in the reference dataset. These cases show that our analysis

275    pipeline can be used to identify new substrates. Once experimentally verified, the new A

276    domains and their substrates can expand the reference dataset, allowing targeted

277    improvement of our phylogeny-focused prediction technique. Finally, there were only two

278    cases (0.8%) that represented true mismatches between observed and predicted amino acids.

279    Altogether, our phylogeny-focused method is highly accurate in predicting pyoverdine peptide

280    structures and in identifying unknown substrates in *Pseudomonas*.
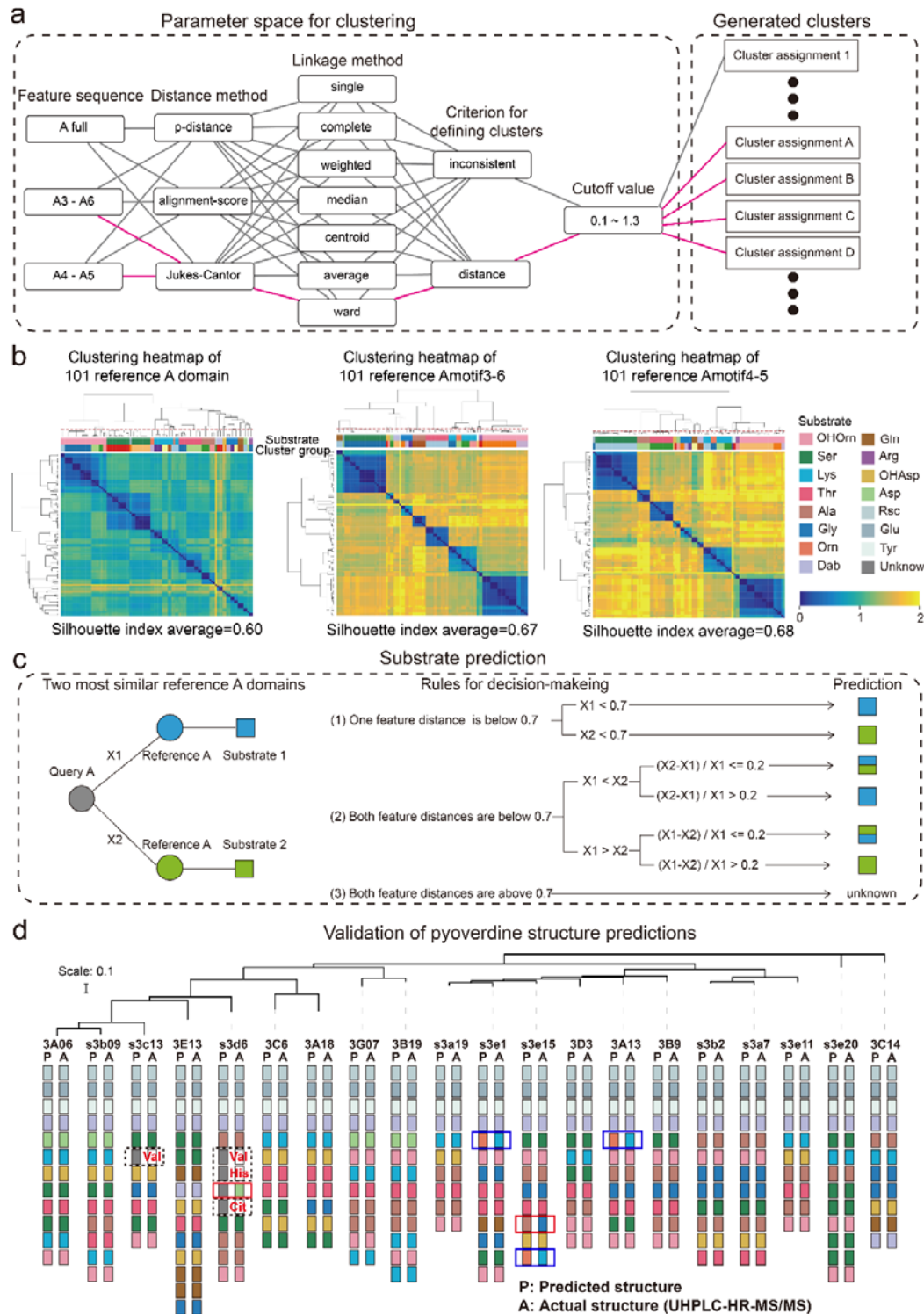
281

**Figure 3 Phylogeny-focused substrate prediction for pyoverdine synthetase assembly**

282     **lines. a.** Information from 101 reference A domains with known amino acid substrates were used to

283     develop an algorithm that predicts substrates from A domain sequence data with high accuracy. The
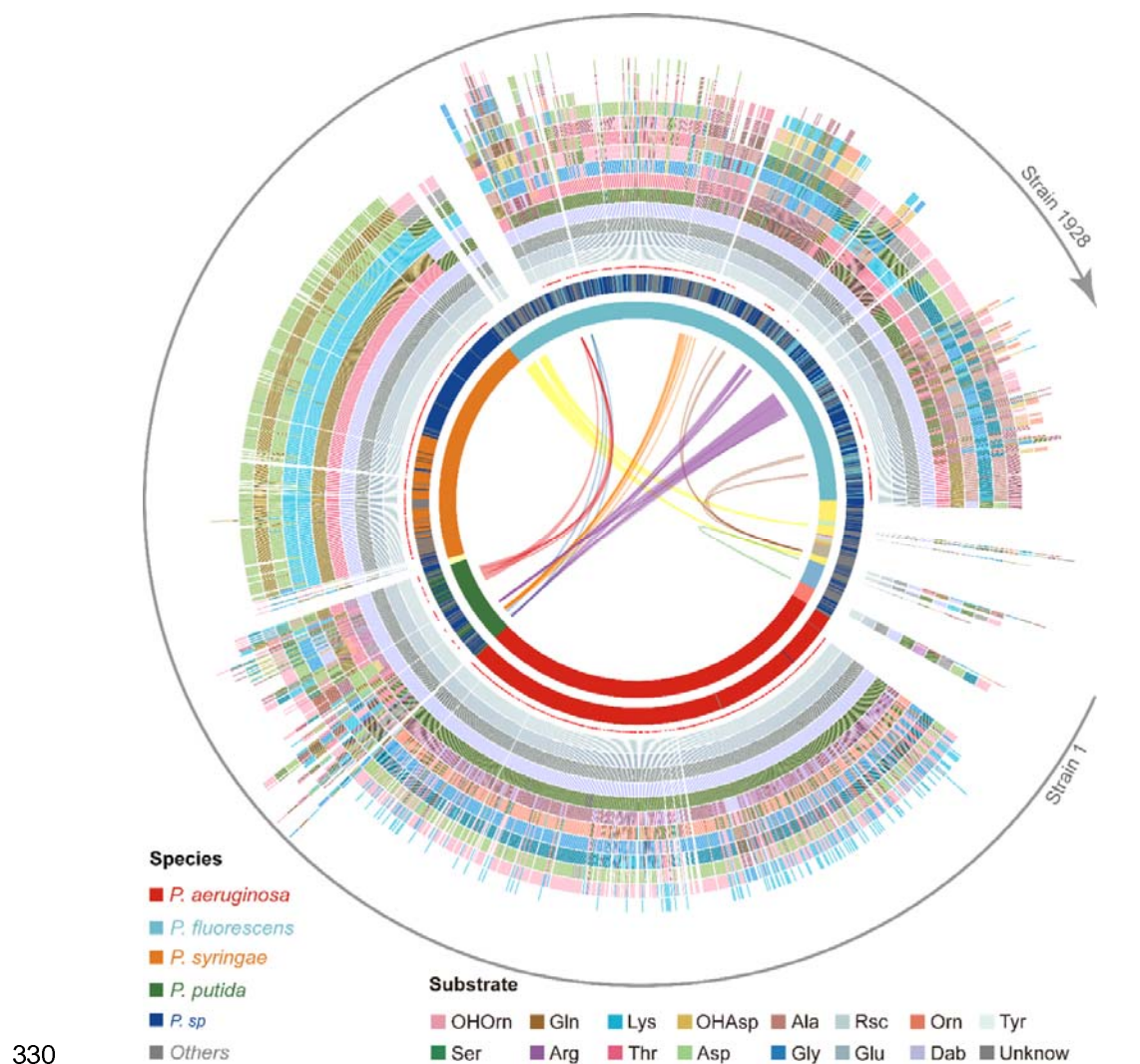
285     challenge is to group the variable A domains into clusters that predict the same substrate (captured by

286     the silhouette index). To find the most distinctive algorithm, we combined different feature sequences of

287     A domains (Amotif) with different distance and linkage methods in our hierarchical clustering analyses.

288     The best performing path is shown in pink. **b.** Heatmap showing the hierarchically clustered distances of

289     the 101 reference A domains as a function of the feature sequence used. Left panel: complete A domain

290     sequences. Middle panel: Amotif3-6 sequences. Right panel: Amotif4-5 sequences. The experimentally

291     validated substrates are shown on top of the heatmaps. The heatmaps show that hierarchical clustering,

292     reliably associating sequence distances with substrate, worked best with the Amotif4-5. **c.** Phylogeny-

293     focused substrate prediction pipeline for query A domains (grey circle) based on Amotif4-5 feature

294     sequence comparisons. X1 and X2 represent the feature distance between the query A domain and two

295     closest reference A domains (blue and green circles), respectively. Three rules are used, based on the

296     feature distances X1 and X2 and a threshold value of 0.7 (50% similarity), to make substrate predictions

297     for the query A domain. There are three possible outcomes: unambiguous substrate prediction (blue or

298     green squares), ambiguous substrate prediction (dual-colored squares), and no prediction ("unknown). **d.**

299     Phylogenetic tree of 20 *Pseudomonas* strains and visualization of their predicted and actual pyoverdine

300     structures to validate our phylogeny-focused substrate prediction pipeline. 228 out of the 237 substrates

301     (96.2%) were correctly predicted. The nine inconsistencies are boxed in blue (Lysine and Ornithine are

302     indistinguishable), in dashed black (correct detection of "unknown" substrates), and in red (true

303     mismatches). Note that our prediction pipeline (as any other pipeline) cannot distinguish between

304     modified variants of the same amino acid.

305     **Section 4: Application of the annotation and prediction pipelines to a full dataset**

306     After successful validation, we applied our bioinformatic pipeline to the 1664 complete NRPS

307 assembly lines annotated in our genome analysis (Figure 2). Across all assembly lines, we

308 were able to predict the substrates of 17,880 A domains (97.75%) without ambiguity, whereas

309 133 A domains (0.73%) were associated with two different substrates, and 279 A domains

310 (1.52%) predicted an unknown substrate (similar to the case of valine above). After

311 considering the presence/absence of an E domain in each module, we derived the structure

312 of 1664 pyoverdines according to method at section 2 (Figure 4). Our prediction yielded 188

313 different pyoverdine molecules, out of which only 37 structures had been previously reported.

314 However, these 37 reported structures were highly abundant across strains (1103 out of

315 1664). Agreeing with previous studies, we observed that the fluorophore is highly conserved

316 among the 188 predicted structures. Moreover, our analysis confirmed that 13 amino acid

317 substrates form the core of all the 188 pyoverdine structures, with most of the variation being

318 attributable to different substrate combinations, peptide lengths, and substrate chirality

319 (Figure 4). In addition, the 279 unknown substrates will significantly increase the repertoire of

320 pyoverdine amino acids if could be characterized by future experiments, despite that they

321 were much rarer than the 13 main substrates. Notably, pyoverdine structural diversity was not

322 strongly linked to phylogeny because the same pyoverdine structure could be found in

323 completely unrelated species, while closely related species often had different pyoverdine

324 structures (Figure 4). These observations suggest that there may be both frequent

325 recombination and horizontal gene transfer of pyoverdine synthetase clusters between

326 species. Taken together, the bioinformatics methods developed in our study can predict a suit

327 of secondary metabolites (pyoverdines) from sequence data with high accuracy, revealing an

328 unprecedented richness and evolutionary history of siderophores within pseudomonads and

329     the discovery of 151 putative novel pyoverdine variants.



**Figure 4 Predicted pyoverdine structural diversity based on our developed algorithm mapped onto the phylogenetic tree comprising all 1928 (non-redundant) *Pseudomonas* strains.** The stacked boxes in the outermost circle show the predicted structure of pyoverdines, whereby each color represents a specific amino acid substrate. Strains without boxes represent non-producers (n = 264). Boxes with two colors indicate cases of ambiguous (dual) substrate prediction. The red dots at the basis of the stacked boxes indicate experimentally validated pyoverdine structures. The inner circle shows the taxonomic species classification following Figure 2e. Because the allocation of strains to species names is often imprecise, we divided the 1928 strains by their

339    phylogenetic distance into 18 clades (color shadings in inner-most circle), out of which 13 contained

340    more than 1 strain. Lines within the inner-most circle link strains from different clades that share the

341    same pyoverdine structures, whereby line colors represent the shared unique pyoverdines. The bending

342    of the lines represents the phylogenetic sequence distances of the connected strain pairs.

343

344    **Section 5: Development of a region-based identification method for annotation of the**

345    **FpvA receptors**

346    In pseudomonads, iron-loaded pyoverdines are recognized by FpvA, a TonB-dependent

347    receptor, that transports the ferri-siderophore into the periplasm[19,32,33] . The protein structure

348    of characterized FpvA variants consists of three domains: The Secretin and TonB N-terminus

349    short domain (STN), the Plug domain (Plug), and the TonB dependent receptor domain

350    (TonB)[19]. While these domains are conserved across FpvA variants and other siderophore

351    receptors, there is substantial variation at the sequence level. This makes it challenging to

352    reliably identify FpvA receptors from sequence data by homologous search. As an example,

353    we were unable to find FpvA genes (with a 60% identity threshold) by homologous search in

354    several genomes although they had complete pyoverdine synthesis machineries. Moreover,

355    there are many other TonB-dependent receptors with fairly high sequence identity to FpvA but

356    that transport other siderophores than pyoverdine (e.g. FpvB, 55% identity, transporting

357    pyoverdine, ferrichrome and ferrioxamine B[34]. Therefore, it is imperative to develop a new

358    comprehensive method for identifying FpvA receptors in *Pseudomonas* genomes.

**Figure 5 A sequence-region-based identification pipeline for annotating FpvA receptors.**

**a.** Heatmap displaying the hierarchically clustered sequence distances (p-distance calculation method,

identity (%) = (1-sequence distance) * 100) of 35 reference siderophore receptors identified in

*Pseudomonas* spp., based on full sequences. No clear discrimination between FpvA, FpvB and other

364    receivers is possible. The order of receptors is consistent across panels (**b**), (**e**), and (**f**). **b.** The pHMM

365    scores of the three standard receptor domains (STN, Plug, and TonBDR) vary across the 35 reference

366    sequences (A: FpvA, B: FpvB and NA: others), but do not allow to distinguish between receptor groups.

367    **c.** FpvA region-based conservation scores from a multi-alignment of the 35 reference sequences

368    mapped to the FpvA sequence of strain *P. aeruginosa* PAO1. All residues within the top 10% of the

369    conservation score denoted with black dots. For each region flanked by two black dots, we calculated

370    the FpvA identification score (heatmap), representing the ability to distinguish FpvA from non-FpvA

371    receptors. **d.** Mapping of the two regions with the highest FpvA identification scores R1(dark red) and

372    R2 (orange) to the crystal structure of FpvA from PAO1 conjugated with pyoverdine (PDB 2IAH). **e.**

373    Heatmap showing the hierarchically clustered sequence distances of 35 reference siderophore

374    receptors based on the R1 sequence region. A clear discrimination between FpvA/FpvB and other

375    receptors emerges. **f.** Heatmap showing the hierarchically clustered sequence distances of 35

376    reference siderophore receptors based on the R2 sequence region. A clear discrimination between

377    FpvA and FpvB receptors emerges. **e.** The pHMM scores of regions R1 and R2 for the 35 siderophore

378    reference receptors are contrasted against each other, yielding a clear separation between FpvA, FpvB

379    and other receptors. Dashed lines indicate the pHMM threshold scores used for later analysis. **f.**

380    Flowchart showing all steps involved in the FpvA annotation from genome sequence data. The red star

381    indicates the start of the workflow.

382       We started our approach by comparing the sequences of 35 reported siderophore

383    receptors, including 21 FpvA, 6 FpvB, and 8 TonB-dependent siderophore receptor

384    sequences often found in *Pseudomonas* genomes, encoding receptors for the uptake of

385    heterologous siderophores (Supplementary_table5). We found that all receptor sequences

386    share a similar length of around 800 amino acids (FpvA and FpvB sequences: 809 ± 10

387    amino acids). We then used the complete sequences to calculate the pair-wise distances by

388    global alignment before applying hierarchical clustering (Figure 5a). We found substantial

389    divergence between FpvA variants to an extent that was comparable to the distance between

390    FpvA and other siderophore receptors. Moreover, FpvB variants clustered with FpvA variants,

391    showing that FpvA identification based on full sequence distances is unachievable. We hence

392    focused on the three typical receptor domains (TonB, Plug, and STN, retrieved from the Pfam

393    database) and applied Profile Hidden Markov Models (pHMM) to calculate the pHMM

394    probability scores for each domain and reference sequence. The probability scores

395    (calculated as the log-odd ratios for emission probabilities and log probabilities for state

396    transitions) had reasonably high scores but no distinction was apparent between the three

397    receptor classes (Figure 5b).

398        We next asked whether there are specific regions within the receptor sequences that are

399    characteristic of FpvA. To address this, we conducted a multiple sequence alignment (MSA)

400    with all 35 reference receptor sequences and mapped them onto the sequence of the well-

401    characterized FpvA of *P. aeruginosa* PAO1 (Figure 5c). MSA allows to identify conserved sites

402    (Figure 5c, black dots representing the top 10% most conserved sites) that are shared by the

403    majority of the reference sequences. We then used these conserved sites to partition the

404    MSA into variable regions which were flanked by two conserved sites. For each variable

405    region, we assessed its predictive power to differentiate FpvA from non-FpvA sequences. For

406    this we defined the "FpvA identification score" analogous to the intercluster-vs-intracluster

407    Calinski-Harabasz variance ratio, as

$$I_{\mathrm{FpvA}} = d_{\mathrm{A:non}}/d_{\mathrm{A:A}}$$

408     where $d_{\mathrm{A:A}}$ is the sequence distance among all 21 FpvA sequences, and $d_{\mathrm{A:non}}$ is the

409     sequence distance between all 21 FpvA and the 14 non-FpvA sequences.

410     Our analysis yielded two locations with noticeably high FpvA identification scores (Figure

411     5c). The region with the highest FpvA identification score (referred to as R1) locates at the

412     intersection of the Plug domain and the barrel structure of the TonB domain (Figure 5d,

413     between 258 Gly and 309 Gly in the PAO1 FpvA). According to the sequence distance matrix,

414     the R1 region allows to distinguish heterologous siderophore receptors from FpvA and FpvB

415     receptors (Figure 5e). The region with the second highest FpvA identification score (referred

416     to as R2) was located in the C-terminal signaling domain (Figure 5d, between 59 Leu and 86

417     Lys in the PAO1 FpvA). The sequence distance matrix revealed that R2 allows to distinguish

418     FpvB from FpvA receptors (Figure 5f).

419     We then constructed two pHMM by (i) the alignment of the 21 FpvA sequences in the R1

420     region, termed R1(FpvA), and (ii) the alignment of the 6 FpvB sequences in the R2 region,

421     termed R2(FpvB). Running R1(FpvA) and R2(FpvB) against all 35 reference sequences

422     revealed a clear separation between the three receptor categories (Figure 5g). Along the

423     R1(FpvA) axis, FpvA and FpvB reference sequences have high R1 scores (minimal score

424     77.0) that separate them from other siderophore receptors (maximal score 38.1), whereas

425     FpvAs references have substantially lower R2 scores (maximal 20.2) than FpvBs (minimal

426     49.0) along the R2(FpvB) axis.

427     Based on these insights, we developed a decision flow chart for annotating FpvAs in

428     *Pseudomonas* genomes (Figure 5h): First, we considered sequences as Fpv-like receptors

429     that share similar properties to the ones identified in our reference database. Particularly,

430     protein coding sequence (CDS) length has to be between 750 and 850 amino acids and the

431     phMM scores for the three typical receptor domains STN, Plug, and TonB have to be greater

432     than 25, 50, and 80, respectively (Figure 5b, red dashed lines). Second, we used the pHMM

433     threshold scores obtained for R1(FpvA) and R2(FpvB) (Figure 5g) to differentiate other

434     siderophore receptors (R1(FpvA) score < 50) from FpvB receptors (R1(FpvA) score > 50 and

435     R2(FpvB) score > 30) and FpvA receptors (R1(FpvA) score > 50 and R2(FpvB) score < 30).

436     Our method effectively identifies FpvA receptors from sequence data and can be readily

437     applied to the entire *Pseudomonas* dataset.

438

439     **Section 6: Application of the receptor annotation pipeline to the full dataset.**
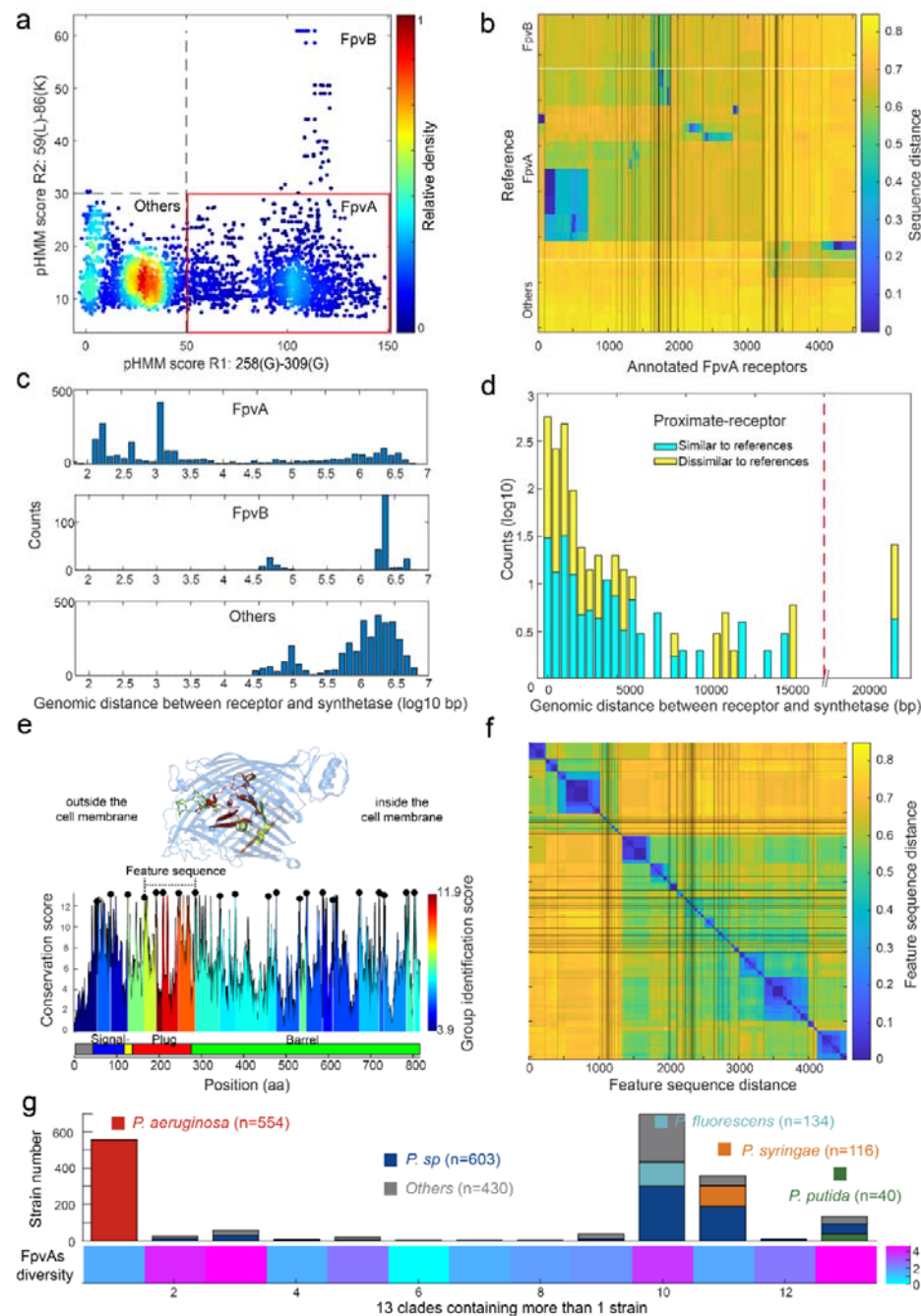
440     The region-based receptor identification pipeline was applied to all 1928 *Pseudomonas*

441     genomes. The analysis identified 4547 FpvAs, 615 FpvBs, and 9139 other TonB-dependent

442     Fpv-like receptors across the dataset (Figure 6a). The 4547 FpvA sequences clustered

443     hierarchically into 114 groups, defined by an identity threshold of 60%. When comparing to

444     the 21 reference FpvAs (Figure 6b), we found that 2293 FpvA sequences have close

445     homologues in the reference data base, while 2254 FpvA sequences lack such close

446     homologues (sequence identity < 50%). These latter sequences, termed as "dissimilar to

447     reference", may represent novel subtypes of FpvA receptors that could not be found by

448     simple homology search. Our analysis further shows that many strains have more than one

449     FpvA receptor.

450        We then asked whether the 4547 FpvAs are found in proximity of pyoverdine Pep

451  synthetase genes as it is commonly the case for cognate FpvA receptors [35]. We thus

452  calculated the proximity between pyoverdine Pep genes and the Fpv-like receptor genes by

453  counting the number of base pairs between the two coding regions. All TonB-dependent

454  receptors that have not been classified as FpvAs were more than 20 kb away from the Pep

455  genes (Figure 6c). In contrast, 92% of the nearest FpvA genes were indeed located within 20

456  kb of their pyoverdine Pep genes (Figure 6d, called proximate receptors). These proximate

457  receptors encompassed both those with close (66%) and more distant (34%) resemblance to

458  the reference receptor types. Overall, this proximity analysis confirmed that our region-based

459  gene identification method can reliably identify FpvA receptors.

460      We next explored the diversity among FpvA receptors in more detail by focusing on the

461  1534 strains that had proximate-receptors within 20 kb of the pyoverdine Pep genes (Figure

462  6d) and using high-confidence FpvAs for sequence feature extraction. When considering the

463  whole gene sequences, these receptors segregated into 44 groups according to single-

464  linkage clustering with an identity threshold of 60% (Figure S2a). To investigate which

465  sequence regions were the most informative for reliable clustering, we used a similar

466  approach as with FpvAs detection by quantifying the "group identification score" for variable

467  regions flanked by highly conserved sites. The higher the score, the stronger a region's

468  capacity to discriminate between FpvA groups. We found that the four regions with the top

469  discrimination capacities all located near the Plug domain surrounding the pyoverdine

470  transmission channel (Figure 6e). The plug domain is known to undergo conformational

471  changes and is involved in pyoverdine selectivity and import[36,37], suggesting that the four

472  high-score regions are responsible for pyoverdine specificity.

473



474

**Figure 6 Application of the receptor annotation pipeline to the full database. a.** Applying

the receptor annotation pipeline to the genomes of the 1928 non-redundant *Pseudomonas* strains yields

14301 Fpv-like receptors, which segregate into 4547 FpvA receptors (red box), 615 FpvB receptors, and

9139 other receptors, based on the pHMM score thresholds for regions R1 and R2. The heatmap

479    indicates receptor density. **b.** Sequence distance matrix between the 35 reference sequences (y-axis)

480    and the 4547 annotated FpvA sequences in the full database (x-axis). Database sequences were

481    ordered by hierarchically clustering and segregated into 114 groups. 2254 of the annotated FpvA

482    sequences have sequence identity < 60% compared to the reference receptors, pointing at novel

483    subtypes of FpvA receptors. **c.** Genomic distance (in base pairs) between each Fpv-like receptor

484    sequence and its pyoverdine peptide synthetase gene (Pep) for annotated FpvA receptors (upper panel),

485    FpvB receptors (middle panel) and other receptors (lower panel). **d.** Distribution of the genomic distance

486    between each FpvA receptor and its nearest pyoverdine peptide synthetase depending on whether the

487    annotated FpvA receptor has high sequency similarity (blue, ≥ 50%) or low sequence similarity (yellow,

488    < 50%) with at least one of the 21 reference FpvAs. **e.** FpvA region-based conservation scores from a

489    multi-alignment of all the annotated FpvA receptors that are proximate (< 20 kbp) to the pyoverdine

490    synthetase cluster mapped to the FpvA sequence of strain *P. aeruginosa* PAO1. All residues within the

491    top 10% of the conservation score denoted with black dots. For each region flanked by two black dots,

492    we calculated the group identification score (heatmap, lower panel), representing the ability of the region

493    to distinguish between different groups of FpvA receptors. Four regions in the plug domains had a

494    particularly high group identification score (called the feature sequence). They are mapped to the crystal

495    structure of FpvA from PAO1 conjugated with pyoverdine (PDB 2IAH, up panel). All four regions

496    surround the pyoverdine transmission channel and are shown in the respective heatmap color. **f.**

497    Heatmap showing the hierarchically clustered distances between the 4547 annotated FpvA receptors

498    based on the feature sequence (comprising the four groups with the highest identification scores). The

499    analysis identifies 94 receptor groups with a 70% identity threshold. **e.** The diversity of FpvA receptors

500    along the 13 phylogeny clades containing more than 1 strain. Receptor diversity was calculated by the

501    Shannon entropy, similar to the alpha-diversity in microbial community.

502        Based on the above insights, we concatenated the four high-score regions (from 168 Pro

503    to 295 Ala in PAO1) into a single "feature sequence". The feature sequence could

504    characterize 98% of the distance matrix compared to the whole sequence (1534 FpvAs,

505    r=0.98, Figure S2a-b) and substantially reduced within-group distance. We applied the

506    concatenated feature sequence approach to all the 4547 annotated FpvAs to calculate the

507    sequence distance matrix. Single-linkage clustering with an identity threshold of 70% revealed

508    a total of 94 groups, out of which 43 groups contained more than 10 members (Figure 6f).

509    The diversity of receptors is hence much larger than currently anticipated as only 3 groups of

510    FpvAs have previously been reported. Finally, we calculated the diversity of receptor FpvAs

511    for each of the 13 phylogenetic clades with more than one strain by the Shannon entropy,

512    which is similar to the alpha-diversity in microbial community (Figure 6g). We noticed large

513    differences in FpvA diversity across the clades and species: clades with *P. aeruginosa* and *P.*

514    *syringae* species had lower FpvAs diversity (1.55 and 1.60) than clades containing *P. putida*

515    and *P. fluorescens* species (4.82 and 3.77). Taken together, the region-based identification

516    method developed in our study can reliably mine the FpvAs (pyoverdines receptors) from

517    genome data, revealing undiscovered diversity of FpvA pyoverdine receptors that are

518    unequally distributed across the different phylogenetic clades of pseudomonads.

519

520    **Discussion**

521    The rapid expansion of sequencing data offers exciting opportunities for microbiology[38-40].

522    One key challenge of current research in the field is to infer biological functions of microbial

523 communities from genome sequence data[41-43]. While this endeavor is increasingly successful

524 for biological functions involving the primary metabolism and the associated complex

525 metabolic flux, reconstructing aspects of the secondary metabolism is much more challenging.

526 The main issue is that neither the function of a secondary metabolite enzyme nor the resulting

527 metabolite can be precisely predicted from gene sequence data. In our study, we tackled this

528 challenge and developed a bioinformatic pipeline to reconstruct the complete secondary

529 metabolism pathway of pyoverdines, a class of iron-scavenging siderophores produced by

530 *Pseudomonas* spp. These secondary metabolites are synthetized by a series of non-

531 ribosomal peptide synthetases and require a specific receptor (FpvA) for uptake. We

532 combined knowledge-guided learning with phylogeny-based methods to predict with high

533 accuracy: (i) the full pyoverdine assembly line, (ii) the substrate specificity for each enzyme

534 within the assembly lines, (iii) the complete chemical structure of pyoverdines, and (iv) the

535 FpvA receptors from genome sequences. After validation, we tested our pipeline with

536 sequence data from 1664 phylogenetically distinct *Pseudomonas* strains and were able to

537 determine 18,292 enzymatic A domains involved in pyoverdine synthesis, reliably predicted

538 97.8% of their substrates, identified 188 different pyoverdine molecule structures and 4547

539 FpvA receptor variants belonging to 94 distinct groups. The uncovered diversity is stunning

540 and goes far beyond currently known levels of variation (73 pyoverdines and 3 FpvA groups).

541 The molecular diversity of iron scavenging capacity highlights its importance among

542 pseudomonads.

543 We show that knowledge-guided learning is an extremely powerful tool to predict enzyme,

544 metabolite, and receptor properties. The establishment of our entire pipeline is based on only

545    101 previously known enzymatic A domains (from 13 known pyoverdine assembly lines) and

546    21 FpvA receptor sequences. Even with this limited amount of information, we were able to

547    predict the substrates of almost all the 18,292 enzymatic A domains and to identify 4547 FpvA

548    receptors from the sequence data. A key insight from our knowledge-guided learning is that

549    comparisons based on the full gene sequences (e.g., for pyoverdine synthetase or receptor)

550    are likely non-informative and unsuitable for obtaining functional information. This is because

551    overall diversity does not stand for functional diversity, meaning that A domains recognizing

552    the same substrate can diverge substantially in their full sequences. The same holds true for

553    receptor sequences: whole-sequence alignments can neither accurately identify FpvA

554    receptors nor reliably separate them into functional groups. Instead, it is imperative to extract

555    informative feature sequences that are defined as sequence stretches within a gene whose

556    diversity is tightly linked to variation in its functioning. We successfully extracted and applied

557    feature sequence comparisons for both A domain substrate prediction and FpvA identification.

558    It is important to note that a knowledge-guided pipeline does not have to be perfect right from

559    the start. For example, our pipeline for pyoverdine structure prediction returned unknowns for

560    several amino acid positions within the PEP. Our experimental verifications then revealed

561    indeed new substrates such as valine and citrulline. This information can then be used to

562    refine our prediction algorithm in a feedback loop.

563        Another main advantage of our bioinformatic pipeline is that it can be applied to draft

564    genomes. This reflects a major improvement compared to existing annotation tools such as

565    antiSMASH[25], which typically has difficulties in recognizing NRPS structures in fragmented

566    genome assemblies. However, draft genomes are the most common data source in

567    microbiology. While our pipeline shows high performance, we need to acknowledge that we

568    still lose many genomes (6087 out of 9599, 63.4%). The reason for the loss is that the

569    pyoverdine synthesis machinery is large, which increases the probability that it is positioned

570    at the end of a contig. We decided to exclude those cases because the annotated synthesis

571    machinery might be truncated and thus incomplete. Thus, the high loss rate of draft genomes

572    is rather due to limitations in sequence quality (too many short contigs) and not due to a

573    limitation of our bioinformatic pipeline. We believe that this limitation will be lesser of a proble

574    in the future as long-read sequencing technologies are quickly becoming cheaper and more

575    reliable.

576        We further show that knowledge-guided learning combined with a phylogeny-focused

577    approach is a powerful tool for predicting the substrate specificity of A domains of synthetases.

578    It outperforms currently known bioinformatics prediction tools of NRPS substrates such as

579    antiSMASH[25]. Most current algorithms[44-48] perform poorly when applied to pyoverdines,

580    particularly when encountering non-proteogenic amino acids. The high accuracy of our

581    algorithm can largely be attributed to our reference set, composing only 13 pyoverdines from

582    *Pseudomonas spp.*, yet capturing most of the substrate diversity. Similarly accurate

583    predictions based on a handful of known substrates among closely related species were

584    observed in several fungal NRPS systems[49]. It is worth noting that when the algorithm output

585    is "unknown," it actually signifies uncharacterized A domains not yet incorporated into the

586    reference data set. This should prompt researchers to pay attention to these A domains, and

587    like in our case, subject them to further experimental investigation. This approach helped us

588    discover new substrates (valine, histidine, citrulline), which had not been previously

589    documented in pyoverdines and were therefore absent from the reference A domains. The

590    novel substrates identified through our structural assessment and mass spectrometry

591    experiment can subsequently be used to enhance the precision of our phylogeny-centered

592    substrate prediction technique in the future, creating a progressive feedback loop of

593    expanding knowledge. Taken together, supervised learning based on a few known

594    compounds produced by species from the same genus probably outperforms generalized

595    prediction algorithms trained on many products from a diverse set of microbes for NRPS

596    substrate predictions.

597    Our results show that both pyoverdine and receptor diversity has been vastly

598    underestimated. While considerable pyoverdine diversity (n=73) has already been captured in

599    previous studies, here we discovered 151 new variants. On the receptor side, the uncovered

600    novel diversity is more dramatic. One reason for this is that research on receptors has mainly

601    focused on the pathogen *P. aeruginosa*[19-22,50]. For this species, three different pyoverdine

602    types were described[21] together with three structurally different FpvA receptor types that each

603    recognize one of the pyoverdine types[22]. While our study confirmed that *P. aeruginosa* strains

604    (n = 554) indeed have only 3 pyoverdine-receptor systems, we also discovered 91 new FpvA

605    groups among environmental *Pseudomonas* spp. Our findings raise the question why there

606    are so many different pyoverdine and receptor variants. One potential explanation is that the

607    benefit of specific siderophores could be context-dependent and locally adapted to multitude

608    of different environmental conditions pseudomonads are exposed to. For example,

609    experimental work has revealed that pyoverdines can be cooperatively shared among strains

610    with matching receptors[33,51], or conversely, pyoverdines can serve as competitive agents by

611    locking away iron from species that have non-matching receptors[52]. Given that bioavailable

612    iron is limited in most natural and host-associated habitats[53-55], the unraveled functional

613    diversity is likely a direct evolutionary consequence of the struggle and competition of

614    microbes for iron. While experimental work is often restricted to a low number of strains, we

615    propose that our bioinformatic pipeline can be used to predict pyoverdine-mediated

616    interaction networks across thousands of strains and across different habitats. We will

617    address this point in a future study.

618    We believe our pipeline could be easily expanded to study iron competition in multi-

619    species communities in the future and perhaps in plant-microbe ecosystems, as siderophores

620    exist ubiquitously and are shared among microbes[56]. To move further, a key question is

621    whether our knowledge-guided approach can be applied to other important secondary

622    metabolites, such as antibiotics, toxins, biosurfactants and pigments? This answer is: not

623    directly but the pipeline development strategies are translational between different types of

624    compounds. As soon as sufficient case-by-case knowledge on a specific system is available,

625    the annotation strategies together with the feature sequence extraction and the phylogeny-

626    focused approach developed in our paper can be applied. For most of the secondary

627    metabolites listed above, there are no receptors as the compounds have purely extra-cellular

628    functions, which substantially simplifies the development of bioinformatic pipelines. In the long

629    run, it will certainly be possible to automate the steps implemented in our workflow so that the

630    algorithms can be applied to a large set of secondary metabolites when fed with an

631    appropriate training set.

632

633    **Data Availability**

634     The source code and parameters used are available in the supplementary material.

635

636     **Acknowledgements**

639

640     **Funding**

650

651     **Author contributions**

652     Shaohua Gu performed the majority of computational analysis in this research and drafted the

653     manuscript. Yuanzhe Shao built the pipeline for retrieving synthetase sequence information

654     and developed the phylogeny-centered method for predicting substrates. Karoline Rehm and

655     Laurent Bigler performed the experiment of pyoverdine structure elucidation. Di Zhang

656    depicted the circular plot. Ruolin He standardized of all NRPS structures into motif-intermotif

657    format. Jiqi Shao assisted in revising the manuscript. Alexandre Jousset and Ville-Petri

658    Friman offered insightful comments and assisted in revising and writing of the manuscript.

659    Rolf Kümmerli and Zhong Wei oversaw the project, designed experiments and revised the

660    manuscript. Zhiyuan Li conceptualized and oversaw the project, conducted the analysis of the

661    receptor annotated method, and revised the manuscript.

662

663    **Competing interests**

664    The authors declare no competing interests.

665

666    **References**

667    1    Zengler, K. & Palsson, B. O. A road map for the development of community systems

668         (CoSy) biology. *Nature Reviews Microbiology* **10**, 366-372, doi:10.1038/nrmicro2763

669         (2012).

670    2    Gu, C., Kim, G. B., Kim, W. J., Kim, H. U. & Lee, S. Y. Current status and applications

671         of genome-scale metabolic models. *Genome Biology* **20**, 121, doi:10.1186/s13059-

672         019-1730-3 (2019).

673    3    García-Jiménez, B., Torres-Bacete, J. & Nogales, J. Metabolic modelling approaches

674         for describing and engineering microbial communities. *Computational and Structural*

675         *Biotechnology Journal* **19**, 226-246, doi:10.1016/j.csbj.2020.12.003 (2021).

676    4    Colarusso, A. V., Goodchild-Michelman, I., Rayle, M. & Zomorrodi, A. R.

677         Computational modeling of metabolism in microbial communities on a genome-scale.

678      *Current Opinion in Systems Biology* **26**, 46-57, doi:10.1016/j.coisb.2021.04.001

679      (2021).

680   5   Scherlach, K. & Hertweck, C. Mediators of mutualistic microbe–microbe interactions.

681      *Natural Product Reports* **35**, 303-308, doi:10.1039/C7NP00035A (2018).

682   6   Trivedi, P., Leach, J. E., Tringe, S. G., Sa, T. & Singh, B. K. Plant–microbiome

683      interactions: from community assembly to plant health. *Nature Reviews Microbiology*

684      **18**, 607-621, doi:10.1038/s41579-020-0412-1 (2020).

685   7   Price-Whelan, A., Dietrich, L. E. P. & Newman, D. K. Rethinking 'secondary'

686      metabolism: physiological roles for phenazine antibiotics. *Nature Chemical Biology* **2**,

687      71-78, doi:10.1038/nchembio764 (2006).

688   8   Thirumurugan, D., Cholarajan, A., Raja, S. & Vijayakumar, R. An introductory chapter:

689      secondary metabolites. (2018).

690   9   Quinn, G. A., Banat, A. M., Abdelhameed, A. M. & Banat, I. M. Streptomyces from

691      traditional medicine: sources of new innovations in antibiotic discovery. *J Med*

692      *Microbiol* **69**, 1040-1048, doi:10.1099/jmm.0.001232 (2020).

693   10  Durand, G. A., Raoult, D. & Dubourg, G. Antibiotic discovery: history, methods and

694      perspectives. *International Journal of Antimicrobial Agents* **53**, 371-382,

695      doi:10.1016/j.ijantimicag.2018.11.010 (2019).

696   11  Penn, K. *et al.* Genomic islands link secondary metabolism to functional adaptation in

697      marine Actinobacteria. *The ISME Journal* **3**, 1193-1203, doi:10.1038/ismej.2009.58

698      (2009).

699   12  Andryukov, B., Mikhailov, V. & Besednova, N. The Biotechnological Potential of

700     Secondary Metabolites from Marine Bacteria. *Journal of Marine Science and*

701     *Engineering* **7** (2019).

702  13  Kautsar, S. A., Blin, K., Shaw, S., Weber, T. & Medema, M. H. BiG-FAM: the

703     biosynthetic gene cluster families database. *Nucleic Acids Research* **49**, D490-D497,

704     doi:10.1093/nar/gkaa812 (2021).

705  14  He, R. *et al.* Knowledge-guided data mining on the standardized architecture of

706     NRPS: Subtypes, novel motifs, and sequence entanglements. *PLOS Computational*

707     *Biology* **19**, e1011100, doi:10.1371/journal.pcbi.1011100 (2023).

708  15  Xu, Z., Park, T. J. & Cao, H. Advances in mining and expressing microbial

709     biosynthetic      gene      clusters.      *Crit      Rev      Microbiol*,      1-20,

710     doi:10.1080/1040841x.2022.2036099 (2022).

711  16  Keller, N. P. Fungal secondary metabolism: regulation, function and drug discovery.

712     *Nature Reviews Microbiology* **17**, 167-180, doi:10.1038/s41579-018-0121-1 (2019).

713  17  Ringel, M. T. & Brüser, T. The biosynthesis of pyoverdines. *Microb Cell* **5**, 424-437,

714     doi:10.15698/mic2018.10.649 (2018).

715  18  Hopkinson, B. M. & Morel, F. M. M. The role of siderophores in iron acquisition by

716     photosynthetic marine microorganisms. *BioMetals* **22**, 659-669, doi:10.1007/s10534-

717     009-9235-2 (2009).

718  19  Cobessi, D. *et al.* The Crystal Structure of the Pyoverdine Outer Membrane Receptor

719     FpvA from Pseudomonas aeruginosa at 3.6Å Resolution. *Journal of Molecular*

720     *Biology* **347**, 121-134, doi:10.1016/j.jmb.2005.01.021 (2005).

721  20  Diggle, S. P. & Whiteley, M. Microbe Profile: Pseudomonas aeruginosa: opportunistic

pathogen and lab rat. *Microbiology (Reading)* **166**, 30-33, doi:10.1099/mic.0.000860

(2020).

21    Meyer, J.-M. *et al.* Use of Siderophores to Type Pseudomonads: The Three

Pseudomonas Aeruginosa Pyoverdine Systems. *Microbiology* **143**, 35-43,

doi:10.1099/00221287-143-1-35 (1997).

22    Bodilis, J. *et al.* Distribution and evolution of ferripyoverdine receptors in

Pseudomonas aeruginosa. *Environmental Microbiology* **11**, 2123-2135,

doi:10.1111/j.1462-2920.2009.01932.x (2009).

23    Kümmerli, R. Iron acquisition strategies in pseudomonads: mechanisms, ecology, and

evolution. *BioMetals*, doi:10.1007/s10534-022-00480-8 (2022).

24    Meyer, J. M. Pyoverdines: pigments, siderophores and potential taxonomic markers

of fluorescent Pseudomonas species. *Arch Microbiol* **174**, 135-142,

doi:10.1007/s002030000188 (2000).

25    Blin, K. *et al.* antiSMASH 5.0: updates to the secondary metabolite genome mining

pipeline. *Nucleic Acids Research* **47**, W81-W87, doi:10.1093/nar/gkz310 (2019).

26    Bateman, A. *et al.* The Pfam protein families database. *Nucleic Acids Research* **32**,

D138-D141, doi:10.1093/nar/gkh121 (2004).

27    Winsor, G. L. *et al.* Enhanced annotations and features for comparing thousands of

Pseudomonas genomes in the Pseudomonas genome database. *Nucleic Acids

Research* **44**, D646-D653, doi:10.1093/nar/gkv1227 (2016).

28    Felnagle, E. A. *et al.* Nonribosomal peptide synthetases involved in the production of

medically relevant natural products. *Mol Pharm* **5**, 191-211, doi:10.1021/mp700137g

744        (2008).

745    29    Süssmuth, R. D. & Mainz, A. Nonribosomal Peptide Synthesis—Principles and

746        Prospects. *Angewandte Chemie International Edition* **56**, 3770-3821,

747        doi:10.1002/anie.201609079 (2017).

748    30    Butaitė, E., Kramer, J., Wyder, S. & Kümmerli, R. Environmental determinants of

749        pyoverdine production, exploitation and competition in natural Pseudomonas

750        communities. *Environmental Microbiology* **20**, 3629-3642, doi:10.1111/1462-

751        2920.14355 (2018).

752    31    Rehm, K., Vollenweider, V., Kümmerli, R. & Bigler, L. A comprehensive method to

753        elucidate pyoverdines produced by fluorescent Pseudomonas spp. by UHPLC-HR-

754        MS/MS. *Analytical and Bioanalytical Chemistry* **414**, 2671-2685, doi:10.1007/s00216-

755        022-03907-w (2022).

756    32    Butaitė, E., Baumgartner, M., Wyder, S. & Kümmerli, R. Siderophore cheating and

757        cheating resistance shape competition for iron in soil and freshwater Pseudomonas

758        communities. *Nature Communications* **8**, 414, doi:10.1038/s41467-017-00509-4

759        (2017).

760    33    Kramer, J., Özkaya, Ö. & Kümmerli, R. Bacterial siderophores in community and host

761        interactions. *Nature Reviews Microbiology* **18**, 152-163, doi:10.1038/s41579-019-

762        0284-4 (2020).

763    34    Chan, D. C. K. & Burrows, L. L. &lt;em&gt;Pseudomonas aeruginosa&lt;/em&gt; FpvB

764        is a high-affinity transporter for xenosiderophores ferrichrome and ferrioxamine B.

765        *bioRxiv*, 2022.2009.2020.508722, doi:10.1101/2022.09.20.508722 (2022).

766    35    González, J. *et al.* Loss of a pyoverdine secondary receptor in Pseudomonas

767          aeruginosa results in a fitter strain suitable for population invasion. *The ISME Journal*

768          **15**, 1330-1343, doi:10.1038/s41396-020-00853-2 (2021).

769    36    Schalk, I. J., Mislin, G. L. & Brillet, K. Structure, function and binding selectivity and

770          stereoselectivity of siderophore-iron outer membrane transporters. *Curr Top Membr*

771          **69**, 37-66, doi:10.1016/B978-0-12-394390-3.00002-1 (2012).

772    37    Greenwald, J. *et al.* FpvA bound to non-cognate pyoverdines: molecular basis of

773          siderophore recognition by an iron transporter. *Molecular Microbiology* **72**, 1246-1259,

774          doi:10.1111/j.1365-2958.2009.06721.x (2009).

775    38    Almeida, A. *et al.* A new genomic blueprint of the human gut microbiota. *Nature* **568**,

776          499-504, doi:10.1038/s41586-019-0965-1 (2019).

777    39    Schloss, P. D. & Handelsman, J. Metagenomics for studying unculturable

778          microorganisms: cutting the Gordian knot. *Genome Biology* **6**, 229, doi:10.1186/gb-

779          2005-6-8-229 (2005).

780    40    Handelsman, J. Metagenomics: Application of Genomics to Uncultured

781          Microorganisms. *Microbiology and Molecular Biology Reviews* **68**, 669-685,

782          doi:10.1128/MMBR.68.4.669-685.2004 (2004).

783    41    Tsilimigras, M. C. B. & Fodor, A. A. Compositional data analysis of the microbiome:

784          fundamentals, tools, and challenges. *Annals of Epidemiology* **26**, 330-335,

785          doi:10.1016/j.annepidem.2016.03.002 (2016).

786    42    Weiss, S. *et al.* Correlation detection strategies in microbial data sets vary widely in

787          sensitivity and precision. *ISME J* **10**, 1669-1681, doi:10.1038/ismej.2015.235 (2016).

43    Faust, K. & Raes, J. Microbial interactions: from networks to models. *Nature Reviews Microbiology* **10**, 538-550, doi:10.1038/nrmicro2832 (2012).

44    Khayatt, B. I., Overmars, L., Siezen, R. J. & Francke, C. Classification of the Adenylation and Acyl-Transferase Activity of NRPS and PKS Systems Using Ensembles of Substrate Specific Hidden Markov Models. *PLOS ONE* **8**, e62136, doi:10.1371/journal.pone.0062136 (2013).

45    Minowa, Y., Araki, M. & Kanehisa, M. Comprehensive Analysis of Distinctive Polyketide and Nonribosomal Peptide Structural Motifs Encoded in Microbial Genomes. *Journal of Molecular Biology* **368**, 1500-1517, doi:10.1016/j.jmb.2007.02.099 (2007).

46    Prieto, C., García-Estrada, C., Lorenzana, D. & Martín, J. F. NRPSsp: non-ribosomal peptide synthase substrate predictor. *Bioinformatics* **28**, 426-427, doi:10.1093/bioinformatics/btr659 (2012).

47    Röttig, M. *et al.* NRPSpredictor2—a web server for predicting NRPS adenylation domain specificity. *Nucleic Acids Research* **39**, W362-W367, doi:10.1093/nar/gkr323 (2011).

48    Zierep, P. F., Ceci, A. T., Dobrusin, I., Rockwell-Kollmann, S. C. & Günther, S. SeMPI 2.0-A Web Server for PKS and NRPS Predictions Combined with Metabolite Screening in Natural Product Databases. *Metabolites* **11**, doi:10.3390/metabo11010013 (2020).

49    Fan, J. *et al.* Biosynthetic diversification of peptaibol mediates fungus-mycohost interactions. *bioRxiv* (2022).

810  50  Smith, E. E., Sims, E. H., Spencer, D. H., Kaul, R. & Olson, M. V. Evidence for

811       diversifying selection at the pyoverdine locus of Pseudomonas aeruginosa. *J*

812       *Bacteriol* **187**, 2138-2147, doi:10.1128/jb.187.6.2138-2147.2005 (2005).

813  51  Gu, S. *et al.* Competition for iron drives phytopathogen control by natural rhizosphere

814       microbiomes. *Nature Microbiology* **5**, 1002-1010, doi:10.1038/s41564-020-0719-8

815       (2020).

816  52  Figueiredo, A. R. T., Özkaya, Ö., Kümmerli, R. & Kramer, J. Siderophores drive

817       invasion dynamics in bacterial communities through their dual role as public good

818       versus public bad. *Ecology Letters* **25**, 138-150, doi:10.1111/ele.13912 (2022).

819  53  Andrews, S. C., Robinson, A. K. & Rodríguez-Quiñones, F. Bacterial iron homeostasis.

820       *FEMS Microbiology Reviews* **27**, 215–237, doi:10.1016/s0168-6445(03)00055-x

821       (2003).

822  54  Boyd, P. W. & Ellwood, M. J. The biogeochemical cycle of iron in the ocean. *Nat*

823       *Geosci* **3**, 675-682, doi:10.1038/ngeo964 (2010).

824  55  Emerson, D., Roden, E. & Twining, B. The microbial ferrous wheel: iron cycling in

825       terrestrial, freshwater, and marine environments. *Frontiers in Microbiology* **3**,

826       doi:10.3389/fmicb.2012.00383 (2012).

827  56  Ruolin, H. *et al.* SIDERITE: Unveiling Hidden Siderophore Diversity in the Chemical

828       Space Through Digital Exploration. *bioRxiv*, 2023.2008.2031.555687,

829       doi:10.1101/2023.08.31.555687 (2023).

830