

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36

Title: ActDES – a Curated Actinobacterial Database for Evolutionary Studies

Jana K. Schniete^{1,3}, Nelly Selem-Mojica², Anna S. Birke¹, Pablo Cruz-Morales², Iain S. Hunter¹, Francisco Barona-Gómez², & Paul A. Hoskisson^{1*}

Affiliations: ¹Strathclyde Institute of Pharmacy and Biomedical Sciences, University of Strathclyde, 161 Cathedral Street, Glasgow, G4 0RE, United Kingdom

²Evolution of Metabolic Diversity Laboratory, Langebio, Cinvestav-IPN, Libramiento Norte Carretera Leon Km 9.6, 36821, Irapuato, Guanajuato, México.

³Biology Department, Edge Hill University, St Helens Road, Ormskirk, Lancashire, L39 4QP, UK.

***Corresponding Author:** ¹Strathclyde Institute of Pharmacy and Biomedical Sciences, University of Strathclyde, 161 Cathedral Street, Glasgow, G4 0RE, United Kingdom, Tel.: +44 (0)141 548 2819; Fax: +44 (0)141 548 4124; E-mail: Paul.hoskisson@strath.ac.uk

Key words: *Streptomyces*, evolution, primary metabolism, specialised metabolism, antibiotics, natural product, biosynthetic gene cluster.

Repositories: All data, databases, files and scripts can be found available on Microbiology Society Figshare [link to be added by the Society], <https://github.com/nselem/ActDES> and <https://hub-binder.mybinder.ovh/user/nselem-actdes-2p4esa8n/tree#notebooks>

Abbreviations: BGC - Biosynthetic gene cluster; CCR – Carbon-catabolite repression; CDS – Coding sequence; HGT – Horizontal gene transfer; PTS - phosphoenolpyruvate-dependent phosphotransferase system; WGS - whole genome sequence

37 **Abstract**

38 *Actinobacteria* is a large and diverse phylum of bacteria that contains medically and
39 ecologically relevant organisms. Many members are valuable sources of bioactive natural
40 products and chemical precursors that are exploited in the clinic and made using the enzyme
41 pathways encoded in their complex genomes. Whilst the number of sequenced genomes has
42 increased rapidly in the last twenty years, the large size, complexity and high G+C content of
43 many Actinobacterial genomes means that the sequences remain incomplete and consist of
44 large numbers of contigs with poor annotation, which hinders large scale comparative genomic
45 and evolutionary studies. To enable greater understanding and exploitation of Actinobacterial
46 genomes, specialised genomic databases must be linked to high-quality genome sequences.
47 Here we provide a curated database of 612 high-quality actinobacterial genomes from 80
48 genera, chosen to represent a broad phylogenetic group with equivalent genome re-
49 annotation. Utilising this database will provide researchers with a framework for evolutionary
50 and metabolic studies, to enable a foundation for genome and metabolic engineering, to
51 facilitate discovery of novel bioactive therapeutics and studies on gene family evolution.

52

53

54 **Significance as a bioresource to the community**

55 The Actinobacteria are a large diverse phylum of bacteria, often with large, complex genomes
56 with a high G+C content. Sequence databases have great variation in the quality of
57 sequences, equivalence of annotation and phylogenetic representation, which makes it
58 challenging to undertake evolutionary and phylogenetic studies. To address this, we have
59 assembled a curated, taxa-specific, non-redundant database to aid detailed comparative
60 analysis of Actinobacteria. ActDES constitutes a novel resource for the community of
61 Actinobacterial researchers that will be useful primarily for two types of analyses: (i)
62 comparative genomic studies - facilitated by reliable identification of orthologs across a set of
63 defined phylogenetically-representative genomes and (ii) phylogenomic studies which will be
64 improved by identification of gene subsets at specified taxonomic level. These analyses can
65 then act as a springboard for the studies of the evolution of virulence genes, the evolution of
66 metabolism and identification of targets for metabolic engineering.

67

68 **Data summary**

69 All genome sequences used in this study can be found in the NCBI taxonomy browser
70 <https://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi> and are summarised along with
71 Accession numbers in Table S1

72 All other data are available on Microbiology Society Figshare [link to be added by the Society]
73

74 a. Perl script files available on GitHub <https://github.com/nselem/ActDES> including
75 details of how to batch annotate genomes in RAST from the terminal
76 <https://github.com/nselem/myrast>

77 b. **Supp. Table S1** List of genomes from NCBI (Actinobacteria database.xlsx)
78 available on Microbiology Society Figshare
79 <https://doi.org/10.6084/m9.figshare.13143407.v1>

80 c. CVS genome annotation files including the FASTA files of nucleotide and amino
81 acids sequences (individual .cvs files) available on Microbiology Society Figshare
82 <https://doi.org/10.6084/m9.figshare.13143407.v1>

83 d. BLAST nucleotide database (.fasta file) available on Microbiology Society Figshare
84 <https://doi.org/10.6084/m9.figshare.13143407.v1>

85 e. BLAST protein database (.fasta file) available on Microbiology Society Figshare
86 <https://doi.org/10.6084/m9.figshare.13143407.v1>

87 f. Supp. Table S2 Expansion table genus level (Expansion table.xlsx Tab Genus
88 level) available on Microbiology Society Figshare [link to be added by the Society]

89 g. Supp. Table S2 Expansion table species level (Expansion table.xlsx Tab species
90 level) available on Microbiology Society Figshare
91 <https://doi.org/10.6084/m9.figshare.13143407.v1>

92 h. All GlcP and Glk data – blast hits from ActDES database, MUSCLE Alignment files
93 and .nwk tree files can be found at available on Microbiology Society Figshare
94 <https://doi.org/10.6084/m9.figshare.13143407.v1>

95 i. Interactive trees in Microreact for Glk tree
96 https://microreact.org/project/w_KDfn1xA/5a178533 and associated files can be
97 found at Microbiology Society Figshare
98 <https://doi.org/10.6084/m9.figshare.13143407.v1>

99 j. Interactive trees in Microreact for GlcP tree
100 https://microreact.org/project/VBUdiQ5_k/045c95e1 and associated files can be
101 found at Microbiology Society Figshare
102 <https://doi.org/10.6084/m9.figshare.13143407.v1>

103 k. Jupyter Notebook for exploring ActDES in MyBinder can be found
104 <https://github.com/nselem/ActDES>

105 **Introduction**

106 The increase in availability of bacterial whole genome sequencing (WGS) provides large
107 amounts of data for evolutionary and phylogenetic analysis. However, there is great variation
108 in the quality, annotation and phylogenetic skew of the data available in large universal
109 databases, meaning that evolutionary and phylogenetic studies can be challenging. To
110 address this variation, curated, high-level, taxa-specific, non-redundant sub-databases need
111 to be assembled to aid detailed analysis. Given that there is a direct correlation between
112 phylogenetic distance and the discovery of novel function [1–3], it is imperative that any
113 derived databases must be phylogenetically representative and non-redundant to enable
114 insight into the evolution of genes, proteins and pathways within a given group of taxa [1].

115 The phylum *Actinobacteria* is a major taxon amongst the *Bacteria*, which includes
116 phenotypically and morphologically diverse organisms found on every continent and in
117 virtually every ecological niche [4]. They are particularly common in soils, yet within their ranks
118 are potential human and animal pathogens such as *Corynebacterium*, *Mycobacterium*,
119 *Nocardia* and *Tropheryma*, inhabitants of the gastrointestinal tract (*Bifidobacterium* and
120 *Scardovia*) as well as plant commensals and pathogens such as *Frankia*, *Leifsonia* and
121 *Clavibacter* [4, 5]. Perhaps the most notable trait of the phylum is the renowned ability to
122 produce bioactive natural products such as antibiotics, anti-cancer agents and immuno-
123 suppressive agents, with genera such as *Amycolatopsis*, *Micromonospora* and *Streptomyces*
124 being particularly prominent [6]. As a result, computational ‘mining’ of Actinobacterial genomes
125 has become an important part of the drug discovery pipeline, with increasing numbers of online
126 resources and software devoted to identification of natural product biosynthetic gene clusters
127 (BGCs)[7–9]. It is important to move beyond approaches that rely on similarity searches of
128 known BGCs and to expand searches to identify hidden chemical diversity within the genomes
129 [6, 7, 10–13].

130 A recent study of 830 Actinobacterial genomes found >11,000 BGCs comprising 4,122
131 chemical families, indicating that there is a vast diversity of strains and chemistry to exploit
132 [14], yet within each of these strains there will be hidden diversity in the form of cryptic BGCs.
133 To exploit this undiscovered diversity as the technology develops and databases expand, new
134 biosynthetic logic will emerge, yet we know little of how natural selection shapes the evolution
135 of BGCs and how biosynthetic precursors are supplied to gene products of BGCs from primary
136 metabolism and to identify targets for metabolic engineering of industrially relevant strains.
137 Such logic will expedite industrial strain improvement processes, enabling titre increases and
138 development of novel molecules, as well as the engineering of strains to use more sustainable
139 feedstocks.

140 To aid this process we have created an Actinobacterial metabolism database including
141 functional annotations for enzymes from 612 species to enable phylum-wide interrogation of

142 gene expansion events that may indicate adaptive evolution, help shape metabolic robustness
143 for antibiotic production [15] or enable the identification of targets for metabolic engineering.
144 **Actinobacterial Database for Evolutionary Studies (ActDES)** provides a curated list of high-
145 quality, phylum specific genomes and data to help researchers navigate the redundancy and
146 inconsistency in sequence databases in a simplified format that enables researchers with little
147 taxonomic knowledge to develop testable evolutionary hypotheses. To demonstrate the utility
148 of ActDES, we have detailed its construction and used it to investigate the glucose
149 permease/glucokinase system phylogeny across the Actinobacteria.

150 **Methods**

151 We generated ActDES, a database for evolutionary analysis of Actinobacterial genomes, in
152 two formats - a database for interrogation by BLASTn or BLASTp for phylogenetic analysis
153 and a primary metabolic gene expansion table, which can be mined at different taxonomic
154 levels (**Supp. Table S1 and S2**) for specific metabolic functions from primary metabolism. A
155 schematic overview of generation of the dataset is shown in **Fig.1**.

156 The database was generated via the NCBI taxonomy browser
157 (<https://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi>) to identify Actinobacterial
158 genome sequences. The quality of the genome sequences was filtered by the number of
159 contigs (<100 contigs per 2Mb of genome sequence) and the genomes were downloaded from
160 the NCBI WGS repository (<https://www.ncbi.nlm.nih.gov/Traces/wgs/>). These genomes were
161 then dereplicated to ensure that the database comprised a wide taxonomic range of the
162 phylum, resulting in 612 species from 80 genera within 13 suborders of the Actinobacteria
163 (**Supp. Table S1**).

164 Each of these 612 genomes was reannotated using RAST. Default settings were used to
165 ensure equivalence of annotation across the database and the annotation files of each
166 genome were downloaded (Data File: cvs files). These annotation files were subsequently
167 used to extract all protein and nucleotide sequences into two files. Each of these files was
168 subsequently converted into BLAST databases (a protein database and a nucleotide database
169 <https://doi.org/10.6084/m9.figshare.12167724>) to facilitate phylogenetic analysis. Sequences
170 of interest can be aligned using MUSCLE [16] and phylogenetic trees constructed using a
171 range of tree construction software such as QuickTree [17], IQ tree [18] or MrBayes [19].
172 Subsequent trees may be visualized in software such as FigTree v1.4.2
173 (<http://tree.bio.ed.ac.uk/software/figtree/>).

174 The RAST annotation files were also used to extract the functional roles of each CDS per
175 genome and the level of gene expansion was assessed for each genome by counting the
176 number of genes per species per functional category (gene function annotation). The dataset
177 was then curated manually for central carbon metabolism and amino acid biosynthesis
178 pathways to create the gene expansion table (**Supp. Table S2**) with the organisms grouped
179 according to their taxonomic position. The quality of the data was checked at each step for
180 duplicates and inconsistencies and was curated manually to exclude faulty entries.

181 As the NCBI taxonomy browser database is overrepresented in *Streptomyces* genomes due
182 to the number of species that have been sequenced relative to other Actinobacteria, this is
183 also reflected in the ActDES database (288 *Streptomyces* genomes from a total of 612
184 genomes). However this was addressed in the expansion table (**Supp. Table S2**) by
185 calculating the mean occurrence of each functional category within each genus and then
186 calculating an overall mean for the phylum to compensate. The mean occurrence of each

187 functional category per genus plus the standard deviation was also calculated and this was
188 used to analyse the occurrence of each functional gene category per species within **Supp.**
189 **Table S2**. A gene function annotation with a gene copy number value above the mean plus
190 the standard deviation for each genus, indicated that there had been a gene expansion event
191 in that species and this was noted. The gene expansion table (**Supp. Table S2**) enables
192 researchers to identify groups of genes of interest for subsequent phylogenetic and
193 evolutionary analysis, which can be performed with confidence due to the highly curated
194 nature of the data included in the database.

195

196 **Results**

197 The gene expansion table (**Supp. Table S2**) lists 612 species of 80 genera within the
198 Actinobacteria with data that provides an extensive analysis at the phylum level, which is the
199 starting point for detailed phylogenomic studies. Gene expansions were identified in separate
200 datasets at the genus and species levels, along with details of the numbers of genes in each
201 functional category per species and the average numbers of genes in each functional category
202 per genus expanded within the genomes. These data can be used subsequently in
203 phylogenomic analyses to identify targets for metabolic engineering and gene function studies.
204 Identification of expanded gene families may also facilitate the recognition of novel natural
205 product biosynthetic gene clusters, for which gene expansion events of primary metabolic
206 genes have been classified to be associated within BGCs as biosynthetic enzymes or through
207 provision of additional copies of antibiotic targets that may subsequently function as resistance
208 mechanisms [6, 11, 20–24].

209 This database has found utility for studying primary metabolic gene expansions in
210 *Streptomyces*. It enabled a detailed *in silico* analysis of the duplication event leading to the
211 two pyruvate kinases in the genus of *Streptomyces* subsequently enabling the functional
212 characterisation of the two isoenzymes to reveal how they contribute to metabolic robustness
213 [15]. ActDES may also be useful for investigating the distribution of primary metabolic genes
214 across the phylum to link phenotype to genotype and phylogenetic position. An initial RpoB
215 phylogeny has been constructed previously using this database [15] which provided a robust
216 universal phylogeny for comparison of individual protein trees [25].

217 To demonstrate the utility of ActDES, the glucose permease/glucokinase system of the
218 Actinobacteria was investigated. The role of nutrient-sensing in regulation of antibiotic
219 biosynthesis is well known [26] with the enzyme glucokinase (Glc) playing a central role in
220 carbon-catabolite repression in *Streptomyces* [27]. In most bacteria, CCR is mediated by the
221 phosphoenolpyruvate-dependent phosphotransferase system (PTS), yet in *Streptomyces*,
222 glucose uptake is mediated by the Major-Facilitator Superfamily (MFS) transporter, glucose
223 permease (GlcP), and there is evidence for direct interaction between Glc and GlcP which

224 may mediate CCR [28]. Understanding the nature and distribution of these enzymes will play
225 a key role in developing industrial fermentations with glucose as major carbon source.
226 Investigating the distribution of the glucose permease/glucokinase system across the phylum
227 shows that GlcP and Glk have been the subject of gene expansion events in some members
228 of the Streptomycetales, most notably the *Streptomyces*, with a patchy distribution of the
229 Glk/GlcP system across the remainder of the phylum (**Table S2; Genus tab**). However, where
230 the Glk/GlcP system is found, the number of expansion events observed is greater for Glk
231 than for GlcP (**Fig. 2 A & B**). The phylogenetic trees (**Fig. 2A & 2B**) clearly show two clades
232 for Glk and GlcP within the Streptomycetales (Interactive trees are available via Microreact
233 [29]: Glk https://microreact.org/project/w_KDfn1xA/5a178533 and GlcP
234 https://microreact.org/project/VBUdiQ5_k/045c95e1). However, these clades differ in the
235 number of sequences, with the Glk clades being equal in number, suggesting that a duplication
236 event has occurred within the Streptomycetales (**Fig. 2A**). This is consistent throughout the
237 order, with the patterns largely the same as that observed for *S. coelicolor*. This species has
238 two ROK-family ATP-dependent glucokinases, SCO2126 (*glkA*) and SCO6260, that share
239 around 50% amino acid sequence identity and each is found in one of the distinct clades
240 (permease-associated kinases and orphan kinases; **Fig. 2A**). Whilst SCO2126 is a GlcP-
241 associated kinase, the gene encoding SCO6260 is located in an operon including genes
242 encoding a putative carbohydrate ABC-transporter system, which has been reported
243 previously [30]. SCO6260 appears to be the only glucokinase in the database that is
244 associated with an ABC-transporter. This may suggest that expansion of the Glk gene family
245 in Streptomycetales might have occurred to extend the number of CCR-mediating kinases in
246 the genome, adding increased regulatory complexity to carbohydrate metabolism in this group
247 of organisms that use CCR as a major regulator of specialised metabolism.

248 The two clades for GlcP within the Streptomycetales differ in size suggesting either gene
249 duplication followed by gene loss, or an expansion through horizontal gene transfer (HGT)
250 has occurred. A detailed examination of these clades by species (**Table S2; Species Tab**)
251 shows the presence of both scenarios. There are duplicated enzymes located within the same
252 clade (as observed in *S. coelicolor*; Group I) or additional copies of the permease which are
253 located in a phylogenetically distinct clade, which lacks congruence with the RpoB tree [15]
254 and remarkably consists entirely of sequences from the genus *Streptomyces* (Group II; **Fig.**
255 **2B**). This suggests that they may have been acquired via HGT. The expansive nature of the
256 duplicated Glk enzymes compared to GlcP may be due to the role played in CCR by the GlkA
257 enzymes [27] and the different transcriptional activities under glycolytic and gluconeogenic
258 conditions [31], yet quite how these different Glk enzymes interact with the permease(s) under
259 various conditions requires further experimental investigation to understand their exact

260 physiological role, and how this may be translated in to industrial strain improvement
261 processes.

262

263 **Discussion**

264 Large scale WGS and phylogenomic analysis is increasingly used for identifying targets for
265 genome and metabolic engineering, studies of metabolic capabilities, pathogen
266 phylogenomics and evolutionary studies. These studies are often complicated by the large
267 number of sequences in the databases, database redundancy and the poor quality of some
268 genome sequence data. The development of the high-quality, curated ActDES database,
269 reported here, enables phylum-wide taxonomic representation of the Actinobacteria coupled
270 with quality-filtered genome data and equivalent annotation for each CDS.

271 The intended primary use of ActDES will be in the study of primary metabolism, but it is not
272 limited. It can also inform the development and evolution of metabolism in strains that produce
273 bioactive metabolites, given the high representation of genera renowned for their ability to
274 produce natural products such as *Streptomyces* and *Micromonospora*. Due to a greater
275 understanding of BGC evolution and genome organisation in Actinobacteria it is becoming
276 increasingly clear that genes whose functions are in primary metabolism may actually
277 contribute directly to the biosynthesis of specialised metabolites and, hence, the identification
278 of duplicates may indicate the presence of cryptic BGCs [6, 11] or, when associated with
279 precursor biosynthetic genes, provide the raw material for the enzymes across multiple BGCs
280 [32–34].

281 ActDES may also find utility in evolutionary studies of expanded gene families across the
282 Actinobacterial phylum that contribute to virulence such as the *mce* locus which is known to
283 facilitate host survival in Mycobacteria [35], but also facilitates xenobiotic substrate uptake in
284 *Rhodococcus* [36] and enables root colonization and survival in *Streptomyces* [37]. With
285 phylum-wide taxonomic representation of established Actinobacterial animal and plant
286 pathogens, the scope for evolutionary studies using these data is enormous.

287

288 **Usage Notes**

289 The CVS files of each genome contains the RAST annotation details in addition to the DNA
290 and protein sequences for each annotated CDS
291 (<https://doi.org/10.6084/m9.figshare.12167880>). The Genome list contains the RAST ID
292 (which is equivalent to the name of the .cvs file) along with the NCBI ID (sequence ID; **Table**
293 **S1**) plus the species name which are included in the dataset. Further details of annotating
294 batches of genomes in RAST can be found here <https://github.com/nselem/myrast>.

295 The primary metabolism expansion tables (**Supp. Table S2**) are organised by metabolic
296 pathway along the top row with the Enzyme Commission (EC) number and functional

297 annotation, with the first column being the taxonomic assignment. The genus table shows the
298 average number of genes of the annotated function. Highlighted cells reflect gene expansion
299 events, i.e. those genes that are present in higher number than the overall mean across the
300 database plus the standard deviation.

301 It is suggested that the gene expansion table (**Supp. Table S2**) is searched in the first instance
302 (either by species or genus of interest or by a specific enzymatic function). This can be carried
303 out by a simple text search. This will then allow the identification of a query sequence from a
304 species or gene of interest (either nucleotide or amino acid sequence) which can then be
305 searched against the curated BLAST database allowing a detailed phylogenetic analysis of a
306 gene/protein of interest by using standard alignment and tree building software tools. These
307 data can also be used in detailed evolutionary analysis of selection, mutation rates etc. We
308 have set up a Jupyter Notebook through the MyBinder (<https://mybinder.org/>) project here to
309 enable ease of use of the code <https://github.com/nselem/ActDES> with a tutorial to enable use
310 of the database <https://github.com/nselem/ActDES>.

311

312

313 **Funding Information**

314 This work was funded through PhD studentships from the Scottish University Life Science
315 Alliance (SULSA) to JKS and by the Industrial Biotechnology Innovation Centre (IBioIC) with
316 GlaxoSmithKline funded PhD studentship to ASB.

317

318 **Author Contributions**

319 Conceptualization: JKS, NS-M, PAH, FB-G

320 Data curation: JKS & NS-M

321 Formal analysis: JKS, N.S-M, PC-M, & ASB

322 Funding acquisition: PAH, FB-G

323 Methodology: JKS, NS-M, PC-M, FB-G & PAH

324 Project administration: PAH & FB-G

325 Supervision: PAH, ISH & F.B-G

326 Writing – original draft: JKS, ASB & PAH

327 Writing – review and editing: JKS, N.S-M, PC-M, ASB, ISH, FB-G & PAH

328

329 **Acknowledgements**

330 We thank the Scottish Universities Life Science Alliance (SULSA) for BioScape PhD funding
331 to JKS, and a Mac Robertson Travelling Scholarship awarded to JKS to visit the laboratory of
332 FB-G., an Industrial Biotechnology Innovation Centre (IBioIC) and GlaxoSmithKline funded
333 PhD studentship to ASB, and a NERC (grant NE/M001415/1), BBSRC (grants BB/N023544/1
334 and BB/T001038/1) and BBSRC/NPRONET (grant NPRONET POC045) to PAH. PAH would
335 also like to acknowledge the support of the Royal Academy of Engineering for the award of a
336 Research Chair in Engineering Biology of Antibiotic Production. Work in the laboratory of FB-
337 G. was funded by CONACyT, Mexico Metabolic Robustness in Streptomyces, and Langebio
338 institutional funds to support PC-M. as a postdoctoral fellow and Royal Society Newton
339 Advanced Fellowship (NAF\R2\18063).

340 **Conflicts of interest**

341 The authors declare that there are no conflicts of interest

342 **Ethical statement**

343 No ethical approval was required.

344 **References**

- 345 1. **Mukherjee S, Seshadri R, Varghese NJ, Eloe-Fadrosh EA, Meier-Kolthoff JP, et al.**
346 1,003 reference genomes of bacterial and archaeal isolates expand coverage of the tree of
347 life. *Nat Biotechnol* 2017;35:676–683.
- 348 2. **Wu D, Hugenholtz P, Mavromatis K, Pukall R, Dalin E, et al.** A phylogeny-driven genomic
349 encyclopaedia of Bacteria and Archaea. *Nature* 2009;462:1056–60.
- 350 3. **Kunin V, Cases I, Enright AJ, Lorenzo V de, Ouzounis CA.** Myriads of protein families,
351 and still counting. *Genome Biol* 2003;4:401.
- 352 4. **Goodfellow M.** Bergey’s Manual of Systematics of Archaea and Bacteria. 2015; Springer.
- 353 5. **Ventura M, Canchaya C, Tauch A, Chandra G, Fitzgerald GF, et al.** Genomics of
354 Actinobacteria: Tracing the Evolutionary History of an Ancient Phylum. *Microbiology and*
355 *Molecular Biology Reviews* 2007;71:495–548.
- 356 6. **Chevrette MG, Gutiérrez-García K, Selem-Mojica N, Aguilar-Martínez C, Yañez-Olvera**
357 **A, et al.** Evolutionary dynamics of natural product biosynthesis in bacteria. *Nat Prod Rep.*
358 Epub ahead of print 2019. DOI: 10.1039/c9np00048h.
- 359 7. **Adamek M, Alanjary M, Ziemert N.** Applied evolution: phylogeny-based approaches in
360 natural products research. *Nat Prod Rep* 2019;36:1295–1312.
- 361 8. **Ziemert N, Alanjary M, Weber T.** The evolution of genome mining in microbes – a review.
362 *Nat Prod Rep* 2016;33:988–1005.
- 363 9. **Medema MH, Fischbach MA.** Computational approaches to natural product discovery. *Nat*
364 *Chem Biol* 2015;11:639–648.
- 365 10. **Adamek M, Alanjary M, Sales-Ortells H, Goodfellow M, Bull AT, et al.** Comparative
366 genomics reveals phylogenetic distribution patterns of secondary metabolites in
367 *Amycolatopsis* species. *BMC Genomics* 2018;19:426.
- 368 11. **Cruz-Morales P, Kopp JF, Martínez-Guerrero C, Yañez-Guerra LA, Selem-Mojica N,**
369 **et al.** Phylogenomic Analysis of Natural Products Biosynthetic Gene Clusters Allows
370 Discovery of Arseno-Organic Metabolites in Model Streptomyces. *Genome Biol Evol*
371 2016;8:1906–1916.
- 372 12. **Navarro-Muñoz JC, Selem-Mojica N, Mullaney MW, Kautsar SA, Tryon JH, et al.** A
373 computational framework to explore large-scale biosynthetic diversity. *Nat Chem Biol* 2019;1–
374 9.

- 375 13. **Sélem-Mojica N, Aguilar C, Gutiérrez-García K, Martínez-Guerrero CE, Barona-**
376 **Gómez F.** EvoMining reveals the origin and fate of natural product biosynthetic enzymes.
377 *Microb Genom*;5. Epub ahead of print 2019. DOI: 10.1099/mgen.0.000260.
- 378 14. **Doroghazi JR, Metcalf WW.** Comparative genomics of actinomycetes with a focus on
379 natural product biosynthetic genes. *BMC Genomics* 2013;14:611.
- 380 15. **Schniete JK, Cruz-Morales P, Selem-Mojica N, Fernández-Martínez LT, Hunter IS, et**
381 **al.** Expanding Primary Metabolism Helps Generate the Metabolic Robustness To Facilitate
382 Antibiotic Biosynthesis in *Streptomyces*. *Mbio* 2018;9:e02283-17.
- 383 16. **Edgar RC.** MUSCLE: a multiple sequence alignment method with reduced time and space
384 complexity. *BMC Bioinformatics* 2004;5:113.
- 385 17. **Howe K, Bateman A, Durbin R.** QuickTree: building huge Neighbour-Joining trees of
386 protein sequences. *Bioinformatics* 2002;18:1546–1547.
- 387 18. **Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, et al.** IQ-TREE 2:
388 New models and efficient methods for phylogenetic inference in the genomic era. *Mol Biol*
389 *Evol* 2020;37:1530–1534.
- 390 19. **Ronquist F, Teslenko M, Mark P van der, Ayres DL, Darling A, et al.** MrBayes 3.2:
391 efficient Bayesian phylogenetic inference and model choice across a large model space.
392 *Systematic Biol* 2012;61:539–42.
- 393 20. **Tang X, Li J, Millán-Aguiñaga N, Zhang JJ, O’Neill EC, et al.** Identification of
394 Thiotetronic Acid Antibiotic Biosynthetic Pathways by Target-directed Genome Mining. *Acs*
395 *Chem Biol* 2015;10:2841–9.
- 396 21. **Schmidt KL, Peterson ND, Kustus RJ, Wissel MC, Graham B, et al.** A Predicted
397 ABC Transporter, FtsEX, Is Needed for Cell Division in *Escherichia coli*. *J Bacteriol*
398 2004;186:785–793.
- 399 22. **Steffensky M, Mühlenweg A, Wang Z-X, Li S-M, Heide L.** Identification of the Novobiocin
400 Biosynthetic Gene Cluster of *Streptomyces spheroides* NCIB 11891. *Antimicrob Agents Ch*
401 2000;44:1214–1222.
- 402 23. **Kling A, Lukat P, Almeida DV, Bauer A, Fontaine E, et al.** Targeting DnaN for
403 tuberculosis therapy using novel griselimycins. *Science* 2015;348:1106–1112.
- 404 24. **Peterson RM, Huang T, Rudolf JD, Smanski MJ, Shen B.** Mechanisms of Self-
405 Resistance in the Platensimycin- and Platencin-Producing *Streptomyces platensis* MA7327
406 and MA7339 Strains. *Chem Biol* 2014;21:389–397.

- 407 25. **Case RJ, Boucher Y, Dahllöf I, Holmström C, Doolittle WF, et al.** Use of 16S rRNA and
408 rpoB Genes as Molecular Markers for Microbial Ecology Studies. *Appl Environ Microb*
409 2007;73:278–288.
- 410 26. **Fernández-Martínez LT, Hoskisson PA.** Expanding, integrating, sensing and
411 responding: the role of primary metabolism in specialised metabolite production. *Curr Opin*
412 *Microbiol* 2019;51:16–21.
- 413 27. **Gubbens J, Janus M, Florea BI, Overkleeft HS, Wezel GP.** Identification of glucose
414 kinase-dependent and -independent pathways for carbon control of primary metabolism,
415 development and antibiotic production in *Streptomyces coelicolor* by quantitative proteomics.
416 *Molecular Microbiology* 2012;86:1490–1507.
- 417 28. **Wezel GP van, König M, Mahr K, Nothaft H, Thomae AW, et al.** A New Piece of an Old
418 Jigsaw: Glucose Kinase Is Activated Posttranslationally in a Glucose Transport-Dependent
419 Manner in *Streptomyces coelicolor* A3(2). *J Mol Microb Biotech* 2006;12:67–74.
- 420 29. **Argimón S, Abudahab K, Goater RJE, Fedosejev A, Bhai J, et al.** Microreact:
421 visualizing and sharing data for genomic epidemiology and phylogeography. *Microb Genom*;2.
422 Epub ahead of print 2016. DOI: 10.1099/mgen.0.000093.
- 423 30. **Bertram R, Schlicht M, Mahr K, Nothaft... H.** In silico and transcriptional analysis of
424 carbohydrate uptake systems of *Streptomyces coelicolor* A3 (2). *Journal of Bacteriology*
425 2004;186:1362–1373.
- 426 31. **Schniete JK, Reumerman R, Kerr L, Tucker NP, Hunter IS, et al.** Differential
427 transcription of expanded gene families in central carbon metabolism of *Streptomyces*
428 *coelicolor* A3(2). *Access Microbiol*. Epub ahead of print 2020. DOI: 10.1099/acmi.0.000122.
- 429 32. **Chan YA, Podevels AM, Kevany BM, Thomas MG.** Biosynthesis of polyketide synthase
430 extender units. *Nat Prod Rep* 2009;26:90–114.
- 431 33. **Pfeifer BA, Khosla C.** Biosynthesis of Polyketides in Heterologous Hosts. *Microbiol Mol*
432 *Biol R* 2001;65:106–118.
- 433 34. **Zhang G, Li Y, Fang L, Pfeifer BA.** Tailoring pathway modularity in the biosynthesis of
434 erythromycin analogs heterologously engineered in *E. coli*. *Science Advances*
435 2015;1:e1500077e1500077.
- 436 35. **Arruda S, Bomfim G, Knights R, Huima-Byron T, Riley L.** Cloning of an *M. tuberculosis*
437 DNA fragment associated with entry and survival inside cells. *Science* 1993;261:1454–1457.

- 438 36. **Mohn WW, Geize R van der, Stewart GR, Okamoto S, Liu J, et al.** The Actinobacterial
439 *mce4* Locus Encodes a Steroid Transporter. *Journal of Biological Chemistry* 2008;283:35368–
440 35374.
- 441 37. **Clark LC, Seipke RF, Prieto P, Willemsse J, Wezel GP van, et al.** Mammalian cell entry
442 genes in *Streptomyces* may provide clues to the evolution of bacterial virulence. *Scientific*
443 *Reports* 2013;3.
- 444 38. **Aziz RK, Bartels D, Best AA, DeJongh M, Disz T, et al.** The RAST Server: Rapid
445 Annotations using Subsystems Technology. *Bmc Genomics* 2008;9:75.
- 446 39. **Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ.** Basic local alignment search
447 tool. *J Mol Biol* 1990;215:403–410.
448

449 **Figures**

450 **Fig. 1. Schematic workflow for the creation of the ActDES database.** Genomes were
451 selected from NCBI Taxonomy browser and uploaded for annotation to RAST[38]. The
452 annotated genomes were the processed for two different analyses. Firstly, the functional roles
453 were downloaded and for each functional role the numbers of occurrences per genome were
454 counted in order to obtain an expansion table (**Supp. Table S2**) by comparing the mean of
455 each genus to the overall mean of all genera. Secondly, the genomes were used to extract all
456 nucleotide and protein sequences in FASTA format which could then be queried by sequence
457 using BLAST [39]. The hits were aligned in MUSCLE [16] and after refinement the alignment
458 was used to construct phylogenetic trees in Quicktree [17].

459

460 **Fig. 2A)** Actinobacterial-wide phylogenetic tree of glucose permeases (GlcP) **B)**
461 Actinobacterial-wide phylogenetic tree of Glucokinases (GlcK). Trees are colour-coded
462 according to the NCBI Taxonomy browser
463 (<https://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi>). Interactive trees are also
464 available via Microreact [29]: GlcP https://microreact.org/project/VBUdiQ5_k/045c95e1 and
465 GlcK https://microreact.org/project/w_KDfn1xA/5a178533.

466

Fig. 1

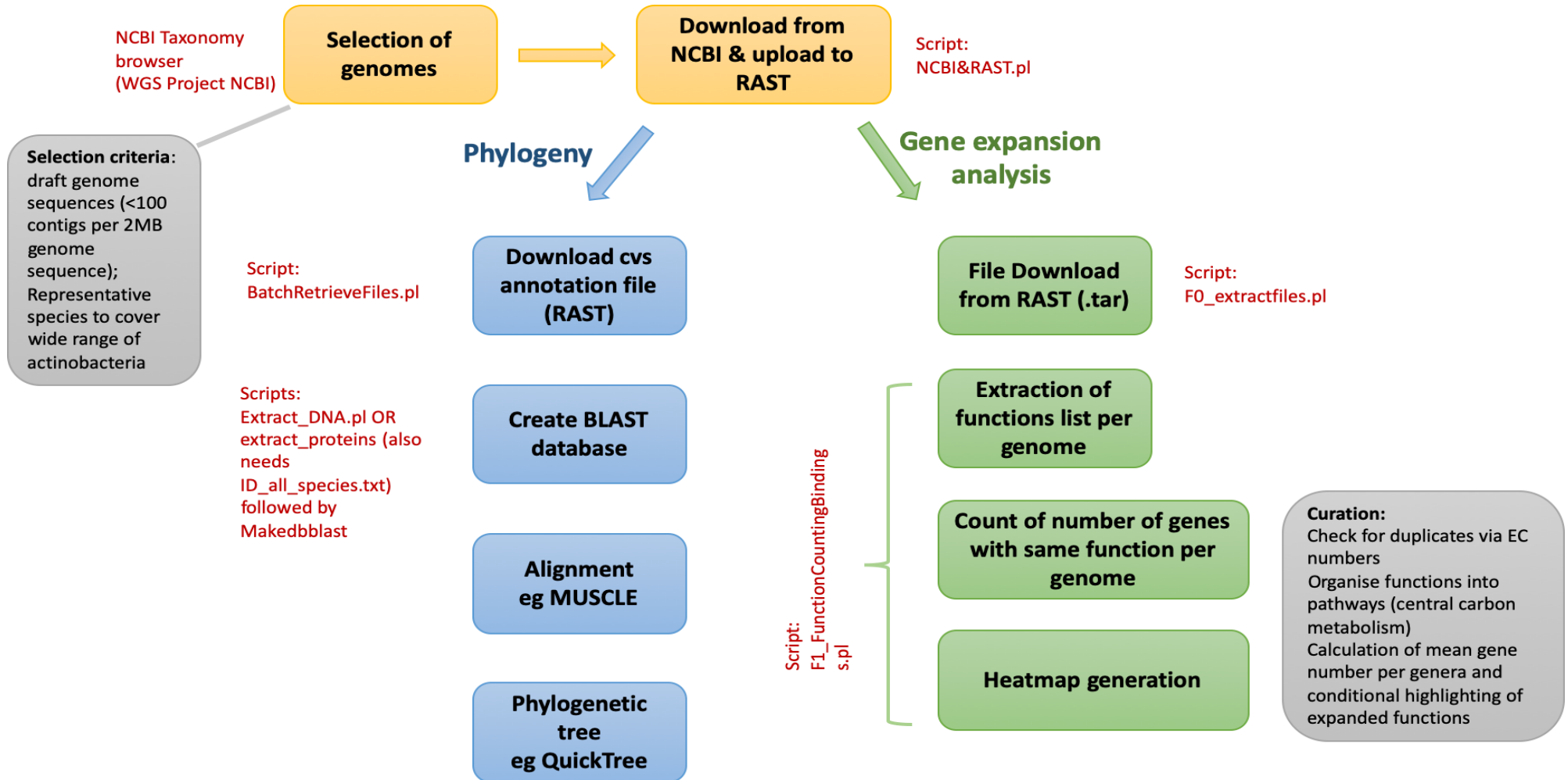
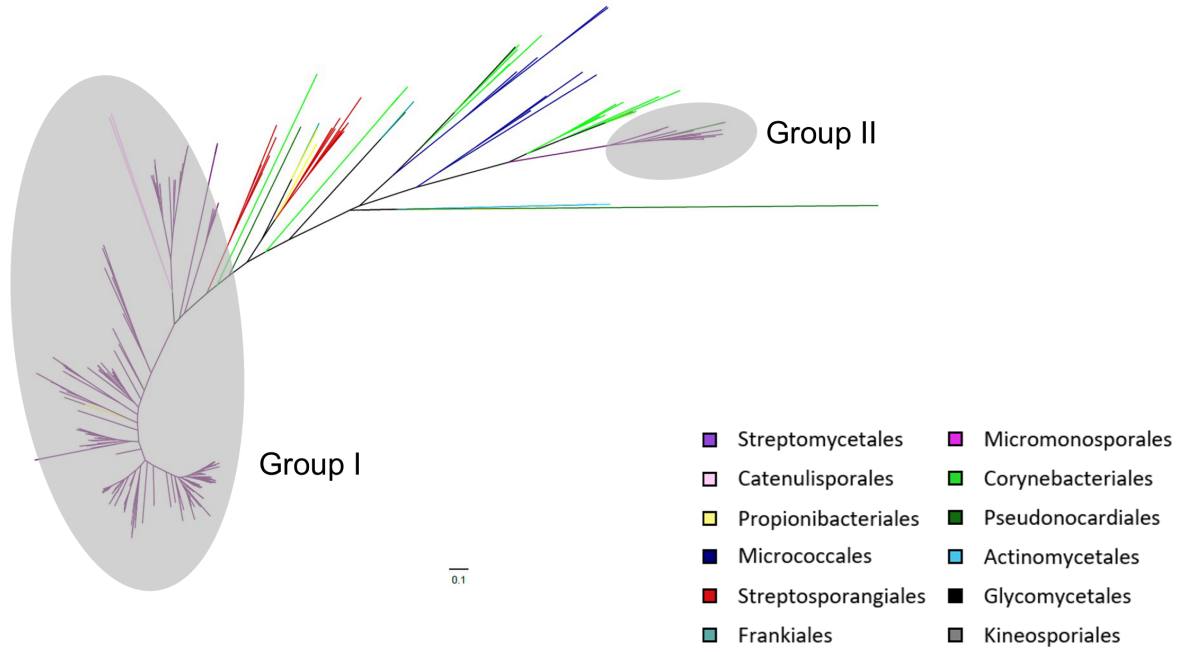


Fig. 2

A. Glucose permease (GlcP)



B. Glucose kinase (GK)

