

1 Application of a Machine Learning Approach Towards the Targeted 2 Identification of Phage Depolymerases

3 Damian J. Magill^{1*} & Timofey A. Skvortsov²

4

5 Affiliations

6 ¹Independent researcher, 14 Chemin du Jard, 37550, Saint-Avertin, France

7 ²School of Pharmacy, Queen's University Belfast, 97 Lisburn Road, Belfast

8 *Corresponding Author: damianjmagill@gmail.com

9

10 Keywords: bacteriophage, depolymerase, machine-learning

11

12 Abstract

13

14 Biofilm production plays a clinically significant role in the pathogenicity of many bacteria,
15 limiting our ability to apply antimicrobial agents and contributing in particular to the
16 pathogenesis of chronic infections. Bacteriophage depolymerases, leveraged by these
17 viruses to circumvent biofilm mediated resistance, represent a potentially powerful weapon
18 in the fight against antibiotic resistant bacteria. Such enzymes are able to degrade the
19 extracellular matrix that is integral to the formation of all biofilms and as such would allow
20 complementary therapies or disinfection procedures to be successfully applied. In this
21 manuscript, we describe the development and application of a machine learning based
22 approach towards the identification of phage depolymerases. We demonstrate that on the
23 basis of a relatively limited number of experimentally proven enzymes and using an amino
24 acid derived feature vector that the development of a powerful model with an accuracy on
25 the order of 90% is possible, showing the value of such approaches in the discovery of novel
26 therapeutic agents.

27

28

29

30

31

32

33

34 **Background**

35

36 Biofilms are the most common form of bacterial lifestyle in nature (1). Biofilm formation by
37 pathogenic bacteria allows for the establishment of a multicellular consortium of clinical
38 significance due to the role such communities play in the persistence of bacterial infection
39 and their resistance to various modes of treatment and disinfection. Indeed, such
40 assemblages confer antimicrobial resistance on multiple levels including limiting the
41 penetrability of antimicrobial compounds, the presence of metabolically inactive persister
42 cells exhibiting intrinsic resistance, and the internal structure of such communities providing
43 an optimal environment facilitating horizontal gene transfer (HGT) of resistance
44 determinants (2). A critical component for the establishment of biofilms and a significant
45 contributor to the resistant phenotype they exhibit is the production of a matrix embedding
46 the biofilm cells consisting of various polymeric compounds, including proteins, extracellular
47 DNA, and polysaccharides. The latter can be broadly categorised as lipopolysaccharides
48 (LPS), which are integral components of cell walls of Gram-negative bacteria, capsular
49 polysaccharides (CPS), loosely associated with bacterial surface, and exopolysaccharides
50 (EPS), released by bacteria into the surrounding environment (3). The ability to remove such
51 polymeric barriers in order to expose the underlying community of cells is a desirable one
52 from the practical point of view, be it for the purposes of surface disinfection, de-fouling, or
53 to improve the biocidal effects of antibiotic treatment.

54 Barrier properties of bacterial biofilms also pose a problem for bacterial viruses
55 (bacteriophages) whose diffusion and ability to infect host cells is reduced within biofilms
56 (4). Targeted degradation of biofilm polysaccharides is a feature of many bacteriophages
57 (phages) which increases the probability of successful infection; this is the result of
58 enzymatic activity of a class of phage-encoded enzymes called depolymerases (DP). The
59 majority of DPs are phage-associated enzymes and belong to lyase and hydrolase classes,
60 with the former constituting a large majority of the well characterised and experimentally
61 validated DPs (5 - 7). Given the global antibiotic crisis we now face, there is a resurgence of
62 interest in both phage and phage-derived therapeutic agents as alternatives. Several
63 recently published reviews describe the structural and functional characteristics of phage
64 DPs and outline their potential applications as biotechnological tools and therapeutic agents
65 (8 - 11).

66 The therapeutic potential of phage DPs was recognised more than 60 years ago (12). Phage
67 DPs are of particular interest due to their potential use in combinatorial therapies with
68 antibiotics or other antimicrobial agents and in the removal of biofilms from medical devices
69 most notably catheters (13; 14). Moreover, as the depolymerases do not kill bacteria, it is
70 posited that they could be employed on their own as anti-virulence agents, decreasing
71 bacterial fitness and facilitating the clearance of the bacteria by the human immune system
72 (10). Therefore, any approach that enhances our ability to identify novel DPs is of great
73 value, especially since it is not always trivial to attribute depolymerase activity to a specific
74 gene. As the polysaccharides produced by even closely related bacterial species may have
75 subtle but significant structural differences, phage DPs acting on them also demonstrate

76 high variability, to the point that the depolymerase domains will sometimes be among the
77 only genomic DNA fragments showing no conservation between phages of the same species
78 (14). Although the majority of known DPs are parts of phage receptor-binding proteins
79 (RBPs) such as tail spikes and thus have conserved N-terminal domains responsible for virion
80 attachment, some depolymerases can be encoded as truncated RBPs (presumably acting as
81 diffusible DPs), further complicating their reliable prediction (15).

82 Machine learning based approaches are proving to be an extremely valuable avenue in all
83 realms of science and this is no less true of phage biology whereby success has been
84 demonstrated through the application of such techniques towards the identification of
85 phage structural proteins (16), host-phage pairs (17), RBPs (18) and lifecycle (19) amongst
86 others (20). Recently published papers expand this list to include endolysins (21) and
87 depolymerases (22). Nevertheless, the ultimate success or failure of machine learning
88 algorithms depends on many factors, including but not limited to the size and composition
89 of training sets, the algorithm used for the problem at hand, and the careful construction of
90 a vector capturing adequately discriminant features (23). Therefore, more ML solutions are
91 needed to expand the computational phage characterisation toolkit and allow for a series of
92 complementary approaches to be available.

93 In this manuscript we describe the development and application of a machine learning
94 approach towards the identification of phage DPs, highlighting that such models should
95 form an integral part of our toolkit enabling the discovery of novel therapeutics. We
96 demonstrate that even a relatively small training set is sufficient to produce a highly
97 generalizable machine learning model capable of accurately predicting DPs in a multitude of
98 phages infecting vastly different bacteria. Indeed, an accuracy of 90% was attained on the
99 test data set and a similar result for genome context predictions that detected the DP within
100 the top 10 predictions.

101

102

103

104

105

106

107

108

109

110

111

112

113 **Methods**

114

115 **Data Set Preparation**

116

117 In order to establish a database of DP sequences that would ultimately fuel our model, we
118 focussed our attention on publications within which depolymerase activity had been
119 experimentally demonstrated. A comprehensive literature search was conducted and a
120 database established consisting of 50 depolymerase sequences. Table S1 presents an
121 overview of this sequence database including the phages and references from which they
122 were found. 28 of the sequences exclusively state CPS as an enzymatic target, 20 target EPS,
123 and the remaining two target LPS and a combination of targets. The vast majority of
124 sequences were *Podoviridae* derived and the database concerned phages infecting Gram
125 negative bacteria. The size range of sequences varied from 150 amino acids to 1267 amino
126 acids in length.

127 To complete this dataset, we required 50 sequences that would serve as the negative non-
128 depolymerase set and thus provide a 1:1 positive to negative sequence set. To do this, we
129 randomly extracted 50 sequences from a soil metagenome (SRR15048733) that were
130 sampled across the size distribution of sequences so as to avoid the introduction of
131 sequence size biases. BLAST searches were conducted with these sequences against the
132 positive depolymerases to ensure the absence of homology followed by HHPred analysis to
133 confirm the absence of domains known to be associated with depolymerase activity.

134 To highlight the dissimilarity in the dataset, we calculated pairwise similarity scores across
135 the entire dataset and represented this as a heatmap (figure 1.).

136

137 **Feature Extraction and Selection**

138

139 A diverse range of features were generated which were derived solely from the amino acid
140 sequences. Eleven of these features were directly calculated using the ProteinAnalysis
141 feature from the BioPython (version 1.73) ProtParam module (24). These were the MW,
142 aromaticity, predicted instability and isoelectric point, GRAVY score, predicted secondary
143 structure (sequence proportion engaged in helices, strands, and turns), extinction
144 coefficients (ox/red), and a combined flexibility score. Beyond this the relative abundance of
145 each amino acid and the total sequence length were also taken into account. As a final set of
146 features, we considered dipeptides and tripeptides as a function of conserved
147 physicochemical properties. Seven groups were established consisting of amino acids with a
148 hydrocarbon R group, those with an uncharged aromatic side chain, sulphur containing,
149 positively charged, negatively charged, polar uncharged, and proline. According to this
150 schema, the dipeptides AE and LD were considered as both belonging to group 15. Whilst
151 allowing us to incorporate dipeptide and tripeptide properties into the model, this also

152 reduced the overall feature set compared to using all possible combinations of amino acids.
153 This was carried out using in-house scripts.

154

155 **Model Selection, Training, and Evaluation**

156

157 With respect to the appropriate choice of machine learning algorithm, we decided to test
158 both support vector machine (SVM) and random forest (RF) approaches (25; 26). This was
159 due to the fact that our data set constituted a small number of samples exhibiting a high
160 feature space. In both cases, we leveraged a grid search in order to assess the
161 hyperparameter space and find the best model configuration for both algorithms. This was
162 conducted using the scikit-learn library (version 0.23.2) (27). We opted for a 5-fold cross
163 validation using an 80/20 split of the dataset.

164 To evaluate model performance, we particularly focussed on the overall accuracy and recall
165 on the cross-validations defined as follows:

166

$$167 \quad Accuracy = \frac{TP+TN}{TP+FP+TN+FN} \quad (1)$$

168

$$169 \quad Recall = \frac{TP}{TP+FN} \quad (2)$$

170

171

172 Where TP, FP, TN, and FN correspond to true positive, false positive, true negative, and false
173 negative respectively regarding the classification performed on the test data. All scores
174 reported are the average obtained following the cross-validation.

175 With respect to the hyperparameters tuned, for SVM both linear and RBF kernels were
176 evaluated along with cost and gamma functions when applicable. For RF, differing numbers
177 of estimators were evaluated using a step size of 100 along with total tree depth, and the
178 minimum samples supporting a branch and split of the tree. In addition to this, we also
179 integrated a two-degree polynomial feature transformation, min/max scaling, and applied
180 entropy-based impurity.

181 Once optimal parameters were determined for the model following evaluation, the final
182 version was created incorporating the entirety of the training set.

183

184

185

186 **Software Package Depolymerase Predict**

187

188 Both the source code and a standalone ready to use version of the application are available
189 as detailed in the “Availability of data and materials” section. A simple user-friendly GUI has
190 been developed through which users can step-by-step upload their sequences, generate the
191 feature vector, and carry out predictions and view the output.

192

193 **Results and Discussion**

194

195 **Feature Generation**

196

197 The application of the feature generation script was carried out on the 100 amino acid
198 sequence input data set. This resulted in the construction of a feature vector with 424
199 descriptors for each of the sequences. An additional column was added to distinguish the
200 depolymerases from the negative cases. This entire training set is presented in Table S2. and
201 can be used directly in the reproduction of our analysis with the parameters outlined below.

202

203 **Model Evaluation and Final Selection**

204

205 The SVM approach was initially applied to the dataset with no hyperparameter tuning, with
206 the application of a linear kernel. This resulted in a model exhibiting an overall accuracy
207 score of 0.70 across all folds. As presented in the normalised confusion matrix in figure 2.
208 this model performed extremely poorly with respect to true and false positives but handled
209 negative cases well. Indeed, hyperparameter tuning did nothing to resolve this problem. The
210 overall accuracy remained unchanged, but the model improved in its ability to correctly
211 identify non-depolymerase sequences with 100% success rate. This was at the cost of
212 decreased performance on positive cases with only 45% of true depolymerases being
213 correctly identified as such.

214 Subsequent application of an RF approach yielded more promising results. This is an
215 ensemble machine learning method that leverages multiple decision trees in order to
216 reduce variance and provide better model generalization. It performs especially well with
217 small sample sets and large feature spaces and so it was expected it may be the best
218 approach to this problem. Application of a tuned RF model indeed showed a much higher
219 level of performance (figure 2). An overall accuracy score of 0.90 was obtained across all
220 folds with similar performance observed with respect to the correct classification of positive
221 and negative cases. It was found that for this case, the following parameters provided
222 optimal performance of the model: use of 1500 estimators with automatic definition of

223 maximum features to be used by each tree. A maximum depth of 30 was applied with a
224 minimum sample support of 3 required for each leaf. Each tree split was evaluated using the
225 entropy-based criterion. The pipeline also integrated a two-degree polynomial feature
226 transformation along with application of min/max scaling.

227

228 **Application of Model Towards Depolymerase Identification**

229

230 In order to further assess the performance of our model, we decided to apply it within the
231 context of whole phage genomes to see whether it could correctly identify depolymerases
232 amongst the other genes. Due to the fact that our model leverages experimentally active
233 depolymerases and we have thus exhausted this option, we were limited to performing this
234 test on computationally predicted enzymes. The first such case was *Pseudomonas* phage
235 pf16; a phage previously characterised by our group (28). Depolymerase activity was
236 previously observed in this phage and extensive computational analysis identified gp215 as
237 the likely candidate with probable pectate lyase activity. We proceeded to analyse pf16
238 gene products using our model and ranked the probabilities accordingly. These results are
239 presented in figure 3. We immediately observed that the predicted depolymerase was
240 ranked 4th by our model in the context of the whole genome. This in itself is a reasonable
241 result however, further analysis of the higher ranked candidates revealed that they possess
242 domains not unrelated to what is observed in depolymerases including endosialidase and
243 VrlC-like domains, the latter speculated to have sialidase activity (29).

244 We further tested the performance of our model in the context of whole genomes by
245 directing our attention towards computationally predicted depolymerases described by
246 Pires *et al.* (11). This provided a good opportunity to test the generalizability of our model as
247 the sequences described in this paper exhibit significant diversity in terms of the domains
248 present and nature of the hosts infected by the phages. We downloaded the genomes of
249 the associated phages, removing some for which the records no longer exist. This resulted in
250 155 genomes on which we applied our model. Predictions were performed, the probabilities
251 ranked, and the position of the putative depolymerase identified. Table S3 presents all of
252 the genomes, the depolymerases and the associated ranking provided by our model. Across
253 all sequences the depolymerase featured as the first prediction 40.6% of the time. This
254 increased to 69.7% and 78.1% for top 3 and top 5 predictions respectively. When
255 considering top 10 and top 20 this grows to a large majority with 87.1% and 94.8%. Most
256 poorly predicted sequences were those containing domains that did not feature in our
257 model, especially DUF867. When we look closer at the distribution of these results we
258 observe a good level of model generalizability in a number of aspects (Figure 4). Despite
259 being fuelled by depolymerases in phages infecting Gram-negative bacteria, the model
260 performs equally well for phages infecting both types. This fact also holds when considering
261 the family of phage and the genus of the host. This implies that the model is leveraging
262 features that are common to a large majority of known depolymerase enzymes.

263

264 **Conclusion**

265 Bacteriophage depolymerases offer a host of promising clinical and biotechnological
266 applications, including the synergistic treatment of infections via biofilm removal. There is
267 however, a need for rapid and accurate identification of such enzymes. In this work we have
268 described the development and application of a machine learning approach that allows for
269 depolymerase prediction with an overall accuracy of 90% using a sequence-derived feature
270 vector. We demonstrated that this model was generalizable to depolymerases from a
271 variety of phages, robustly predicting them in the context of the genomes across several
272 hosts and enzyme classes. This highlights the power that such approaches can offer in the
273 identification of industrially and/or clinically useful enzymes.

274

275 **Declarations**

276

277 **Ethics approval and consent to participate**

278 Not applicable

279

280 **Consent for publication**

281 Not applicable

282

283 **Availability of data and materials**

284

285 The source code for the application can be found via the following URL:

286 <https://github.com/DamianJM/Depolymerase-Predict.git>

287 In addition, a standalone version of the application is available with all dependencies and
288 training set compiled within at the following address:

289 <https://sourceforge.net/projects/depolymerase-predict>

290 The training dataset has been provided as part of the supplementary data which allows for
291 our work to be reproduced. Depolymerases used in the development of the model are
292 detailed in Supplementary table 1, with all accession numbers and associated literature
293 references provided.

294

295 **Competing interests**

296 Not applicable

297

298 **Funding**

299 The authors declare that they received no specific funding for this work

300

301 **Authors' contributions**

302 DM developed the software; DM and TS analysed the data; DM and TS wrote the
303 manuscript and approved the final version along with figures.

304

305 **Acknowledgements**

306 Not applicable

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328 **Figure Legends**

329

330 **Figure 1. Heatmap of pairwise similarity scores calculated for the training dataset**

331 Grayscale colours correspond to percentage identity as provided in the associated legend. The
332 negative and positive components of the dataset are highlighted with braces and associated labels.
333 As highlighted by the scale of the legend, the global identities of the matrix are rather low, showing a
334 high level of dissimilarity between the sequences.

335

336 **Figure 2. Normalised confusion matrices summarising model performance on test data**

337 Matrices give the proportion of depolymerase (DP) and non-depolymerase (Not DP) that are
338 correctly identified by the model, corresponding thus to the true/false positive and true/false
339 negative proportions. Matrices are shown for non-optimised SVM (a), optimised SVM (b), and
340 optimised RF (c) models.

341

342 **Figure 3. Top Predictions of *Pseudomonas* phage pf16 depolymerases**

343 The graph highlights that probability reported by the model of the gene product being a
344 depolymerase. Gene products are labelled accordingly. The putative depolymerase previously
345 reported is highlighted on the graph and the modelling of this protein shown with respect to a
346 known EPS depolymerase and endopolygalacturonase as reported in *Magill et al. (2017)*.

347

348 **Figure 4. Graphs showing ranking of depolymerases predicted by the model**

349 Rankings performed on depolymerase predictions from genomes described by Pires *et al. (2016)*.
350 Rankings are coloured by depolymerase domains (a), family of the phage described (b), whether the
351 host is Gram-positive or negative (c), and by the host genus (d).

352

353

354

355

356

357

358

359

360

361

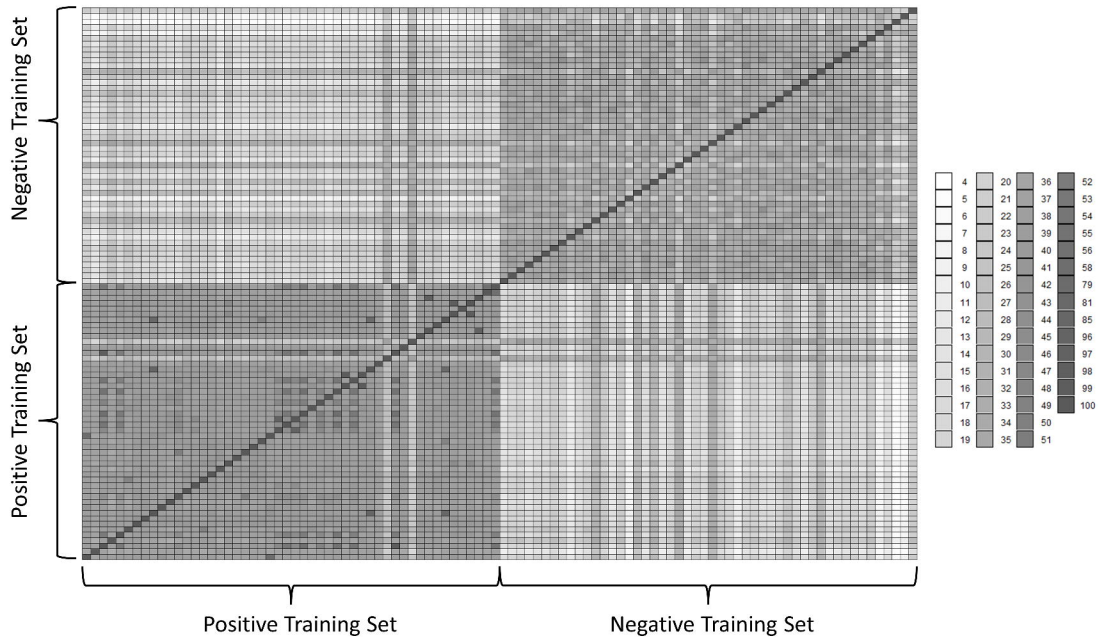
362 References

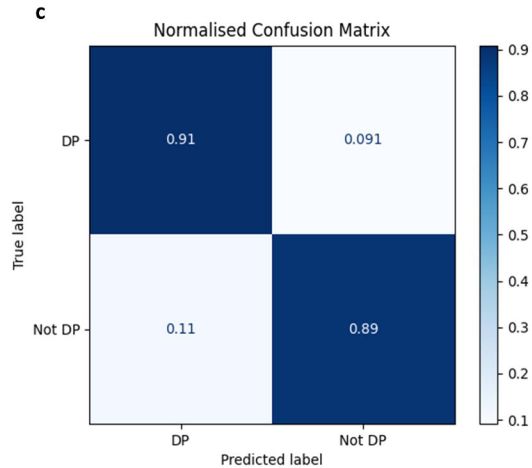
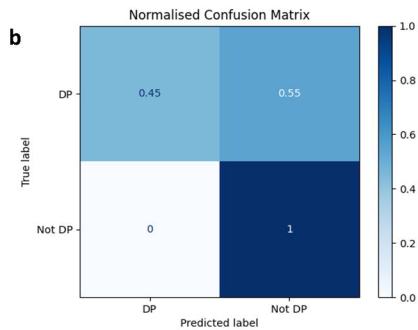
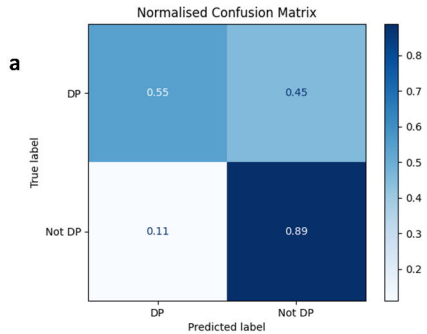
363

- 364 1. Flemming, H.C. and Wuertz, S., 2019. Bacteria and archaea on Earth and their
365 abundance in biofilms. *Nature Reviews Microbiology*, 17(4), pp.247-260.
366
- 367 2. Uruén, C., Chopo-Escuin, G., Tommassen, J., Mainar-Jaime, R.C. and Arenas, J., 2020.
368 Biofilms as promoters of bacterial antibiotic resistance and tolerance. *Antibiotics*,
369 10(1), p.3.
370
- 371 3. Mostowy, R.J. and Holt, K.E., 2018. Diversity-generating machines: genetics of
372 bacterial sugar-coating. *Trends in microbiology*, 26(12), pp.1008-1021.
373
- 374 4. Simmons, M., Drescher, K., Nadell, C.D. and Bucci, V., 2018. Phage mobility is a core
375 determinant of phage–bacteria coexistence in biofilms. *The ISME journal*, 12(2),
376 pp.531-543.
377
- 378 5. Majkowska-Skrobek, G., Łatka, A., Berisio, R., Maciejewska, B., Squeglia, F., Romano,
379 M., Lavigne, R., Struve, C. and Drulis-Kawa, Z., 2016. Capsule-targeting
380 depolymerase, derived from Klebsiella KP36 phage, as a tool for the development of
381 anti-virulent strategy. *Viruses*, 8(12), p.324.
382
- 383 6. Olszak, T., Shneider, M. M., Latka, A., Maciejewska, B., Browning, C., Sycheva, L. V.,
384 et al. (2017). The O-specific polysaccharide lyase from the phage LKA1 tailspike
385 reduces *Pseudomonas* virulence. *Sci. Rep.* 7:16302. doi: 10.1038/s41598-017-16411-
386 4
387
- 388 7. Thompson, J.E., Pourhossein, M., Waterhouse, A., Hudson, T., Goldrick, M., Derrick,
389 J.P. and Roberts, I.S., 2010. The K5 lyase KflA combines a viral tail spike structure
390 with a bacterial polysaccharide lyase mechanism. *Journal of Biological Chemistry*,
391 285(31), pp.23963-23969.
392
- 393 8. Knecht, L.E., Veljkovic, M. and Fieseler, L., 2020. Diversity and function of phage
394 encoded depolymerases. *Frontiers in Microbiology*, 10, p.2949.
395
- 396 9. Latka, A., Maciejewska, B., Majkowska-Skrobek, G., Briers, Y. and Drulis-Kawa, Z.,
397 2017. Bacteriophage-encoded virion-associated enzymes to overcome the
398 carbohydrate barriers during the infection process. *Applied microbiology and*
399 *biotechnology*, 101(8), pp.3103-3119.
400
- 401 10. Oliveira, H., Drulis-Kawa, Z. and Azeredo, J., 2022. Exploiting phage-derived
402 carbohydrate depolymerases for combating infectious diseases. *Trends in*
403 *Microbiology*.
404

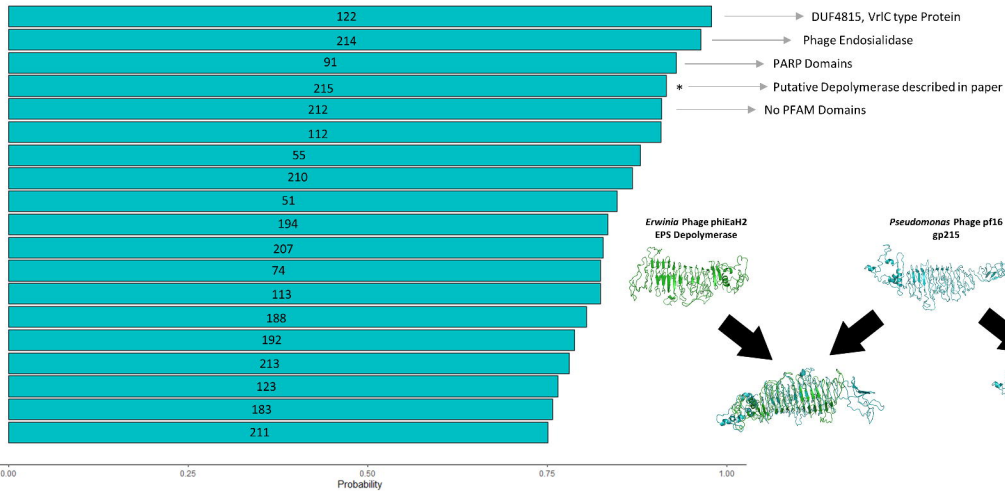
- 405 11. Pires, D.P., Oliveira, H., Melo, L.D., Sillankorva, S. and Azeredo, J., 2016.
406 Bacteriophage-encoded depolymerases: their diversity and biotechnological
407 applications. *Applied microbiology and biotechnology*, 100(5), pp.2141-2151.
408
- 409 12. Adams, M.H. and Park, B.H., 1956. An enzyme produced by a phage-host cell
410 system: II. The properties of the polysaccharide depolymerase. *Virology*, 2(6),
411 pp.719-736.
412
- 413 13. Shahed-Al-Mahmud, M., Roy, R., Sugiokto, F.G., Islam, M.N., Lin, M.D., Lin, L.C. and
414 Lin, N.T., 2021. Phage ϕ AB6-borne depolymerase combats *Acinetobacter baumannii*
415 biofilm formation and infection. *Antibiotics*, 10(3), p.279.
416
- 417 14. Rice, C.J., Kelly, S.A., O'Brien, S.C., Melaugh, E.M., Ganacias, J.C., Chai, Z.H., Gilmore,
418 B.F. and Skvortsov, T., 2021. Novel phage-derived depolymerase with activity against
419 *Proteus mirabilis* biofilms. *Microorganisms*, 9(10), p.2172.
420
- 421 15. Latka, A., Leiman, P.G., Drulis-Kawa, Z. and Briers, Y., 2019. Modeling the
422 architecture of depolymerase-containing receptor binding proteins in *Klebsiella*
423 phages. *Frontiers in microbiology*, 10, p.2649.
424
- 425 16. Cantu, V.A., Salamon, P., Seguritan, V., Redfield, J., Salamon, D., Edwards, R.A. and
426 Segall, A.M., 2020. PhANNs, a fast and accurate tool and web server to classify phage
427 structural proteins. *PLoS computational biology*, 16(11), p.e1007845.
428
- 429 17. Boeckaerts, D., Stock, M., Criel, B., Gerstmans, H., De Baets, B. and Briers, Y., 2021.
430 Predicting bacteriophage hosts based on sequences of annotated receptor-binding
431 proteins. *Scientific reports*, 11(1), pp.1-14.
432
- 433 18. Boeckaerts, D., Stock, M., De Baets, B. and Briers, Y., 2022. Identification of Phage
434 Receptor-Binding Protein Sequences with Hidden Markov Models and an Extreme
435 Gradient Boosting Classifier. *Viruses*, 14(6), p.1329.
436
- 437 19. Hockenberry, A.J. and Wilke, C.O., 2021. BACPHLIP: predicting bacteriophage
438 lifestyle from conserved protein domains. *PeerJ*, 9, p.e11396.
439
- 440 20. Nami, Y., Imeni, N. and Panahi, B., 2021. Application of machine learning in
441 bacteriophage research. *BMC microbiology*, 21(1), pp.1-8.
442
- 443 21. Criel, B., Taelman, S., Van Criekeing, W., Stock, M. and Briers, Y., 2021. PhaLP: A
444 Database for the Study of Phage Lytic Proteins and Their Evolution. *Viruses*, 13(7),
445 p.1240.
446
- 447 22. Duarte, J.A.G., 2021. PhageDPO: phage depolymerase finder (Doctoral dissertation).
448

- 449 23. Greener, J.G., Kandathil, S.M., Moffat, L. and Jones, D.T., 2022. A guide to machine
450 learning for biologists. *Nature Reviews Molecular Cell Biology*, 23(1), pp.40-55.
451
- 452 24. Cock, P.J., Antao, T., Chang, J.T., Chapman, B.A., Cox, C.J., Dalke, A., Friedberg, I.,
453 Hamelryck, T., Kauff, F., Wilczynski, B. and De Hoon, M.J., 2009. Biopython: freely
454 available Python tools for computational molecular biology and bioinformatics.
455 *Bioinformatics*, 25(11), pp.1422-1423.
456
- 457 25. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel,
458 M., Prettenhofer, P., Weiss, R., Dubourg, V. and Vanderplas, J., 2011. Scikit-learn:
459 Machine learning in Python. *the Journal of machine Learning research*, 12, pp.2825-
460 2830.
461
- 462 26. Breiman, L., 2001. Random forests. *Machine learning*, 45(1), pp.5-32.
463
- 464 27. Boser, B.E., Guyon, I.M. and Vapnik, V.N., 1992, July. A training algorithm for optimal
465 margin classifiers. In *Proceedings of the fifth annual workshop on Computational*
466 *learning theory* (pp. 144-152).
467
- 468 28. Magill, D.J., Krylov, V.N., Shaburova, O.V., McGrath, J.W., Allen, C.C., Quinn, J.P. and
469 Kulakov, L.A., 2017. Pf16 and phiPMW: Expanding the realm of *Pseudomonas putida*
470 bacteriophages. *PLoS One*, 12(9), p.e0184307.
471
- 472 29. Billington SJ, et al. Complete nucleotide sequence of the 27-kilobase virulence
473 related locus (vrl) of *Dichelobacter nodosus*: evidence for extrachromosomal origin.
474 *Infect Immun.* 1999;67:1277–1286.
475
476
477
478
479





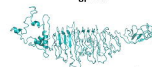
Top pf16 Depolymerase Predictions



Erwinia Phage phiEaH2
EPS Depolymerase



Pseudomonas Phage pf16
gp215



Fusarium moniliforme
Endopolygalacturonase

