1    Alignment-free identification of COI DNA barcode data with the Python package Alfie

2

3    Cameron M. Nugent[1,2,*], Sarah J. Adamowicz[1]

4

5    [1] Department of Integrative Biology, University of Guelph. Guelph, Ontario, Canada

6    [2] Biodiversity Institute of Ontario, University of Guelph. Guelph, Ontario, Canada

7    [*] Corresponding author: nugentc@uoguelph.ca

8

9    **Abstract**

10

11    Characterization of biodiversity from environmental DNA samples and bulk metabarcoding data

12    is hampered by off-target sequences that can confound conclusions about a taxonomic group of

13    interest. Existing methods for isolation of target sequences rely on alignment to existing

14    reference barcodes, but this can bias results against novel genetic variants. Effectively parsing

15    targeted DNA barcode data from off-target noise improves the quality of biodiversity estimates

16    and biological conclusions by limiting subsequent analyses to a relevant subset of available data.

17    Here, we present Alfie, a Python package for the alignment-free classification of cytochrome c

18    oxidase subunit I (COI) DNA barcode sequences to taxonomic kingdoms. The package

19    determines $k$-mer frequencies of DNA sequences, and the frequencies serve as input for a neural

20    network classifier that was trained and tested using ~58,000 publicly available COI sequences.

21    The classifier was designed and optimized through a series of tests that allowed for the optimal

22    set of DNA $k$-mer features and optimal machine learning algorithm to be selected. The neural

23    network classifier rapidly assigns COI sequences to kingdoms with greater than 99% accuracy

24    and is shown to generalize effectively and make accurate predictions about data from previously

25    unseen taxonomic classes. The package contains an application programming interface that

26    allows the Alfie package's functionality to be extended to different DNA sequence classification

27    tasks to suit a user's need, including classification of different genes and barcodes, and

28    classification to different taxonomic levels. Alfie is free and publicly available through GitHub

29    (https://github.com/CNuge/alfie) and the Python package index (https://pypi.org/project/alfie/).

30

31    **Keywords:** eDNA, environmental DNA, metabarcoding, COI, machine learning, neural

32    network, alignment-free, classification

**Introduction**

Biodiversity is declining across the globe. Millions of species face the threat of extinction, and ecosystems are being irreversibly altered due to loss of biomass and changes in species composition (Barnosky *et al.* 2011; Ceballos *et al.* 2015). To maintain the health of ecosystems and curb biodiversity loss, informed conservation and management practices are required. Achievement of conservation goals is limited by a lack of fundamental information about species composition for many of the world's ecosystems. It is therefore imperative that technological solutions are developed to enable the accurate and efficient characterization of the world's biodiversity, so that existing species can be catalogued, and informed conservation strategies can be developed to protect the planet's ecosystems.

The field of DNA barcoding offers a technological solution to the problem of taxonomically classifying organismal specimens (Hebert *et al.* 2003). Instead of relying on laborious and error-prone phenotypic classifications, sequence diversity within standardized gene regions is used to enable both specimen identification and species discovery (Hebert *et al.* 2003; Ratnasingham & Hebert 2007; Hubert & Hanner 2015). The field has advanced from the barcoding of single specimens to the bulk analysis of samples, known as metabarcoding (Hajibabaei *et al.* 2011, 2016; Taberlet *et al.* 2012; Cristescu 2014), as well as multi-marker (Stefanni *et al.* 2018) and metagenomics approaches (Cuvelier *et al*. 2010). These methods have been applied in environmental biomonitoring, where multiple species are identified at once through the collection of environmental DNA (eDNA) (Taberlet *et al.* 2012). Despite the widespread adoption of these techniques, a fundamental problem persists: the accurate and repeatable characterization of biodiversity from eDNA and bulk-sample metabarcoding data is difficult, and conclusions drawn from analyses are strongly affected by methodological decisions (Clare *et al.* 2016; Braukmann *et al.* 2019).

Environmental biomonitoring often aims to answer ecological questions through the targeted examination of a taxonomic group of interest. DNA barcodes from a group of focus are targeted using group-specific PCR primers for one or more selected marker genes in the PCR amplification step that precedes high-throughput sequencing (Braukmann *et al.* 2019; Wilson *et al.* 2019). Some commonly used primers are overly general, which results in the amplification of non-target barcodes, introducing noise into data and confounding efforts to characterize true species composition for targeted taxonomic groups (Brandon-Mong *et al.* 2015; Zinger *et al.* 2019). Additionally, intra-group PCR bias can further confound the characterization of biodiversity. The over representation of certain taxa within the target group can result in other taxa being overlooked due to poorer amplification and sequencing coverage (Elbrecht & Leese 2015).

Shotgun sequencing of eDNA overcomes the primer issues of eDNA metabarcoding but also produces substantial sequencing noise and sequences from non-standardized genomic regions (Stat *et al.* 2017; Wilson *et al.* 2019). A trade off therefore exists; shotgun sequencing overcomes the amplification bias associated with PCR, but the majority of shotgun sequencing outputs cannot be assigned even high-level taxonomic classifications with confidence (Stat *et al.* 2017; Singer *et al.* 2020). Despite present technical limitations, eDNA shotgun sequencing and other next-generation biomonitoring techniques are seeing increased adoption thanks to their potential to characterize biodiversity more broadly (Makiola *et al.* 2020). Within this next generation of biomonitoring methodologies, tools leveraging machine-learning algorithms and

2

78  available data will be essential to overcoming the limitations associated with existing methods
79  (Cordier *et al.* 2019).
80      The detection of the presence and abundance of species from a specific group is
81  hampered by off-target barcodes that are amplified and sequenced in metabarcode analysis. The
82  failure to parse target sequences effectively from off-target noise can result in erroneously
83  inflated estimates of biodiversity (Bengtsson *et al.* 2011). Currently, the characterization of
84  biodiversity via metabarcode samples is primarily dependent on the alignment of sequences
85  against a pre-defined set of reference barcodes or comparison of sequences against taxon-specific
86  models (Altschul *et al.* 1990; Wang *et al.* 2007; Bengtsson *et al.* 2011; Bengtsson-Palme *et al.*
87  2015). These processes limit comparison to previously characterized barcode sequences,
88  potentially exhibiting bias against novel genetic variants. The methods are also computationally
89  intensive, often requiring each novel variant to be compared to each reference entry. These
90  methods would therefore be improved through the incorporation of an alignment-free pre-
91  filtering step that allowed for target sequences to be rapidly and accurately isolated from the
92  whole set of metabarcode output sequences using algorithms with lower computational
93  complexity (Zielezinski *et al.* 2017). This would reduce the number of spurious barcodes and
94  improve inflated biodiversity estimates. Additionally, the speed of analyses would be improved
95  by limiting subsequent alignment-based analyses to the isolated target sequences.
96      Alignment-free methods have been widely applied in biological sequence annotation and
97  classification problems (Zielezinski *et al.* 2017). Alignment-free comparison is defined as any
98  method of quantifying sequence similarity that does not produce an alignment; these methods are
99  generally less computationally intensive and can be as effective as conventional alignments
100  (Bonham-Carter *et al.* 2014; Zielezinski *et al.* 2017). To compare sequences without alignment,
101  features must be extracted from sequences in order to characterize their structure. One common
102  set of alignment-free features is $k$-mer counts, where the number of occurrences of fixed length
103  DNA words of length $k$ are quantified (Crusoe *et al.* 2015). These features can be used as inputs
104  for machine learning models trained to predict classifications such as the taxonomic designation
105  associated with sequences (Solis-Reyes *et al.* 2018). Machine learning models that operate on $k$-
106  mer input features have previously been applied in DNA barcode sequence classification and
107  other predictive tasks (Kuksa & Pavlovic 2009; Langenkämper *et al.* 2014; Ainsworth *et al.*
108  2016; Cordier *et al.* 2017). The application of these tools is often limited to specific taxonomic
109  classification tasks (Kuksa & Pavlovic 2009), or they rely on user-provided sets of sequence data
110  for model training (Langenkämper *et al.* 2014).
111      The goals of this study were to: (1) develop a high-level alignment-free taxonomic
112  classification tool for metabarcoding and environmental DNA marker gene data. This tool was
113  initially designed for the kingdom-level classification of barcode sequences from the most
114  common animal barcode, a region of the mitochondrial cytochrome c oxidase subunit I (COI)
115  gene. (2) To achieve this, we explore different feature sets ($k$-mer sizes) and machine learning
116  algorithms to determine the optimal machine learning architecture for alignment-free barcode
117  classification. (3) To make the tool accessible to other researchers, we develop a Python package
118  and command line interface to allow the alignment-free classifier to be easily deployed in future
119  research applications. (4) Within the Python package, we also develop an application
120  programming interface (API) to facilitate the construction of customized alignment-free
121  classifiers for any barcode, gene, or taxonomic group of interest. Addressing these goals led to
122  the creation of the Python package Alfie, which contains a kingdom-level alignment-free DNA
123  barcode classifier, as well as an API to aid users in custom alignment-free classifier construction.

3

124 Alfie is free and publicly available through GitHub (https://github.com/CNuge/alfie) and the
125 Python package index (https://pypi.org/project/alfie/).
126
127 **Methods**
128
129       **Data acquisition**
130
131 The Barcode Of Life Data system (BOLD) (Ratnasingham & Hebert 2007) was queried to obtain
132 all publicly available sequences for the DNA barcode: cytochrome c oxidase subunit I (COI)
133 (https://github.com/CNuge/data-alfie). Sequences were filtered to ensure a minimum length of
134 300 base pairs (bp). The five kingdom-level classifications used by the BOLD database (Animal,
135 Bacteria and Archaea, Fungi, Plant, Protist) were maintained and utilized as the labels in
136 subsequent classifier development. As a result of BOLD's mandate to catalogue animal
137 biodiversity, the database displays a significant sampling bias towards the animal kingdom. To
138 ensure that models could be trained effectively and not be biased towards animal classification,
139 down sampling of the animal data was performed to ensure more even representation of
140 sequences among kingdoms. Stratified sampling of animal sequences was performed to obtain a
141 representative subsample of 0.2% of the total set of sequences available (sequences were
142 sampled proportionally on the taxonomic level: class; a sample size of 0.2% was chosen as this
143 yielded a set of animal sequences roughly equal to the kingdom with the second highest number
144 of available COI barcodes, plants) (Table 1). To train models robust to variable data quality and
145 barcode sequence coverage, each individual barcode sequence was randomly subsampled, with a
146 200-600 base pair subsection of the complete barcode being retained at random and subsequently
147 utilized in model training and testing.
148       Prior to splitting the data into a train and test set, a validation set was created to provide a
149 stringent test of the final models' ability to make external predictions. From each kingdom, a
150 complete taxonomic class was withheld to create the validation set and simulate rare or
151 previously unseen sequences. The class withheld from each kingdom was chosen manually, with
152 selection being based on the distribution of barcodes across the taxonomic classes of the given
153 kingdom. Barcode distribution was variable across kingdoms, so no suitable rule-based selection
154 method was found; classes with intermediate levels of representation within their kingdom were
155 selected. Classes with intermediate representation levels were chosen to provide good sample
156 sizes for subsequent classification tests without grossly detracting from the size of available
157 training data. For the protist kingdom, two classes were selected for inclusion in the validation
158 set due to small intra-class barcode counts. The composition of the final validation set is
159 described in Table 2. After the validation set was withheld, the remaining data were split into a
160 train and test (stratified split on level: kingdom), with 80% of data comprising the training set,
161 and the other 20% being withheld as the test set (Table 2; Supplementary File S1).
162
163       **Feature set evaluation – *k*-mer size**
164
165 Following the train-test split, different sets of alignment-free features were generated, and the
166 accuracy of kingdom-level classifications by the resulting models were tested. For barcode
167 sequences in the training set, *k*-mer frequencies were generated for values of *k* from 1 to 6.
168 *K*-mer frequencies (count of a given *k*-mer divided by the total number of *k*-mers counted in a
169 given barcode) were used as model inputs, so as to standardize the scale of input values and also

4

170 ensure the models were robust to inputs of different lengths. For each *k*-mer feature set, deep
171 neural networks with five hidden neuron layers were trained and evaluated through 5-fold cross
172 validation (neural networks implemented using the package Tensorflow Version 2.1.0, Abadi *et*
173 *al.* 2016). The choice of deep neural network-based classifiers with five hidden neuron layers
174 was based on exploratory data analysis and preliminary model construction that showed this
175 architecture to produce effective classifiers. The number of neurons in the hidden layers of the
176 neural network were adjusted according to the size of the input feature set (Table 3). The 5-fold
177 loss and accuracy metrics for the neural networks with different *k*-mer inputs were compared via
178 a one-factor analysis of variance (ANOVA) to determine if there were significant differences in
179 classification accuracy for different feature sets (*k*-mer sizes) and to select an optimal value of *k*
180 for further model testing.
181
182 **Algorithm evaluation**
183
184 After selection of the optimal *k*-mer size, a series of different machine learning models were fit
185 using the training set and optimized through a grid search of hyperparameters. Five classification
186 algorithms were utilized: *k* nearest neighbour (KNN), support vector machine (SVM), random
187 forest (RF), extreme gradient boosting (XGB), and deep neural network (DNN). All models were
188 deployed using the Python programming language (Version 3.7.4). The KNN, SVM, and RF
189 models were implemented using the package scikit-learn (Version 0.21.3, Pedregosa *et al.* 2011),
190 the XGB model was implemented using the package XGBoost (Version 0.90, Chen & Guestrin
191 2016), and the DNN was implemented using the package Tensorflow (Version 2.1.0, Abadi *et al.*
192 2016). In order to select optimal hyperparameters and optimize performance, for each algorithm
193 a grid search was performed using scikit-learn's GridSearchCV function to train a series of
194 models on the training data set using 5-fold cross validation (Supplementary File S2). Optimal
195 hyperparameters were selected based on the highest classification accuracy. For the DNN, a
196 custom grid search script was used, with 5-fold cross validation and several potential values for
197 each of the models' respective hyperparameters (Supplementary File S3).
198 Following the selection of optimal hyperparameter sets through the grid searches, a final
199 version of each model was trained using the optimal set of hyperparameters and the complete
200 training data set. Final trained models were then used to make predictions for the previously
201 withheld test and validation sets (Table 1; Table 2). Predicted classifications were compared to
202 true values to determine the model with the highest classification accuracy. A single optimal
203 alignment-free kingdom-level classifier was selected for inclusion in the Alfie package based on
204 the accuracy of predictions made on the test and validation data. Several secondary classifier
205 characteristics were also considered to ensure model reusability. Specifically, the file size of the
206 trained models and the time required to make predictions were quantified to ensure that the
207 package's memory and time requirements were not prohibitive. The Alfie package was then
208 constructed to allow for the model to be reused in external analyses.
209
210
211 **Results and Discussion**
212
213 ***K*-mer size**
214

5

215    The cross-validation accuracy scores for the different neural networks and corresponding $k$-mer
216    feature sets were compared to determine an optimal $k$-mer feature size. The results showed that
217    the accuracy of models improved with the $k$-mer feature size, with diminishing improvements
218    beyond $k = 3$ (Table 3; Figure 1). A one-factor ANOVA revealed the differences to be significant
219    ($p < 2e\text{-}16$, F statistic = 318.3, $DF_{1,2} = 5, 24$), and a subsequent Tukey's HSD test showed the
220    accuracy of both $k = 1$ and $k = 2$ to differ significantly from all larger values of $k$ but no
221    significant differences in the performance of pairwise comparisons between $k$ 3-6. A final $k$
222    value of 4 was selected for subsequent tests, due to the insignificant differences between the
223    values of $k = 3$ to $k = 6$ and the conservative choice to select a $k$-mer size one larger than the
224    apparent minimal effective feature set.
225
226    **Training and validation**
227
228    For each of the machine learning algorithms, a grid search was used to obtain an optimal
229    hyperparameter set (Supplementary File S3). Final models were trained using the complete
230    training data set and then used to make predictions for the test and validation sets (Table 1; Table
231    4). Performance on the test data (withheld barcodes from taxonomic groups otherwise
232    represented in the training data) was strong for all models, with the lowest classification
233    accuracy exceeding 98% (RF), and all other models exceeding 99.5% accuracy (Table 4). All
234    models made less accurate kingdom-level predictions on the validation data (barcodes from
235    taxonomic classes that were completely withheld during training) (Table 5). The accuracy was
236    more variable across models as well. On the validation data, the accuracy score of the RF model
237    was 0.861, and accuracy for the KNN model was 0.927, indicating poorer generalization for
238    these methods to previously unseen data. Each of the DNN, SVM, and XGB models had
239    accuracy >97% on the validation data, and the most accurate model was the DNN (0.976).
240
241    **Final model**
242
243    The DNN (operating on 4-mer input features) was selected as the final default kingdom-level
244    classification model for the Alfie package. The DNN provided the highest accuracy on the
245    validation data, as well as high accuracy on the test dataset. These results indicated that the
246    model was not likely to be over fit to the training data and that it was able to generalize
247    effectively and make predictions about data from previously unseen taxonomic classes. This
248    generalizability of the model to rare or unseen taxa is an important feature that indicates the Alfie
249    package can likely be used effectively in the analysis of under-studied environments where
250    uncharacterized biodiversity is more likely to be present. The 4-mer DNN's high accuracy on the
251    test and validation data also indicated that the features and model can effectively capture a
252    taxonomic signal despite no alignment being performed and variable input sequence length. The
253    model was robust to sequences of variable lengths that spanned various subsections of the COI
254    barcode region (variable start and stop positions in the COI barcode region, as opposed to
255    primer-standardized sub-regions). This indicates that the alignment-free classification by Alfie is
256    an effective method for processing DNA barcoding, metabarcoding (specific subsections of the
257    barcode region in a given study), and potentially even applied in analysis of metagenomics data
258    (non-standardized fragments from shotgun sequencing).
259
260

261        **Alignment-free model framework**

262

263   The design and testing of the Alfie package presented here focuses on high-level (kingdom)
264   classification for the most common animal barcode, COI. However, the Alfie package provides a
265   robust framework that a user can easily apply to produce and test alignment-free classification
266   tools for any taxonomic distinction, DNA barcode, or combination thereof (Supplementary File
267   S4). As a kingdom-level classifier, Alfie acts as an effective data filter, allowing the barcode
268   sequences from a kingdom of interest to be separated from the large amount of off-target noise
269   common in metabarcode or metagenomics data. The alignment-free methods can be reapplied to
270   further home in on taxonomic targets; for example, using publicly available data
271   (https://github.com/CNuge/data-alfie) a binary classifier can be trained and subsequently
272   deployed with Alfie to allow for any taxonomic group of interest to be separated from a complete
273   set of COI metabarcode sequences. Using other publicly available data (i.e. Pruesse *et al.* 2007;
274   Banchi *et al.* 2020), the same custom model construction and training tools in Alfie can be used
275   to construct binary or multiclass alignment-free classification tools for other DNA barcodes or
276   genes.
277        Although the Alfie package is an effective alignment-free classification framework at
278   high taxonomic levels, traditional alignments are likely more effective for lower-level
279   classification tasks (i.e. classification to genus or species level). The *k*-mer frequency method
280   used by Alfie is not likely to be effective for resolving differences between closely related
281   species with more subtle genetic differences than those seen at higher taxonomic levels.
282   Similarly, for taxonomic groups with few representatives and no closely related outgroups,
283   available training data may be scant, providing a limitation in training of DNNs or other machine
284   learning models which rely on abundant training data. The integration of alignment-based and
285   alignment-free methods for biological sequence classification has been shown to leverage the
286   strengths of the individual approaches to yield an efficient and accurate classification method
287   (Borozan *et al.* 2015).
288        A similar hybrid approach using the Alfie package for filtration of sequences and
289   subsequent alignment of sequences for a group of interest can narrow the scope of the
290   application of alignment methods and thereby improve both analysis speed and accuracy. The
291   alignment-free model construction framework of Alfie can allow for multiple models to be
292   trained with relative ease and applied in conjunction with one another to isolate barcode
293   sequences of interest from large and messy inputs such as metagenomics data. Models could be
294   trained and applied to: (a) separate sequences from key mitochondrial genes from other
295   sequences, (b) assign sequences to a barcode or gene of origin, (c) conduct kingdom-level
296   classification for different barcode genes, and (d) conduct classification at lower taxonomic
297   levels. All this could be accomplished using the same 4-mer frequency data and would allow for
298   messy inputs to be filtered and categorized. Processing of metagenomics data in this manner
299   would allow subsequent alignment effort to be more strategically targeted, improving analysis
300   speed and accuracy.

301

302

303        **Conclusions**

304

305   We have developed and tested the Python package Alfie, which extracts *k*-mer features and uses
306   a neural network to make kingdom-level classifications of COI DNA barcode fragments with

307  greater than 99% accuracy. The Alfie package can therefore be used to separate barcode data for
308  a kingdom of interest from off-target noise, narrowing the scope of subsequent analyses to only
309  relevant data. The model is robust to full-length barcodes and short sequence fragments and is
310  therefore an effective classifier for use in both barcode and metabarcode analyses. The Alfie
311  package can be incorporated into broader analyses pipelines (Elbrecht *et al*. 2018; Cordier *et al.*
312  2019) and paired with tools that conduct quality control (Callahan *et al.* 2016; Nugent *et al.*
313  2020) and taxonomic annotation (Altschul *et al.* 1990; Wang *et al.* 2007) to characterize
314  biodiversity from large and complex data sets. The default model of Alfie is limited to kingdom-
315  level classification for the most common animal barcode, COI. Researchers may expand upon
316  this narrow scope to fit custom research needs by using the training module of Alfie. This allows
317  Alfie to be applied in different taxonomic classification tasks or for the classification of data
318  from different DNA barcodes (where labelled training data are available). The generalized and
319  customized nature of the Alfie package will allow for it to adapt along with the field of
320  biodiversity genomics. As metagenomics becomes more prevalent, the Alfie package can be
321  expanded with additional default models for tasks such as the isolation of mitochondrial DNA or
322  sequences from specific mitochondrial genes from large, messy shotgun sequencing datasets.
323
324

**References**

Abadi M, Barham P, Chen J, Chen Z, Davis A, Dean J, Devin M, Ghemawat S, Irving G, Isard M, Kudlur M (2016) Tensorflow: A system for large-scale machine learning. In 12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16) 265-283.

Agarap AF (2018) Deep learning using rectified linear units (relu). arXiv preprint arXiv:1803.08375.

Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990) Basic local alignment search tool. Journal of Molecular Biology 215(3):403-10 DOI: https://doi.org/10.1016/S0022-2836(05)80360-2

Ainsworth D, Sternberg MJ, Raczy C, Butcher SA (2017) k-SLAM: accurate and ultra-fast taxonomic classification and gene identification for large metagenomic data sets. Nucleic Acids Research 45(4):1649-56. DOI: https://doi.org/10.1093/nar/gkw1248

Banchi E, Ametrano CG, Greco S, Stanković D, Muggia L, Pallavicini A (2020) PLANiTS: a curated sequence reference dataset for plant ITS DNA metabarcoding. Database. DOI: https://doi.org/10.1093/database/baz155

Barnosky AD, Matzke N, Tomiya S, Wogan GO, Swartz B, Quental TB, Marshall C, McGuire JL, Lindsey EL, Maguire KC, Mersey B (2011) Has the Earth's sixth mass extinction already arrived?. Nature, 471(7336), 51-57. DOI: https://doi.org/10.1038/nature09678.

Bengtsson J, Eriksson KM, Hartmann M, Wang Z, Shenoy BD, Grelet GA, Abarenkov K, Petri A, Rosenblad MA, Nilsson RH (2011) Metaxa: a software tool for automated detection and discrimination among ribosomal small subunit (12S/16S/18S) sequences of archaea, bacteria, eukaryotes, mitochondria, and chloroplasts in metagenomes and environmental sequencing datasets. Antonie Van Leeuwenhoek 100(3):471. DOI: https://doi.org/10.1007/s10482-011-9598-6

Bengtsson-Palme J, Hartmann M, Eriksson KM, Pal C, Thorell K, Larsson DG, Nilsson RH (2015) METAXA2: improved identification and taxonomic classification of small and large subunit rRNA in metagenomic data. Molecular Ecology Resources (6):1403-14. DOI: https://doi.org/10.1111/1755-0998.12399

Bonham-Carter O, Steele J, Bastola D (2014) Alignment-free genetic sequence comparisons: a review of recent approaches by word analysis. Briefings in Bioinformatics 15(6):890-905. DOI: https://doi.org/10.1093/bib/bbt052

Borozan I, Watt S, Ferretti V (2015) Integrating alignment-based and alignment-free sequence similarity measures for biological sequence classification. Bioinformatics 31(9):1396-404. DOI: https://doi.org/10.1093/bioinformatics/btv006

371  Braukmann TW, Ivanova NV, Prosser SW, Elbrecht V, Steinke D, Ratnasingham S, de Waard JR,
372      Sones JE, Zakharov EV, Hebert PD (2019) Metabarcoding a diverse arthropod mock
373      community. Molecular Ecology Resources 19(3):711-27. DOI:
374      https://doi.org/10.1111/1755-0998.13008
375
376  Brandon-Mong GJ, Gan HM, Sing KW, Lee PS, Lim PE, Wilson JJ (2015) DNA metabarcoding
377      of insects and allies: an evaluation of primers and pipelines. Bulletin of Entomological
378      Research 105(6):717-27. DOI: https://doi.org/10.1017/S0007485315000681
379
380  Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP (2016). DADA2:
381      high-resolution sample inference from Illumina amplicon data. Nature Methods, 13(7),
382      581. DOI: https://doi.org/10.1038/nmeth.3869
383
384  Ceballos G, Ehrlich, PR, Barnosky AD, García A, Pringle RM, Palmer TM (2015) Accelerated
385      modern human–induced species losses: Entering the sixth mass extinction. Science
386      advances, 1(5), e1400253. DOI:  https://doi.org/10.1126/sciadv.1400253
387
388  Chen T, Guestrin C (2016) Xgboost: A scalable tree boosting system. In Proceedings of the 22nd
389      acm sigkdd international conference on knowledge discovery and data mining 2016 pp.
390      785-794. DOI: https://doi.org/10.1145/2939672.2939785
391
392  Cordier T, Esling P, Lejzerowicz F, Visco J, Ouadahi A, Martins C, Cedhagen T, Pawlowski J
393      (2017) Predicting the ecological quality status of marine environments from eDNA
394      metabarcoding data using supervised machine learning. Environmental Science &
395      Technology 51(16):9118-26. DOI: https://doi.org/10.1021/acs.est.7b01518
396
397  Cordier T, Lanzén A, Apothéloz-Perret-Gentil L, Stoeck T, Pawlowski J (2019) Embracing
398      environmental genomics and machine learning for routine biomonitoring. Trends in
399      Microbiology 27(5):387-97. DOI: https://doi.org/10.1016/j.tim.2018.10.012
400
401  Clare EL, Chain FJ, Littlefair JE, Cristescu ME (2016) The effects of parameter choice on
402      defining molecular operational taxonomic units and resulting ecological analyses of
403      metabarcoding data. Genome 59(11):981-90. DOI: https://doi.org/10.1139/gen-2015-
404      0184
405
406  Cristescu ME (2014) From barcoding single individuals to metabarcoding biological
407      communities: towards an integrative approach to the study of global biodiversity. Trends
408      in Ecology & Evolution 29(10):566-71. DOI: https://doi.org/10.1016/j.tree.2014.08.001
409
410  Crusoe MR, Alameldin HF, Awad S, Boucher E, Caldwell A, Cartwright R, Charbonneau A,
411      Constantinides B, Edvenson G, Fay S, Fenton J (2015) The khmer software package:
412      enabling efficient nucleotide sequence analysis. F1000Research 4. DOI:
413      https://doi.org/10.12688/f1000research.6924.1
414
415  Cuvelier ML, Allen AE, Monier A, McCrow JP, Messié M, Tringe SG, Woyke T, Welsh RM,
416      Ishoey T, Lee JH, Binder BJ (2010) Targeted metagenomics and ecology of globally

417    important uncultured eukaryotic phytoplankton. Proceedings of the National Academy of
418    Sciences 107(33):14679-84. DOI: https://doi.org/10.1073/pnas.1001665107
419
420 Elbrecht V, Leese F (2015) Can DNA-based ecosystem assessments quantify species abundance?
421    Testing primer bias and biomass—sequence relationships with an innovative
422    metabarcoding protocol. PLoS ONE 10(7). DOI:
423    https://doi.org/10.1371/journal.pone.0130324
424
425 Elbrecht V, Vamos EE, Steinke D, Leese F (2018) Estimating intraspecific genetic diversity from
426    community DNA Metabarcoding Data. PeerJ, 6, e4644. DOI:
427    https://doi.org/10.7717/peerj.4644
428
429 Hajibabaei M, Shokralla S, Zhou X, Singer GA, Baird DJ (2011) Environmental barcoding: a
430    next-generation sequencing approach for biomonitoring applications using river benthos.
431    PLoS ONE 6(4). DOI: https://doi.org/10.1371/journal.pone.0017497
432
433 Hajibabaei M, Baird DJ, Fahner NA, Beiko R, Golding GB (2016) A new way to contemplate
434    Darwin's tangled bank: how DNA barcodes are reconnecting biodiversity science and
435    biomonitoring. Philosophical Transactions of the Royal Society B: Biological Sciences
436    371(1702):20150330. DOI: https://doi.org/10.1098/rstb.2015.0330
437
438 Hebert PDN, Cywinska A, Ball SL, Dewaard JR (2003) Biological identifications through DNA
439    barcodes. Proceedings of the Royal Society of London. Series B: Biological Sciences
440    270(1512):313-21. DOI: https://doi.org/10.1098/rspb.2002.2218
441
442 Hubert N, Hanner R (2015) DNA barcoding, species delineation and taxonomy: a historical
443    perspective. DNA Barcodes 3(1):44-58. DOI: https://doi.org/10.1515/dna-2015-0006
444
445 Kingma DP, Ba J (2014) Adam: A method for stochastic optimization. arXiv preprint
446    arXiv:1412.6980.
447
448 Kuksa P, Pavlovic V (2009) Efficient alignment-free DNA barcode analytics. BMC
449    Bioinformatics 10(S14):S9. DOI: https://doi.org/10.1186/1471-2105-10-S14-S9
450
451 Langenkämper D, Goesmann A, Nattkemper TW (2014) Ake-the accelerated *k*-mer exploration
452    web-tool for rapid taxonomic classification and visualization. BMC Bioinformatics
453    15(1):384. DOI: https://doi.org/10.1186/s12859-014-0384-0
454
455 Makiola A, Compson ZG, Baird DJ, Barnes MA, Boerlijst SP, Bouchez A, Brennan G, Bush A,
456    Canard E, Cordier T, Creer S (2020) Key questions for next-generation biomonitoring.
457    Frontiers in Environmental Science. DOI: https://doi.org/10.3389/fenvs.2019.00197
458
459 Nugent CM, Elliott TA, Ratnasingham S, Adamowicz SJ (2020) Coil: an R package for
460    cytochrome C oxidase I (COI) DNA barcode data cleaning, translation, and error
461    evaluation. Genome. 63(6):291-305. DOI: https://doi.org/10.1139/gen-2019-0206
462

11

463 Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer
464        P, Weiss R, Dubourg V, Vanderplas J (2011) Scikit-learn: Machine learning in Python.
465        Journal of Machine Learning Research 12(Oct):2825-30.
466
467 Pruesse E, Quast C, Knittel K, Fuchs BM, Ludwig W, Peplies J, Glöckner FO (2007) SILVA: a
468        comprehensive online resource for quality checked and aligned ribosomal RNA sequence
469        data compatible with ARB. Nucleic Acids Research 35(21):7188-96. DOI:
470        https://doi.org/10.1093/nar/gkm864
471
472 Ratnasingham S, Hebert PDN (2007) BOLD: The Barcode of Life Data System (http://www.
473        barcodinglife.org). Molecular Ecology Notes. 7(3):355-64. DOI:
474        https://doi.org/10.1111/j.1471-8286.2007.01678.x
475
476 Singer GA, Shekarriz S, McCarthy A, Fahner N, Hajibabaei M (2020) The utility of a
477        metagenomics approach for marine biomonitoring. bioRxiv. DOI:
478        https://doi.org/10.1101/2020.03.16.993667
479
480 Solis-Reyes S, Avino M, Poon A, Kari L (2018) An open-source *k*-mer based machine learning
481        tool for fast and accurate subtyping of HIV-1 genomes. PLoS ONE 13(11). DOI:
482        https://doi.org/10.1371/journal.pone.0206409
483
484
485 Stat M, Huggett MJ, Bernasconi R, DiBattista JD, Berry TE, Newman SJ, Harvey ES, Bunce M
486        (2017) Ecosystem biomonitoring with eDNA: metabarcoding across the tree of life in a
487        tropical marine environment. Scientific Reports 7(1):1-1. DOI:
488        https://doi.org/10.1038/s41598-017-12501-5
489
490 Stefanni S, Stanković D, Borme D, de Olazabal A, Juretić T, Pallavicini A, Tirelli V (2018)
491        Multi-marker metabarcoding approach to study mesozooplankton at basin scale.
492        Scientific Reports 8(1):1-3. DOI: https://doi.org/10.1038/s41598-018-30157-7
493
494 Taberlet P, Coissac E, Hajibabaei M, Rieseberg LH (2012) Environmental DNA. Molecular
495        Ecology 21(8):1789-93. DOI: https://doi.org/10.1111/j.1365-294X.2012.05542.x
496
497 Wang Q, Garrity GM, Tiedje JM, Cole JR (2007) Naive Bayesian classifier for rapid assignment
498        of rRNA sequences into the new bacterial taxonomy. Applied and Environmental
499        Microbiology 73(16):5261-7. DOI: https://doi.org/10.1128/AEM.00062-07
500
501 Wilson JJ, Brandon-Mong GJ, Gan HM, Sing KW (2019) High-throughput terrestrial
502        biodiversity assessments: mitochondrial metabarcoding, metagenomics or
503        metatranscriptomics?. Mitochondrial DNA Part A 30(1):60-7. DOI:
504        https://doi.org/10.1080/24701394.2018.1455189
505
506 Zielezinski A, Girgis HZ, Bernard G, Leimeister CA, Tang K, Dencker T, Lau AK, Röhling S,
507        Choi JJ, Waterman MS, Comin M (2019) Benchmarking of alignment-free sequence

508    comparison methods. Genome Biology 20(1):144. DOI: https://doi.org/10.1186/s13059-
509        019-1755-7
510

511    Zinger L, Bonin A, Alsos IG, Bálint M, Bik H, Boyer F, Chariton AA, Creer S, Coissac E,
512        Deagle BE, De Barba M (2019) DNA metabarcoding—Need for robust experimental
513        designs to draw sound ecological conclusions. Molecular Ecology 28(8):1857-62. DOI:
514        https://doi.org/10.1111/mec.15060
515

**Supplementary Files**

**Supplementary File S1 –** Training, test, and validation data sets used in model training and analysis

**Supplementary File S2 –** Python script for custom grid search of hyperparameters for optimization of the neural network.

**Supplementary File S3 –** The parameters utilized in the grid search for each of the five machine learning algorithms tested in the design of the Alfie package.

**Supplementary File S4 –** Jupyter notebook with tutorial demonstrating how to apply the Alfie classifier in the Python programming language, and how to train custom alignment-free classifiers using the Alfie training module.

531     **Tables and Figures**

532

533     **Table 1.** The numbers of COI barcode sequences obtained from BOLD for each kingdom and
534     the number of sequences retained within different data sets used in development of the Alfie
535     package. The raw barcode counts represent the complete set of publicly available sequences for
536     the given kingdom. The 'Barcodes utilized' column is the total number of sequences used in the
537     analysis for the given kingdoms after filtering based on minimum sequence length and down
538     sampling to decrease imbalanced representation of the different kingdoms. The breakdown of
539     these sequences between the train, test, and validation data sets is also shown.

| Kingdom | Raw barcode count | Barcodes utilized | Train data set size | Test data set size | Validation data set size (see Table 2) |
|---|---|---|---|---|---|
| **Animal** | 1,137,552 | 23,493 | 18,189 | 4,547 | 757 |
| **Bacteria and Archaea** | 5,565 | 5,547 | 4,380 | 1,095 | 72 |
| **Fungi** | 1,407 | 1,368 | 1,038 | 260 | 70 |
| **Plant** | 22,638 | 22,599 | 18,017 | 4,505 | 77 |
| **Protist** | 5,029 | 5,026 | 4,014 | 1,003 | 9 |
| **Total** | 1,172,191 | 58,033 | 45,638 | 11,410 | 985 |

540

15

541   **Table 2**. The taxonomic breakdown of the validation data set. For each kingdom, a taxonomic
542   class with a near average number of sequences in the kingdom's whole data set was chosen for
543   exclusion from the training set and inclusion in the validation data set. The names of the
544   taxonomic classes and the numbers of barcode sequences withheld from training and testing for
545   subsequent validation are shown.

| Kingdom | Withheld class | Sequence count |
|---|---|---|
| **Animal** | Diplopoda | 757 |
| **Bacteria and Archaea** | Flavobacteria | 72 |
| **Fungi** | Leotiomycetes | 70 |
| **Plant** | Liliopsida | 77 |
| **Protist** | Heterotrichea and Colpodea | 9 |

546

16

547    **Table 3.** The architectures of the neural networks tested in conjunction with the different $k$-mer
548    feature sets. For each $k$-mer feature set and corresponding neural network, the average loss and
549    accuracy scores from 5-fold cross validation on the training data are presented. Each neural
550    network was comprised of a dense input layer (neuron number = number of unique $k$-mers, or
551    $4^k$), five hidden layers of neurons (neuron counts for each layer given in table), and a dense
552    output layer (neuron size equal to number of classes). The input and hidden layers utilized a
553    rectified linear unit (relu) activation function (Agarap 2018), and the hidden layers had dropout
554    rates of 0.3. The final output layer utilized a softmax activation function, and the models were
555    trained using an Adam optimizer (Kingma & Ba 2014), minimizing sparse categorical cross
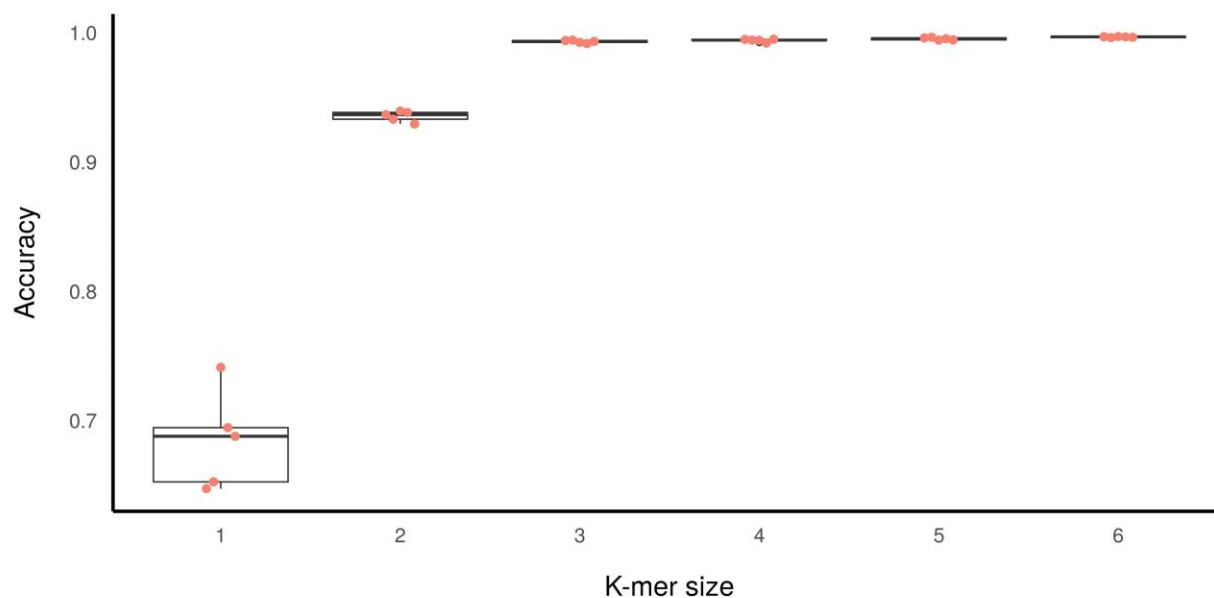556    entropy.

557

| $K$-mer size | NN hidden layers sizes | Average accuracy | Average loss |
|:---:|:---:|:---:|:---:|
| 1 | [4,64,128,32,16] | 0.684 | 0.899 |
| 2 | [16,64,128,64,16] | 0.935 | 0.216 |
| 3 | [64,128,64,32,16] | 0.993 | 0.038 |
| 4 | [256,128,64,32,16] | 0.994 | 0.033 |
| 5 | [1024,512,256,64,16] | 0.995 | 0.047 |
| 6 | [2080,1040,520,260,130] | 0.997 | 0.023 |

558

559 **Table 5.** The accuracy scores for the predictions made by the five different machine learning
560 models (trained on 4-mer frequency features and the complete training data set). Accuracy on the
561 test and validation data sets (Table 1) are shown.

562

| Algorithm | Test accuracy | Validation accuracy |
|---|---|---|
| DNN | 0.996 | 0.976 |
| Support Vector Machine | 0.996 | 0.974 |
| K Nearest Neighbors | 0.997 | 0.927 |
| Random Forest | 0.983 | 0.861 |
| XGBoost | 0.998 | 0.972 |

563

18

564
565   **Figure 1.** Boxplot of the 5-fold cross validation accuracy results for the training of models of
566   different *k*-mer feature sets and corresponding neural network architectures on the training data.
567   Each dot represents an accuracy score for one of the individual fold in the cross-validation
568   corresponding to the given *k*-mer feature set.
569

**Acknowledgements**

571

**Competing Interests**

586

The authors have declared that no competing interests exist.

588

**Author Contributions**

590

The study was conceived and designed by CMN and SJA. Development of the Alfie package
was performed by CMN. The initial draft of the manuscript was written by CMN. CMN and SJA
contributed to the editing of the manuscript.