1	Computational antigen discovery for eukaryotic pathogens using
2	Vacceed
3	
4	
5	
6	Stephen J. Goodswen <sup>1</sup> , Paul J. Kennedy <sup>2</sup> , John T. Ellis <sup>1,*</sup>
7	
8	<sup>1</sup> School of Life Sciences, University of Technology Sydney (UTS), Ultimo, NSW, Australia.
9	<sup>2</sup> School of Software, Faculty of Engineering and Information Technology and the Centre for
10	Artificial Intelligence, University of Technology Sydney (UTS), Ultimo, NSW, Australia.
11	
12	Correspondence: John.Ellis@uts.edu.au
1.2	
13	

**Running Head: Eukaryotic pathogen** antigen discovery using Vacceed

#### 15 Abstract

16 Bioinformatics programs have been developed that exploit informative signals encoded 17 within protein sequences to predict protein characteristics. Unfortunately, there is no program 18 as yet that can predict whether a protein will induce a protective immune response to a 19 pathogen. Nonetheless, predicting those pathogen proteins most likely from those least likely 20 to induce an immune response is feasible when collectively using protein characteristics that 21 can be predicted. Vacceed is a computational pipeline that manages different standalone 22 bioinformatics programs to predict various protein characteristics, which offer supporting 23 evidence on whether a protein is secreted or membrane associated. A set of machine learning 24 algorithms predicts the most likely pathogen proteins to induce an immune response given 25 the supporting evidence. This chapter provides step by step descriptions of how to configure 26 and operate *Vacceed* for a eukaryotic pathogen of the user's choice.

27

28

Key Words *Vacceed*, machine learning, *in silico* vaccine discovery, computational antigen
discovery, eukaryotic pathogen

#### 32 **1. Introduction**

33 Proteins sequences are not random assemblies of amino acids. There is a precise biological 34 reason why one particular amino acid is connected to another, which ultimately contributes to 35 a protein's distinctive characteristics [1]. Researchers, over the last two decades, have 36 developed bioinformatics programs that exploit informative signals or patterns encoded 37 within these amino acid sequences to predict protein characteristics. Examples of these 38 characteristics are subcellular localization [2], presence and location of signal peptide 39 cleavage sites [3], and transmembrane topology[4]. With respect to discovering protein 40 vaccine candidates, no signal has yet been detected that helps predict a characteristic 41 signifying a protein's contributing capacity to a protective immune response in a host. 42 Consequently, the current computational antigen discovery aspiration is to distinguish those 43 pathogen proteins most likely (referred to henceforth as positives) from those least likely 44 (referred to henceforth as negatives) to induce an immune response.

45 *Vacceed* is the collective name for a configurable pipeline of linked bioinformatics programs, 46 Perl scripts, R functions and Linux shell scripts [5]. It was inspired by the principles of 47 reverse vaccinology [6], whereby antigen discovery starts *in silico* using the pathogen 48 genome rather than the traditional culture-based method of cultivating and dissecting the 49 pathogen itself. Vacceed has been designed to facilitate an automated, high-throughput 50 computational approach to predict vaccine candidates against eukarvotic pathogens given 51 protein sequences [7]. The pipeline uses various standalone bioinformatics programs to 52 predict various protein characteristics. Vacceed is grounded on the underlying premise that there is an expected difference between the set of characteristics defining positives to those of 53 54 negatives. These differences are typically not apparent to an observer and hence applying a rule-based approach to distinguish proteins is not feasible. Conversely, machine learning 55

(ML) has the capacity to detect obscure differences. *Vacceed* uses a set of ML algorithms trained on protein characteristics of known positives and negatives to distinguish if a yet to be classified protein is a positive or negative [8]. So far, *Vacceed* has been used in studies to predict vaccine candidates for *Neospora caninum* [9] and *Cystoisospora suis* [10].

60 This chapter provides step by step descriptions of how to configure and operate *Vacceed* for a 61 eukaryotic pathogen of the user's choice. A prerequisite for pathogen choice, nonetheless, is a 62 substantial representation of the pathogen's proteome in the form of quality protein 63 sequences.

64 2. Vacceed core background information

65 Vacceed can be downloaded from: https://github.com/goodswen/vacceed/releases. The
66 download package includes a comprehensive Vacceed User Guide and sample data. Note that
67 Vacceed has been designed for a Linux operating system and has only been tested on Red Hat
68 Enterprise Linux 7.5, but is expected to work on most Linux distributions.

69 Each data processing stage in the *Vacceed* pipeline is an independent resource, which is built 70 from a central Linux shell script encapsulating all programs needed to perform specific but 71 related tasks. Typical tasks include predicting a particular protein characteristic. By default, 72 Vacceed uses seven bioinformatics programs to predict protein characteristics: Signal 5.0 [11] (predicts presence and location of signal peptide cleavage sites using deep neural 73 74 networks); WoLF PSORT 0.2 [12] and TargetP 1.1 [2] (predict subcellular localization); 75 TMHMM 2.0 [4] (predicts transmembrane domains in proteins); Phobius 1.01 [13] (predicts 76 transmembrane topology and signal peptides); DeepLoc 1.0 [14] (predicts eukaryotic protein 77 subcellular localization using deep learning); and IEDB peptide-MHC binding predictors 78 (MHCI version 2.17 and MHCII version 2.16.3) [15] (see Note 1). Observe that each of the

seven programs have specific version numbers on which *Vacceed* has been tested. There is no
assurance older or newer program versions will work.

The most pertinent file from a user's perspective is a species configuration file in a headerkey format (see Fig. 1). For example, [Resources] is the header, and 'name' is the key. Text following the '=' sign is configurable. A suggested convention is to have one configuration file for each target pathogen. *Vacceed* is started by entering only one command in a Linux Shell (or terminal) e.g. perl startup *xx*, where *xx* is a user specified code that links *Vacceed* to the target pathogen configuration file. No other commands are required.

Once *Vacceed* is started, each resource listed after the 'name' key is consecutively executed.
Resource names can be in any order or even excluded with the exception of VALIDATE (*see* **Note 2**) and EVIDENCE (*see* **Note 3**), which must always be the first and last in the list,
respectively. Any key in the configuration file can be used as a variable replacement in the
rest of the configuration file. That is, a '\$' character preceding a word denotes a variable e.g.
\$work dir is replaced by '\$HOME/vacceed' throughout the configuration file on execution.

93 Typical Vacceed run times are dependent on various factors including numbers of proteins to 94 process, programs to execute (resources), computer processors (cores), and the amount of 95 memory. For example, a test with 500 proteins processed through all resources completed in 96 3 hours, 21 minutes, and 17 seconds using Red Hat Enterprise Linux Workstation release 7.5, 97 64 bit kernel, and 32 MB memory with 8 cores; however, the same test without the resources 98 MHCI and MHCII completed in 23 minutes and 54 seconds. Vacceed takes advantage of 99 multi-core processors. By default, the proteins to process are split into subsets by the number 100 of cores and then each subset is processed in parallel.

# **3. Methods**

demonstration purposes.

102	3.1 Running Vacceed with sample data		
103	The Vacceed installation provides sample data comprising a small collection of Toxoplasma		
104	gondii proteins as input. The purpose of this section is to test the Vacceed installation.		
105	1. Install <i>Vacceed</i> (see Note 4).		
106	2. Edit the species configuration file 'toxoplasma.ini' located in the directory		
107	<install_dir>/vacceed/start/config_dir (where <install_dir> is the directory in which</install_dir></install_dir>		
108	Vacceed was installed). Under the [Resources] header, remove MHCI and MHCII		
109	( <i>see</i> <b>Note 5</b> ).		
110	3. Under the [Main] header, change the current path assigned to work_dir to		
111	install_dir/vacceed/.		
112	4. Under the [Main] header, assign an appropriate e-mail address to email_url		
113	5. In a command-line terminal, change directory to <i>install_dir</i> /vacceed/start.		
114	6. Enter the command: <b>perl startup tg</b>		
115	7. An e-mail will automatically be sent either when the pipeline is successfully		
116	completed or immediately when an error occurs. A log file is attached to the e-mail		
117	providing details of success or failure (see Note 6).		
118	8. If successful, the main output file called 'vaccine_candidates' is created in the		
119	directory install_dir/vacceed/toxoplasma/proteome. This file contains a list of all		
120	processed proteins ranked on average ML scores (see Fig 2. and Note 7).		
121	3.2 Running <i>Vacceed</i> with user provided data		
122	Once the Vacceed installation has been successfully tested, Vacceed can be configured and		
123	operated for a eukaryotic pathogen of the user's choice. Neospora caninum is used here for		

125	1. Collect <i>all</i> known protein sequences of the target pathogen into one file ( <i>see</i> <b>Note 8</b> ).
126	The sequences must be in a FASTA format with a sequence identifier in the following
127	layout: $>xx$  protein Identifier (ID)  text (optional), where xxx can be any characters
128	e.g. 'tr' or 'sp' as per UniProt identifiers.
129	2. Copy file from step #1 into <i>install_dir</i> /vacceed/start/proteome
130	3. Copy the entire template_species directory to a user-named directory, e.g., neospora.
131	4. Copy the species configuration file 'toxoplasma.ini' located in the directory
132	install_dir/vacceed/start/config_dir to 'neospora.ini'.
133	5. Add a new line to startup.ini located in <i>install_dir</i> /vacceed/start/:
134	nc< Neospora caninum <pipeline<neospora.ini< config_dir<="" install_dir="" start="" td="" vacceed=""></pipeline<neospora.ini<>
135	6. Edit neospora.ini to match the following:
136	work_dir="install_dir/vacceed"
137	species_dir="neospora"
138	email_url="your_email@address" (user e-mail address)
139	proteome_fasta="proteome.fasta" (protein sequence file as per step #1)
140	prot_id_prefix="xxx" (needs to match the sequence identifier as per step #1)
141	7. Modify the [Resources] in neospora.ini, if required. That is, remove any resource
142	names between VALIDATE and EVIDENCE that are not required e.g. MHCI and
143	MHCII.
144	8. Change directory to <i>install_dir</i> /vacceed/start in a command-line terminal.
145	9. Enter the command: <b>perl startup nc</b> (where 'nc' is as per step #5).
146	10. Check results in 'vaccine_candidates' in <i>install_dir</i> /vacceed/neospora/proteome
147	3.3 Creating pathogen specific training data
148	Training data here is essentially the collection of predicted evidence (referred to henceforth
149	as evidence profiles) from the seven bioinformatics programs for those proteins known to be

150	positive or negative. A training data file called 'train_profiles' is provided with the Vacceed				
151	package as part of the <i>T. gondii</i> sample data (see Note 9). A previous study [8] tested				
152	Vacceed with different evidence profiles compiled from different eukaryotic species. It				
153	concluded that there is no fundamental difference in evidence profile patterns, e.g., a model				
154	trained on one species can be used to classify proteins from another. This is because the				
155	bioinformatics programs are designed or ML trained for eukaryotes in general. Therefore, the				
156	creation of a pathogen specific training dataset is not a mandatory step. However, an ideal				
157	training dataset is one that contains the greatest variety of evidence profiles (see Note 10)				
158	irrespective of the source species, e.g., quality and variety are indisputably the most				
159	important factors that impact the accuracy of ML algorithms [8]. A new or amended training				
160	file is recommended under any of the following circumstances: a bioinformatics program is				
161	upgraded, i.e., it has improved accuracy; experimentally proved immunogenic proteins				
162	become available; and a new prediction program is added (see Section 3.5).				
163	1. Collect as many proteins as possible for the target species that are known to induce an				
164	immune response in the relevant host. The proteins will represent the 'positives' for				
165	the training file (see Note 11).				
166	2. Collect proteins that <i>do not</i> induce an immune response. These proteins will represent				

167 the 'negatives' (*see* **Note 12**).

- 168 3. Create a file (e.g., positives.fasta) containing the positive sequences in a FASTA
  169 format.
- 4. Create a file (e.g., negatives.fasta) containing the negative sequences in a FASTA
  format.
- 172 5. Copy both FASTA files into *install\_dir*/vacceed/start/proteome

173 6. Copy the entire template\_species directory to a user-named directory, e.g., training.

174 7. Copy 'toxoplasma.ini' to 'train.ini'.

175	8. Add a new line to startup.ini located in <i>install_dir</i> /vacceed/start/:
176	train< Neospora caninum <pipeline<train.ini< config_dir<="" install_dir="" start="" td="" vacceed=""></pipeline<train.ini<>
177	9. Edit train.ini to match the following:
178	work_dir="install_dir/vacceed"
179	species_dir="training"
180	email_url="your_email@address" (user e-mail address)
181	proteome_fasta="positives.fasta" (as per step #3)
182	prot_id_prefix="xxx" (needs to match the sequence identifier)
183	11. Modify [Resources] in train.ini if required e.g. remove any resource not installed or
184	required.
185	12. Change directory to <i>install_dir</i> /vacceed/start in a command-line terminal.
186	13. Enter the command: <b>perl startup train</b> (where 'train' is as per step #8).
187	14. Copy the file 'evidence_profiles' from
188	install_dir/vacceed/training/pipeline/evidence/output to
189	install_dir/vacceed/training/pipeline/evidence/training_files
190	15. Rename evidence_profiles to a user-defined name e.g. neospora_profiles
191	16. Add ',YES' to the end of each row in the new training file (exclude the first row). The
192	'YES' is the required target label for the positives.
193	17. Edit train.ini to match the following:
194	proteome_fasta="negatives.fasta" (as per step #4)
195	18. Change directory to <i>install_dir</i> /vacceed/start in a command-line terminal.
196	19. Enter the command: perl startup train
197	20. Add ',NO' (i.e., the required target label for the negatives) to the end of each row in
198	evidence_profiles in <i>install_dir</i> /vacceed/training/pipeline/evidence/output

- 199 21. Append the entire contents of the amended evidence\_profiles (except first row) to the200 new training file, e.g., neospora profiles.
- 201 22. Copy new training file to *install\_dir*/vacceed/<new species>/evidence/training\_files
   202 where <new species> is the directory created for the target species, e.g., neospora.
- 203 23. Edit the species configuration file, e.g., neospora.ini and change the value of the
- 204 train\_file key under header [EVIDENCE] to the new training file, e.g.,

205 neospora\_profiles (*see* Note 13).

- 206 24. The new training data should be evaluated with techniques such as k-fold cross
- 207 validation (*see* **Note 14**) and the ML algorithm parameters tweaked to improve
- 208 performance (*see* **Note 15**).

#### 209 **3.4 Creating MHCI and MHCII training data**

- 210 This section is only applicable when using resources MHCI and/or MHCII *and* the target
- 211 pathogen host is **not** human. By default, *Vacceed* uses human alleles (e.g. HLA-A\*01:01) for
- 212 peptide-MHC binding predictions. The following describes steps required to setup MHCI for
- a host other than human, e.g., mouse.
- 1. Follow steps #1 to #5 from Section 3.3.
- 215 2. Create a file (e.g. mouse\_mchI\_alleles) in a comma delimited format containing all
- 216 required mouse alleles and peptide lengths e.g. H-2-IAb,8 where each 'allele, length'
- is on a separate line (*see* **Note 16**).
- 218 3. Copy the entire template\_species directory to a user-named directory, e.g., mouse.
- 4. Copy mouse\_mchI\_alleles to < *install\_dir*/Vacceed/mouse/pipeline/mhci/alleles
- 220 5. Copy 'toxoplasma.ini' to 'mouse.ini'.
- 221 6. Add a new line to startup.ini located in *install\_dir*/vacceed/start/:
- 222 m<mouse\_pipeline<mouse.ini<*install\_dir*/Vacceed/start/config\_dir
- 223 7. Edit mouse.ini to match the following:

224	work_dir="install_dir/vacceed"
225	species_dir="mouse"
226	email_url="your_email@address" (user e-mail address)
227	proteome_fasta="positives.fasta"
228	prot_id_prefix="xxx" (needs to match the sequence identifier)
229	allele_file="mouse_mchI_alleles" (located under the resource [MHCI_files])
230	8. Modify [Resources] in mouse.ini to 'name=VALIDATE,MHCI'
231	9. Change directory to <i>install_dir</i> /vacceed/start in a command-line terminal.
232	10. Enter the command: <b>perl startup m</b> (where 'm' is as per step #6).
233	11. Copy mhci_ml.txt from install_dir/vacceed/mouse/pipeline/mhci/output from
234	install_dir/vacceed/mouse/pipeline/mhci/training_files
235	12. Rename mhci_ml.txt to a user-defined name e.g. mouse_mhci_ml.txt
236	13. Add ',YES' to the end of each row in the new training file (exclude the first row)
237	14. Edit mouse.ini to match the following:
238	proteome_fasta="negatives.fasta"
239	15. Change directory to <i>install_dir</i> /vacceed/start in a command-line terminal.
240	16. Enter the command: perl startup m
241	17. Add ',NO' to the end of each row in mhci_ml.txt in
242	install_dir/vacceed/mouse/pipeline/mhci/output
243	18. Append the entire contents of the amended mhci_ml.txt (except first row) to the new
244	training file, e.g., mouse_mhci_ml.txt.
245	19. Copy new training file to <i>install_dir</i> /vacceed/ <new species="">/mhci/training_files</new>
246	where <new species=""> is the directory created for the target species, e.g., new_mouse.</new>

247 20. Edit the species configuration file, e.g., new\_mouse.ini, and change the value of the
248 train\_file key under header [MHCI\_files] to the new training file, e.g.,

249 mouse\_mhci\_ml.txt

250 21. Repeat the steps above to create a MHCII training file, but change mhci to mhcii (*see*251 Note 17).

## 252 **3.5 Add a new resource**

253 New programs to predict protein characteristics will inevitably be developed in the future.

254 This section describes how to incorporate a new program into *Vacceed*, which essentially is

- adding a new resource with the goal of extracting relevant evidence from the new program
- 256 output to append to evidence profiles.
- Install and test new program with sample data to ensure it runs successfully from any
   directory (*see* Note 18).
- 259 2. Determine the input requirements and the output format of new program.
- 260 3. Add a new resource name, e.g., program\_Z in an appropriate configuration file:
- 261 [Resources]

# 262 name=VALIDATE,WOLF,TMHMM,**PROGRAM\_Z**,EVIDENCE

- 4. Add a new section to the same configuration file. The easiest way to do this is to copy
  an existing resource and amend accordingly (see Fig. 3). The texts highlighted in red
  are the only parts expected to be changed.
- 5. Create a new directory in *install\_dir*/vacceed/*new\_species*/pipeline using the same
- 267 name as the new resource (but in lowercase), e.g., program\_z.
- 268 6. Create two directories called 'output' and 'scripts' in the program\_z directory.
- 269 7. Copy 'template\_resource\_script' from
- 270 *install\_dir*/vacceed/*new\_species*/pipeline/common\_programs to
- 271 *install\_dir*/vacceed/*new\_species*/pipeline/program\_z

- 8. Rename 'template\_resource\_script' to a user-named file e.g. program\_z\_script (*see*Note 19)
- 274 9. Amend program z script where it states << Add new programs here >>, e.g., 275 echo "script step=\">> executing program z\"" >> \$script dir/script\$chr no 276 echo "program z \$required input \$out dir" >> \$script dir/script\$chr no || error exit 277 Where \$required input is the input as determined in step #2. 278 10. A generic Perl script called 'get evidence.pl' (located in: 279 install dir/vacceed/new species/pipeline/common programs) can be amended accordingly to extract the relevant evidence from the program z output file (see Note 280 281 **20**). Alternatively, any programming language can be used to write a program to 282 extract evidence. In such a case, the program name would need to replace 283 'get evidence.pl' in program z script. Regardless of the extraction program, 284 evidence needs to be saved in a user-named file with the suffix ' evd', e.g., 285 programz evd in the directory 286 install dir/vacceed/new species/pipeline/evidence/output.

#### **287 4. Notes**

- The bioinformatics programs are third-party and are not part of the *Vacceed* package.
   Furthermore, installation steps for the third-party programs are not described in this
   chapter. Most of the programs provide a **ReadMe** file with instructions. Even so,
   these installations are still a challenging aspect to preparing *Vacceed* ready for use. It
   is highly recommended to seek the help of an administrator or an experienced Linux
   user.
- 294 2. *Vacceed* checks to see if a protein sequence contains invalid letters, e.g., characters
  295 other than [ACDEFGHIKLMNPQRSTVWY].

296	3.	Vacceed collates relevant, predicted protein characteristics (typically in the form of
297		numerical values) into one file called evidence_profiles, i.e., contents of all files with
298		the extension '_evd' in the evidence/output directory are combined as columns into
299		evidence_profiles.
300	4.	Ensure that each third-party program runs successfully before testing Vacceed.
301	5	The computation of peptide-MHC predictions takes less than a few minutes
302		depending on the computer environment to run the test when both MHCI and MHCII
303		are removed. Furthermore, MHCI and MHCII predictions on the whole are not
304		accurate [16] (particularly MHCII [17]) and only marginally contributed to the
305		Vacceed end result when tested with the T. gondii sample data [8].
306	6	If Vacceed fails with the test data then it will inevitably fail with any other data. The
307		expected reason for the failure is installation issues of one or more of the third-party
308		programs (see Note 4). The log file may give clues as to which third-party program(s)
309		is the culprit.
310	7	The ML algorithms used are listed in the configuration file under the header
311		[Evidence] and the key 'algorithms'.
312	8	It is recommended that all known pathogen proteins are processed irrespective of
313		protein name or expected function. This allows for an unbiased approach.
314	9	This training file contains 475 positives of mainly <i>T. gondii proteins</i> (nine are <i>N</i> .
315		caninum). A small selection of these proteins are known to induce an immune
316		response, but most are proteins predicted to be membrane-associated or secreted, i.e.,
317		proteins exposed to the immune system. There are 501 T. gondii proteins representing
318		negatives, which were defined by the protein's predicted sub-cellular location, i.e.,
319		neither membrane-associated nor secreted.

10 Variety, in this instance from a ML perspective, is having a generalised selection of
proteins in the training file that are representative of all conceivable types of positive
and negative proteins, e.g., with a limited selection, a ML algorithm may not
generalize to evidence profiles not seen when it was learning (i.e., poorly predicts
when given new data).

325 11 Finding training proteins for most species is not a trivial task. The expectation is that a 326 thorough search of the literature will be required. Even then, there may still be an 327 inadequate number of examples to create a training file. A suggested compromise is to 328 use positive proteins from a closely related organism or proteins 'expected' to induce 329 or not induce an immune response. For instance, use proteins known to be exposed to 330 the immune system (e.g., membrane associated or secreted proteins) for positives and 331 non-exposed proteins (e.g., proteins normally located in the interior of the organism) 332 as negatives.

A drawback for collecting negative examples is that a protein cannot definitively be
 defined a negative unless it has been explicitly tested in a laboratory.

The same proteins should never be used for training and evaluation. This would
introduce biased results. Typically, the proteins are randomly divided into two sets.
One set containing the majority of data e.g. 80% for training. The other set (e.g., 20%)
used to evaluate the trained model's performance.

14 *k*-Fold cross-validation is a resampling statistical method used to estimate the
performance of ML models. The '*k*' refers to the number of groups that a given data
sample is to be split, e.g., 10-fold cross-validation indicates the sample data is split
into 10 groups. One group in turn is used as a test dataset and the remaining groups
used for training. The average of the *k* evaluation scores provides an indication of

how the model is expected to perform when used to make predictions on data not usedduring model training.

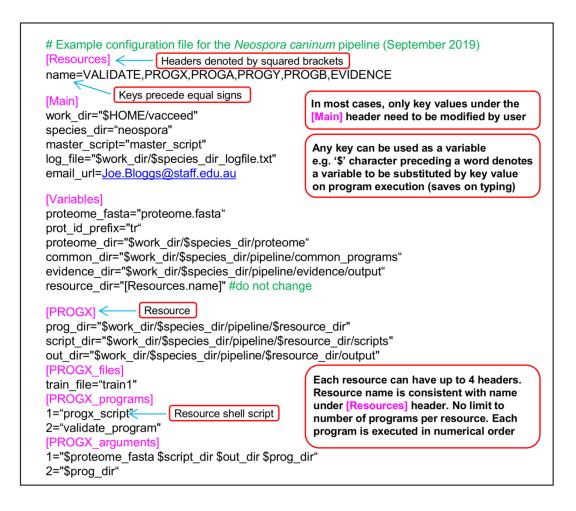
15 The distributed version of *Vacceed* is configured to run ML algorithms via R 346 347 functions contained in packages. The algorithms are executed using Rscript. There are 348 three R functions in *install dir*/vacceed/<new species>/evidence/ that encapsulate the 349 relevant command for each algorithm: <al> wrapper.R, <al> runPred.R, and <al> makePred.R, where <al> is the algorithm abbreviation. Parameters to fine tune 350 351 the algorithms can be modified in <al> makePred.R e.g parameters 'ntree' and/or 'mtry' in rf makePred.R, where rf = random forest, ntree = number of decision trees, 352 353 and 'mtry' = number of variables to try at each split in the decision tree. 354 16 Run the following command to see available class I alleles: ./src/predict binding.py 355 IEDB recommended mhc (only listed alleles can be used). 356 17 Run the following command to see available class II alleles: 357 python mhc II binding.py allele (only listed alleles can be used). 358 18 May need to append new program location to the PATH variable. 359 19 This is a template script only and will need to be edited appropriately to suit the new 360 program. There are user comments denoted by a '#' symbol, but a familiarity with 361 Linux scripting is expected. 362 20 Amending get evidence.pl requires experience in writing Perl scripts. Reading step 363 #8 under the section 'Adding a new resource' in the Vacceed User Guide may prove 364 useful when amending get evidence.pl.

365 Acknowledgements

- 366 SJG gratefully acknowledges Zoetis (Pfizer) Animal Health for funding the development of
- 367 *Vacceed* through a PhD scholarship.
- 368 References

- 369 1. Lee D, Redfern O, Orengo C (2007) Predicting protein function from sequence and structure.
- 370 Nature Reviews Molecular Cell Biology 8 (12):995-1005. doi:10.1038/nrm2281
- 2. Emanuelsson O, Brunak S, von Heijne G, Nielsen H (2007) Locating proteins in the cell using
- TargetP, SignalP and related tools. Nature Protocols 2 (4):953-971. doi:10.1038/nprot.2007.131
- 373 3. Petersen TN, Brunak S, von Heijne G, Nielsen H (2011) SignalP 4.0: discriminating signal peptides
- from transmembrane regions. Nature Methods 8 (10):785-786. doi:10.1038/nmeth.1701
- 4. Krogh A, Larsson B, von Heijne G, Sonnhammer ELL (2001) Predicting transmembrane protein
- topology with a hidden Markov model: Application to complete genomes. Journal of Molecular
- 377 Biology 305 (3):567-580. doi:10.1006/jmbi.2000.4315
- 378 5. Goodswen SJ, Kennedy PJ, Ellis JT (2014) Vacceed: a high-throughput in silico vaccine candidate
- discovery pipeline for eukaryotic pathogens based on reverse vaccinology. Bioinformatics 30
- 380 (16):2381-2383. doi:10.1093/bioinformatics/btu300
- 381 6. Rappuoli R (2000) Reverse vaccinology. Current Opinion in Microbiology 3 (5):445-450.
- 382 doi:10.1016/s1369-5274(00)00119-3
- 383 7. Goodswen SJ, Kennedy PJ, Ellis JT (2013) A guide to in silico vaccine discovery for eukaryotic
   384 pathogens. Briefings in Bioinformatics 14 (6):753-774. doi:10.1093/bib/bbs066
- 385 8. Goodswen SJ, Kennedy PJ, Ellis JT (2013) A novel strategy for classifying the output from an in
- 386 silico vaccine discovery pipeline for eukaryotic pathogens using machine learning algorithms. BMC
- 387 Bioinformatics 14. doi:10.1186/1471-2105-14-315
- 388 9. Goodswen SJ, Kennedy PJ, Ellis JT (2017) On the application of reverse vaccinology to parasitic
- 389 diseases: a perspective on feature selection and ranking of vaccine candidates. International Journal
- 390 for Parasitology 47 (12):779-790. doi:10.1016/j.ijpara.2017.08.004
- 391 10. Palmieri N, Shrestha A, Ruttkowski B, Beck T, Vogl C, Tomley F, Blake DP, Joachim A (2017)
- 392 The genome of the protozoan parasite Cystoisospora suis and a reverse vaccinology approach to
- 393 identify vaccine candidates. International Journal for Parasitology 47 (4):189-202.
- 394 doi:10.1016/j.ijpara.2016.11.007
- 395 11. Armenteros JJA, Tsirigos KD, Sonderby CK, Petersen TN, Winther O, Brunak S, von Heijne G,
- Nielsen H (2019) SignalP 5.0 improves signal peptide predictions using deep neural networks. Nature
   Biotechnology 37 (4):420-+. doi:10.1038/s41587-019-0036-z
- 398 12. Horton P, Park KJ, Obayashi T, Fujita N, Harada H, Adams-Collier CJ, Nakai K (2007) WoLF
- 399 PSORT: protein localization predictor. Nucleic Acids Research 35:W585-W587.
- 400 doi:10.1093/nar/gkm259
- 401 13. Kall L, Krogh A, Sonnhammer ELL (2004) A combined transmembrane topology and signal
- 402 peptide prediction method. Journal of Molecular Biology 338 (5):1027-1036.
- 403 doi:10.1016/j.jmb.2004.03.016
- 404 14. Armenteros JJA, Sonderby CK, Sonderby SK, Nielsen H, Winther O (2017) DeepLoc: prediction
- 405 of protein subcellular localization using deep learning. Bioinformatics 33 (21):3387-3395.
- 406 doi:10.1093/bioinformatics/btx431
- 407 15. Vita R, Zarebski L, Greenbaum JA, Emami H, Hoof I, Salimi N, Damle R, Sette A, Peters B
- 408 (2010) The Immune Epitope Database 2.0. Nucleic Acids Research 38:D854-D862.
- 409 doi:10.1093/nar/gkp1004
- 410 16. Bui HH, Sidney J, Peters B, Sathiamurthy M, Sinichi A, Purton KA, Mothe BR, Chisari FV,
- 411 Watkins DI, Sette A (2005) Automated generation and evaluation of specific MHC binding predictive
- 412 tools: ARB matrix applications. Immunogenetics 57 (5):304-314. doi:10.1007/s00251-005-0798-y
- 413 17. Wang P, Sidney J, Dow C, Mothe B, Sette A, Peters B (2008) A systematic assessment of MHC
- 414 class II peptide binding predictions and evaluation of a consensus approach. Plos Computational
- 415 Biology 4 (4). doi:10.1371/journal.pcbi.1000048
- 416

# 418 Figures



419

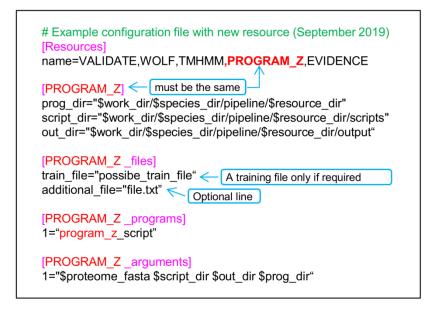
420 Fig. 1 Extract from a species configuration file

**#ID,ada,knn,nb,nn,rf,svm,average\_ML\_score** BBOV\_IV006420,1.000,1.000,1.000,1.000,0.0994,0.999 BBOV\_II001970,1.000,1.000,1.000,1.000,0.983,0.997 BBOV\_III005590,0.984,0.667,0.883,0.800,0.604,0.872,0.801 BBOV\_I004490,0.403,0.333,0.980,0.200,0.391,0.202,0.418 BBOV\_IV002290,0.427,0.667,0.000,0.000,0.463,0.141,0.283 BBOV\_III001400,0.015,0.333,0.764,0.200,0.255,0.121,0.281 BBOV\_III011900,0.000,0.000,0.000,0.001,0.007,0.001

422

# 423 Fig. 2 Extract from main *Vacceed* output file 'vaccine candidates'

424 Where ID = protein identifier, ada = adaptive boosting, knn = k-nearest neighbour classifier, 425 nb = Naive Bayes classifier, nn = neural network, rf = random forest, and svm = support 426 vector machines. vaccine candidates is a comma delineated file containing an ordered list of 427 all machine learning (ML) algorithm scores for each protein processed (seven in this 428 instance). Each ML algorithm generates probabilities that the YES and NO classifications are 429 correct, but only YES probabilities are displayed in the output. The 'average ML score' for each protein is the average probabilities of the YES classifications. The list order is 430 431 descending based on 'average ML score' value. An appropriate threshold value (e.g., 0.5) can 432 be compared to the average ML score to determine the relevant class, positive or negative.



#### 435 Fig. 3 Example of new resource added to species configuration file