

Machine learning as an effective method for identifying true SNPs in polyploid plants

Walid Korani¹, Josh P. Clevenger¹, Ye Chu¹ and Peggy Ozias-Akins^{1*}

Institute of Plant Breeding, Genetics and Genomics, University of Georgia, Tifton, Georgia,
United States of America

*Corresponding author

pozias@uga.edu

Abstract

Single Nucleotide Polymorphisms (SNPs) have many advantages as molecular markers since they are ubiquitous and co-dominant. However, the discovery of true SNPs especially in polyploid species is difficult. Peanut is an allopolyploid, which has a very low rate of true SNP calling. A large set of true and false SNPs identified from the *Arachis* 58k Affymetrix array was leveraged to train machine learning models to select true SNPs straight from sequence data. These models achieved accuracy rates of above 80% using real peanut RNA-seq and whole genome shotgun (WGS) re-sequencing data, which is higher than previously reported for polyploids. A 48K SNP array, Axiom *Arachis*2, was designed using the approach which revealed 75% accuracy of calling SNPs from different tetraploid peanut genotypes. Using the method to simulate SNP variation in peanut, cotton, wheat, and strawberry, we show that models built with our parameter sets achieve above 98% accuracy in selecting true SNPs. Additionally,

23 models built with simulated genotypes were able to select true SNPs at above 80% accuracy
24 using real peanut data, demonstrating that our model can be used even if real data are not
25 available to train the models. This work demonstrates an effective approach for calling highly
26 reliable SNPs from polyploids using machine learning. A novel tool was developed for
27 predicting true SNPs from sequence data, designated as SNP-ML (SNP-Machine Learning,
28 pronounced “snip mill”), using the described models. SNP-ML additionally provides
29 functionality to train new models not included in this study for customized use, designated SNP-
30 MLer (SNP-Machine Learner, pronounced “snip miller”). SNP-ML is freely available for public
31 use.

32

33 **Introduction**

34 Single Nucleotide Polymorphisms (SNPs) are a major source of variation across plant genotypes.
35 Therefore, the demand for discovery of a large number of SNPs increased after the advent of
36 Next-Generation Sequencing (NGS). However, the extraction of true SNPs in polyploid
37 organisms is challenging. Cultivated peanut is an allotetraploid, which poses an exceptional
38 challenge for the discovery of true SNPs since the two parental diploid genomes (A and B) are
39 very similar and the natural polymorphisms among peanut genotypes are very low [1,2].

40

41 Using Restriction-site-Associated DNA (RAD) sequencing, a large number of SNPs were
42 identified in peanut diploid species; however, very few SNPs were discovered in cultivated
43 peanuts [3]. Generally, the true SNP discovery in tetraploid peanut using NGS data is very low
44 [4-6]. Sliding Window Extraction of Explicit Polymorphisms (SWEEP) was developed to
45 improve the SNP calling by filtering out the polymorphisms between the two parental

46 subgenomes [7]. However, SNP calling in tetraploid peanut requires additional improvement. An
47 Affymetrix SNP array was designed using the SWEEP pipeline and showed promising
48 genotyping results among cultivated peanuts [8]. The chip showed that SWEEP identified ~40%
49 true SNPs in tetraploid peanut genotypes. The array provided an unprecedented number of
50 validated true and false SNP calls that can be leveraged with machine learning to increase the
51 accuracy of selection of true SNPs straight from sequence data. The ability to have confidence in
52 *in silico* SNP calls gives researchers access to all avenues of sequence-based genotyping
53 methods.

54

55 Machine learning applies sets of different algorithms that facilitate pattern recognition and
56 classification leading to prediction by creating models using existing data [9]. Machine learning
57 algorithms are divided into two major classes; *i.e.* supervised and unsupervised. Supervised
58 algorithms train previously well classified existing objects to predict the classes of new objects
59 based on available features. Unsupervised algorithms cluster objects depending on their features
60 without providing pre-defined classes. Both algorithms are used widely in different biological
61 fields; *e.g.* coding region recognition, signal peptide prediction, biomarker identification, disease
62 gene recognition, metabolic network detection, and protein-protein interaction [10-15]. For SNP
63 calling, neural networks were used to differentiate between true SNPs and sequence errors and
64 this method showed promising results for human SNPs [16]. In plants, neural networks also were
65 used to classify called SNPs as true or false positives and the approach showed a positive
66 prediction rate of 84.8% on the testing sets of soybean [17]. However, there has been little
67 application of machine learning in polyploid organisms where the occurrence of more than one

68 subgenome with high similarity to each other increases the complexity of read mapping and
69 confounds the calling of true SNPs.

70

71 In this study, different supervised machine learners were used to improve the discovery of
72 tetraploid peanut SNPs, utilizing the information of sequencing features and mapping data of the
73 validated true- and false-positive SNP data sets extracted from analysis of the *Arachis*
74 Affymetrix array. A new 48K SNP array was designed and validated based on the analysis of
75 this method. Simulated SNP variant data from peanut, cotton, wheat, and strawberry also were
76 used to extend the functionality of machine learning to other allopolyploids. Models trained with
77 simulated data then were used to select SNPs from real peanut data with an accuracy exceeding
78 80%. This result has implications for using machine learning to select true SNPs in polyploid
79 crops where no large validation sets are available. A tool was created, SNP-MLer (SNP-Machine
80 Learner), which allows users to train models for use in selecting true SNPs from sequence data.
81 The user can completely customize parameter sets used in training the models or default to the
82 complete set used to train the peanut models. The models then can be implemented in SNP-ML
83 (SNP-Machine Learning) to select true SNPs in new data sets.

84

85 **Materials and methods**

86 **Data sets**

87 The re-sequencing data set was created using 21 tetraploid *A. hypogaea* genotypes described in
88 Clevenger et al. (2017) [8] and deposited publically at ncbi.nlm.nih.gov (Bio Project
89 PRJNA340877 and Bio Samples SAMN05721179 to SAMN05721198). The RNA-seq data set

90 has information from nine tetraploid peanut genotypes described in Clevenger et al. (2015,
91 2016a, 2016b) [18-20]. Validated true and false-positive SNP sets were based on testing the
92 *Arachis* Affymetrix array with 384 peanut genotypes [8]. Mapping parameters were extracted
93 from the vcf files used for the original design of the array. All positions of SNPs and surrounding
94 sequence are based on the *A. duranensis* and *A. ipaensis* v1 pseudomolecules
95 [<https://peanutbase.org/>, 1]

96

97 **Creating and testing a new Affymetrix array based on SNP-ML**

98 A new affymttrix array was designed containing 28,218 SNPs which were extracted by SNP-ML
99 using peanut real data re-sequencing model of neural network and tree bagger (S1 Table). The
100 previously described 21 genotypes alongside 8 more genotypes and 103 minicore peanut lines
101 [8] were assayed on the array and all 28,218 SNP-ML-derived markers were manually curated
102 for polymorphism. A total of 21,112 markers were validated as polymorphic between genotypes
103 (75%).

104

105 **Creating and testing the machine learning models**

106 The data sets were prepared by R statistical software, *e.g.* extracting the attributes, randomly
107 created training and testing sets and preparing fasta files for SNP flanking segments. Various
108 toolboxes in MATLAB R2015b (the University of Georgia campus-wide site licensing
109 agreement) were used for different purposes. Bioinformatics Toolbox was used for calculating
110 the thermodynamic parameters, molecular weights and GC contents, Statistics and Machine
111 Learning Toolbox was used for creating and testing the different models of supervised machine

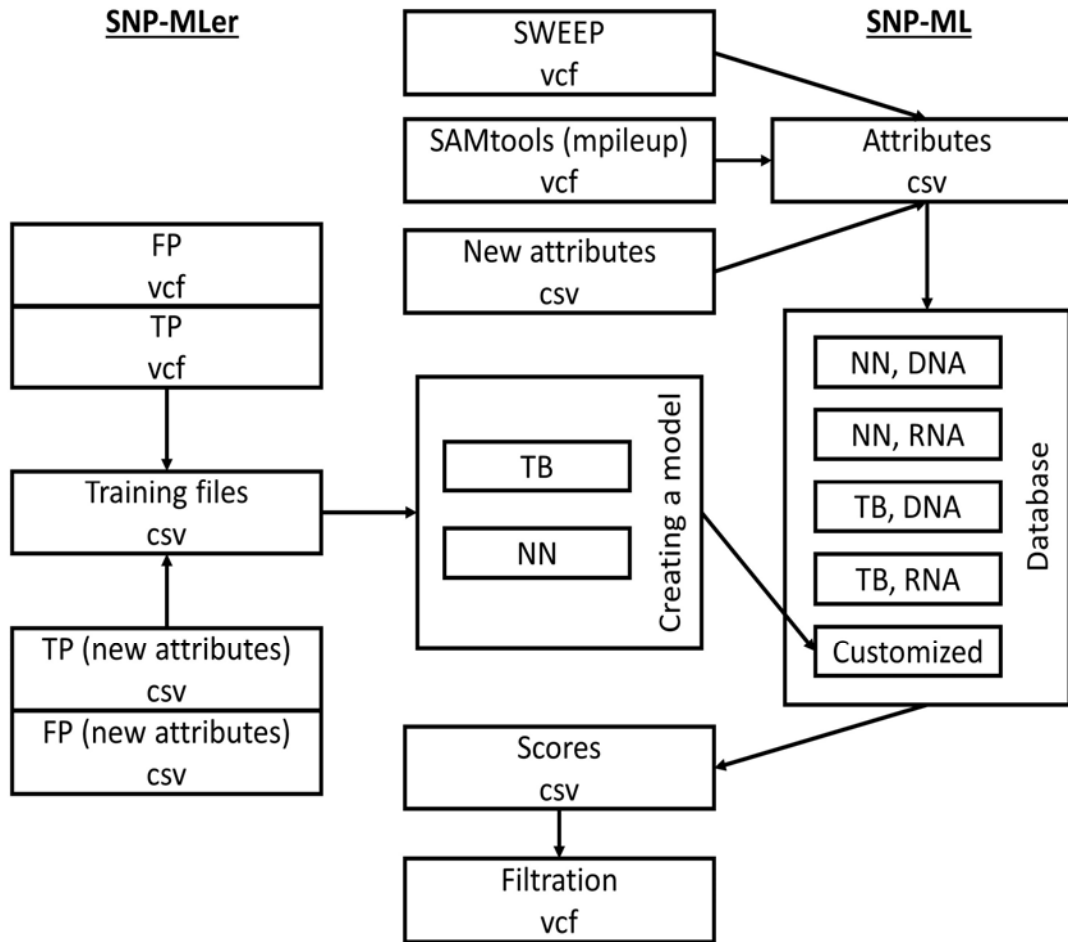
112 learning and Graphics functions were used for producing the bar plots and ROC (Receiver
113 Operating Characteristic) graphs. The specific arguments of the different machine learning
114 models are given in S2 Table.

115

116 **SNP-ML construction**

117 We built paired (neural network and TB) specific trainer models for the two data types, WGS re-
118 sequencing and RNA-seq. The models were built and stored in four files by a python script. In
119 addition, three C++ classes were built, vcf.h, csv_write.h and csv.h, to process vcf and csv files.
120 The SNP-ML main steps are illustrated in Fig 1. It uses C++ class vcf.h to extract the eight
121 selected attributes from the input file, which is a vcf file of the output of SNP calling by mpileup
122 of SAMtools, either directly or after primary filtration by SWEEP. The output is saved using the
123 C++ class csv_write.h into a csv file, which is read by a python script to be applied to one pair of
124 stored models (two files, one for neural network and the other for TB) depending on the data
125 type. The two score sets are saved to a csv file, which is read by C++ class csv.h. The scores are
126 filtered by passing only SNPs that have a value higher than the cutoff of neural network, which
127 can be selected by the user (the default is 0.5), and occurred in the two score sets (shared SNPs
128 in the output of the neural network and TB score file), in case the user selects that option. The
129 scores are stored in csv files and the corresponding SNPs are stored in a vcf file.

130



144

145

146

147 **Fig 1: SNP-ML/SNP-MLer infrastructure.**

148

149 To extend the program applications, a second tool was designed, designated SNP-MLer
 150 (pronounced ‘snip miller’) to allow users to create predictors that are suitable for interested
 151 species/experimental conditions. SNP-MLer uses reading/writing approach as described above, it
 152 takes validated true-positive and false-positive vcf files as input and generates predictor models
 153 as outputs.

154 Both tools, SNP-ML and SNP-MLer allow the user to skip or select some of the eight attributes,
155 and to apply new user defined attributes as csv files.

156

157 **SNP-ML requirements**

158 The script was written by C++ and python 2.7.1 (S1 file). C++ was used for processing the data,
159 input, output and filtering. The binary file was created by GCC 4.1.2 that was run on Red Hat
160 4.1.2-55 linux system. Python was used for creating the neural network and bagging machine
161 learning models and applying the prediction using them. Different python packages were used
162 for these purposes, i.e. numpy-1.11.0 (SciPy.org), scipy-0.17.1 (SciPy.org), pandas-0.18.1
163 (pandas.pydata.org), python-dateutil-2.0 (pypi.python.org), pytz-2016.4 (pypi.python.org),
164 scikit-learn-0.17.1 (scikit-learn.org) and pyrenn 0.1 (pyrenn.readthedocs.io).

165

166 **Creating and testing models using simulated data**

167 The pseudo molecule assembly AD1_BGI of cotton [<https://www.cottongen.org/>,21], the
168 pseudomolecule assembly of the 3B chromosome of wheat [22], the contigs of TGACv1 wheat
169 genome [<https://plants.ensembl.org/index.html>], the pseudomolecule assembly of *Fragaria vesca*
170 Genome v1.1 [<https://www.rosaceae.org/>, 23], and the contigs of *F. nipponica* Genome v1.0
171 (FNI_r1.1), *F. nubicola* Genome v1.0 (FNU_r1.1) and *F. orientalis* Genome v1.0 (FOR_r1.1)
172 [<https://www.rosaceae.org/>,24] were downloaded.

173 10,000 random loci were assigned in Chromosomes Aradu.A01, At_chr1, 3B and LG1, of
174 peanut, cotton, wheat and strawberry, respectively. The loci were randomly mutated five times to
175 form five synthetic genotypes using ART tool [25]. HiSeq 125 bp paired end sequences with
176 different depths, 10x to 50x, were generated. The fastq produced files were mapped using BWA

177 0.7.10 [26] with default parameters to synthetic references as follows: a synthetic tetraploid
178 reference containing Aradu.A01 and Araip.B01 chromosomes for peanut, a synthetic tetraploid
179 reference containing At_chr1 and Dt_chr1 for cotton, a synthetic hexaploid reference containing
180 3B chromosome and the contigs of A and D genomes for wheat, and a synthetic octoploid
181 reference containing LG1 chromosome and the contigs of FNI, FNU and FOR genomes for
182 strawberry. SNPs were called using samtools mpileup 1.2 and bcftools 1.2.1 with default
183 parameters without filtration. The SNP calling was carried out twice for every species. SNPs
184 between two genotypes were called in the first instance and SNPs among the five genotypes
185 were called in the second.

186

187 For each species, the SNPs located among the 10,000 loci were extracted in a separate vcf file,
188 and considered to be True-positive (TP) SNPs. Any others identified by the program were
189 extracted in another vcf file, and considered to be False-positive (FP) SNPs. Seventy percent of
190 each one were randomly selected, and combined to be used as training sets, and the remaining
191 30% were used as testing sets for Neural Network models using Matlab R2015b (the University
192 of Georgia campus-wide site licensing agreement).

193

194 **Testing simulated data against the real data:**

195 For peanut, 21 synthetic genotypes with 10X depth were generated and SNPs were called in four
196 batches (three with five and one with six genotypes). The simulated data were used to train the
197 model to mimic the conditions of the real data.

198 All sets of the TP and FP simulated data were used to train the models, to increase the strength,
199 and the testing sets of the real data were re-applied to these simulated models. The generation of

200 synthetic genotypes and carrying out the machine learning (training and testing) were applied as
201 described above.

202

203 **Results and discussion**

204 **Identification and evaluation of attributes for the model**

205 A set of 18,057 validated true-positive SNPs and 26,050 false-positive SNPs were collected from
206 the Axiom *Arachis* 58K SNP array [8]. These SNPs had been identified using SWEEP from 21
207 tetraploid peanut genotypes. The true-positive rate achieved was 40%, which was higher than
208 previous efforts in peanut, but still inadequate. All of the mapping data in vcf form was available
209 from the initial SNP calling, which provided the ability to test the hypothesis that machine
210 learning would increase the accuracy of true SNP selection.

211 Seventy percent of the array-validated true- and false-positive SNPs (12,640 and 18,235,
212 respectively) were randomly selected to train the machine learning model. Seventeen different
213 attributes to be used in the model were calculated from sequences surrounding these SNPs (Table
214 1). These attributes were categorized into two groups, *i.e.* sequence and map features. The first
215 machine learning approach used in biological applications was neural networks where it was
216 used for recognizing the transcriptional start sites in *Escherichia coli* [9]. Since that time, it has
217 become one of the most common machine learning approaches. In addition, neural networks
218 have many advantages such as detection of all possible interactions between predictor variables,
219 the ability to detect complex nonlinear relationships between independent and dependent
220 variables, and applicability for different types of data sets [27]. Therefore, we used neural
221 networks to build our first model and to select the most effective attributes.

222

224 **Table 1:** The attributes that were used for building the machine learning models.

Attribute abbreviation	Attribute description	Group
gc	Lowest GC contents of the segment of SNP and seven flanking nucleotides	1
mw	Highest molecular weight of the segment of SNP and seven flanking nucleotides	1
tm	Highest melting temperature of the segment of SNP and seven flanking nucleotides [32]	1
dh	Highest enthalpy (in kilocalories per mole) of the segment of SNP and seven flanking nucleotides [32]	1
ds	Highest entropy (in calories per mole-degrees Kelvin) of the segment of SNP and seven flanking nucleotides [32]	1
dg	Highest free energy (in kilocalories per mole) of the segment of SNP and seven flanking nucleotides [32]	1
dp	The number of reads that cover the SNP	2
n1	The number of reads with the reference nucleotide	2
n2	The number of reads with the alternate nucleotide	2
mq	The Root Mean Square (RMS) mapping quality	2
af	EM estimate of the site allele frequency of the strongest non-reference allele	2
qual	Phred-scaled probability of all samples being homozygous reference	2
no	SNP counts in the segment of SNP and 150 flanking nucleotides	2
lg	The mean of middle phred-scaled data likelihoods of all homozygous reference genotypes	2
n1/n2	The ratio of the number of reads with the reference nucleotide to the alternate one	2
freq1	Frequency of the reference nucleotide in the segment of SNP and 150 flanking nucleotides	1
freq2	Frequency of the alternate nucleotide in the segment of SNP and 150 flanking nucleotides	1

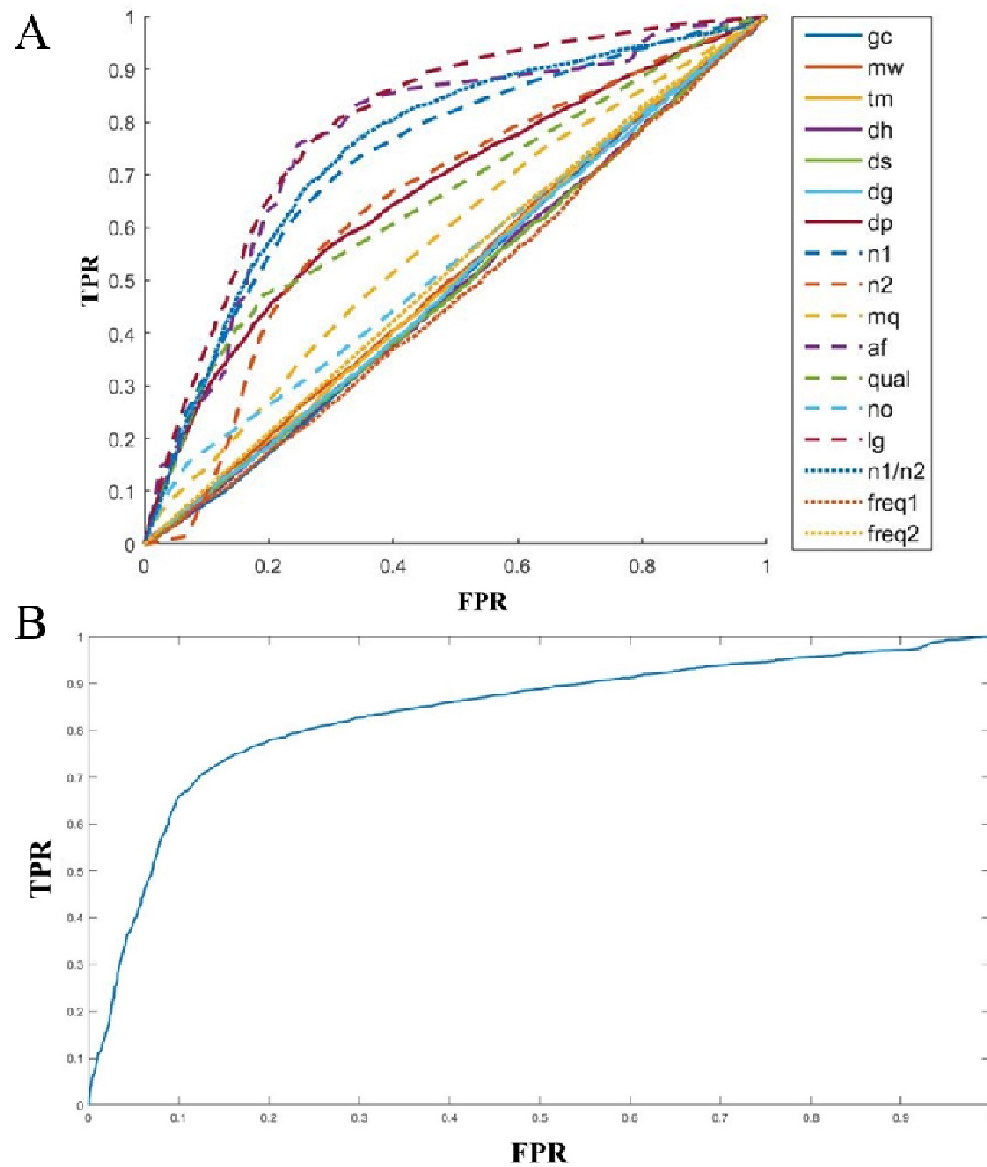
225 *group1: sequence features, group2: map features, bold records: the selected attributes.

226

227 Sequence features previously were used for genome wide *de novo* prediction purposes such as
 228 the prediction of coding regions, and to build a reliable neural network model for SNP calling in
 229 humans [16]. Thermodynamics of nucleic acids are important for diagnostic genetic markers for
 230 diseases, SNP sequencing on a genome-wide scale, designing PCR primers and creating probes

231 for cloning and hybridization experiments [28]. Since thermodynamic parameters give
232 indications for DNA molecule stability, they were used widely to predict the DNA secondary
233 structure [29]. Therefore, we calculated the thermodynamic parameters ΔH , ΔS and ΔG
234 for the SNP locations and flanking seven nucleotides (15 nucleotide segments) and incorporated
235 the highest values from each pair of alternate SNP segments into the model. The higher value is
236 associated with less stable states. Melting temperature (T_m) also was used in the same manner as
237 it shares the primary components of ΔH and ΔS . Molecular weight was included since the
238 change of a nucleotide affects the molecular weight of the DNA molecule. Lower GC contents
239 decrease the stability of the DNA molecule. Therefore, we used the lower GC percentage of the
240 two 15 nucleotide segments (the one with reference nucleotide versus the one with alternative
241 nucleotide). In addition, frequency of the reference and alternate nucleotides in the sequences
242 adjacent to the SNP location were calculated (for the seven nucleotides before and after the SNP
243 location) and included in the model. We hypothesized that higher abundance of a particular
244 nucleotide (reference or alternate nucleotide) would lower the probability of a true SNP.
245 The map features represent the quality of the mapping process and sequence data. Nine mapping
246 parameters were selected to be used in the training model, namely quality features, *i.e.* mq
247 (mapping quality) and qual (SNP quality); read abundance features, *i.e.* dp (depth of reads
248 covering the SNP), af (minor allele frequency), n1 (reads with a reference base), n2 (reads with
249 an alternate base) and n1/n2 (ratio of reference reads to alternate reads). In addition, a probability
250 feature of homozygous reference genotypes (lg) was included. Some of these attributes, *i.e.* dp,
251 n1, n2 and qual, were successfully used to create a neural network classifier for SNP calling for
252 soybean [17]. Therefore, we assumed that these attributes and related features are good
253 candidates for building a classifier in polyploids.

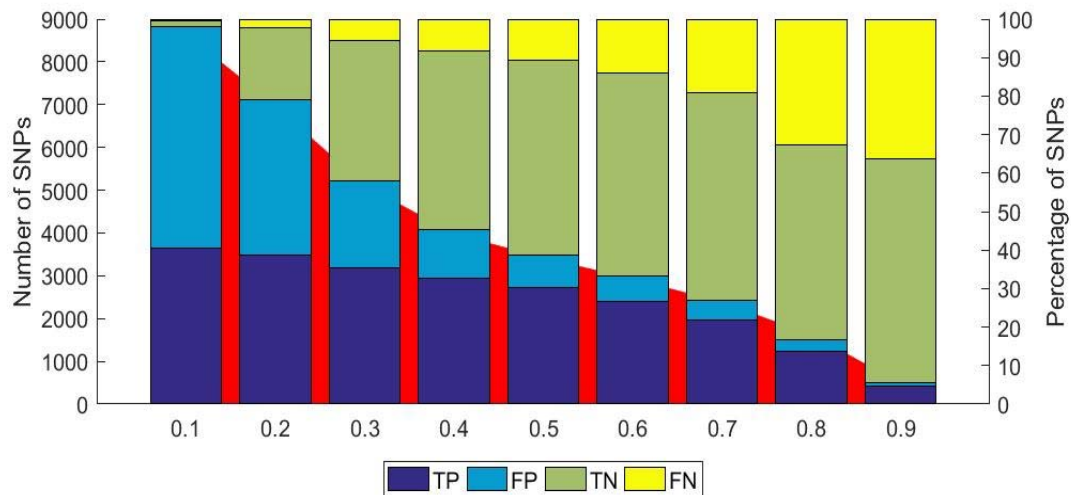
254 Twenty percent of the array-validated SNPs (3,611 true- and 5,210 false-positive SNPs) were
255 used to test the model. Neural network models were applied to every one of the seventeen
256 attributes independently and the relationship between false positives to false negatives was
257 plotted for every model (Fig 2A). Interestingly, eight out of 17 attributes, all eight being map
258 attributes, strongly affected the trainer (Fig 2A). These eight attributes were used for building
259 one model, which showed a high reliability in classification of true- and false-positive SNPs (Fig
260 2B). The neural network score output of the testing data was applied to different neural network
261 score cutoffs, from 0.1 to 0.9 by 0.1 intervals. The confusion matrices (predicted vs. actual)
262 showed a gradual increase in the percentage of true negative (TN; false-positive SNP on the
263 array and not called by SNP-ML) and decrease in the percentage of true positives (TP; true-
264 positive SNP on the array and called by SNP-ML) as the cutoff increased (Fig 3). Increasing the
265 cutoff over 0.5 dramatically decreased the percentage of TP SNPs, and also led to loss of a large
266 number of valid SNPs (FN; true-positive SNP on the array but missed by SNP-ML). On the other
267 hand, decreasing the cutoff below 0.5 increased the occurrence of a large number of false-
268 positive SNPs (FP; false-positive SNP on the array and called by SNP-ML), an undesirable
269 result. The cutoff of 0.5 showed a reasonable trade-off for recovery of the largest possible
270 number of TP while minimizing FP and FN SNPs. These confusion matrices confirmed the
271 efficiency of the eight selected attributes to build a reliable classifier.



288

289 **Fig 2: Receiver Operating Characteristic (ROC) curve of the attributes used in the neural**
290 **network model trainer** A. Independent applications of the 17 attributes. B. The combined
291 application of the selected eight attributes.

292



301

302 **Fig 3: Bar plots representing the confusion matrices of the testing data using different**

303 **cutoffs in neural network model, TP: True Positive (validated as a true SNP on the array and**

304 **called by SNP-ML), FP: False Positive (not a true SNP according to array data but called by**

305 **SNP-ML), TN: True Negative (not a true SNP according to array data and not called by SNP-**

306 **ML), FN: False Negative (validated as a true SNP on the array and not called by SNP-ML). The**

307 **red area shows the number of SNPs which are recognized by the model. The left Y scale presents**

308 **the number of SNPs within every class and the right Y scale presents the percentage SNPs of**

309 **every class to the total SNPs.**

310

311 **Comparison among different supervised machine learning models**

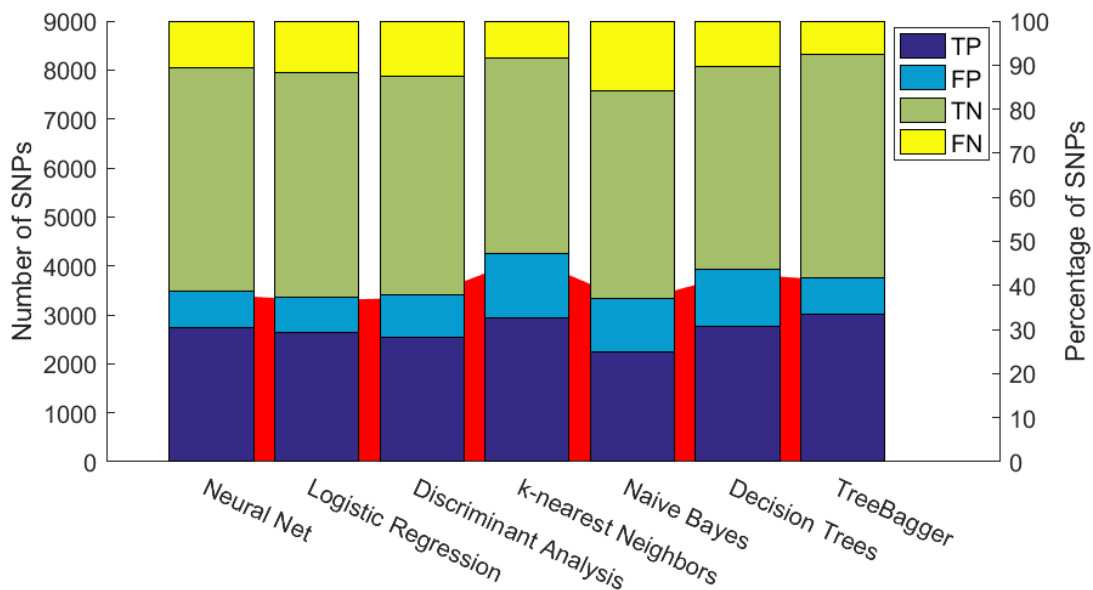
312 **using the selected attributes**

313 The training data set was used to build training models by applying different supervised

314 algorithms, *i.e.* Logistic Regression (LR), Discriminant Analysis (DA), K-nearest Neighbors

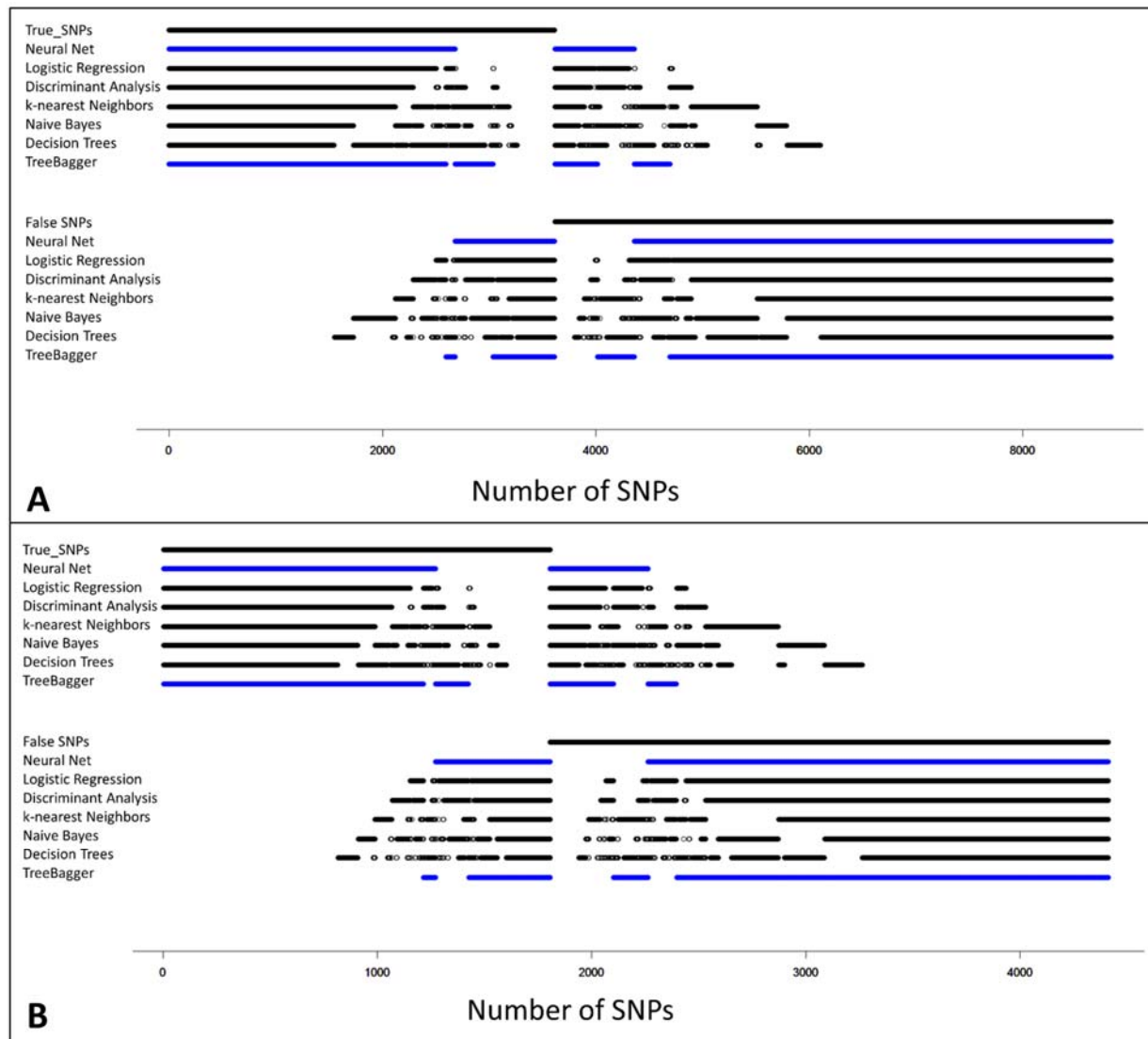
315 (KN), Naïve Bayes (NB), Decision Trees (DT) and TreeBagger (TB). The testing data set was

316 applied for these trainers along with the neural network output of 0.5 cutoff (Fig 4). All models
317 showed 60 to 80% true-positive rates relative to the number of SNPs extracted by a respective
318 model, or between 25.0 to 33.4% of the total number of SNPs in the testing set. KN showed the
319 highest false-positive rate and the neural network gave the lowest rate. Conversely, NB showed
320 the lowest true-positive rate while TB produced the best rate. However, both TB and neural
321 network showed the best trade-off between the two rates (Fig 5 and S3 Table). Therefore, we
322 combined these two models to increase accuracy. TB was first described around 50 years after
323 the first neural network approach was proposed [30]. It reduces the variance among observations
324 and avoids overfitting, which are two limitations for neural network, thus it works as a
325 complementary model to neural network to overcome its drawbacks.



326
327 **Fig 4: Bar plots represented the confusion matrices of the testing data using supervised**
328 **machine learning algorithms**, TP: True Positive, FP: False Positive, TN: True Negative, FN:
329 False Negative. The red area shows the number of SNPs which are recognized by the model

330 (positive total). The left Y scale presents the number of SNPs within every class and the right Y
331 scale presents the percentage SNPs of every class to the total SNPs.



332
333 **Fig 5: A dot plot of the trade-off combination between the different machine learning**
334 **algorithms on the testing set (A) and validation set (B), every dot shows a single true or false**
335 **SNP (upper lines) and corresponding dots shows if this SNP was called using different machine**
336 **learning algorithms.**

337

338 To further test the model, the remaining 10% of the original data set, 1,806 validated true-
339 positive and 2,605 false-positive SNPs, was used as a validation set. This data set was applied to
340 the combined NN + TB model. A total of 1,510 SNPs was extracted by the model and 1,214 of
341 those were true-positive SNPs. Therefore, the combined model efficiency increased to 80%
342 versus 73% (1,271 out of 1,792) and 76% (1,369 out of 1,797) of using only neural network or
343 TB, respectively. However, 33% of validated SNPs were lost through the prediction process
344 using the combined model. The validation set of SNPs called using SWEEP and identified as
345 true or false using the chip, is provided in S4 Table, along with detection state using only neural
346 network, only TB, or the combined model.

347

348 **Model validation on Axiom *Arachis2* 48K SNP array:**

349 To validate the model for further real world analysis, 28,218 markers were selected to be
350 included in a newly designed SNP array. A set of 133 tetraploid peanut genotypes and lines was
351 genotyped on the chip. Polymorphisms were found in 21,112 SNPs among the tested genotypes,
352 which revealed an accuracy of 75%. This represents the largest validation experiment to test a
353 bioinformatics method developed to identify SNPs in polyploid species and provides the highest
354 true positive validation rate reported in polyploids.

355

356 **Building models for RNA-seq**

357 Unlike the re-sequencing data, RNA-seq provides data that measure gene expression and can
358 produce a very high depth at specific loci [31]. The values of the attributes are different from the
359 genomic re-sequencing data. For this reason, a specific model was built for RNA-seq data using

360 sequence from nine tetraploid peanut genotypes. The analysis of this data set with SWEEP
361 produced 3,525 SNP-chip overlapped SNPs, 2,143 true and 1,382 false SNPs.
362 Eighty percent of the array-validated SNPs were used for training the models, 1,714 true- and
363 1,104 false-positive SNPs, and the remaining 20% of SNPs were used as a testing set, 429 true-
364 and 278 false-positive SNPs. Two models were built, *i.e.* neural network and TB, and the scored
365 results were combined. The combined model extracted 371 SNPs (using the cutoff of 0.5 for
366 neural network model). Of the SNPs extracted, 328 of them were true SNPs. The accuracy of
367 true SNP discovery was raised to 88%. However, 101 SNPs were lost (~24%).

368

369 **Application in other polyploids**

370 In the absence of validation SNP sets for allotetraploid cotton (*Gossypium hirsutum*),
371 allohexaploid wheat (*Triticum aestivum*), or allo-octoploid strawberry (*Fragaria x ananassa*), a
372 simulation experiment was carried out to generate allelic variation. Genome sequence for each
373 species was downloaded and five genotypes were simulated in one of the subgenomes while
374 keeping the other subgenomes constant. The locations of true-positive SNPs thus were known
375 due to the *in silico* mutation of the sequence and any other SNPs called by the program were
376 considered false-positive. Because only one subgenome was mutated to derive the genotypes, all
377 true SNPs were subgenome-specific. The true and false SNP calls were randomly categorized as
378 training set (70%) and testing set (30%). The training set was used to train neural network
379 models which were then used to select SNPs from the testing set. Simulations for all three
380 species achieved accuracy of greater than 99% at five different sequence coverage depths (10X,
381 20X, 30X, 40X and 50X) (Table 2 and S5 Table). A peanut simulation was also included for
382 comparison.

383

384 **Table 2:** The SNP-ML calling accuracy on different polyploid simulated data.

Depth	Cotton	Wheat	Strawberry	Peanut
	True positive %	True positive %	True positive %	True positive %
10X	100.00	99.64	99.72	99.85
20X	99.85	99.65	99.49	99.96
30X	100.00	99.88	99.73	99.96
40X	99.93	100.00	99.57	99.92
50X	99.96	99.88	98.15	99.96

385

386

387 **Application of simulation trained models on real data**

388 Next, it was tested if models trained with simulated genotypes could achieve high accuracy in
389 predicting true SNPs from real data, using the validation SNP sets available for peanut. Models
390 that were trained in the simulations discussed above were used to select SNPs from the 21
391 genotypes of peanut (S6 Table). Each run of SNP-ML was performed three times to show
392 variation between runs. For peanut, the models trained with simulated data were able to select
393 true SNPs with accuracy on average of 78%. This result strongly suggests that this method can
394 be used effectively in species where there are no large validation sets to train the models, but
395 some reference sequence is available. This result, combined with the simulation results and
396 results on real peanut data led us to construct a novel tool, SNP-ML, to carry out these analyses.
397 The tool is designed to be highly flexible so that it can be used effectively in the broadest sense.

398

399 **SNP-MLer**

400 All of the models discussed in this work are provided in the SNP-ML subdirectory “/db”. They
401 include the peanut WGS and RNA-seq-trained models from real data and the models trained
402 from simulated cotton, wheat, and strawberry data. The binary executable tool, SNP-MLer, will
403 take two files as input, a vcf file containing true-positive SNP calls and a vcf file containing
404 false-positive SNP calls. By default, SNP-MLer will train a neural network model using these
405 sets of SNP calls and the eight parameters used in this work. The user has the ability with ‘-skip’
406 to not use one or more (up to seven) of the parameters if they wish or use ‘-custom’ to specify
407 selected parameters in a comma-delimited sequence. Additionally, the user can use ‘-m’ to train
408 a treebagging model as well. Most importantly, the user can add customized parameters to
409 include in the model training by invoking ‘-addnew1’ and ‘-addnew2’. These options take csv
410 files that include one or more new parameter lists for the true-positive SNP calls (-addnew1) and
411 the false-positive calls (-addnew2). The user also needs to add the prefix name for the new model
412 using ‘-o’.

413

414 **SNP-ML**

415 If the user has trained new models using SNP-MLer or will use the models trained in this study,
416 all models are located in the ‘\db’ folder for use with SNP-ML. SNP-ML is the tool that will take
417 as input (-i) a vcf file of the SNP calls of interest. It is recommended to first use SWEEP to filter
418 most of the false-positive SNP calls, but it is not required. The name of model to be used for
419 SNP selection (-iM) should also be given as input to SNP-ML. The program contains currently
420 two models, “peanut_DNA” for use with WGS data, and “peanut_RNA” for use with RNA-seq
421 data. Any new models trained with SNP-MLer by the user will be included in this folder as well.

422 Users can submit any newly trained models to be included in new versions of SNP-ML by
423 emailing the author. SNP-ML has similar options as SNP-MLer to skip (-skip) or customize (-
424 custom) parameter sets for SNP prediction, and to invoke the treebagging model (-m) or add new
425 parameters (-addnew; for custom trained models). An additional option (-c) allows the user to
426 increase or decrease the stringency of true-positive selection from the default of 0.5. As
427 discussed above and in Fig 3, increasing this cutoff will decrease false positives (decreasing
428 selection of false SNPs) while increasing false negatives (limiting recovery of validated true
429 SNPs) while decreasing the cutoff has the opposite effect.

430 The program is freely available for public use under MIT license and can be downloaded from
431 <https://github.com/w-korani/SNP-ML>. A help file containing detailed information about using
432 the program can be accessed by typing SNP-ML -h.

433

434 **Conclusions**

435 We introduce a highly reliable method for calling SNPs for polyploid species using machine
436 learning. To have a good classifier, the most effective attributes should be determined. Many
437 attributes were tested and the best were selected for creating the model. In addition, different
438 supervised machine learning algorithms were tested and the best ones for the data sets, neural
439 network and bagging, were combined. We built and tested our method on peanut, an
440 allotetraploid for which identifying true SNPs has been difficult. In addition, a 48K SNP array
441 was designed using SNP-ML was created and showed high accuracy. The method was then used
442 on simulated data from three other allopolyploids with different ploidy levels and achieved high
443 accuracy. Most importantly, we showed that simulated data can be used to train models that
444 achieve similar accuracy in selecting true SNPs using real data as do models trained with real

445 data. The implication is that for species where there are no large validation sets available, our
446 method can still be used to efficiently select true SNPs. With this important result in mind, SNP-
447 MLer was developed; a tool that will train new neural network or treebagging models with user
448 inputted data. Subsequently, SNP-ML can be used with newly trained models or included peanut
449 models to select true SNPs for two different data set types, re-sequencing and RNA-seq. The
450 flexibility and functionality of these tools allow the user a completely customizable experience,
451 giving the ability to use the power of machine learning to researchers of all expertise levels.

452

453 **Acknowledgments**

454 This work was supported by the Peanut Foundation, the Agriculture and Food Research Initiative
455 competitive grant 2012-85117-19435 of the USDA National Institute of Food and Agriculture
456 and the Feed the Future Innovation Lab for Collaborative Research on Peanut Productivity and
457 Mycotoxin Control (Peanut and Mycotoxin Innovation Lab), supported by funding from the
458 United States Agency for International Development (USAID).

459

460 **Availability of data and materials**

461 SNP-ML, SNP-MLer and extendable database are freely available for public use under MIT
462 license and can be downloaded from <https://github.com/w-korani/SNP-ML>

463

464 **Authors' contributions**

465 WK collected sequence attributes, applied the models, programmed SNP-ML/SNP-MLer and
466 drafted the manuscript. JC collected map attributes, designed SNP array, edited and revised the

467 manuscript. CY collected the validation data of the new SNP array. PO conceived and supervised
468 the project, secured funding, and revised and submitted the manuscript.

469

470 **Competing interests**

471 The authors declare that they have no competing interests.

472

473 **Supporting information**

474 **S1 Table: Affymetrix SNP array structure.**

475 **S2 Table: The specific arguments of the different machine learning models.**

476 **S3 Table: The efficiency of the different machine learning models.**

477 **S4 Table: The validation set of SNPs called using SWEEP, neural network, TB, and the
478 combined model.**

479 **S5 Table: The efficiency of SNP-ML SNP calling for simulated data of cotton, wheat,
480 strawberry and peanut.**

481 **S6 Table: The efficiency of using simulated data model to call SNPs from real data in
482 peanut.**

483 **S1 File: The source code of C++ and python of SNP-ML (SNP-ML_source_code.zip).**

484

485 **References**

486

487 1. Bertoli DJ, Cannon SB, Froenicke L, Huang G, Farmer AD, Cannon EK, et al. The
488 genome sequences of *Arachis duranensis* and *Arachis ipaensis*, the diploid ancestors of

- 489 cultivated peanut. *Nature Genetics* 2016;48(4):438-446.
- 490 2. Kochert G, Halward T, Branch WD, Simpson CE: RFLP variability in peanut (*Arachis*
491 *hypogaea* L.) cultivars and wild species. *TAG Theoretical and applied genetics*
492 *Theoretische und Angewandte Genetik* 1991;81(5):565-570.
- 493 3. Gupta SK, Baek J, Carrasquilla-Garcia N, Penmetsa RV: Genome-wide polymorphism
494 detection in peanut using next-generation restriction-site-associated DNA (RAD)
495 sequencing. *Molecular Breeding* 2015;35(7):145.
- 496 4. Zhou X, Xia Y, Ren X, Chen Y, Huang L, Huang S, Liao B, Lei Y, Yan L, Jiang H:
497 Construction of a SNP-based genetic linkage map in cultivated peanut based on large
498 scale marker development using next-generation double-digest restriction-site-associated
499 DNA sequencing (ddRADseq). *BMC Genomics* 2014;15:351.
- 500 5. Khera P, Upadhyaya HD, Pandey MK, Roorkiwal M, Sriswathi M, Janila P, et al. Single
501 Nucleotide Polymorphism-based genetic diversity in the reference set of peanut (*Arachis*
502 *spp.*) by developing and applying cost-effective kompetitive allele specific polymerase
503 chain reaction genotyping assays. *The Plant Genome* 2013;6(3): 1-11.
- 504 6. Peng Z, Gallo M, Tillman BL, Rowland D, Wang J: Molecular marker development from
505 transcript sequences and germplasm evaluation for cultivated peanut (*Arachis hypogaea*
506 L.). *Molecular genetics and genomics : MGG* 2016, 291(1):363-381.
- 507 7. Clevenger JP, Ozias-Akins P: SWEEP: A tool for filtering high-quality SNPs in
508 polyploid crops. *G3 (Bethesda, Md)* 2015;5(9):1797-1803.
- 509 8. Clevenger J, Chu Y, Chavarro C, Agarwal G, Bertioli DJ, Leal-Bertioli SC, et al.
510 Genome-wide SNP genotyping resolves signatures of selection and tetrasomic
511 recombination in peanut. *Molecular Plant* 2017;10(2):309-322.

- 512 9. Tarca AL, Carey VJ, Chen XW, Romero R, Draghici S: Machine learning and its
513 applications to biology. *PLoS Computational Biology* 2007;3(6):e116.
- 514 10. Bostan B, Greiner R, Szafron D, Lu P: Predicting homologous signaling pathways using
515 machine learning. *Bioinformatics (Oxford, England)* 2009;25(22):2913-2920.
- 516 11. Lingner T, Kataya AR, Antonicelli GE, Benichou A, Nilssen K, Chen XY, Siemsen T,
517 Morgenstern B, Meinicke P, Reumann S: Identification of novel plant peroxisomal
518 targeting signals by a combination of machine learning methods and in vivo subcellular
519 targeting analyses. *The Plant Cell* 2011;23(4):1556-1572.
- 520 12. Swan AL, Mobasher A, Allaway D, Liddell S, Bacardit J: Application of machine
521 learning to proteomics data: classification and biomarker identification in postgenomics
522 biology. *Omics : A Journal of Integrative Biology* 2013;17(12):595-610.
- 523 13. Jowkar GH, Mansoori EG: Perceptron ensemble of graph-based positive-unlabeled
524 learning for disease gene identification. *Computational Biology and Chemistry*
525 2016;64:263-270.
- 526 14. Roche-Lima A: Implementation and comparison of kernel-based learning methods to
527 predict metabolic networks. *Network Modeling and Analysis in Health Informatics and*
528 *Bioinformatics* 2016;5:26.
- 529 15. Melo R, Fieldhouse R, Melo A, Correia JD, Cordeiro MN, Gumus ZH, Costa J, Bonvin
530 AM, Moreira IS: A Machine learning approach for hot-spot detection at protein-protein
531 interfaces. *International Journal of Molecular Sciences* 2016;17(8):1215.
- 532 16. Unneberg P, Stromberg M, Sterky F: SNP discovery using advanced algorithms and
533 neural networks. *Bioinformatics (Oxford, England)* 2005;21(10):2528-2530.
- 534 17. Matukumalli LK, Grefenstette JJ, Hyten DL, Choi IY, Cregan PB, Van Tassell CP:

- 535 Application of machine learning in SNP discovery. *BMC Bioinformatics* 2006;7:4.
- 536 18. Clevenger J, Chavarro C, Pearl SA, Ozias-Akins P, Jackson SA: Single nucleotide
537 polymorphism identification in polyploids: A review, example, and recommendations.
538 *Molecular Plant* 2015;8(6):831-846.
- 539 19. Clevenger J, Chu Y, Scheffler B, Ozias-Akins P: A developmental transcriptome map for
540 allotetraploid *Arachis hypogaea*. *Frontiers in Plant Science* 2016;7:1446.
- 541 20. Clevenger J, Marasigan K, Liakos V, Sobolev V, Vellidis G, Holbrook C, Ozias-Akins P:
542 RNA sequencing of contaminated seeds reveals the state of the seed permissive for pre-
543 harvest aflatoxin contamination and points to a potential susceptibility factor. *Toxins*
544 2016;8(11): 317.
- 545 21. Li F, Fan G, Lu C, Xiao G, Zou C, Kohel RJ, et al. Genome sequence of cultivated
546 Upland cotton (*Gossypium hirsutum* TM-1) provides insights into genome evolution.
547 *Nature Biotechnology* 2015;33(5):524-530.
- 548 22. Choulet F, Alberti A, Theil S, Glover N, Barbe V, Daron J, et al. Structural and
549 functional partitioning of bread wheat chromosome 3B. *Science (New York, NY)*
550 2014;345(6194):1249721.
- 551 23. Shulaev V, Sargent DJ, Crowhurst RN, Mockler TC, Folkerts O, Delcher AL, et al. The
552 genome of woodland strawberry (*Fragaria vesca*). *Nature Genetics* 2011;43(2):109-116.
- 553 24. Hirakawa H, Shirasawa K, Kosugi S, Tashiro K, Nakayama S, Yamada M, et al.
554 Dissection of the octoploid strawberry genome by deep sequencing of the genomes of
555 *Fragaria* species. *DNA Research : An International Journal for Rapid Publication of*
556 *Reports on Genes and Genomes* 2014;21(2):169-181.
- 557 25. Huang W, Li L, Myers JR, Marth GT: ART: a next-generation sequencing read

- 558 simulator. *Bioinformatics (Oxford, England)* 2012;28(4):593-594.
- 559 26. Li H, Durbin R: Fast and accurate long-read alignment with Burrows-Wheeler transform.
560 *Bioinformatics (Oxford, England)* 2010;26(5):589-595.
- 561 27. Tu JV: Advantages and disadvantages of using artificial neural networks versus logistic
562 regression for predicting medical outcomes. *Journal of Clinical Epidemiology*
563 1996;49(12n):1225-1231.
- 564 28. Wu P, Nakano S, Sugimoto N: Temperature dependence of thermodynamic properties for
565 DNA/DNA and RNA/DNA duplex formation. *European Journal of Biochemistry*
566 2002;269(12):2821-2830.
- 567 29. SantaLucia J, Jr., Hicks D: The thermodynamics of DNA structural motifs. *Annual*
568 *Review of Biophysics and Biomolecular Structure* 2004;33:415-440.
- 569 30. Breiman L: Bagging Predictors. *Machine Learning* 1996;24(2):123-140.
- 570 31. Lopez-Maestre H, Brinza L, Marchet C, Kielbassa J, Bastien S, Boutigny M, et al. SNP
571 calling from RNA-seq data without a reference genome: identification, quantification,
572 differential analysis and impact on the protein sequence. *Nucleic Acids Research*
573 2016;44(19):e148.
- 574 32. Sugimoto N, Nakano S, Yoneyama M, Honda K: Improved thermodynamic parameters
575 and helix initiation factor to predict stability of DNA duplexes. *Nucleic Acids Research*
576 1996;24(22):4501-4505.