1    **ADMIXPIPE: Population analyses in ADMIXTURE for non-model**

2    **organisms**

3    Steven M. Mussmann[1], Marlis R. Douglas[1], Tyler K. Chafin[1], and Michael E. Douglas[1]

4    [1]Department of Biological Sciences, University of Arkansas, Fayetteville, AR 72701

5
6    *Corresponding author and person to whom reprint requests should be addressed:*

7    Steven M. Mussmann

8    Department of Biological Sciences

9    University of Arkansas

10   Fayetteville, AR 72701

11   Voice: 479-575-5529

12   e-mail: smussmann@gmail.com

13

14

15   **Abstract**

16   **Background:** Research on the molecular ecology of non-model organisms, while

17   previously constrained, has now been greatly facilitated by the advent of reduced-

18   representation sequencing protocols. However, tools that allow these large datasets to

19   be efficiently parsed are often lacking, or if indeed available, then limited by the

20   necessity of a comparable reference genome as an adjunct. This, of course, can be

21   difficult when working with non-model organisms. Fortunately, pipelines are currently

22   available that avoid this prerequisite, thus allowing data to be *a priori* parsed. An oft-

23   used molecular ecology program (i.e., STRUCTURE), for example, is facilitated by such

24   pipelines, yet they are surprisingly absent for a second program that is similarly popular

25   and computationally more efficient (i.e., ADMIXTURE). The two programs differ in that

26   ADMIXTURE employs a maximum-likelihood framework whereas STRUCTURE uses a

27   Bayesian approach, yet both produce similar results. Given these issues, there is an

28   overriding (and recognized) need among researchers in molecular ecology for

29   bioinformatic software that will not only condense output from replicated ADMIXTURE

30   runs, but also infer from these data the optimal number of population clusters (K).

31

32   **Results:** Here we provide such a program (i.e., ADMIXPIPE) that (a) filters SNPs to allow

33   the delineation of population structure in ADMIXTURE, then (b) parses the output for

34   summarization and graphical representation via CLUMPAK. Our benchmarks effectively

35   demonstrate how efficient the pipeline is for processing large, non-model datasets

36   generated via double digest restriction-site associated DNA sequencing (ddRAD).

37    Outputs not only parallel those from STRUCTURE, but also visualize the variation among

38    individual ADMIXTURE runs, so as to facilitate selection of the most appropriate *K*-value.

39

40    **Conclusions:** ADMIXPIPE successfully integrates ADMIXTURE analysis with popular

41    variant call format (VCF) filtering software to yield file types readily analyzed by

42    CLUMPAK. Large population genomic datasets derived from non-model organisms are

43    efficiently analyzed via the parallel-processing capabilities of ADMIXTURE. ADMIXPIPE is

44    distributed under the GNU Public License and freely available for Mac OSX and Linux

45    platforms at: https://github.com/stevemussmann/admixturePipeline.

46

47    **Keywords:** RADseq, SNP analysis, Population Genomics, Population Structure,

48    ADMIXTURE analysis

49

50    **Background**

51        Advances in genomics during the past decade have accelerated research in

52    molecular ecology by significantly increasing the capacity of researchers to generate

53    vast quantities of data at relatively low cost. These advances largely represent the

54    development of reduced representation genomic libraries [1–3] that identify tens of

55    thousands of SNPs for non-model organisms, coupled with high-throughput sequencing

56    methods that efficiently genotype fewer SNPs for thousands of individuals [4]. However,

57    data generation, particularly through these novel and affordable marker-discovery

58    methods [5], has greatly outpaced analytical capabilities, and especially so with regard

59    to evolutionary and conservation genomics.

60    Here, technological advances have also precipitated a suite of new analytical

61    issues. The thousands of SNPs generated in a typical RADseq project may exhibit

62    biases that impact the inferences that can be drawn from these data [6], and which

63    necessitate careful data filtration to avoid [7]. Yet, the manner by which data are filtered

64    represents a double-edged sword. While it is certainly mandated (as above), the

65    procedures involved must be carefully evaluated in the context of each study, in that

66    downstream analyses can be seriously impacted [8, 9], to include the derivation of

67    population structure [10].

68    For example, the analysis of multilocus codominant markers in evaluation of

69    population structure is frequently accomplished using methods that make no *a priori*

70    assumptions about underlying structure. One of the most popular options to accomplish

71    this is the program STRUCTURE [11–13]. However, it necessitates that users test specific

72    clustering values (K), and conduct *post hoc* evaluation of these results so as to

73    determine an optimal K [14]. This typically involves searching a complicated parameter

74    space using heuristic algorithms for Maximum Likelihood (ML) and Bayesian (BA)

75    methods that, in turn, provide additional complications such as a tendency to sample

76    local optima [15].

77    A common strategy to mitigate this is to sample multiple independent replicates

78    at each K, using different random number seeds for initialization. These results are

79    subsequently collated and evaluated to assess confidence that global rather than local

80    optima have indeed been sampled. Clearly, this procedure must be automated so as to

81    alleviate the onerous task of testing multiple replicates across a range of K-values.

82    Pipelines to do so are available for STRUCTURE, and have been deployed on high-

4

83    performance computing systems via integrated parallelization (STRAUTO,

84    PARALLELSTRUCTURE) [16, 17]. Multiple programs have likewise been developed for

85    handling STRUCTURE output (i.e., CLUMPP, DISTRUCT) [18, 19]; and pipelines constructed

86    to assess the most appropriate K-values (i.e., STRUCTUREHARVESTER, CLUMPAK) [20, 21].

87         Despite the considerable focus on STRUCTURE, few such resources have been

88    developed for a popular alternative program (i.e., ADMIXTURE [22]). The Web of Science

89    indexing service indicates that (as of January, 2020) it has been cited 1,812 times since

90    initial publication (September, 2009). This includes 479 (26.4%) in 2019 alone. Despite

91    its popularity, it has just a single option that promotes the program as part of a pipeline

92    (i.e., SNIPLAY3 [23]), and unfortunately it requires a reference genome as an adjunct for

93    its application. Needless to say, its applicability is thus limited for those laboratories that

94    employ non-model organisms as study species.

95         Options for post-processing of ADMIXTURE results are similarly deficit. However,

96    one positive is that CLUMPAK is flexible enough in its implementation to allow for the

97    incorporation of ADMIXTURE output, as well as that of STRUCTURE. Furthermore, no

98    available software currently exists that can summarize the variation in cross-validation

99    (CV) values, the preferred method for selecting an optimal K-value in ADMIXTURE [24].

100        Here we describe a novel software package that integrates ADMIXTURE as the

101    primary component of an analytical pipeline that also incorporates the filtering of data as

102    part of its procedure. This, in turn, provides a high-throughput capability that not only

103    generates input for ADMIXTURE but also evaluates the impact of filtering on population

104    structure. ADMIXPIPE also automates the process of testing multiple K-values, conducts

105    replicates at each K, and automatically formats these results as input for the CLUMPAK

5

106    pipeline. Optional post-processing scripts are also provided as a part of the toolkit to

107    process CLUMPAK output, and to visualize the variability among CV values for

108    independent ADMIXTURE runs. Sections of the pipeline are specifically designed for use

109    with non-model organisms, as these are increasingly common study species in

110    evolutionary and conservation genomic investigations.

111

112    **Implementation**

113       ADMIXPIPE requires two input files: a population map and a standard VCF file.

114    The population map is a tab-delimited text file with each row representing a sample

115    name/ population pair. The VCF file is filtered according to user-specified command line

116    options that include the following: minor allele frequency (MAF) filter, biallelic filtering,

117    data thinning measured in basepairs (bp), and missing data filtering (for both individuals

118    and loci). Users may also remove specific samples from their analysis by specifying a

119    file of sample names to be ignored. All filtering and the initial conversion to PLINK

120    (PED/MAP) format [25] is handled by VCFTOOLS [26].

121       ADMIXPIPE is intended for use with non-model organisms that lack genomic

122    reference data, and given this, additional conversions are required before the PLINK-

123    formatted files will be accepted by ADMIXTURE. Popular software packages for *de novo*

124    assembly of RADseq data, such as pyRAD [27, 28] produce VCF files with each locus

125    as an individual "chromosome." This, in turn, yields output that exceeds the number of

126    chromosomes in those model organisms for which PLINK was originally designed. The

127    initial MAP file is therefore modified to append a letter at the start of each "chromosome"

6

128  number. PLINK is then executed using the "–allow-extra-chr 0" option that treats loci as

129  unplaced contigs in the final PED/ MAP files submitted to ADMIXTURE.

130      The main element of the pipeline executes ADMIXTURE on the filtered data. The

131  assessment of multiple K values and multiple replicates is automated based upon user-

132  specified command line input. The user defines minimum and maximum K values to be

133  tested, in addition to the number of replicates for each K. Users may also specify the

134  number of processor cores to be utilized by ADMIXTURE, and the cross-validation number

135  which is utilized in determining optimal K. The final outputs of the pipeline include a

136  compressed results file and a population file that are submitted as-is to CLUMPAK for

137  processing and visualization.

138      The pipeline also offers two accessory scripts for processing of CLUMPAK output.

139  The first (i.e., distructRerun.py) compiles the major clusters identified by CLUMPAK,

140  generates DISTRUCT input files, executes DISTRUCT, and extracts CV-values for all major

141  cluster runs. The second script (i.e., cvSum.py) plots the boxplots of CV-values against

142  each K so as to summarize the distribution of CV-values for multiple ADMIXTURE runs.

143  This permits the user to make an informed decision on the optimal K by graphing how

144  these values vary according to independent ADMIXTURE runs.

145      ADMIXTURE is the only component of the pipeline that is natively parallelized.

146  Therefore, we performed benchmarking to confirm that processing steps did not

147  significantly increase runtime relative to that expected for ADMIXTURE. Data for

148  benchmarking were selected from a recently published paper that utilized ADMIXPIPE for

149  data processing [29]. The test data contained 343 individuals and 61,910 SNPs. Four

150  data thinning intervals (i.e.,1, 25, 50, and 100) yielded SNP datasets of variable size for

151    performance testing. All filtering intervals were repeated with variable numbers of

152    processor cores (i.e.,1, 2, 4, 8, and 16). Sixteen replicates of ADMIXTURE were first

153    conducted for each K=1-8 at each combination of thinning interval and number of

154    processor cores, for a total of 20 executions of the pipeline. The process was then

155    repeated for each K=9-16, for an additional 20 runs of the pipeline. Memory profiling

156    was conducted through the python3 'mprof' package at K=16, with a thinning interval of

157    1 as a final test of performance. All tests were completed on a computer equipped with

158    dual Intel Xeon E5-4627 3.30GHz processors, 256GB RAM, and with a 64-bit Linux

159    environment.

160

161    **Results**

162        The filtering intervals resulted in datasets containing 61,910 (interval = 1bp),

163    25,851 (interval = 25bp), 19,140 (interval = 50bp), and 12,527 SNPs (interval = 100bp).

164    Runtime increased linearly with the number of SNPs analyzed, regardless of the

165    number of processors utilized (Figure 1: $R^2$ = 0.975, df = 58). For example, increasing

166    the number of SNPs from 12,527 to 61,910 (494% increase) produced an average

167    increase of 519% in ADMIXPIPE runtime (SD = 41.6%).

168        Little change was observed in response to increasing the numbers of processor

169    cores from K=1-8 (Figure 2A). A slight decrease in performance was observed in some

170    cases, particularly for the largest dataset. This trend changed at higher K-values, as

171    substantial gains were observed at K=9-16 when processors were increased from 1 to

172    4. The most dramatic performance increase was observed for the 61,910 SNP dataset,

173    where a 24.3-hour (34.5%) reduction in computation time occurred when processors

174  increased from 1 to 4. However, only marginal improvements occurred when processors

175  were increased from 1 to 8 (24.5 hours; 34.7%) or 16 (26.2 hours; 37.7%).

176      Profiling also revealed efficient and consistent memory usage. The greatest

177  memory spike occurred during the initial filtering steps, when peak memory usage

178  reached approximately 120 MB. All subsequent usage held constant at ~60 MB as

179  ADMIXTURE runs progressed.

180

181  **Discussion**

182      The performance of ADMIXPIPE improved with the number of processor cores

183  utilized at higher K-values. However, it did not scale at the rate suggested in the original

184  ADMIXTURE publication. We have been unable to attribute the difference in performance

185  to any inherent property of our pipeline. Filtering and file conversion steps at the

186  initiation of ADMIXPIPE are non-parallel sections. Reported times for completion of these

187  steps were approximately constant across runs, with the maximum reported time being

188  eight seconds. This indicates that ADMIXTURE itself is the main driver of performance, as

189  it comprises the vast majority of system calls made by ADMIXPIPE.

190      The original performance increase documented for ADMIXTURE was 392% at K=3,

191  utilizing four processor cores [24]. Unfortunately, we could not replicate this result with

192  our benchmarking data [29], or the original test data (i.e.,324 samples; 13,928 SNPs)

193  [24] which parallels our own. When we attempted to replicate the original benchmark

194  scores, we found that it also failed to scale as the number of processor cores increased

195  (1-core $\bar{x}$ = 40.63 seconds, $\sigma$ = 0.90; 4-core $\bar{x}$ = 47.46 seconds, $\sigma$ = 4.71). Furthermore,

196  we verified that performance did increase with up to four processor cores at higher K

197    values (K≥9). We therefore view this as 'expected behavior' for ADMIXTURE, and find no

198    reason to believe that ADMIXPIPE has negatively impacted the performance of any

199    individual program.

200         Results of ADMIXPIPE were similar to those found by STRUCTURE for the test

201    dataset, as evaluated in an earlier publication [29], and gauged for the optimum K=8.

202    This is not surprising, given that ADMIXTURE implements the same likelihood model as

203    does STRUCTURE [22]. However, minor differences have previously been noted for both

204    programs in the assignment probabilities [29, 30].

205         Memory usage was efficient and constant, with the greatest increase occurring

206    when PLINK was executed. Thus, users will be able to execute ADMIXPIPE on their

207    desktop machines for datasets sized similarly to that evaluated herein. Performance

208    gains were minimal with >4 processors, and this (again) reduces the necessity for

209    supercomputer access, since desktop computers with ≥4 processor cores are now

210    commonplace. However, given the built-in parallelization capabilities of ADMIXPIPE, its

211    application on dedicated high-performance computing clusters will be beneficial when

212    runtime considerations are necessary, such as when evaluating K>8, or SNPs≥20,000.

213         Finally, our integration of common SNP filtering options provides the flexibility to

214    quickly filter data and assess the manner by which various filtering decisions impact

215    results. A byproduct of the filtering process is the production of a STRUCTURE-formatted

216    file that will facilitate comparisons with other popular algorithms that assess population

217    structure. These options are important tools, particularly given recent documentation

218    regarding of the impacts of filtering on downstream analyses. We thus suggest that

10

219 users implement existing recommendations on filtering RAD data, and use these to

220 investigate subsequent impacts on their own data [7–10].

221

222 **Conclusions**

223   Benchmarking has demonstrated that the benefits of ADMIXPIPE (e.g., low

224 memory usage and performance scaling with low numbers of processor cores at high K-

225 values) will prove useful for researchers with limited access to advanced computing

226 resources. ADMIXPIPE also allows the effects of common filtering options to be assessed

227 on population structure of study species by coupling this process with the determination

228 of population structure. Integration with CLUMPAK, and our custom options that allow

229 plotting of data, to include variability in CV-values and customization of population-

230 assignment plots, will facilitate the selection of appropriate K-values and allow variability

231 to be assessed across runs. These benefits thus allow researchers to implement

232 recommendations regarding assignment of population structure in their studies, and to

233 accurately report the variability found in their results [31]. In conclusion, ADMIXPIPE is a

234 new tool that successfully fills a contemporary gap found in pipelines that assess

235 population structure. It is our hope that ADMIXPIPE, and its subsequent improvements

236 will greatly facilitate the analysis of SNP data in non-model organisms.

237

238

239

240 **Acknowledgements**

11

241 Computational resources were provided by the Arkansas High Performance Computing

242 Center (AHPCC) and the NSF Jetstream XSEDE Resource (XSEDE Allocation: TG-

243 BIO160065). This research represents partial fulfillment of the Ph.D. degree (SMM) in

244 Biological Sciences at University of Arkansas.

245

246 **Funding**

247 We acknowledge indirect financial support from the University of Arkansas in the form of

248 university endowments. These include the Bruker Professorship in Life Sciences

249 (MRD), the 21st Century Chair in Global Change Biology (MED), a Doctoral Academy

250 Fellowship (SMM), and a Distinguished Doctoral Fellowship (TKC). Funding agencies

251 played no role in the design and/or conclusions of this study.

252

253 **Availability of data and materials**

254 Data utilized for benchmarking was part of an earlier publication, and is available on

255 Data Dryad (https://datadryad.org/stash/dataset/doi:10.5061/dryad.d3q3220). Source

256 code for ADMIXPIPE is released under the GNU General Public License v3.0 at

257 https://github.com/stevemussmann/admixturePipeline. The pipeline will run on Unix-

258 based operating systems such as Mac OSX and Linux. It is compatible with Python 2.7+

259 and Python 3.5+. Dependencies include other freely available software packages

260 (ADMIXTURE, DISTRUCT, PLINK, and VCFTOOLS).

261

262

263 **Authors' Contributions**

264    SMM, MRD, and MED designed the study; SMM and TKC authored the Python code for

265    ADMIXPIPE; TKC and SMM completed data analyses and program testing; all authors

266    contributed in drafting the manuscript, and all approved the final version.

267

268    **Competing interests**

269    The authors declare that they have no competing interests.

270

271    **Consent for publication**

272    Not applicable.

273

274    **Ethics approval and consent to participate**
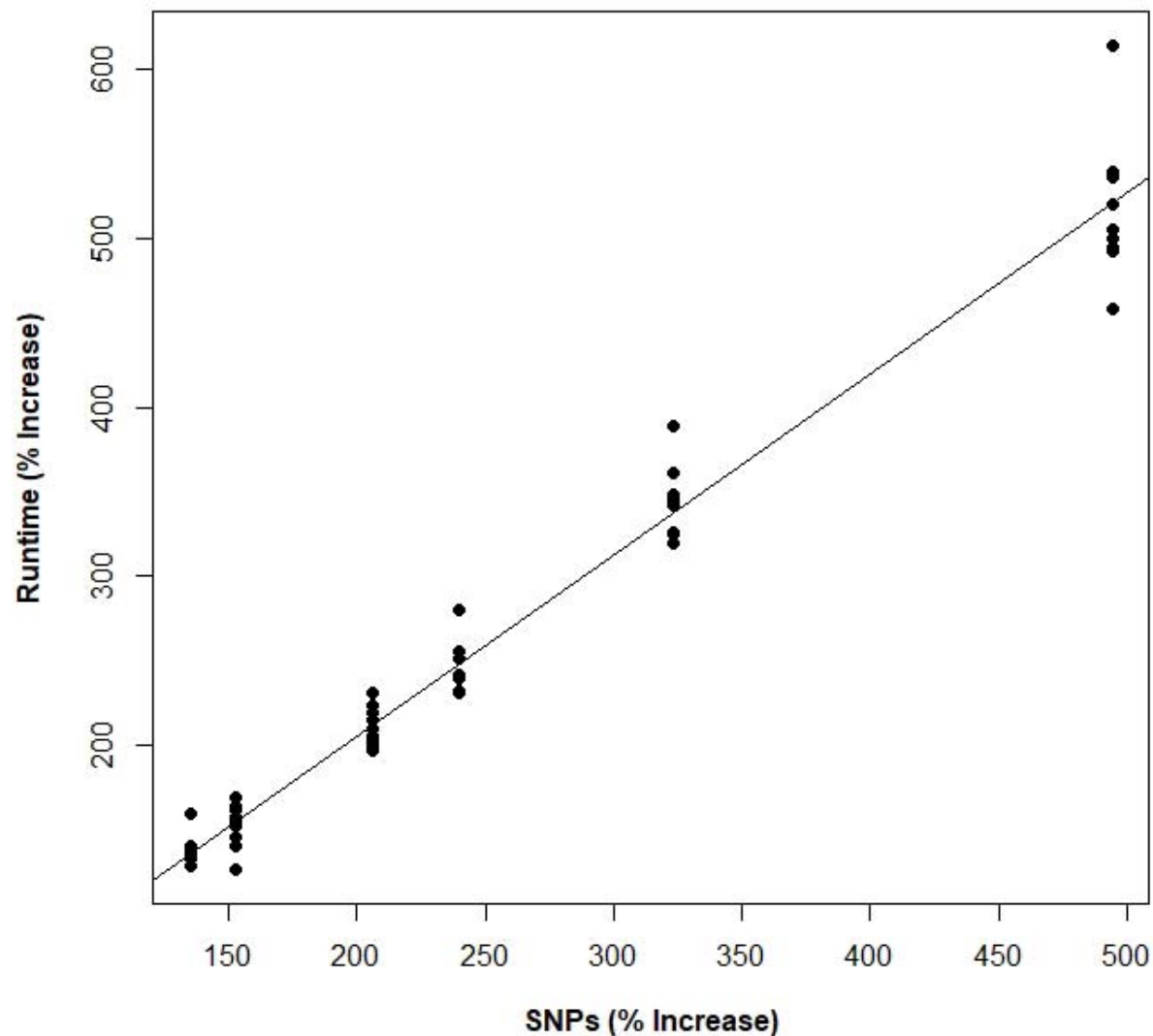
275    Not applicable.

276

277

## References

278

279    1. Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HE. Double digest RADseq:
280    an inexpensive method for *de novo* SNP discovery and genotyping in model and non-
281    model species. PLoS ONE. 2012;7:1–11. doi:10.1371/journal.pone.0037135.

282    2. Ali OA, O'Rourke SM, Amish SJ, Meek MH, Luikart G, Jeffres C, et al. RAD capture
283    (rapture): flexible and efficient sequence-based genotyping. Genetics. 2016;202:389.
284    doi:10.1534/genetics.115.183665.

285    3. Miller MR, Dunham JP, Amores A, Cresko WA, Johnson EA. Rapid and cost-effective
286    polymorphism identification and genotyping using restriction site associated DNA (RAD)
287    markers. Genome Res. 2007;17:240–8. doi:10.1101/gr.5681207.

288    4. Campbell NR, Harmon SA, Narum SR. Genotyping-in-Thousands by sequencing
289    (GT-seq): A cost effective SNP genotyping method based on custom amplicon
290    sequencing. Mol Ecol Resour. 2015;15:855–67. doi:10.1111/1755-0998.12357.

291    5. Benestan LM, Ferchaud A-L, Hohenlohe PA, Garner BA, Naylor GJP, Baums IB, et
292    al. Conservation genomics of natural and managed populations: building a conceptual
293    and practical framework. Mol Ecol. 2016;25:2967–77. doi:10.1111/mec.13647.

294    6. DaCosta JM, Sorenson MD. Amplification biases and consistent recovery of loci in a
295    double-digest RAD-seq protocol. PLoS ONE. 2014;9:1–14.
296    doi:10.1371/journal.pone.0106713.

297    7. O'Leary SJ, Puritz JB, Willis SC, Hollenbeck CM, Portnoy DS. These aren't the loci
298    you're looking for: Principles of effective SNP filtering for molecular ecologists. Mol Ecol.
299    2018;27:3193–206. doi:10.1111/mec.14792.

300    8. Shafer ABA, Peart CR, Tusso S, Maayan I, Brelsford A, Wheat CW, et al.
301    Bioinformatic processing of RAD-seq data dramatically impacts downstream population
302    genetic inference. Methods Ecol Evol. 2017;8:907–17. doi:10.1111/2041-210X.12700.

303    9. Linck E, Battey CJ. Minor allele frequency thresholds strongly affect population
304    structure inference with genomic data sets. Mol Ecol Resour. 2019;19:639–47.
305    doi:10.1111/1755-0998.12995.

306    10. Díaz-Arce N, Rodríguez-Ezpeleta N. Selecting RAD-Seq Data Analysis Parameters
307    for Population Genetics: The More the Better? Front Genet. 2019;10:533.
308    doi:10.3389/fgene.2019.00533.

309    11. Pritchard JK, Stephens M, Donnelly P. Inference of population structure using
310    multilocus genotype data. Genetics. 2000;155:945–59.
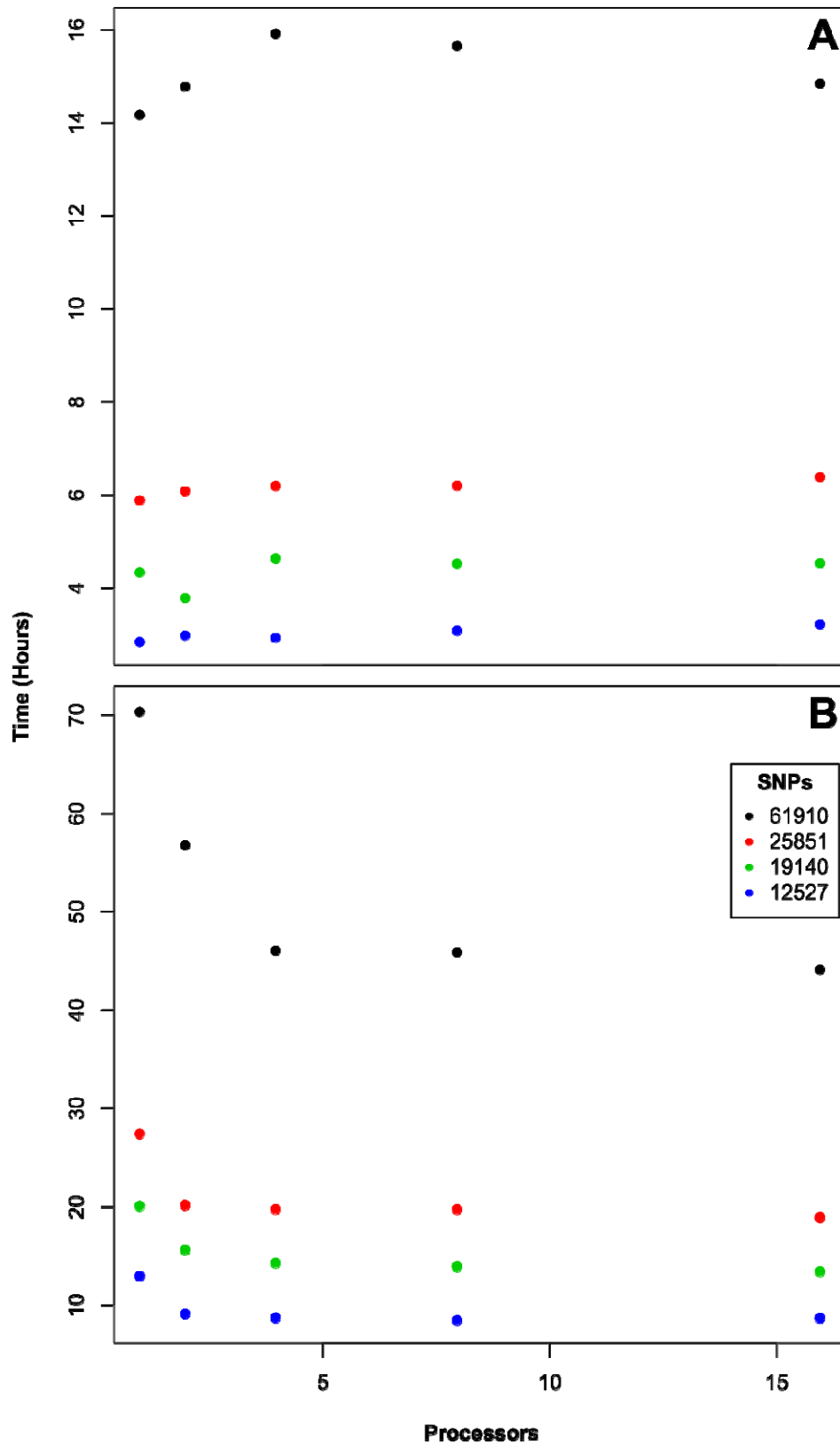311    http://www.genetics.org/content/155/2/945.abstract.

14

312  12. Falush D, Stephens M, Pritchard JK. Inference of Population Structure Using
313  Multilocus Genotype Data: Linked Loci and Correlated Allele Frequencies. Genetics.
314  2003;164:1567. http://www.genetics.org/content/164/4/1567.abstract.

315  13. Hubisz MJ, Falush D, Stephens M, Pritchard JK. Inferring weak population structure
316  with the assistance of sample group information. Mol Ecol Resour. 2009;9:1322–32.
317  doi:10.1111/j.1755-0998.2009.02591.x.

318  14. Evanno G, Regnaut S, Goudet J. Detecting the number of clusters of individuals
319  using the software Structure: a simulation study. Mol Ecol. 2005;14:2611–20.

320  15. Verdu P, Pemberton TJ, Laurent R, Kemp BM, Gonzalez-Oliver A, Gorodezky C, et
321  al. Patterns of Admixture and Population Structure in Native Populations of Northwest
322  North America. PLOS Genet. 2014;10:e1004530. doi:10.1371/journal.pgen.1004530.

323  16. Chhatre VE, Emerson KJ. STRAUTO: automation and parallelization of STRUCTURE
324  analysis. BMC Bioinformatics. 2017;18:192. doi:10.1186/s12859-017-1593-0.

325  17. Besnier F, Glover KA. PARALLELSTRUCTURE: A R Package to Distribute Parallel
326  Runs of the Population Genetics Program STRUCTURE on Multi-Core Computers. PLOS
327  ONE. 2013;8:e70651. doi:10.1371/journal.pone.0070651.

328  18. Rosenberg NA. distruct: a program for the graphical display of population structure.
329  Mol Ecol Notes. 2004;4:137–8. doi:10.1046/j.1471-8286.2003.00566.x.

330  19. Jakobsson M, Rosenberg NA. CLUMPP: a cluster matching and permutation program
331  for dealing with label switching and multimodality in analysis of population structure.
332  Bioinformatics. 2007;23:1801–6. doi:10.1093/bioinformatics/btm233.

333  20. Kopelman NM, Mayzel J, Jakobsson M, Rosenberg NA, Mayrose I. CLUMPAK: a
334  program for identifying clustering modes and packaging population structure inferences
335  across K. Mol Ecol Resour. 2015;15:1179–91. doi:10.1111/1755-0998.12387.

336  21. Earl DA, vonHoldt BM. STRUCTURE HARVESTER: a website and program for
337  visualizing STRUCTURE output and implementing the Evanno method. Conserv Genet
338  Resour. 2012;4:359–61. doi:10.1007/s12686-011-9548-7.

339  22. Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in
340  unrelated individuals. Genome Res. 2009;19:1655–64. doi:10.1101/gr.094052.109.

341  23. Dereeper A, Homa F, Andres G, Sempere G, Sarah G, Hueber Y, et al. SNiPLAY3: a
342  web-based application for exploration and large scale analyses of genomic variations.
343  Nucleic Acids Res. 2015;43:W295–300. doi:10.1093/nar/gkv351.

344  24. Alexander DH, Lange K. Enhancements to the ADMIXTURE algorithm for individual
345  ancestry estimation. BMC Bioinformatics. 2011;12:246. doi:10.1186/1471-2105-12-246.

346  25. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK:
347  a tool set for whole-genome association and population-based linkage analyses. Am J
348  Hum Genet. 2007;81:559–75. http://www.ncbi.nlm.nih.gov/pmc/articles/PMC1950838/.

349  26. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA. The variant
350  call format and VCFTOOLS. Bioinformatics. 2011;27. doi:10.1093/bioinformatics/btr330.

351  27. Eaton DA. PyRAD: assembly of *de novo* RADseq loci for phylogenetic analyses.
352  Bioinformatics. 2014;30:1844–9. doi:10.1093/bioinformatics/btu121.

353  28. Eaton DAR, Overcast I. ipyrad: Interactive assembly and analysis of RADseq
354  datasets. Bioinformatics. 2020. doi:10.1093/bioinformatics/btz966.

355  29. Chafin TK, Douglas MR, Martin BT, Douglas ME. Hybridization drives genetic
356  erosion in sympatric desert fishes of western North America. Heredity. 2019;123:759–
357  73. doi:10.1038/s41437-019-0259-2.

358  30. Raj A, Stephens M, Pritchard JK. fastSTRUCTURE: Variational Inference of
359  Population Structure in Large SNP Data Sets. Genetics. 2014;197:573–89.
360  doi:10.1534/genetics.114.164350.

361  31. Janes JK, Miller JM, Dupuis JR, Malenfant RM, Gorrell JC, Cullingham CI, et al. The
362  $K = 2$ conundrum. Mol Ecol. 2017;26:3594–602. doi:10.1111/mec.14187.

363

364

**Figure 1.** The percent increase in runtime for ADMIXPIPE exhibits a nearly 1:1 ratio with respect to percent increase in the number of SNPs. Data is based upon pairwise comparisons in runtime and input size increases for four datasets of varying size (61,910 SNPs, 25,851 SNPs,, 19,140 SNPs, and 12,527 SNPs). $R^2 = 0.975$, degrees of freedom=58.

17

371

372    **Figure 2.** Results of benchmarking ADMIXPIPE for two ranges of population clustering
373    (K) values. Time is presented in hours on the Y-axis. Plot A shows total runtime for 20
374    replicates each of K=1-8. Plot B shows total runtime for 16 replicates each of K=9-16.
375    The number of processor cores (CPU=1, 2, 4, 8, and 16) was varied across runs. Four
376    data thinning intervals (1, 25, 50, and 100) produced variable numbers of SNPs
377    (61,910, 25,851, 19,140, and 12,527 respectively).