

SOFTWARE

Open Access



ADMIXPIPE: population analyses in ADMIXTURE for non-model organisms

Steven M. Mussmann^{1,2*} , Marlis R. Douglas¹, Tyler K. Chafin¹ and Michael E. Douglas¹

* Correspondence: smussmann@gmail.com

¹Department of Biological Sciences, University of Arkansas, Fayetteville, AR 72701, USA

²Present address: Molecular Ecology Laboratory, Southwestern Native Aquatic Resources and Recovery Center (SNARRC), U.S. Fish & Wildlife Service, PO Box 219, Dexter, NM 88230, USA

Abstract

Background: Research on the molecular ecology of non-model organisms, while previously constrained, has now been greatly facilitated by the advent of reduced-representation sequencing protocols. However, tools that allow these large datasets to be efficiently parsed are often lacking, or if indeed available, then limited by the necessity of a comparable reference genome as an adjunct. This, of course, can be difficult when working with non-model organisms. Fortunately, pipelines are currently available that avoid this prerequisite, thus allowing data to be a priori parsed. An oft-used molecular ecology program (i.e., STRUCTURE), for example, is facilitated by such pipelines, yet they are surprisingly absent for a second program that is similarly popular and computationally more efficient (i.e., ADMIXTURE). The two programs differ in that ADMIXTURE employs a maximum-likelihood framework whereas STRUCTURE uses a Bayesian approach, yet both produce similar results. Given these issues, there is an overriding (and recognized) need among researchers in molecular ecology for bioinformatic software that will not only condense output from replicated ADMIXTURE runs, but also infer from these data the optimal number of population clusters (K).

Results: Here we provide such a program (i.e., ADMIXPIPE) that (a) filters SNPs to allow the delineation of population structure in ADMIXTURE, then (b) parses the output for summarization and graphical representation via CLUMPAK. Our benchmarks effectively demonstrate how efficient the pipeline is for processing large, non-model datasets generated via double digest restriction-site associated DNA sequencing (ddRAD). Outputs not only parallel those from STRUCTURE, but also visualize the variation among individual ADMIXTURE runs, so as to facilitate selection of the most appropriate K -value.

Conclusions: ADMIXPIPE successfully integrates ADMIXTURE analysis with popular variant call format (VCF) filtering software to yield file types readily analyzed by CLUMPAK. Large population genomic datasets derived from non-model organisms are efficiently analyzed via the parallel-processing capabilities of ADMIXTURE. ADMIXPIPE is distributed under the GNU Public License and freely available for Mac OSX and Linux platforms at: <https://github.com/stevemussmann/admixturePipeline>.

Keywords: RADseq, SNP analysis, Population genomics, Population structure, ADMIXTURE analysis



Background

Advances in genomics during the past decade have accelerated research in molecular ecology by significantly increasing the capacity of researchers to generate vast quantities of data at relatively low cost. These advances largely represent the development of reduced representation genomic libraries [1–3] that identify tens of thousands of SNPs for non-model organisms, coupled with high-throughput sequencing methods that efficiently genotype fewer SNPs for thousands of individuals [4]. However, data generation, particularly through these novel and affordable marker-discovery methods [5], has greatly outpaced analytical capabilities, and especially so with regard to evolutionary and conservation genomics.

Technological advances have also precipitated a suite of new analytical issues. The thousands of SNPs generated in a typical RADseq project may exhibit biases that impact the inferences that can be drawn from these data [6], and which necessitate careful data filtration to avoid [7]. Yet, the manner by which data are filtered represents a double-edged sword. While it is certainly mandated (as above), the procedures involved must be carefully evaluated in the context of each study, in that downstream analyses can be seriously impacted [8, 9], to include the derivation of population structure [10].

For example, the analysis of multilocus codominant markers in evaluation of population structure is frequently accomplished using methods that make no a priori assumptions about underlying population structure. One of the most popular methods in this regard is the program STRUCTURE [11–13]. However, it necessitates that users test specific clustering values (K), and conduct post hoc evaluation of results so as to determine an optimal K [14]. This typically involves searching a complicated parameter space using heuristic algorithms for Maximum Likelihood (ML) and Bayesian (BA) methods that, in turn, provide additional complications such as a tendency to sample local optima [15].

A common mitigation strategy is to sample multiple independent replicates at each K , using different random number seeds for initialization. These results are subsequently collated and evaluated to assess confidence that global rather than local optima have indeed been sampled. Clearly, this procedure must be automated so as to alleviate the onerous task of testing multiple replicates across a range of K -values. Pipelines to do so are available for STRUCTURE, and have been deployed on high-performance computing systems via integrated parallelization (STRAUTO, PARALLELSTRUCTURE) [16, 17]. Multiple programs have likewise been developed for handling STRUCTURE output (i.e., CLUMPP, DISTRUCT) [18, 19]; and pipelines constructed to assess the most appropriate K -values (i.e., STRUCTUREHARVESTER, CLUMPAK) [20, 21].

Despite the considerable focus on STRUCTURE, few such resources have been developed for a popular alternative program (i.e., ADMIXTURE [22]). The Web of Science indexing service indicates that (as of January, 2020) ADMIXTURE has been cited 1812 times since initial publication (September, 2009). This includes 479 (26.4%) in 2019 alone. Despite its popularity, it has just a single option that promotes the program as part of a pipeline (i.e., SNIPLAY3 [23]), which unfortunately requires a reference genome as an adjunct for its application. Needless to say, its applicability is thus limited for those laboratories that employ non-model organisms as study species.

Options for post-processing of ADMIXTURE results are similarly limited, but some packages do exist. One positive is that CLUMPAK is flexible enough in its

implementation to allow for the incorporation of ADMIXTURE output, as well as that of STRUCTURE. Alternatively, PONG provides options for processing and visualizing ADMIXTURE outputs [24]. However, no available software currently exists to summarize variation in cross-validation (CV) values, the preferred method for selecting an optimal K -value in ADMIXTURE [25].

Here we describe a novel software package that integrates ADMIXTURE as the primary component of an analytical pipeline that also incorporates the filtering of data as part of its procedure. This, in turn, provides a high-throughput capability that not only generates input for ADMIXTURE but also evaluates the impact of filtering on population structure. ADMIXPIPE also automates the process of testing multiple K -values, conducts replicates at each K , and automatically formats these results as input for the CLUMPAK pipeline. Optional post-processing scripts are also provided as a part of the toolkit to process CLUMPAK output, and to visualize the variability among CV values for independent ADMIXTURE runs. Sections of the pipeline are specifically designed for use with non-model organisms, as these are the dominant study species in evolutionary and conservation genomic investigations.

Implementation

The workflow for ADMIXPIPE is presented in Fig. 1. The pipeline requires two input files: a population map and a standard VCF file. The population map is a tab-delimited text file with each row representing a sample name/ population pair. The VCF file is filtered according to user-specified command line options that include the following: minor allele frequency (MAF) filter, biallelic filtering, data thinning measured in base-pairs (bp), and missing data filtering (for both individuals and loci). Users may also remove specific samples from their analysis by designating a file of sample names to be

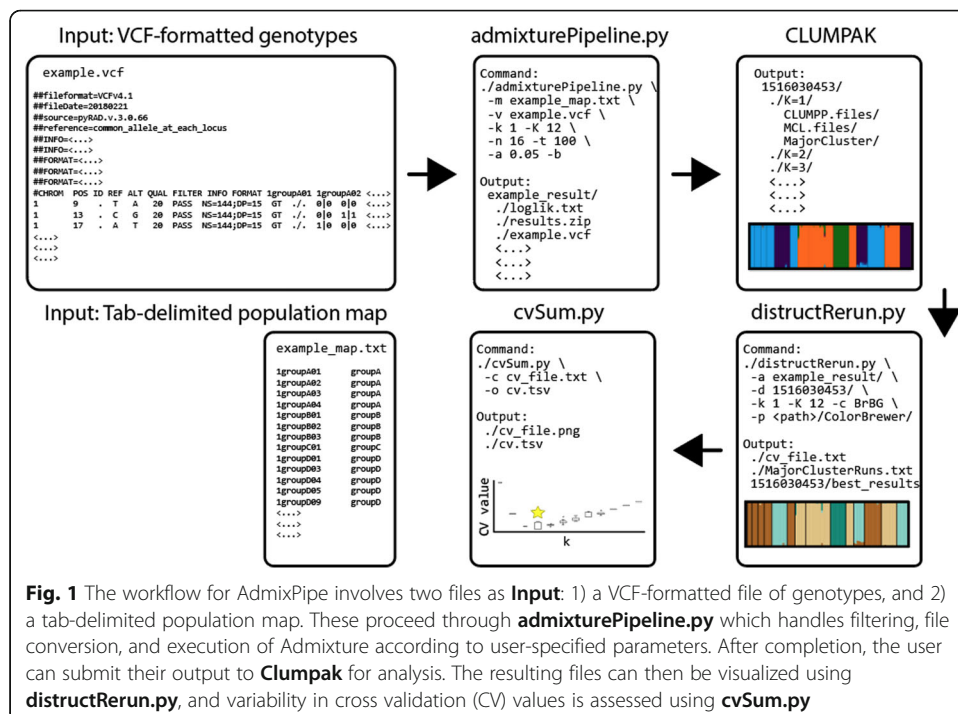


Fig. 1 The workflow for AdmixPipe involves two files as **Input**: 1) a VCF-formatted file of genotypes, and 2) a tab-delimited population map. These proceed through **admixturePipeline.py** which handles filtering, file conversion, and execution of Admixture according to user-specified parameters. After completion, the user can submit their output to **Clumpak** for analysis. The resulting files can then be visualized using **distractRerun.py**, and variability in cross validation (CV) values is assessed using **cvSum.py**

ignored. All filtering and the initial conversion to PLINK (PED/MAP) format [26] is handled by VCFTOOLS [27].

An important consideration in filtering is mitigation of linkage disequilibrium. VCFTOOLS can calculate linkage disequilibrium statistics, however these do not consider population information, thereby increasing the potential for type I error [28]. PLINK not only suffers from these limitations, but also requires a “window size” input that specifies the lengths of genomic regions within which statistical comparisons among loci are conducted. This is typically inappropriate for non-model organisms due to a lack of whole-genome resources. Non-overlapping contigs produced via reduced-representation methods can be short (e.g., 100 bp), making it a reasonable assumption that all SNPs within a contig are linked. Therefore, we suggest specifying a thinning interval in excess of the longest contig length to ensure that ADMIXPIPE samples a single SNP per contig. This method is homologous to solutions implemented in popular RADseq assembly pipelines such as STACKS and IPYRAD to minimize linkage disequilibrium in datasets [29, 30].

Additional conversions following the filtering and initial conversion via VCFTOOLS are required before the PLINK-formatted files will be accepted by ADMIXTURE. Popular software packages for de novo assembly of RADseq data, such as pyRAD [29, 31] produce VCF files with each locus as an individual “chromosome.” As a consequence, these pipelines produce outputs in which the number of “chromosomes” exceeds the number present in the model organisms for which PLINK was originally designed. The initial MAP file is therefore modified to append a letter at the start of each “chromosome” number. PLINK is then executed using the “–allow-extra-chr 0” option that treats loci as unplaced contigs in the final PED/ MAP files submitted to ADMIXTURE.

The main element of the pipeline executes ADMIXTURE on the filtered data. The assessment of multiple K values and multiple replicates is automated, based upon user-specified command line input. The user defines minimum and maximum K values to be tested, in addition to the number of replicates for each K . Users may also specify the number of processor cores to be utilized by ADMIXTURE, and the cross-validation number that is utilized in determining optimal K . The final outputs of the pipeline include a compressed results file and a population file that are ready for direct submission to CLUMPAK for processing and visualization.

The pipeline also offers two accessory scripts for processing of CLUMPAK output. The first (i.e., `distructRerun.py`) compiles the major clusters identified by CLUMPAK, generates DISTRUCT input files, executes DISTRUCT, and extracts CV-values for all major cluster runs. The second script (i.e., `cvSum.py`) plots the boxplots of CV-values against each K so as to summarize the distribution of CV-values for multiple ADMIXTURE runs. This permits the user to make an informed decision on the optimal K by graphing how these values vary according to independent ADMIXTURE runs.

ADMIXTURE is the only component of the pipeline that is natively parallelized. Therefore, we performed benchmarking to confirm that processing steps did not significantly increase runtime relative to that expected for ADMIXTURE. Data for benchmarking were selected from a recently published paper that utilized ADMIXPIPE for data processing [32]. The test data contained 343 individuals and 61,910 SNPs. Four data thinning intervals (i.e., 1, 25, 50, and 100) yielded SNP datasets of variable size for performance testing. All filtering intervals were repeated with variable numbers of processor cores

(i.e., 1, 2, 4, 8, and 16). Sixteen replicates of ADMIXTURE were first conducted for each $K = 1-8$ at each combination of thinning interval and number of processor cores, for a total of 20 executions of the pipeline. The process was then repeated for each $K = 9-16$, for an additional 20 runs of the pipeline. Memory profiling was conducted through the python3 'mprof' package at $K = 16$, with a thinning interval of 1 as a final test of performance. All tests were completed on a computer equipped with dual Intel Xeon E5-4627 3.30GHz processors, 256GB RAM, and with a 64-bit Linux environment.

Results

The filtering intervals resulted in datasets containing 61,910 (interval = 1 bp), 25,851 (interval = 25 bp), 19,140 (interval = 50 bp), and 12,527 SNPs (interval = 100 bp). Runtime increased linearly with the number of SNPs analyzed, regardless of the number of processors utilized (Fig. 2a: $R^2 = 0.975$, $df = 58$). For example, increasing the number of SNPs from 12,527 to 61,910 (494% increase) produced an average increase of 519% in ADMIXPIPE runtime ($SD = 41.6\%$).

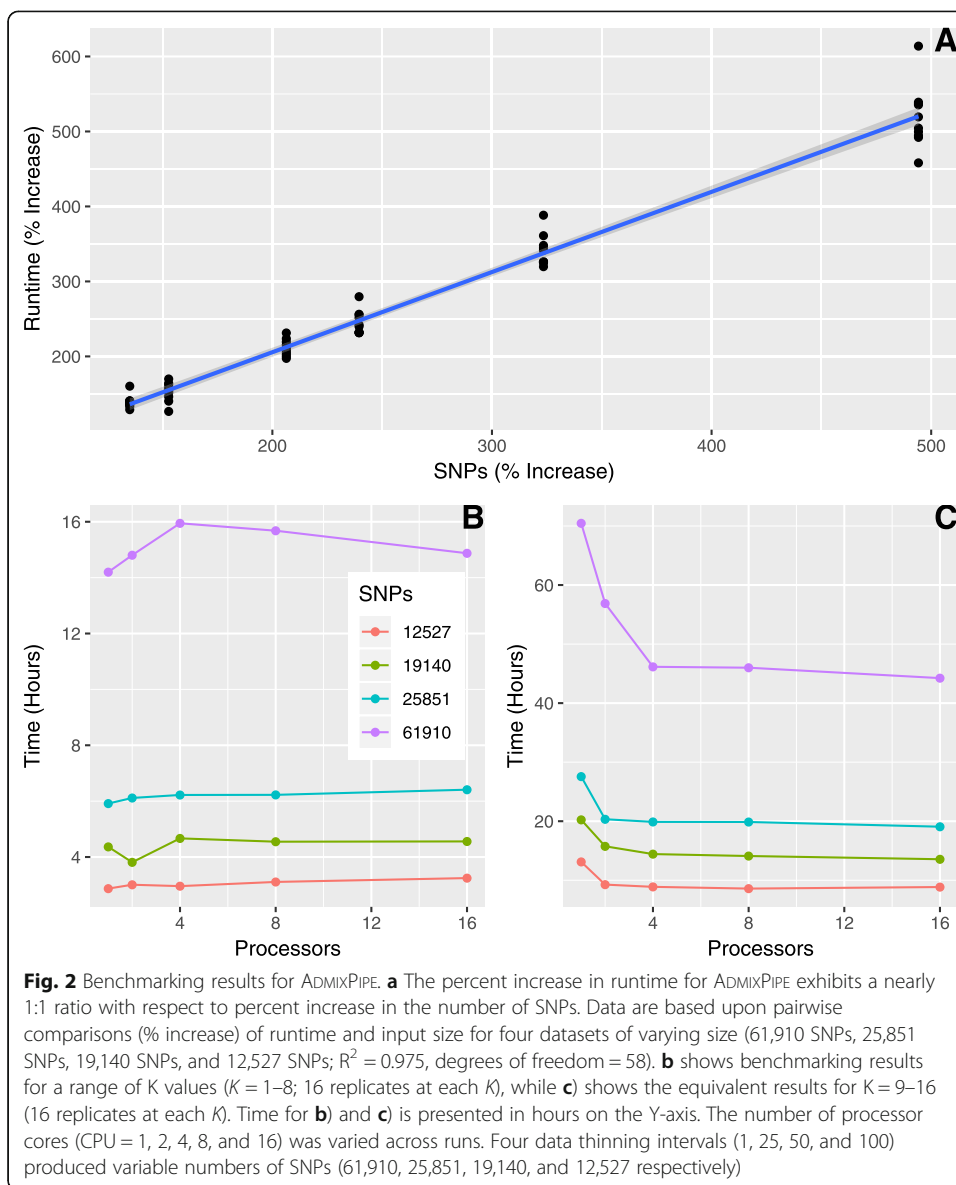
Little change was observed in response to increasing the numbers of processor cores from $K = 1-8$ (Fig. 2b). A slight decrease in performance was observed in some cases, particularly for the largest dataset. This trend changed at higher K -values, as substantial gains were observed at $K = 9-16$ (Fig. 2c) when processors were increased from 1 to 4. The most dramatic performance increase was observed for the 61,910 SNP dataset, where a 24.3-h (34.5%) reduction in computation time occurred when processors increased from 1 to 4. However, only marginal improvements occurred when processors were increased from 1 to 8 (24.5 h; 34.7%) or 16 (26.2 h; 37.7%).

Profiling also revealed efficient and consistent memory usage of ADMIXPIPE. The greatest memory spike occurred during the initial filtering steps, when peak memory usage reached approximately 120 MB. All subsequent usage held constant at ~60 MB as ADMIXTURE runs progressed.

Discussion

The performance of ADMIXPIPE improved with the number of processor cores utilized at higher K -values. However, it did not scale at the rate suggested in the original ADMIXTURE publication. We have been unable to attribute the difference in performance to any inherent property of our pipeline. Filtering and file conversion steps at the initiation of ADMIXPIPE are non-parallel sections. Reported times for completion of these steps were approximately constant across runs, with the maximum being 8 seconds. This indicates that ADMIXTURE itself is the main driver of performance, as it comprises the vast majority of system calls made by ADMIXPIPE.

The original performance increase documented for ADMIXTURE was 392% at $K = 3$, utilizing four processor cores [25]. Unfortunately, we could not replicate this result with our benchmarking data [32], or the original test data (i.e., 324 samples; 13,928 SNPs) [25] which parallels our own. When we attempted to replicate the original benchmark scores, we found that it also failed to scale as the number of processor cores increased (1-core $\bar{x} = 40.63$ s, $\sigma = 0.90$; 4-core $\bar{x} = 47.46$ s, $\sigma = 4.71$). Furthermore, we verified that performance did increase with up to four processor cores at higher K values ($K \geq 9$). We therefore view this as 'expected behavior' for ADMIXTURE, and find



no reason to believe that ADMIXPIPE has negatively impacted the performance of any individual program.

Results of ADMIXPIPE were similar to those estimated by STRUCTURE for the test dataset, as evaluated in an earlier publication [32], and gauged for the optimum $K = 8$. This is not surprising, given that ADMIXTURE implements the same likelihood model as does STRUCTURE [22]. However, minor differences have previously been noted for both programs in the assignment probabilities [32, 33].

Memory usage was efficient and constant, with the greatest increase occurring when PLINK was executed. Thus, users will be able to execute ADMIXPIPE on their desktop machines for datasets sized similarly to those evaluated herein. Performance gains were minimal with > 4 processors, and this (again) reduces the necessity for supercomputer access, since desktop computers with ≥ 4 processor cores are now commonplace. However, given the built-in parallelization capabilities of ADMIXTURE, its application on

dedicated high-performance computing clusters will be beneficial when runtime considerations are necessary, such as when evaluating $K > 8$, or SNPs $\geq 20,000$.

Finally, our integration of common SNP filtering options provides the flexibility to quickly filter data and assess the manner by which various filtering decisions impact results. A byproduct of the filtering process is the production of a STRUCTURE-formatted file that will facilitate comparisons with other popular algorithms that assess population structure. These options are important tools, particularly given recent documentation regarding the impacts of filtering on downstream analyses. We thus suggest that users implement existing recommendations on filtering RAD data, and use these to investigate subsequent impacts on their own data [7–10].

Conclusions

Benchmarking has demonstrated that the benefits of ADMIXPIPE (e.g., low memory usage and performance scaling with low numbers of processor cores at high K -values) will prove useful for researchers with limited access to advanced computing resources. ADMIXPIPE also allows the effects of common filtering options to be assessed on population structure of study species by coupling this process with the determination of population structure. Integration with CLUMPAK, and our custom options that allow plotting of data, to include variability in CV-values and customization of population-assignment plots, will facilitate the selection of appropriate K -values and allow variability to be assessed across runs. These benefits will allow researchers to implement recommendations regarding assignment of population structure in their studies, and to accurately report the variability found in their results [34]. In conclusion, ADMIXPIPE is a new tool that successfully fills a contemporary gap found in pipelines that assess population structure. We anticipate that ADMIXPIPE, and its subsequent improvements, will greatly facilitate the analysis of SNP data in non-model organisms.

Availability and requirements

Project name: AdmixPipe: A Method for Parsing and Filtering VCF Files for Admixture Analysis.

Project home page: <https://github.com/stevemussmann/admixturePipeline>

Operating system(s): Linux, Mac OSX.

Programming language: Python.

Other requirements: Python 2.7+ or Python 3.5+; Python argparse and matplotlib libraries; Dependencies include additional software packages (ADMIXTURE v1.3, DISTRICT v1.1, PLINK 1.9 beta 4.5 or higher, and VCFTOOLS v0.1.16).

License: GNU General Public License v3.0.

Any restrictions to use by non-academics: None.

Abbreviations

BA: Bayesian Analysis; CV: Cross-validation; ddRAD: Double digest Restriction-site Associated DNA; MAF: Minor Allele Frequency; ML: Maximum Likelihood; SNP: Single Nucleotide Polymorphism; VCF: Variant Call Format

Acknowledgements

Computational resources were provided by the Arkansas High Performance Computing Center (AHPCC) and the NSF Jetstream XSEDE Resource (XSEDE Allocation: TG-BIO160065). This research represents partial fulfillment of the Ph.D. degree (SMM) in Biological Sciences at University of Arkansas. We also thank two anonymous reviewers whose comments greatly improved the quality of this manuscript. The use of trade, product, industry, or firm names is for informative purposes only and does not constitute an endorsement by the U.S. Government or the U.S. Fish and Wildlife Service. Links to non-Service websites do not imply any official U.S. Fish and Wildlife Service endorsement of the opinions or

ideas expressed therein or guarantee the validity of the information provided. The findings, conclusions, and opinions expressed in this article represent those of the authors, and do not necessarily represent the views of the U.S. Fish & Wildlife Service.

Authors' contributions

SMM, MRD, and MED designed the study; SMM and TKC authored the Python code for ADMIXPIPE; TKC and SMM completed data analyses and program testing; all authors contributed in drafting the manuscript, and all approved the final version.

Funding

We acknowledge indirect financial support from the University of Arkansas in the form of university endowments. These include the Bruker Professorship in Life Sciences (MRD), the twenty-first Century Chair in Global Change Biology (MED), a Doctoral Academy Fellowship (SMM), and a Distinguished Doctoral Fellowship (TKC). Funding agencies played no role in the design and/or conclusions of this study.

Availability of data and materials

Data utilized for benchmarking was part of an earlier publication, and are available on Data Dryad (<https://datadryad.org/stash/dataset/doi:10.5061/dryad.d3q3220>). Source code for ADMIXPIPE is released under the GNU General Public License v3.0 at <https://github.com/steveusmann/admixturePipeline>. The pipeline will run on Unix-based operating systems such as Mac OSX and Linux. It is compatible with Python 2.7+ and Python 3.5+. Dependencies include other freely available software packages (ADMIXTURE, DISTRUCT, PLINK, and VCFTOOLS).

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Received: 25 February 2020 Accepted: 23 July 2020

Published online: 29 July 2020

References

- Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HE. Double digest RADseq: an inexpensive method for *de novo* SNP discovery and genotyping in model and non-model species. *PLoS One*. 2012;7:1–11. <https://doi.org/10.1371/journal.pone.0037135>.
- Ali OA, O'Rourke SM, Amish SJ, Meek MH, Luikart G, Jeffres C, et al. RAD capture (rapture): flexible and efficient sequence-based genotyping. *Genetics*. 2016;202:389. <https://doi.org/10.1534/genetics.115.183665>.
- Miller MR, Dunham JP, Amores A, Cresko WA, Johnson EA. Rapid and cost-effective polymorphism identification and genotyping using restriction site associated DNA (RAD) markers. *Genome Res*. 2007;17:240–8. <https://doi.org/10.1101/gr.5681207>.
- Campbell NR, Harmon SA, Narum SR. Genotyping-in-thousands by sequencing (GT-seq): a cost effective SNP genotyping method based on custom amplicon sequencing. *Mol Ecol Resour*. 2015;15:855–67. <https://doi.org/10.1111/1755-0998.12357>.
- Benestan LM, Ferchaud A-L, Hohenlohe PA, Garner BA, Naylor GJP, Baums IB, et al. Conservation genomics of natural and managed populations: building a conceptual and practical framework. *Mol Ecol*. 2016;25:2967–77. <https://doi.org/10.1111/mec.13647>.
- DaCosta JM, Sorenson MD. Amplification biases and consistent recovery of loci in a double-digest RAD-seq protocol. *PLoS One*. 2014;9:1–14. <https://doi.org/10.1371/journal.pone.0106713>.
- O'Leary SJ, Puritz JB, Willis SC, Hollenbeck CM, Portnoy DS. These aren't the loci you're looking for: principles of effective SNP filtering for molecular ecologists. *Mol Ecol*. 2018;27:3193–206. <https://doi.org/10.1111/mec.14792>.
- Shafer ABA, Peart CR, Tusso S, Maayan I, Brelsford A, Wheat CW, et al. Bioinformatic processing of RAD-seq data dramatically impacts downstream population genetic inference. *Methods Ecol Evol*. 2017;8:907–17. <https://doi.org/10.1111/2041-210X.12700>.
- Linck E, Batten CJ. Minor allele frequency thresholds strongly affect population structure inference with genomic data sets. *Mol Ecol Resour*. 2019;19:639–47. <https://doi.org/10.1111/1755-0998.12995>.
- Diaz-Arce N, Rodríguez-Ezpeleta N. Selecting RAD-Seq data analysis parameters for population genetics: the more the better? *Front Genet*. 2019;10:533. <https://doi.org/10.3389/fgene.2019.00533>.
- Pritchard JK, Stephens M, Donnelly P. Inference of population structure using multilocus genotype data. *Genetics*. 2000;155:945–59 <http://www.genetics.org/content/155/2/945.abstract>.
- Falush D, Stephens M, Pritchard JK. Inference of population structure using multilocus genotype data: linked loci and correlated allele frequencies. *Genetics*. 2003;164:1567 <http://www.genetics.org/content/164/4/1567.abstract>.
- Hubisz MJ, Falush D, Stephens M, Pritchard JK. Inferring weak population structure with the assistance of sample group information. *Mol Ecol Resour*. 2009;9:1322–32. <https://doi.org/10.1111/j.1755-0998.2009.02591.x>.
- Evanno G, Regnaut S, Goudet J. Detecting the number of clusters of individuals using the software structure: a simulation study. *Mol Ecol*. 2005;14:2611–20.
- Verdu P, Pemberton TJ, Laurent R, Kemp BM, Gonzalez-Oliver A, Gorodetzky C, et al. Patterns of admixture and population structure in native populations of Northwest North America. *PLoS Genet*. 2014;10:e1004530. <https://doi.org/10.1371/journal.pgen.1004530>.

16. Chhatre VE, Emerson KJ. StrAuto: automation and parallelization of STRUCTURE analysis. *BMC Bioinformatics*. 2017;18:192. <https://doi.org/10.1186/s12859-017-1593-0>.
17. Besnier F, Glover KA. ParallelStructure: A R package to distribute parallel runs of the population genetics program STRUCTURE on multi-core computers. *PLoS One*. 2013;8:e70651. <https://doi.org/10.1371/journal.pone.0070651>.
18. Rosenberg NA. Distruct: a program for the graphical display of population structure. *Mol Ecol Notes*. 2004;4:137–8. <https://doi.org/10.1046/j.1471-8286.2003.00566.x>.
19. Jakobsson M, Rosenberg NA. CLUMPP: a cluster matching and permutation program for dealing with label switching and multimodality in analysis of population structure. *Bioinformatics*. 2007;23:1801–6. <https://doi.org/10.1093/bioinformatics/btm233>.
20. Kopelman NM, Mayzel J, Jakobsson M, Rosenberg NA, Mayrose I. CLUMPAK: a program for identifying clustering modes and packaging population structure inferences across K. *Mol Ecol Resour*. 2015;15:1179–91. <https://doi.org/10.1111/1755-0998.12387>.
21. Earl DA, von Holdt BM. STRUCTURE HARVESTER: a website and program for visualizing STRUCTURE output and implementing the Evanno method. *Conserv Genet Resour*. 2012;4:359–61. <https://doi.org/10.1007/s12686-011-9548-7>.
22. Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in unrelated individuals. *Genome Res*. 2009;19:1655–64. <https://doi.org/10.1101/gr.094052.109>.
23. Dereeper A, Homa F, Andres G, Sempere G, Sarah G, Hueber Y, et al. SNIPlay3: a web-based application for exploration and large scale analyses of genomic variations. *Nucleic Acids Res*. 2015;43:W295–300. <https://doi.org/10.1093/nar/gkv351>.
24. Behr AA, Liu KZ, Liu-Fang G, Nakka P, Ramachandran S. Pong: fast analysis and visualization of latent clusters in population genetic data. *Bioinformatics*. 2016;32:2817–23. <https://doi.org/10.1093/bioinformatics/btw327>.
25. Alexander DH, Lange K. Enhancements to the ADMIXTURE algorithm for individual ancestry estimation. *BMC Bioinformatics*. 2011;12:246. <https://doi.org/10.1186/1471-2105-12-246>.
26. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet*. 2007;81:559–75. <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC1950838/>.
27. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA. The variant call format and VCFtools. *Bioinformatics*. 2011;27:2156–8. <https://doi.org/10.1093/bioinformatics/btr330>.
28. Law B, Buckleton JS, Triggs CM, Weir BS. Effects of population structure and admixture on exact tests for association between loci. *Genetics*. 2003;164:381–7. <https://pubmed.ncbi.nlm.nih.gov/12750348>.
29. Eaton DAR, Overcast I. Ipyrad: interactive assembly and analysis of RADseq datasets. *Bioinformatics*. 2020;36:2592–4. <https://doi.org/10.1093/bioinformatics/btz966>.
30. Rochette NC, Rivera-Colón AG, Catchen JM. Stacks 2: analytical methods for paired-end sequencing improve RADseq-based population genomics. *Mol Ecol*. 2019;28:4737–54. <https://doi.org/10.1111/mec.15253>.
31. Eaton DA. PyRAD: assembly of *de novo* RADseq loci for phylogenetic analyses. *Bioinformatics*. 2014;30:1844–9. <https://doi.org/10.1093/bioinformatics/btu121>.
32. Chafin TK, Douglas MR, Martin BT, Douglas ME. Hybridization drives genetic erosion in sympatric desert fishes of western North America. *Heredity*. 2019;123:759–73. <https://doi.org/10.1038/s41437-019-0259-2>.
33. Raj A, Stephens M, Pritchard JK. fastSTRUCTURE: variational inference of population structure in large SNP data sets. *Genetics*. 2014;197:573–89. <https://doi.org/10.1534/genetics.114.164350>.
34. Janes JK, Miller JM, Dupuis JR, Malenfant RM, Gorrell JC, Cullingham CI, et al. The K = 2 conundrum. *Mol Ecol*. 2017;26:3594–602. <https://doi.org/10.1111/mec.14187>.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

