

# ***dDocent*: a RADseq, variant-calling pipeline designed for population genomics of non-model organisms**

Restriction-site associated DNA sequencing (RADseq) has become a powerful and useful approach for population genomics. Currently, no software exists that utilizes both paired-end reads from RADseq data to efficiently produce population-informative variant calls, especially for non-model organisms with large effective population sizes and high levels of genetic polymorphism. *dDocent* is an analysis pipeline with a user-friendly, command-line interface designed to process individually barcoded RADseq data (with double cut sites) into informative SNPs/Indels for population-level analyses. The pipeline, written in BASH, uses data reduction techniques and other stand-alone software packages to perform quality trimming and adapter removal, *de novo* assembly of RAD loci, read mapping, SNP and Indel calling, and baseline data filtering. Double-digest RAD data from population pairings of three different marine fishes were used to compare *dDocent* with *Stacks*, the first generally available, widely used pipeline for analysis of RADseq data. *dDocent* consistently identified more SNPs shared across greater numbers of individuals and with higher levels of coverage. This is due to the fact that *dDocent* quality trims instead of filtering and incorporates both forward and reverse reads in assembly, mapping, and SNP calling, thus enabling use of reads with Indel polymorphisms. The pipeline and a comprehensive user guide can be found at ( <http://dDocent.wordpress.com> ).

1 ***dDocent: a RADseq, variant-calling pipeline designed for population genomics of non-model***  
 2 ***organisms***

3 JONATHAN B. PURITZ<sup>†</sup>, CHRISTOPHER M. HOLLENBECK, AND JOHN R. GOLD

4 *Marine Genomics Laboratory, Harte Research Institute, Texas A&M University-Corpus Christi,*  
 5 *6300 Ocean Drive, Corpus Christi, Texas 78412-5869*

6 <sup>†</sup>Author to whom correspondence should be addressed.

7 Email: jonathan.puritz@tamucc.edu Phone: 361-825-3343 Fax: 361-825-2050

# INTRODUCTION

Next-generation sequencing (NGS) has transformed the field of genetics into genomics by providing DNA sequence data at an ever increasing rate and reduced cost (Mardis, 2008). The nascent field of population genomics relies on NGS coupled with laboratory methods to reproducibly reduce genome complexity to a few thousand loci. The most common approach, restriction-site associated DNA sequencing (RADseq), uses restriction endonucleases to randomly sample the genome at locations adjacent to restriction-enzyme recognition sites that, when coupled with Illumina sequencing, produces high coverage of homologous SNP (Single Nucleotide Polymorphism) loci. As such, RADseq provides a powerful method for population level genomic studies (Ellegren, 2014; Narum et al., 2013; Rowe et al., 2011).

The original RADseq approach (Baird et al., 2008), and initial population genomic studies employing it (Hohenlohe et al., 2010), focused on SNP discovery and genotyping on the first (forward) read only. This is because the original RADseq method (Baird et al., 2008) utilized random shearing to produce RAD loci; paired-end reads were not of uniform length or coverage, making it problematic to find SNPs at high and uniform levels of coverage across a large proportion of individuals. As a result, the most comprehensive and widely used software package for analysis of RADseq data, *Stacks* (Catchen et al., 2013, 2011), provides SNP genotypes based only on first-read data. In contrast, RADseq approaches such as ddRAD (Peterson et al., 2012), 2bRAD (Wang et al., 2012), and ezRAD (Toonen et al., 2013) rely on restriction enzymes to define both ends of a RAD locus, largely producing RAD loci of fixed length (fRAD). Paired-end Illumina sequencing of fRAD fragments provides an opportunity to significantly expand the number of SNPs that can be genotyped from a single RADseq library.

Here, the variant-calling pipeline *dDocent* is introduced as a tool for generating population genomic data; a brief methodological outline of the analysis pipeline also is presented. *dDocent* is a wrapper script designed to take raw fRAD data and produce population informative SNP

calls (SNPs that are shared across the majority of individuals and populations), taking full advantage of both paired-end reads. *dDocent* is configured for organisms with high levels of nucleotide and INDEL polymorphisms, such as are found in many marine organisms (Guo et al., 2012; Keever et al., 2009; Sodergren et al., 2006; Waples, 1998; Ward et al., 1994); however, the pipeline also can be adjusted for low polymorphism species. As input, *dDocent* takes paired FASTQ files for individuals and outputs raw SNP and INDEL calls as well as filtered SNP calls in VCF format. The pipeline and a comprehensive online manual can be found at (<http://dDocent.wordpress.com>). Finally, results of pipeline analyses, using both *dDocent* and *Stacks*, of populations of three species of marine fishes are provided to demonstrate the utility of *dDocent* compared to *Stacks*, the first and most comprehensive, existing software package for RAD population genomics.

## METHODS

### *Implementation and basic usage*

The *dDocent* pipeline is written in BASH and will run using most Unix-like operating systems. *dDocent* is largely dependent on other bioinformatics software packages, taking advantage of programs designed specifically for each task of the analysis and ensuring that each modular component can be updated separately. Proper implementation depends on the correct installation of each third-party packages/tools. A full list of dependencies can be found in the user manual at (<http://ddocent.wordpress.com/ddocent-pipeline-user-guide/>) and a sample script to automatically download and install the packages in a Linux environment can be found at the *dDocent* repository (<https://github.com/jpuritz/dDocent>).

*dDocent* is run by simply switching to a directory containing input data and starting the program. There is no configuration file, and *dDocent* will proceed through a short series of command-line prompts, allowing the user to establish analysis parameters. After all required variables are configured, including an e-mail address for a completion notification, *dDocent*

provides instructions on how to move the program to the background and run, undisturbed, until completion. The pipeline is designed to take advantage of multiple processing-core machines and, whenever possible, processes are invoked with multiple threads or occurrences. For most Linux distributions, the number of processing cores should be automatically detected. If *dDocent* cannot determine the number of processors, it will ask the user to input the value.

There are two distinct modules of *dDocent*: *dDocent.FB* and *dDocent.GATK*. *dDocent.FB* uses minimal, BAM-file preparation steps before calling SNPs and INDELs, simultaneously using FreeBayes (Garrison & Marth, 2012). *dDocent.GATK* uses GATK (McKenna et al., 2010) for INDEL realignment, SNP and INDEL genotyping (using HaplotypeCaller), and variant quality-score recalibration, largely following GATK Best Practices recommendations (Auwera & Carneiro, 2013; DePristo et al., 2011). The modules represent two different strategies for SNP/INDEL calling that are completely independent of one another. Currently, *dDocent.FB* is easier to implement, substantially faster to execute, and depends on software that is commercially unrestricted; consequently, the remainder of this paper focuses on *dDocent.FB*. Additional information on *dDocent.GATK* may be found in the user guide.

### *Data input requirements*

*dDocent* requires demultiplexed forward and paired-end FASTQ files for every individual in the analysis (flRAD data only). A simple naming convention (a single-word locality code/name and a single-word sample identifier separated by an underscore) must be followed for every sample; examples are *LOCA\_IND01.F.fq* and *LOCA\_IND01.R.fq*. A sample script for using a text file containing barcodes and sample names and *process\_radtags* from *Stacks* (Catchen et al., 2013) to properly demultiplex samples and put them in the proper *dDocent* naming convention, can be found at the *dDocent* repository (<https://github.com/jpuritz/dDocent>).

### *Quality trimming*

After *dDocent* checks that it is recognizing the proper number of samples in the current directory, it asks the user if s/he wishes to proceed with quality trimming of sequence data. If directed, *dDocent* can use the program *Trim Galore!* ([http://www.bioinformatics.babraham.ac.uk/projects/trim\\_galore/](http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/)) to simultaneously remove Illumina adapter sequences and trim ends of reads of low quality. By default, *Trim Galore!* looks for double-digest RAD adapters (Peterson et al., 2012) and trims bases with quality scores less than PHRED 10 (corresponding to a 10% chance of error in the base call). The read mapping and variant calling steps of *dDocent* account for base quality, so minimal trimming of the data is needed. Typically, quality trimming only needs to be performed once, so the option exists to skip this step in subsequent *dDocent* analyses.

#### *De novo assembly*

Without reference material, population genomic analyses from RADseq depend on *de novo* assembly of a set of reference contigs. Intrinsically, not all RAD loci appear in all individuals due to stochastic processes inherent in library preparation and sequencing and to polymorphism in restriction-enzyme restriction sites (Catchen et al., 2011). Moreover, populations can contain large levels of within-locus polymorphism, making generation of a reference sequence computationally difficult. *dDocent* minimizes the amount of data used for assembly by taking advantage of the fact that flRAD loci present in multiple individuals should have higher levels of exactly matching reads (forward and reverse) than loci that are only present in a few individuals. Caution is advised for unique reads with low levels of coverage throughout the data set as they likely represent sequencing errors or polymorphisms that are shared only by a few individuals.

In the first step of the assembly process, untrimmed, paired-end reads are reverse complemented and concatenated to forward reads. Unique paired reads are identified and their occurrences are counted in the entire data set. These data are tabulated into the number of unique reads per levels of 1X to 50X coverage; a graph is then generated and printed to the terminal.

The distribution usually follows an asymptotic relationship (Figure 1), with a large proportion of reads only having one or two occurrences, meaning they likely will not be informative on a population scale. Highly polymorphic RAD loci still should have at least one allele present at the level of expected sequence coverage, so this can be used as a guide for informative data. The user chooses a cut-off level of coverage for reads to be used for assembly – note that all reads are still used for subsequent steps of the pipeline.

After a cut-off level is chosen, remaining concatenated reads are divided back into forward- and reverse-read files and then input directly into the RADseq assembly program *Rainbow* (Chong et al., 2012). The default parameters of *Rainbow* are used except that the maximum number of mismatches used in initial clustering is changed from four to six to help account for highly polymorphic species with large effective population sizes. In short, *Rainbow* clusters forward reads based on similarity; clusters are then recursively divided, based on reverse reads, into groups representing single alleles. Reads in merged clusters are then assembled using a greedy algorithm (Pop & Salzberg, 2008). *dDocent* then selects the longest contig for each cluster as the representative reference sequence for that RAD locus. If the forward read does not overlap with the reverse read (almost always the case with flRAD), the forward read is concatenated to the reverse read with ten ‘N’ characters as padding to represent the unknown insert. If forward and reverse reads do overlap, then a full contig is created without N padding. Finally, reference sequences are clustered based on overall sequence similarity (chosen by user, 90% by default), using the program *CD-HIT* (Fu et al., 2012; Li & Godzik, 2006). This final cluster step reduces the data set further, based on overall sequence identity after assembly. Alternatively, *de novo* assembly can be skipped and the user can provide a FASTA file with reference sequences.

*Read mapping*

*dDocent* uses the MEM algorithm (Li, 2013) of *BWA* (Li & Durbin, 2009, 2010) to map quality-trimmed reads to the reference contigs. Users can deploy the default values of *BWA* or set an alternative value for each mapping parameter (match score, mismatch score, and gap-opening penalty). The default settings are meant for mapping reads to the human genome, so users are encouraged to experiment with mapping parameters. *BWA* output is ported to *SAMtools* (Li et al., 2009), saving disk space, and alignments are saved to the disk as binary alignment/Map (BAM). BAM files are then sorted and indexed.

#### *SNP and INDEL discovery and genotyping*

*dDocent* uses a two-step process to optimize the computationally intensive task of SNP/INDEL calling. First, quality-trimmed forward and reverse reads are reduced to unique reads. This data set is then mapped to all reference sequences, using the previously entered mapping settings (see *Read Mapping* above). From this alignment, a set of intervals is created using *BEDtools* (Quinlan & Hall, 2010). The interval set saves computational time by directing the SNP-/INDEL-calling software to examine only reference sequences along contigs that have high quality mappings. Second, the interval list is then split into multiple files, one for each processing core, allowing SNP/INDEL calling to be optimized with a scatter-gather technique. The program *FreeBayes* (Garrison & Marth, 2012) is then executed multiple times simultaneously (one execution per processor and genomic interval). *FreeBayes* is a Bayesian-based, variant-detection software that uses assembled haplotype sequences to simultaneously call SNPs, INDELS, multi-nucleotide polymorphisms (MNPs), and complex events (e.g., composite insertion and substitution events) from alignment files; *FreeBayes* has the added benefit for population genomics of using reads across multiple individuals to improve genotyping (Garrison & Marth, 2012). *FreeBayes* is run with minimal changes to the default parameters; minimum mapping quality score and base quality score are set to PHRED 10. After all executions of

155 *FreeBayes* are completed, raw SNP/INDEL calls are concatenated into a single variant call file  
156 (VCF), using VCFtools (Danecek et al., 2011).

### 157 *Variant Filtering*

158 Final SNP data-set requirements are likely to be highly dependent on specific goals and aims  
159 of individual projects. To that end, *dDocent* uses *VCFtools* (Danecek et al., 2011) to provide only  
160 basic level filtering, mostly for run diagnostic purposes. *dDocent* produces a final VCF file that  
161 contains all SNPs, INDELS, MNPs, and complex events that are called in 90% of all individuals,  
162 with a minimum quality score of 30. Users are encouraged to use VCFtools and *vcflib* (part of  
163 the *FreeBayes* package; <https://github.com/ekg/vcflib>) to fully explore and filter data  
164 appropriately.

### 165 *Comparison between dDocent and Stacks*

166 Two sample localities, each comprising 20 individuals, were chosen randomly from  
167 unpublished RADseq data sets of three different, marine fish species: red snapper (*Lutjanus*  
168 *campechanus*), red drum (*Sciaenops ocellatus*), and silk snapper (*Lutjanus vivanus*). These three  
169 species are part of ongoing RADseq projects in our laboratory, and preliminary analyses  
170 indicated high levels of nucleotide polymorphisms across all populations. Double-digest RAD  
171 libraries were prepared, generally following Peterson *et al.* (2012). Individual DNA extractions  
172 were digested with *EcoRI* and *MspI*. A barcoded adapter was ligated to the *EcoRI* site of each  
173 fragment and a generic adapter was ligated to the *MspI* site. Samples were then equimolarly  
174 pooled and size-selected between 350 and 400 bp, using a Qiagen Gel Extraction Kit. Final  
175 library enhancement was completed using 12 cycles of PCR, simultaneously enhancing properly  
176 ligated fragments and adding an Illumina Index for additional barcoding. Libraries were  
177 sequenced on three separate lanes of an Illumina HiSeq 2000 at the University of Texas Genomic  
178 Sequencing and Analysis Facility. Raw sequence data were archived at NCBI's Short Read  
179 Archive (SRA) under Accession SRP041032.

Demultiplexed individual reads were analyzed with *dDocent* (version 1.0), using three different levels of final reference contig clustering (90%, 96%, and 99% similarity) in an attempt to alter the most comparable analysis variable in *dDocent* to match the maximum distance between stacks parameter and the maximum distance between stacks from different individuals parameter of *Stacks*. The coverage cut-off for assembly was 12 for red snapper, 13 for red drum, and nine for silk snapper. All *dDocent* runs used mapping variables of one, three, and five for match-score value, mismatch score, and gap-opening penalty, respectively. For comparisons, complex variants were decomposed into canonical SNP and INDEL representation from the raw VCF files, using *vcfallelicprimitives* from *vcflib* (<https://github.com/ekg/vcflib>).

For analysis with *Stacks* (version 1.08), reads were demultiplexed and cleaned using *process\_radtags*, removing reads with ‘N’ calls and low-quality base scores. Because *dDocent* inherently uses both reads for SNP/INDEL genotyping, forward reads and reverse reads were processed separately with *denovo\_map.pl*, using three different sets of parameters. The first set had a minimum depth of coverage of two to create a stack, a maximum distance of two between stacks, and a maximum distance of four between stacks from different individuals, with both the deleveraging algorithm and removal algorithms enabled. The second set had a minimum depth of coverage of three to create a stack, a maximum distance of four between stacks, and a maximum distance of eight between stacks from different individuals, with both the deleveraging algorithm and removal algorithms enabled. The third set had a minimum depth of coverage of three to create a stack, a maximum distance of four between stacks, and a maximum distance of 10 between stacks from different individuals, with both the deleveraging algorithm and removal algorithms enabled. SNP calls were output in VCF format.

For both *dDocent* and *Stacks* runs, VCFtools was used to filter out all INDELs and SNPs that had a minor allele count of less than five. SNP calls were then evaluated at different individual-coverage levels: the total number of SNPs; the number of SNPs called in 75%, 90%, and 99% of

individuals at 3X coverage; the number of SNPs called in 75% and 90% of individuals at 5X coverage; the number of SNPs called in 75% and 90% of individuals at 10X coverage; and the number of SNPs called in 75% and 90% of individuals at 20X coverage. Overall coverage levels for red snapper were lower and likely impacted by a few low-quality individuals; consequently, the number of 5X and 10X SNPs shared among 90% of individuals (after removing the bottom 10% of individuals in terms of coverage) were compared instead of SNP loci shared at 20X coverage. Results from two runs of *Stacks* (one using forward and one using reverse reads) were combined for comparison with *dDocent*, which inherently calls SNPs on both reads. All analyses and computations were performed on a 32-core Linux workstation with 128 GB of RAM.

## RESULTS AND DISCUSSION

Results of SNP calling, including run times (in minutes) for each analysis (not including quality trimming), are presented in Table 1. Data from high coverage SNP calls, averaged over all runs for each pipeline, are presented in Figure 1. While *Stacks* called a larger number of low coverage SNPs, limiting results to higher individual coverage and to higher individual call rates revealed that *dDocent* consistently called more high-quality SNPs. Run times were equivalent for both pipelines.

At almost all levels of coverage in three different data sets, *dDocent* called more SNPs across more individuals than *Stacks*. Two key differences between *dDocent* and *Stacks* likely contribute these discrepancies: (i) quality trimming instead of quality filtering, and (ii) simultaneous use of forward and reverse reads by *dDocent* in assembly, mapping, and genotyping, instead of clustering as employed by *Stacks*. As with any data analysis, quality of data output is directly linked to the quality of data input. Both *dDocent* and *Stacks* use procedures to ensure that only high-quality sequence data are retained; however, *Stacks* removes an entire read when a sliding window of bases drops below a preset quality score (PHRED 10, by default), while *dDocent* via *Trim Galore!* trims off low-quality bases, preserving high-quality bases of each read. Filtering

instead of trimming results in fewer reads entering the *Stacks* analysis (between 65%-95% of the data compared to *dDocent*; data not shown), generating lower levels of coverage and fewer SNP calls than *dDocent*.

*dDocent* offers two advantages over *Stacks*: (i) it is specifically designed for paired-end data and utilizes both forward and reverse reads for *de novo* RAD loci assembly, read mapping, variant discovery, and genotyping; and (ii) it aligns reads to reference sequence instead of clustering by identity. Using both reads to cluster and assemble RAD loci helps to ensure that portions of the genome with complex mutational events, including INDELs or small repetitive regions, are properly assembled and clustered as homologous loci. Additionally, using *BWA* to map reads to reference loci enables *dDocent* to properly align reads with INDEL polymorphisms, increasing coverage and subsequent variant discovery and genotyping. Clustering methods employed by *Stacks*, whether clustering alleles within an individual or clustering loci between individuals, effectively remove reads, alleles, and loci with INDEL polymorphisms because the associated frame shift effectively inflates the observed number of base-pair differences. For organisms with large effective population sizes and high levels of genetic diversity, such as many marine organisms (Waples, 1998; Ward et al., 1994), removing reads and loci with INDEL polymorphisms will result in a loss of shared loci and coverage.

*dDocent* is specifically designed to efficiently generate SNP and INDEL polymorphisms that are shared across multiple individuals. To that end, the output reference contigs and variant calls represent a subset of the total, genomic information content of the raw input data; RAD loci and variants present in single individuals are largely ignored. Other analysis software, such as the scripts published by Peterson et al. (2012), represent a more comprehensive alternative for generating for a full *de novo* assembly of RAD loci and would increase the chance of discovering individual level polymorphisms. For population genomics, loci that are not shared by at least 50% of all individuals and/or have minor allele frequencies of less than 5% are often filtered out.

255 *dDocent* saves computational time by ignoring these loci from the outset of assembly; however,  
256 users can pass in a more comprehensive reference (including an entire genome) in order to  
257 include all possible variant calls from the data.

## 258 **CONCLUSION**

259 *dDocent* is an open-source, freely available population genomics pipeline configured for  
260 species with high levels of nucleotide and INDEL polymorphisms, such as many marine  
261 organisms. The *dDocent* pipeline reports more SNPs shared across greater numbers of  
262 individuals and with higher levels of coverage than current alternatives. The pipeline and a  
263 comprehensive online manual can be found at (<http://dDocent.wordpress.com>) and  
264 (<https://github.com/jpuritz/dDocent>).

## 265 **ACKNOWLEDGEMENTS**

266 We thank T. Krabbenhoft for assistance in beta testing, and C. Bird and D. Portnoy for useful  
267 discussions and comments on the manuscript. We also would like to thank the three reviewers  
268 for their substantial help with troubleshooting the user guide and installation process on multiple  
269 computing platforms.

- 271 Auwera G, Carneiro M. 2013. From FastQ Data to High-Confidence Variant Calls: The Genome  
272 Analysis Toolkit Best Practices Pipeline. *Current Protocols in Bioinformatics*: 1–33.
- 273 Baird NA, Etter PD, Atwood TS, Currey MC, Shiver AL, Lewis ZA, Selker EU, Cresko WA,  
274 Johnson EA. 2008. Rapid SNP discovery and genetic mapping using sequenced RAD  
275 markers. *PloS ONE* 3: e3376.
- 276 Catchen JM, Amores A, Hohenlohe P, Cresko W, Postlethwait JH. 2011. Stacks: building and  
277 genotyping Loci de novo from short-read sequences. *G3 (Bethesda, Md.)* 1: 171–182.
- 278 Catchen J, Hohenlohe PA, Bassham S, Amores A, Cresko WA. 2013. Stacks: an analysis tool set  
279 for population genomics. *Molecular ecology* 22: 3124–3140.
- 280 Chong Z, Ruan J, Wu C. 2012. Rainbow : an integrated tool for efficient clustering and  
281 assembling RAD-seq reads. : 1–6.
- 282 Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE, Lunter G,  
283 Marth GT, Sherry ST, McVean G, Durbin R. 2011. The variant call format and VCFtools.  
284 *Bioinformatics (Oxford, England)* 27: 2156–2158.
- 285 DePristo MA, Banks E, Poplin R, Garimella K V, Maguire JR, Hartl C, Philippakis AA, Angel G  
286 del, Rivas MA, Hanna M, McKenna A, Fennell TJ, Kernysky AM, Sivachenko AY,  
287 Cibulskis K, Gabriel SB, Altshuler D, Daly MJ. 2011. A framework for variation discovery  
288 and genotyping using next-generation DNA sequencing data. *Nature genetics* 43: 491–498.
- 289 Ellegren H. 2014. Genome sequencing and population genomics in non-model organisms. *Trends*  
290 *in ecology & evolution* 29: 51–63.
- 291 Fu L, Niu B, Zhu Z, Wu S, Li W. 2012. CD-HIT: accelerated for clustering the next-generation  
292 sequencing data. *Bioinformatics (Oxford, England)* 28: 3150–3152.
- 293 Garrison E, Marth G. 2012. Haplotype-based variant detection from short-read sequencing. : 1–9.
- 294 Guo B, Zou M, Wagner A. 2012. Pervasive indels and their evolutionary dynamics after the fish-  
295 specific genome duplication. *Molecular biology and evolution* 29: 3005–3022.
- 296 Hohenlohe PA, Bassham S, Etter PD, Stiffler N, Johnson EA, Cresko WA. 2010. Population  
297 genomics of parallel adaptation in threespine stickleback using sequenced RAD tags. *PLoS*  
298 *genetics* 6: e1000862.
- 299 Keever CC, Sunday J, Puritz JB, Addison JA, Toonen RJ, Grosberg RK, Hart MW. 2009.  
300 Discordant distribution of populations and genetic variation in a sea star with high dispersal  
301 potential. *Evolution; international journal of organic evolution* 63: 3214–3227.
- 302 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.  
303 *Bioinformatics (Oxford, England)* 25: 1754–1760.
- 304 Li H, Durbin R. 2010. Fast and accurate long-read alignment with Burrows-Wheeler transform.  
305 *Bioinformatics (Oxford, England)* 26: 589–595.
- 306 Li W, Godzik A. 2006. Cd-hit: a fast program for clustering and comparing large sets of protein  
307 or nucleotide sequences. *Bioinformatics (Oxford, England)* 22: 1658–1659.
- 308 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R.  
309 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics (Oxford,*  
310 *England)* 25: 2078–2079.
- 311 Li H. 2013. Aligning sequence reads , clone sequences and assembly contigs with BWA-MEM.  
312 00: 1–3.
- 313 Mardis ER. 2008. Next-generation DNA sequencing methods. *Annual review of genomics and*  
314 *human genetics* 9: 387–402.

315 McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K,  
 316 Altshuler D, Gabriel S, Daly M, DePristo MA. 2010. The Genome Analysis Toolkit: a  
 317 MapReduce framework for analyzing next-generation DNA sequencing data. *Genome*  
 318 *research* 20: 1297–1303.

319 Narum SR, Buerkle CA, Davey JW, Miller MR, Hohenlohe PA. 2013. Genotyping-by-sequencing  
 320 in ecological and conservation genomics. *Molecular ecology* 22: 2841–2847.

321 Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HE. 2012. Double digest RADseq: an  
 322 inexpensive method for de novo SNP discovery and genotyping in model and non-model  
 323 species. *PloS one* 7: e37135.

324 Pop M, Salzberg S. 2008. Bioinformatics challenges of new sequencing technology. *Trends in*  
 325 *Genetics* 24: 142–149.

326 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic  
 327 features. *Bioinformatics (Oxford, England)* 26: 841–842.

328 Rowe HC, Renaut S, Guggisberg A. 2011. RAD in the realm of next-generation sequencing  
 329 technologies. *Molecular ecology* 20: 3499–3502.

330 Sodergren E, Weinstock GM, Davidson EH, Cameron RA, Gibbs RA, Angerer RC, Angerer LM,  
 331 Arnone MI, Burgess DR, Burke RD, Coffman JA, Dean M, Elphick MR, Etensohn CA,  
 332 Foltz KR, Hamdoun A, Hynes RO, Klein WH, Marzluff W, et al. 2006. The genome of the  
 333 sea urchin *Strongylocentrotus purpuratus*. *Science (New York, N.Y.)* 314: 941–952.

334 Toonen RJ, Puritz JB, Forsman ZH, Whitney JL, Fernandez-Silva I, Andrews KR, Bird CE. 2013.  
 335 ezRAD: a simplified method for genomic genotyping in non-model organisms. *PeerJ* 1:  
 336 e203.

337 Wang S, Meyer E, McKay JK, Matz M V. 2012. 2b-RAD: a simple and flexible method for  
 338 genome-wide genotyping. *Nature methods* 9: 808–810.

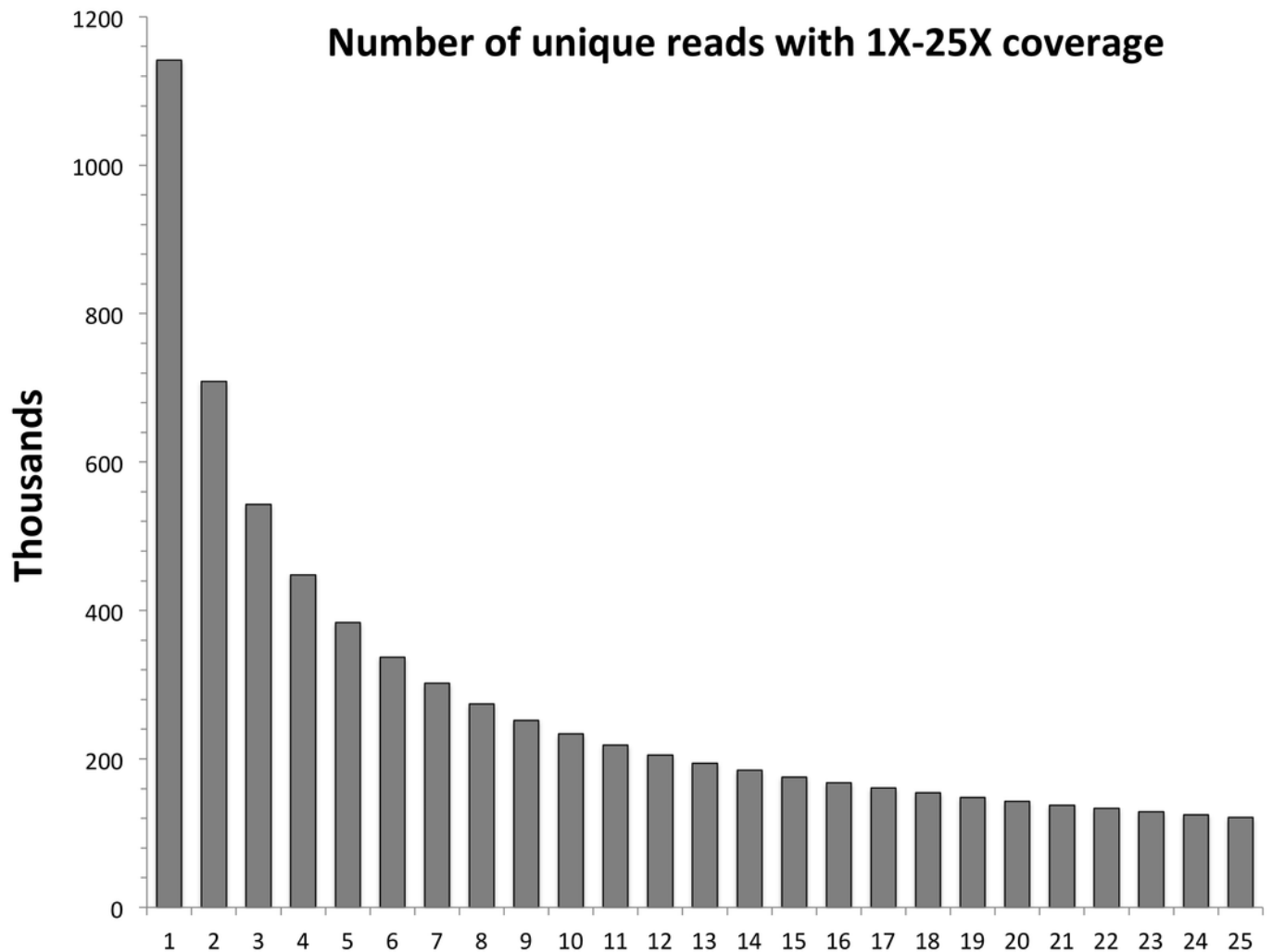
339 Waples RS. 1998. Separating the wheat from the chaff: patterns of genetic differentiation in high  
 340 gene flow species. *Journal of Heredity* 89: 438–450.

341 Ward RD, Woodward M, Skibinski DOF. 1994. A comparison of genetic diversity levels in  
 342 marine, freshwater, and anadromous fishes. *Journal of Fish Biology* 44: 213–232.

# Figure 1

Levels of coverage for each unique read in the red snapper data set.

The horizontal axis represents the minimal level of coverage, while the vertical axis represents the number of unique paired reads in thousands.



## Table 1 (on next page)

Results from individual runs of *dDocent* and *Stacks*.

*dDocent* runs varied in the level of similarity used to cluster reference sequences: A (90%), B (96%), and C (99%). For *Stacks*, forward reads and reverse reads were separately processed with *denovo\_map.pl* (*Stacks* version 1.08), using three different sets of parameters: A, minimum depth of coverage of two to create a stack, a maximum distance of two between stacks, and a maximum distance of four between stacks from different individuals; B, minimum depth of coverage of three to create a stack, a maximum distance of four between stacks, and a maximum distance of eight between stacks from different individuals; and C, minimum depth of coverage of three to create a stack, a maximum distance of four between stacks, and a maximum distance of 10 between stacks from different individuals. For *dDocent*, complex variants were decomposed into canonical SNP and Indel calls and Indel calls were filtered out. SNP calls were evaluated at different individual coverage levels: (i) total number of SNPs; (ii) number of SNPS called in 75%, 90%, and 99% at 3X coverage; (iii) number of SNPS called in 75% and 90% of individuals at 5X coverage; (iv) number of SNPS called in 75% and 90% of individuals at 10X coverage; and, (v) number of SNPS called in 75% and 90% of individuals at 20X coverage. Run times are in minutes. Results from forward and reverse reads of *Stacks* were combined for comparison with *dDocent*, which inherently calls SNPs on both reads.

Table 1. Results from individual runs of *dDocent* and *Stacks*. *dDocent* runs varied in the level of similarity used to cluster reference sequences: A (90%), B (96%), and C (99%). For *Stacks*, forward reads and reverse reads were separately processed with *denovo\_map.pl* (*Stacks* version 1.08), using three different sets of parameters: A, minimum depth of coverage of two to create a stack, a maximum distance of two between stacks, and a maximum distance of four between stacks from different individuals; B, minimum depth of coverage of three to create a stack, a maximum distance of four between stacks, and a maximum distance of eight between stacks from different individuals; and C, minimum depth of coverage of three to create a stack, a maximum distance of four between stacks, and a maximum distance of 10 between stacks from different individuals. For *dDocent*, complex variants were decomposed into canonical SNP and INDEL calls and INDEL calls were filtered out. SNP calls were evaluated at different individual coverage levels: (i) total number of SNPs; (ii) number of SNPS called in 75%, 90%, and 99% at 3X coverage; (iii) number of SNPS called in 75% and 90% of individuals at 5X coverage; (iv) number of SNPS called in 75% and 90% of individuals at 10X coverage; and, (v) number of SNPS called in 75% and 90% of individuals at 20X coverage. Run times are in minutes. Results from forward and reverse reads of *Stacks* were combined for comparison with *dDocent*, which inherently calls SNPs on both reads.

	<i>dDocent</i> A	<i>dDocent</i> B	<i>dDocent</i> C	<i>Stacks</i> A	<i>Stacks</i> B	<i>Stacks</i> C
	Red snapper					
Total 3X SNPS	53,298	53,316	53,361	28,817	33,479	53,298
75% 3X SNPs	21,195	20,990	20,724	4,150	5,735	21,195
90% 3X SNPs	9,102	8,850	8,639	675	987	9,102
99% 3X SNPs	78	47	15	-	-	78

75% 5X SNPs	14,881	14,594	14,339	2,632	4,351	14,881
90% 5X SNPs	5,021	4,925	4,785	179	579	5,021
75% 10X SNPs	7,556	7,318	7,154	783	1,618	7,556
90% 10X SNPs	1,414	1,340	1,286	7	48	1,414
90% IND 90% 5X	10,267	10,026	9,798	806	1,807	10,267
90% IND 90% 10x	4,242	4,093	3,974	129	441	4,242
Run time	41	41	42	70	47	41
Red drum						
Total 3X SNPs	46,378	46,688	46,832	45,792	50,821	46,378
75% 3X SNPs	36,745	36,905	36,900	24,134	28,991	36,745
90% 3X SNPs	32,356	32,424	32,330	13,439	17,946	32,356
99% 3X SNPs	11,906	11,910	11,774	828	1,264	11,906
75% 5X SNPs	34,279	34,393	34,336	21,021	26,526	34,279
90% 5X SNPs	28,532	28,566	28,431	10,494	15,282	28,532
75% 10X SNPs	27,523	27,605	27,488	12,928	17,018	27,523
90% 10X SNPs	19,434	19,442	19,283	4,159	6,734	19,434
75% 20X SNPs	15,080	15,111	14,981	2,276	3,538	15,080
90% 20X SNPs	7,365	7,409	7,304	243	1,974	7,365
Run time	43	45	45	58	55	43
Silk snapper						
Total 3X SNPs	68,668	68,825	68,861	48,742	55,505	68,668
75% 3X SNPs	30,771	30,391	30,051	7,596	9,705	30,771
90% 3X SNPs	14,952	14,673	14,415	2,007	3,439	14,952
99% 3X SNPs	4,294	4,060	3,952	132	527	4,294
75% 5X SNPs	20,534	20,188	19,968	4,789	7,290	20,534
90% 5X SNPs	9,103	8,750	8,533	1,225	2,573	9,103
75% 10X SNPs	9,765	9,400	9,159	2,094	3,547	9,765
90% 10X SNPs	3,923	3,691	3,490	489	1,224	3,923
75% 20X SNPs	4,069	3,832	3,624	703	1,415	4,069
90% 20X SNPs	1,431	1,313	1,228	136	417	1,431
Run time	88	95	59	93	89	88

# Figure 2

SNP results averaged across the three different run parameters for *dDocent* and *Stacks*.

(A) Red snapper, (B) Red drum, (C) Silk snapper (see Methods or Table 1 for SNP categories description). Error bars represent one standard error.

