

Author's accepted manuscript (postprint)

Customize and get the most out of your reduced-representation sequencing experiment with the new simulation software RADinitio

Choquet, M.

Published in: Molecular Ecology Resources

DOI: 10.1111/1755-0998.13218

Available online: 30 Jun 2020

This is the peer reviewed version of the following article: Choquet, M. (2021). Customize and get the most out of your reduced-representation sequencing experiment with the new simulation software RADinitio. *Molecular Ecology Resources*, 21, 351-354. doi: 10.1111/1755-0998.13218, which has been published in final form at <https://onlinelibrary.wiley.com/doi/epdf/10.1111/1755-0998.13218>. This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for Use of Self-Archived Versions.

1 Title:

2
3 **Customize and get the most out of your reduced-representation sequencing experiment**
4 **with the new simulation software *RADinitio***

5
6
7 Author:

8 **Marvin Choquet***

9 *Faculty of Biosciences and Aquaculture, Nord University, Bodø, Norway

10
11
12 ORCID:

13 Marvin Choquet <https://orcid.org/0000-0001-6719-2332>

14
15
16 Correspondence:

17 Dr Marvin Choquet

18 Faculty of Biosciences and Aquaculture, Nord University

19 Campus Mørkved, 8049 Bodø, Nordland, Norway

20 marvin.choquet@nord.no

21
22
23
24 **Whole genome sequencing is still often a difficult, costly and time-consuming task. The**
25 **emergence of various genome reduced-representation sequencing (RRS) protocols such as**
26 **restriction site-associated DNA sequencing (RADseq) has facilitated the access to genome-**
27 **wide information, without the need for whole genome sequencing. Reaching the full**
28 **potential of RRS protocols though, requires adjustments and tailoring to the species under**
29 **investigation. To that end, simulation software have been developed to guide researchers**
30 **in the customization of their RADseq experiment, but the extent to which these tools**
31 **mimic the behavior of a protocol in generating sequencing data is limited. In this current**
32 **issue of *Molecular Ecology Resources*, Rivera-Colón et al. (2020) introduce *RADinitio*, a**
33 **new software for simulating RADseq data designed to perform simulations at the highest**
34 **level of representativeness. By taking into account the effects of library preparation and**
35 **sequencing parameters on the resulting sequences, *RADinitio* allows the precise**
36 **identification of the sources of failure when designing a RADseq experiment. This new**
37 **software represents a considerable advance in RADseq data simulation and will likely lead**
38 **to increased success in RADseq experiments.**

39
40 Keywords: Genotyping-by-sequencing, RADseq, population genomics, next-generation
41 sequencing, non-model species

46 Over the last decade, the rapid development of genome reduced-representation
47 sequencing (RRS) protocols has revolutionized the fields of molecular ecology, evolutionary
48 genetics and conservation genetics. With diverse procedures, RRS protocols aim to reduce
49 the complexity of a genome by sampling a fraction of it, sufficient to address various
50 biological questions, and much easier to sequence and analyze compared to a whole
51 genome. These protocols offer a solution for exploring population genomics at a reasonable
52 cost in model organisms and in non-model organisms often left understudied due to their
53 genome complexity.

54

55 Restriction site-associated DNA sequencing (RADseq) is currently one of the most
56 popular approaches and consists in using restriction enzyme(s) to digest a genome, followed
57 by the sequencing of restriction sites flanking regions. The appeal of this method relies on
58 its applicability to theoretically any type of organism, with or without prior genomic
59 resources associated (Davey & Blaxter, 2010). A variety of RADseq protocols have been
60 described (Andrews et al, 2016) and in spite of reviews assessing pros and cons related to
61 the use of each method (e.g. Andrews & Luikart, 2014), figuring out the most fitting protocol
62 for a new experiment may still be challenging. Standardized protocols exist though, such as
63 the ezRAD (Toonen et al, 2013), and those are often perceived as attractive due to their
64 reported experimental ease of implementation. However, the relevance of data yielded by
65 any protocol, even standardized, to answer a biological question cannot be ensured, except
66 through prior prospective data simulation.

67

68 The article by Rivera-Colón et al. (2020) in the current issue of *Molecular Ecology*
69 *Resources* introduces the new simulation software *RADinitio*. The function of *RADinitio* is to
70 simulate behaviors of different RADseq protocols and parameters of library preparation on
71 a particular species in generating sequencing data. Thereby the user can make informed
72 decisions on how to tailor a specific RADseq experiment. The software *RADinitio* can be
73 downloaded from <https://pypi.org/project/radinitio/> and is used via a command-line
74 interface. Its development emanated from the authors' view that two sources of error often
75 impede researchers from reaching the full potential of success with their RADseq
76 experiment. The first challenge relies in the selection of a protocol (e.g. single-digest or
77 double-digest RADseq) that may be suboptimal for a specific organism. The second
78 challenge lies in the processes of library preparation (i.e. quality of starting template, choice
79 of restriction enzyme) and sequencing (i.e. coverage given). Other simulation software were
80 developed in the past, but *RADinitio* represents an important step forward in RADseq data
81 simulations due to its capacity to generate variants actually relevant for population genetic
82 studies and due to the inclusion of parameters of library preparation and sequencing in its
83 simulations.

84

85 Rivera-Colón et al. (2020) describe *RADinitio* as a three-step pipeline. The user first
86 needs to feed the software with a reference sequence that will be used to simulate a
87 metapopulation. Genetic variants will then be generated via the coalescent simulator
88 *msprime* (Kelleher et al, 2016) following a demographic model that can be defined by the
89 user to match as precisely as possible the study system. In the case when the user does not
90 know the details of the underlying model, *RADinitio* uses reasonable default parameters.
91 The second step consists in digesting *in silico* the reference genome provided by the user
92 with either one or two restriction enzymes, selected by the user, thereby generating a series

93 of RAD loci across the genome. This set of loci is intersected with the set of genetic variants
94 so as to keep only the variants present within RAD loci. At this step, *RADinitio* reproduces a
95 very famous characteristic of RADseq by taking into account the possibility of mutations at
96 restriction sites (following mutation and recombination rates defined by the user), and
97 subsequently saving only sequences with intact cut-sites, thus simulating the allele dropout
98 effect. Thirdly, paired-end sequences are generated from the extracted pool of RAD alleles.
99 Again, *RADinitio* reflects what happens in reality by including read duplicates to imitate the
100 effect of PCR amplification, following the number of PCR cycles defined by the user, which
101 allows exploring the impact of performing more or less PCR cycles. Sequences are produced
102 according to a sequencing coverage determined by the user and random sequencing errors
103 are added to each individual read pair, mirroring Illumina sequencing patterns of error
104 rates.

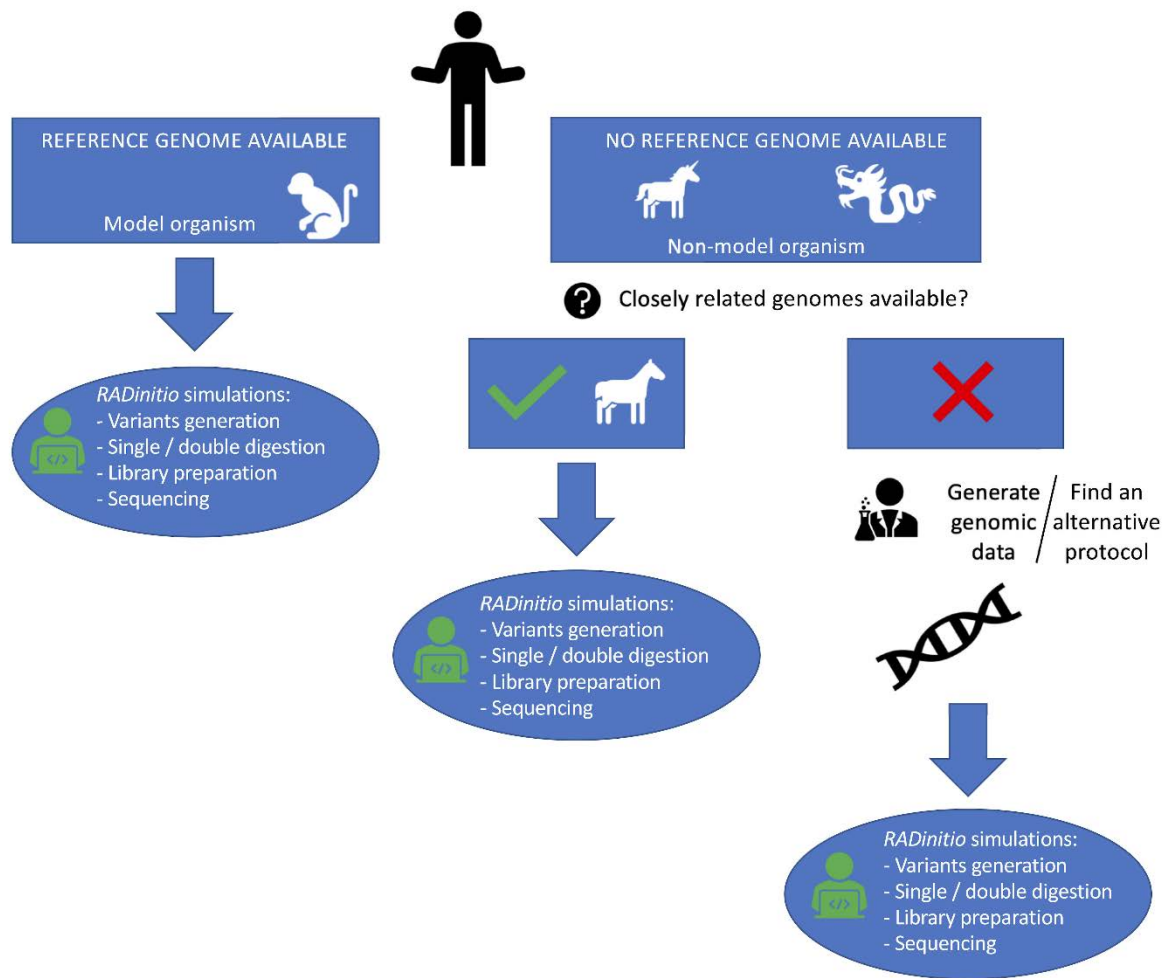
105
106 The level of detail thought by the authors to make *RADinitio* mimic the effects of several
107 variables (via parameters defined by the user) on data generated by a RADseq experiment
108 brings RADseq simulation software to a new level. By considering allele dropout resulting
109 from natural polymorphism in *RADinitio* simulations on a set of species, Rivera-Colón et al.
110 (2020) were able to characterize for the first time the degree of contribution of library
111 preparation and sequencing parameters to total allele dropout. They found that using a bad
112 template quality and a suboptimal sequencing coverage influence much more the amount
113 of allele dropout compared to the effect of natural polymorphism at enzyme restriction
114 sites. The software *RADinitio* is therefore the first of its kind that will guide users in keeping
115 allele dropout to a minimum.

116
117 As mentioned above, the usefulness of *RADinitio* depends on the availability of a
118 reference sequence, which may or may not be available depending on the species under
119 investigation. For a model species where a genome is available, the process is
120 straightforward (Fig. 1). By testing different genomes with different conditions, Rivera-Colón
121 et al. (2020) concluded that *RADinitio* simulations resemble empirical RAD loci. The user can
122 then confidently use the *RADinitio --tally-rad-loci* command to calculate the number of
123 genetic variants expected to be yielded in a RADseq experiment depending on the protocol
124 and parameters chosen. Rivera-Colón et al. (2020) also tested the scenario of a non-model
125 species with no reference genome available, using the example of a salmonid fish. In that
126 situation, the authors concluded that *RADinitio* can still perform informative simulations
127 under the condition that several genomes from a variety of related species of different
128 evolutionary distances are used as input (Fig. 1). There are, however, numerous cases of
129 non-model organisms for which not even one closely related genome is available, making
130 tailoring of a RADseq experiment very difficult. This is particularly striking in the field of
131 marine zooplankton, left largely ignored by genomics, where despite the huge diversity of
132 organisms reported, only a very few genomes are published (Bucklin et al, 2018). One
133 alternative may thus be to generate sequencing data from the species of interest prior to
134 the RADseq experiment in order to feed *RADinitio* with at least some genomic data (Fig. 1).
135 The inconvenient of such practice is that it may be difficult to assess how much data will be
136 necessary for the simulations to be accurate in regard to the actual genome of the target
137 species, and *RADinitio* was not tested for this. Besides, the success of a RADseq experiment
138 cannot be guaranteed if data simulation is not feasible or not accurate, in which case other
139 RRS protocols may need to be considered or developed. A case study performed on the

140 non-model zooplankton species *Calanus finmarchicus*, known for its large genome,
 141 illustrates the potential challenges linked to the absence of a reference genome when trying
 142 to simulate RADseq data (although *RADinitio* was not available back then) (Choquet et al,
 143 2019). In that study, a specific RRS protocol relying on target capture had to be developed
 144 to achieve generation of variants despite starting with no reference genome, via the
 145 sequencing of a draft-transcriptome instead.

147 To conclude, the new simulation software *RADinitio* is a promising tool that should be
 148 used before starting any RADseq experiment in species where a reference genome is
 149 available, or at least several closely related genomes. The level of representativeness
 150 implemented in *RADinitio* simulations will help users customize their experiment to get the
 151 most out of it. This pipeline represents a substantial advance for the field of RRS and
 152 particularly for RADseq users.

153
 154



155
 156 **Figure 1:** Different scenarios when starting a RADseq experiment

157 References:

158

159 Andrews KR, Good JM, Miller MR, Luikart G, Hohenlohe PA (2016) Harnessing the power of
160 RADseq for ecological and evolutionary genomics. *Nature Reviews Genetics*, **17**, 81.

161 Andrews KR, Luikart G (2014) Recent novel approaches for population genomics data
162 analysis. *Molecular Ecology*, **23**, 1661-1667.

163 Bucklin A, DiVito KR, Smolina I *et al.* (2018) Population genomics of marine zooplankton. In:
164 *Population Genomics: Marine Organisms* (ed. Rajora Om P & Oleksiak MF), pp. 61-102.
165 *Springer*.

166 Choquet M, Smolina I, Dhanasiri AK *et al.* (2019) Towards population genomics in non-
167 model species with large genomes: a case study of the marine zooplankton *Calanus*
168 *finmarchicus*. *Royal Society Open Science*, **6**, 180608.

169 Davey JW, Blaxter ML (2010) RADSeq: next-generation population genetics. *Briefings in*
170 *functional genomics*, **9**, 416-423.

171 Kelleher J, Etheridge AM, McVean G (2016) Efficient coalescent simulation and genealogical
172 analysis for large sample sizes. *PLoS computational biology*, **12**.

173 Toonen RJ, Puritz JB, Forsman ZH *et al.* (2013) ezRAD: a simplified method for genomic
174 genotyping in non-model organisms. *PeerJ*, **1**, e203.

175