# Fantastic beasts and how to sequence them:
## genomic approaches for obscure model organisms

Mikhail V. Matz

Department of Integrative Biology, University of Texas at Austin, Austin, TX 78712, USA

matz@utexas.edu

*Summary:*

Application of genomic approaches to "obscure model organisms" (OMOs), meaning species with little or no genomic resources, enables increasingly sophisticated studies of genomic basis of evolution, acclimatization and adaptation in real ecological contexts. Here, I highlight sequencing solutions and data handling techniques most suited for genomic analysis of OMOs.

---

*Glossary:*

- **Allele Frequency Spectrum, AFS** (same as Site Frequency Spectrum, SFS): histogram of the number of segregating variants depending on their frequency in one or more populations.
- **Restriction site-Associated DNA (RAD) sequencing**: family of diverse genotyping methods that sequence short fragments of the genome adjacent to recognition site(s) for specific restriction endonuclease(s).
- **Linkage Disequilibrium (LD)**: in this review, correlation of genotypes at a pair of markers across individuals.
- **LD block**: typical distance between markers in the genome across which their genotypes remain correlated.
- **Genome scan:** profiling of genotypes along the genome looking for unusual patterns. Often used to look for signatures of natural selection or introgression.
- **"Denser-than-LD" genotyping**: genotyping of several polymorphic markers per LD block.
- **Highly contiguous reference**: genome or transcriptome reference sequence containing the least amount of fragmentation.
- **Phased data**: data showing which SNP alleles belong to the same homologous chromosome copy.
- **Cross-tissue gene expression analysis**: looking for individual-specific shifts in gene expression detectable across multiple tissues. Such shifts are predominantly genetic in nature.

---

The focus of this review is mainly on the type of sequencing data required and how to obtain it in the most cost-efficient way rather than on analytical approaches. That said, I could not help but mention some highly promising analytical methods that are not yet broadly adopted by OMO researchers, such as demographic inference based on allele frequency spectra and annotation-independent analyses of gene expression data.

44

45 I will start with the summary of general types of questions in OMO studies and corresponding

46 data types required. We might be interested in the following four layers of genomic information,

47 each requiring a specific type of experimental and reference data:

48

49 1. Genome-wide patterns of neutral variation. This data can elucidate population structure,

50 population sizes, and migration rates, as well as changes of these parameters through time. This

51 analysis benefits from high quality genotype calls but does not require dense genome coverage; it

52 can even be performed in the absence of reference genome.

53

54 2. Regions in the genome particularly affected by non-drift processes (natural selection,

55 introgression, etc). This type of analysis, typically referred to as "genome scanning", takes

56 genome-wide neutral variation as baseline and looks for regions in the genome exhibiting highly

57 dissimilar patterns. It requires "denser-than-LD" genotyping and a highly contiguous reference

58 (see Glossary) to make sure no signal is overlooked.

59

60 3. Genome-wide gene expression, an extremely information-rich resource reflecting both

61 environmental and genetic variation. Streamlined transcript counting methods represent a cost-

62 efficient alternative to the industry-standard RNA-seq for generating quantitative data. Analysis

63 of gene expression does not require a genome reference, although a transcriptome reference must

64 be generated at some point. The reference does not have to be highly contiguous.

65

66 4. Epigenetics, here limited to DNA methylation. A variety of methods have been recently

67 developed that can generate data for DNA methylation analysis. For vertebrates, genome

68 reference is needed, but for other animals or plants, in which DNA methylation is much less

69 prevalent and predominantly occurs in exons, transcriptome or exome presents a good cost-

70 efficient alternative. The reference does not have to be highly contiguous.

71

72 **Genome-wide neutral variation**

73

74 *Allele Frequency Spectrum analysis*

75

76 Neutral genetic markers are traditionally analyzed using summaries of allele frequency

77 differences between populations, such as $F_{ST}$. The large amount of markers accessible through

78 next-generation sequencing opened up the possibility to dramatically enhance this approach by

79 modeling the evolution of the whole allele frequency spectrum (AFS, see Glossary). AFS

80 represents a rich source of information to fit alternative models with time-resolved population

81 sizes and migration rates as parameters (Box 1) based on coalescent simulations (*fastsimcoal2*,

82 [1]), diffusion approximation (*dadi*, [2]), or ordinary differential equations (*moments*, [3]). Model

83 selection is then based on likelihood ratio tests or Akaike information criterion. The new *moments*

84 method is particularly promising, as it is substantially faster than its predecessors and includes

85 built-in bootstrap, demographic model plotting, and capacity to analyze up to five populations

86 simultaneously. It is also very helpful that *moments* inherits the python code structure well

87 familiar to *dadi* practitioners.

2

88
89    *Experimental data*
90
91    The data required for AFS analysis is several thousand biallelic neutral single nucleotide
92    polymorphisms (SNPs). Ideally, SNPs must not be closely physically linked in the genome to
93    represent independent data points, although it is fully appropriate to analyze linked SNPs with
94    AFS methods. The lack of requirement for contiguous SNP coverage makes various flavors of
95    restriction site-associated DNA (RAD) sequencing (recently reviewed in [4,5], see Glossary) well
96    suited for this analysis. In our experience, *dadi* [2] and *moments* [3] work robustly with 5-10
97    thousand SNPs (a typical RAD output) when analyzing populations individually or in pairs.
98    Fitting models with three (*dadi*) or more (*moments, fastsimcoal2*) populations might be
99    problematic with this relatively low number of SNPs but is usually not required for OMOs (Box
100   1). Recent population size changes are often of special interest in OMOs; since they
101   predominantly affect rare alleles, their robust detection requires 20 or more high-quality
102   genotypes per population [6]. This preference for more individuals rather than more SNPs per
103   individual is an additional factor that makes cost-efficient RAD the approach of choice for AFS-
104   based analysis. That said, relatively low number of independent (unlinked) SNPs generated by
105   some RAD protocols might limit the power of the AFS analysis, and a good subject for a future
106   study would be the effect of the number of unlinked SNPs on AFS model selection and
107   uncertainties of parameter estimates. In this regard it is worth noting that RAD flavors differ
108   considerably in the number of unlinked loci in the genome that they interrogate [4,5].
109
110   For demographic inference, the AFS data must be filtered to exclude potential sites under
111   selection. Whichever test is used to identify such sites (for example, Bayescan, [7]), for their
112   removal the false discovery rate should be set as high as 0.5 to ensure purging of the majority of
113   non-neutral sites. Although under this setting half of the removed sites would be neutral, their
114   removal will not affect the overall AFS as long as the removed fraction does not comprise more
115   than 1-2% of the total number of sites.
116
117   *Genotyping quality*
118
119   In diploids, the most common genotyping error is missing one of the alleles in a heterozygote (i.e.,
120   a false homozygote call); and the next most common error is missing the whole SNP locus
121   entirely. Both these "missing data" errors are due to insufficient sequencing coverage, the
122   problem that is pervasive in today's OMO studies. Such errors strongly affect AFS in the region
123   of rare alleles, which is unfortunate since rare alleles are the most informative about recent
124   population history [6,8]. A telltale sign of poor heterozygote calling is under-representation of
125   singletons, but frequencies of doubletons and higher-order frequency bins are also distorted,
126   which has strong effect on AFS itself and inferred demographic parameters until mean
127   sequencing coverage approaches ~10x [9]. When coverage is 10x or higher a good way to filter
128   data is to select SNPs genotyped in >90-95% of samples [10]; importantly for RAD approach,
129   this would select SNPs that are unlikely to be affected by null alleles due to mutation in the
130   restriction endonuclease recognition site [4]. For obvious reasons, for AFS analysis genotype
131   calls should never be quality-filtered based on allele frequencies (for example, retaining only

3

132     variants that are detected in a minimum of two individuals or requiring minor allele frequency to
133     exceed some cutoff). A robust empirical way to evaluate the consistency of genotype calls is to
134     compare results for independently processed biological samples of the same genotype [11]. Such
135     genotyping replicates are quite feasible in RAD and are also useful to identify true SNPs for
136     training variant quality score recalibration model of the GATK pipeline [12]. For low-coverage
137     data (<10x), a general solution is provided by the *ANGSD* package [13], which generates AFS as
138     well as other population genetic statistics based on genotype likelihoods without actually calling
139     genotypes [14]. This method generates unbiased single-population AFS even with 2x coverage
140     [9]. Still, there is a concern that high variation in coverage across samples and populations might
141     affect *ANGSD* statistics; to avoid this potential issue it is recommended to discard the lowest-
142     coverage outliers and down-sample reads from highest-coverage outliers (J. Ross-Ibarra, pers.
143     comm.).
144
145     *PCR duplicates*
146
147     Presence of PCR duplicates in many early RAD applications has been repeatedly highlighted as a
148     source of genotyping errors [4,15] due to induced over-dispersion of read counts among alleles
149     and loci. Interestingly, the proportion of PCR duplicates does not depend on the number of PCR
150     cycles performed during library preparation. Instead, it depends on the ratio between the number
151     of reads sequenced ($N_r$) and the number of unique fragments present in the sample prior to PCR
152     ($N_o$): the fraction of duplicates is the same as expected when sampling $N_r$ from $N_o$ with
153     replacement. Fortunately, PCR duplicates are easy to identify and remove using degenerate tags
154     ligated to RAD fragments prior to amplification [16]. Most present-day RAD protocols now
155     implement this simple deduplication procedure, including the current version of 2bRAD [11].
156
157     *Genome reference for AFS analysis*
158
159     A great advantage of RAD-based AFS analysis for OMOs is that SNPs can be called based on
160     RAD reads themselves, without the need for genome reference. Several *de novo* RAD genotyping
161     pipelines have been developed, such as STACKS, pyRAD, and UNEAK (see references in [4])
162     that work for most RAD flavors, plus a similarly structured pipeline for 2bRAD
163     (https://github.com/z0on/2bRAD_denovo) that takes into account the fact that in 2bRAD either
164     strand of the locus can be sequenced. Still, using a reference genome to call RAD genotypes
165     provides three important advantages. First, it identifies physically linked (and thus potentially
166     non-independent) groups of SNPs, to be resampled as units during AFS bootstrap. The second
167     advantage is particularly important for OMOs sampled in the field: mapping to reference genome
168     automatically discards reads from contaminant DNA sources (viruses, bacteria, ingested food,
169     symbionts etc). To be able to discard such contaminants in *de novo* RAD pipeline the experiment
170     must include at least one sample generated from a clean source and consider only the RAD loci
171     observed in that sample.
172
173     The third advantage of reference-based genotyping is the possibility to discriminate between
174     ancestral and derived SNP alleles, to attain the best power of AFS-based inference. Counter-
175     intuitively, the best reference for AFS analysis is not a genome of the species under investigation

176    but a genome of a related outgroup species, separated from the focal one by a few million years
177    of evolution, because the SNP state as in the outgroup can be assumed to represent the ancestral
178    state (e.g., [17]). If the reads are mapped to the same-species genome, to identify ancestral states
179    of the variants a single well-sequenced RAD sample of an outgroup taxon could be included. The
180    analysis will then be limited to sites that can be successfully genotyped both in ingroup and
181    outgroup; in effect, the result is going to be the same as when mapping the reads from whole
182    project to an outgroup's genome. Although some proportion of ancestral states will be
183    misidentified due to incomplete lineage sorting, convergence or technical artifacts, this error is
184    easy to account for by including a single additional parameter into the model, specifying the
185    proportion of the AFS that needs to be flipped when predicting the data (e.g., [18]). The reference
186    for AFS does not have to be highly contiguous; the contigs should be just long enough to cover a
187    typical LD block for meaningful bootstrapping.
188
189    **Genome scanning**
190
191    Since outlier regions by definition occupy only a small portion of the genome and typically do not
192    form a single cluster, their confident detection requires "denser-than-LD" genotyping (see
193    Glossary). It has been argued that in most situations, RAD-like approaches would sample the
194    genome too sparsely to satisfy this requirement [19,20]. Although many successful genome scans
195    based on RAD have been published [21], RAD cannot be recommended for genome scanning
196    since it inevitably leaves considerable fraction of the genome unexplored. Even when LD is
197    known to be extensive enough for RAD to produce "denser-than-LD" genotyping, a better
198    solution might be to take full advantage of the extended LD and go instead for low coverage
199    whole-genome sequencing (WGS) followed by imputation, to obtain full-genome phased data
200    (Table 1).
201
202    The types of sequencing approaches for genome scanning with their pros and cons are
203    summarized in Table 1. Importantly, all of them require highly accurate reads mapped to a
204    reference for confident SNP detection, making short Illumina reads the genotyping data type of
205    choice. Some of the very promising approaches that have not yet been fully adopted for OMOs
206    are exome-seq and ultra-low whole genome sequencing (WGS) with imputation. Exome-seq used
207    to be a prerogative of model organisms because of the need for exome-capture platform
208    development, but it has recently been shown that OMO exome can be captured just as efficiently
209    using bead-bound normalized cDNA obtained from the OMO itself (EecSeq Puritz 2017). Such
210    "home-made exome" sequencing could become an excellent alternative to RAD since it would
211    interrogate essentially all the interpretable genetic variation for a comparably low cost.  Ultra-low
212    WGS with imputation used to require extensive reference haplotype panels available only for
213    well-established model organisms. However, several methods have been recently developed
214    (most notably STITCH, [22]) that can impute phased genotypes and correct genotyping errors in
215    ultra-low coverage data without relying on reference panels. Still, their applicability for each new
216    OMO must be experimentally confirmed because the success of imputation critically depends on
217    multiple polymorphisms occurring within a typical LD block, and whether this is so is not known
218    for OMOs *a priori*. Demographic events such as strong recent bottleneck, domestication, or
219    recent colonization would make imputation more efficient because of more extensive LD and

220     small number of founding haplotypes [22], and conversely, in large outbred populations
221     imputation will be less accurate and might require sequencing of a very large number (thousands)
222     of individuals. The accuracy of imputation can be evaluated by sequencing a few individuals at
223     high coverage (>10x) to generate high-confidence genotype calls and then attempting to impute
224     them based on sub-sampled read sets to emulate low coverage. It must be noted that it is
225     inappropriate to measure imputation accuracy by imputing genotype calls masked in high-quality
226     datasets (as in, for example, [23]): masked data do not contain false homozygote calls and
227     therefore do not correctly represent the real-life situation.
228
229     **Gene expression**
230
231     There are many aspects to gene expression, of which I here focus on just one: abundance
232     or protein-coding (polyadenylated) transcripts. The reason is that transcript abundance is by far
233     the most interpretable and it can be very easily analyzed in OMOs.
234
235     *Counting transcripts instead of resequencing them*
236
237     Typical RNA-seq [24] resequences the whole transcriptome in each sample, but there is a
238     much more economic way to count abundances of protein-coding transcripts: sequence just a
239     single fragment per each transcript molecule and count reads corresponding to each gene. TagSeq
240     [25], for example, sequences a single randomly generated fragment near the 3'-end of the
241     transcript, which is the most economic use of sequencing effort and removes bias towards longer
242     transcripts. In a recent benchmarking study TagSeq was actually more accurate than the standard
243     RNA-seq in measuring transcript abundances, despite nearly tenfold lower cost [26]. More
244     recently introduced QuantSeq [27] is conceptually similar to TagsSeq: it also sequences a single
245     randomly generated fragment near the 3'-end of each transcript but has a different library
246     preparation procedure, implemented as a kit from Lexogen (https://www.lexogen.com/quantseq-
247     3mrna-sequencing/). Bioinformatics analysis for both TagSeq ad QuantSeq is highly simplified
248     compared to typical RNA-seq. TagSeq was originally designed for OMOs and so its pipeline uses
249     transcriptome rather than genome as a reference to attribute reads to genes
250     (https://github.com/z0on/tag-based_RNAseq). One notable feature of the current version of
251     TagSeq pipeline is that it includes removal of PCR duplicates based on adaptor-derived
252     degenerate tags [11], similarly to 2bRAD and for the same reason – to avoid PCR-associated
253     over-dispersion or read counts.
254
255     *Analysis of gene expression "beyond gene lists"*
256
257     The unfortunate tradition that OMO research inherits from the biomedical field is putting
258     too much emphasis on possible functional implications of expression changes of specific genes.
259     For OMOs, this is bound to remain inconclusive because gene annotations are often absent,
260     tentative or based predominantly on similarity to human genes, which may or may not serve the
261     same function in the OMO. Even greater problem is interpretation bias: too often researchers
262     focus primarily on genes that "make sense" and ignore the rest. This leads to conclusions

263  reflecting predominantly the researchers' idea of what *should* be going on rather than what is
264  actually happening.
265
266       Table 2 lists alternative ways of objective analysis of gene expression data that are
267  enabled by the large sample sizes feasible with TagSeq or QuantSeq. They either do not require
268  gene annotations or rely sufficiently general functional summaries to be robust to occasional
269  missing or mis-annotations. Particularly useful for OMOs are analyses that use gene expression
270  patterns as anonymous multivariate readouts to compare and classify samples, such as principal
271  coordinate analysis (PCoA) or differential analysis of principal components  (DAPC). Related
272  multivariate analyses to visualize and classify genome-wide gene expression data, recently
273  reviewed in [28], have become the mainstream tool of single-cell RNA-seq, where they are used
274  to discover cell types and quantify differences between them. With appropriate experimental
275  design, in OMOs these analyses can lead to much more definitive biological conclusions than
276  studies scrutinizing long lists of differentially expressed genes passing a certain significance
277  cutoff.
278
279  *Gene expression as functional summary of genotype*
280
281  Gene expression is best known for its context-dependence reflecting phenotypic plasticity, which
282  is the view inherited from biomedical research dealing with genetically uniform models. In
283  natural populations, one of the most important sources of gene expression variation is genetic
284  difference among individuals, manifested as context-*in*dependent, individual-specific deviations
285  in gene expression. This is easy to demonstrate in clonally replicated organisms such as corals. In
286  two reciprocal transplantation experiments performed on different coral species in different
287  oceans, stable between-genotype differences accounted for more than 50% of total gene
288  expression variation despite transplantation of clonal fragments for up to a year to highly
289  dissimilar sites [29,30]. In non-clonal model organisms such as mice or humans, the best
290  demonstration of the effect of genetic variation on gene expression are abundant differences in
291  expression between alleles of the same gene [31,32]. In humans, fixed between-population
292  differences are exemplified by hundreds of genes that are differentially expressed between
293  African and European Americans [33]. All this suggests that gene expression can be a proxy of
294  not only phenotypic plasticity and acclimatization, but of genetic variation and adaptation. A
295  major advantage of the use of gene expression for these types of studies is that gene expression
296  integrates over many functionally relevant variants in the genome and thus represents a
297  condensed functional summary of the genotype [34].
298
299  In humans, nearly half of all genetic variants affecting gene expression have detectable effects in
300  all tissues [32], and so one feasible way to separate genotype-specific gene expression from
301  context-dependent variation might be to perform "cross-tissue" comparison (see Glossary) to
302  isolate body-wide expression shifts [35]. In the coming years, cross-tissue or similar analysis is
303  likely to become a major approach to study functional genetic variation in natural populations.
304
305

**Epigenetics**

Among many covalent chromatin modifications I will discuss DNA methylation since it currently receives the most attention in OMOs. Still it must be mentioned that in plants histone methylation appears be no less and perhaps even more involved in acclimatization and transgenerational plasticity [36]. While vertebrates show high methylation throughout the genome, invertebrates and plants methylate their genomes sparsely and mostly in protein-coding regions (so-called gene body methylation, GBM, [37]). The function of his ubiquitous and evolutionarily ancient DNA modification remains unclear [38,39] and the greatest challenge in the next few years will be to decipher it. The most important questions are: (i) Does GBM affect gene expression? (ii) Can it be modified on ecological timescale, to achieve acclimatization to a novel environment? (iii) Can acquired changes in GBM be transmitted across generations? If the answers to all three questions are "yes", then we have a mechanism for transgenerational inheritance of acquired traits, which is an exciting (albeit tentative, [40]) possibility. Table 3 summarizes the methods for generating DNA methylation data. If every gene in the genome has to be interrogated, MBD-seq and meDIP provide the best resolution for sequencing effort [38]. If the goal of the study is to characterize general methylation patterns rather than identify specific genes, highly cost-efficient solutions are provided by RRBS-seq and methylRAD. For studies requiring single base resolution, the best approach appears to be direct detection by PacBio or ONT – however, these exciting developments still require validation in complex genomes.

**Generating a reference sequence**

For all approaches described here, the accuracy of the reference sequence in terms of per-base error rate must only be high enough to allow unambiguous mapping of high-accuracy (Illumina) reads. The gold standard of genome sequence quality, Q30 or 99.9% accuracy, would not provide any benefit compared to a rough draft accuracy of 99%. Occasional errors in the reference would manifest themselves as SNPs that are not polymorphic in the analyzed samples and therefore irrelevant for analysis. This is the same reason why it is possible to use a genome of a related species as a reference.

For AFS analysis, which does not require highly contiguous reference, even a rough genome draft that can be assembled from a single lane worth of 150b paired-end reads from Illumina HiSeq would be suitable. However, substantially better options are now becoming available for a comparably low price tag. The technology offered by 10x Genomics [41] attaches specific barcodes to short reads originating from the same long DNA fragment, which allows assembling Illumina HiSeq data into very long haplotypes. The two single-molecule long-read "third-generation sequencing" methods, Single Molecule Real Time (SMRT) sequencing by PacBio and nanopore sequencing by ONT, produce reads with broad length distribution, including exceedingly long ones (tens to hundreds of kilobases) resulting in a qualitatively more contiguous genome assemblies [42–45] (Table 4, see [43] for recent benchmarking study of assembly pipelines). At the moment of this writing, read accuracy and cost of data for PacBio (Sequel system) and ONT (R9 flow cell) were equivalent; PacBio generated higher proportion of long reads than ONT; however, PacBio's library prep required ten fold more high-quality DNA than

350  ONT. Both for PacBio and ONT it is critically important to obtain high molecular weight DNA in
351  fragments exceeding 20kb in length. For new OMOs, it is also essential to confirm that the DNA
352  is accessible to enzymatic modifications by trying to digest it with a frequent-cutting restriction
353  endonuclease.
354
355  For genome scanning, gene expression, or invertebrate DNA methylation analyses targeting
356  protein-coding sequences (exome) genome sequence might not be the best reference; instead, a
357  highly contiguous transcriptome assembly would be preferable. Until now the standard way to
358  generate a *de novo* transcriptome was to perform high-coverage RNA-seq and assemble the
359  results with Trinity [46]. In the coming years, it is expected that even higher-quality and lower-
360  cost OMO transcriptomes would be generated by PacBio or ONT sequencing of full-length
361  cDNA (or, for ONT, direct mRNA sequencing). The long-read capacity of these technologies
362  would essentially obviate the need for assembly, leaving only the sequence correction procedure
363  to be performed.
364
365  Finally, which tissue or body part to sample for sequencing? For genome sequencing it does not
366  matter much as long as contamination by other DNA sources can be kept to a minimum, but for
367  *de novo* transcriptomics it is not a trivial question, as gene expression varies dramatically across
368  tissues and life cycle stages. In mammals, there is definitely an organ of choice that expresses
369  nearly all genes in the genome: testis. Rather unexpected transcriptome complexity in the testis is
370  putatively due to chromatin re-packaging during spermatogenesis, which results in genome-wide
371  transcription leakage [47]. If so, testis might be a good choice for *de novo* transcriptomics not
372  only for mammals but for any organism that produces compact sperm.
373
374  **Note on data sharing**
375
376  As we have seen, the best power of ecological genomics in OMOs is achieved using a genome or
377  transcriptome reference. Every new reference dataset enables new biological questions, and the
378  whole OMO field will get a great boost if these resources are promptly shared. Please consider
379  rapidly sharing your reference data, at least as soon as the initial preprint of your paper is posted
380  to bioRxiv and ideally sooner, by distributing the link to data through research-related email list
381  or professional twitter feed.
382

383  **References**

384  1    Excoffier, L. *et al.* (2013) Robust demographic inference from genomic and SNP data.
385       *PLoS Genet.* 9, e1003905
386  2    Gutenkunst, R.N. *et al.* (2009) Inferring the joint demographic history of multiple
387       populations from multidimensional SNP frequency data. *PLoS Genet.* 5, e1000695
388  3    Jouganous, J. *et al.* (2017) Inferring the Joint Demographic History of Multiple
389       Populations: Beyond the Diffusion Approximation. *Genetics* DOI:
390       10.1534/genetics.117.200493
391  4    Andrews, K.R. *et al.* (2016) Harnessing the power of RADseq for ecological and
392       evolutionary genomics. *Nat. Rev. Genet.* 17, 81–92
393  5    Puritz, J.B. *et al.* (2014) Demystifying the RAD fad. *Mol. Ecol.* 23, 5937–5942
394  6    Robinson, J.D. *et al.* (2014) Sampling strategies for frequency spectrum-based population
395       genomic inference. *Bmc Evol. Biol.* 14, 254

396   7    Günther, T. and Coop, G. (2013) Robust identification of local adaptation from allele
397        frequencies. *Genetics* 195, 205–20
398   8    Schiffels, S. *et al.* (2016) Iron Age and Anglo-Saxon genomes from East England reveal
399        British migration history. *Nat. Commun.* 7, 1–9
400   9    Han, E. *et al.* (2014) Characterizing bias in population genetic inferences from low-
401        coverage sequencing data. *Mol. Biol. Evol.* 31, 723–35
402   10   Matz, M. V *et al.* (2017) Potential for rapid genetic adaptation to warming in a Great
403        Barrier Reef coral. *bioRxiv* at <http://www.biorxiv.org/content/early/2017/06/18/114173>
404   11   Dixon, G.B. *et al.* (2015) Genomic determinants of coral heat tolerance across latitudes.
405        *Science (80-. ).* 348, 1460–1462
406   12   Mckenna, A. *et al.* (2010) The Genome Analysis Toolkit: A MapReduce framework for
407        analyzing next-generation DNA sequencing data. *Genome Res.* 20, 1297–1303
408   13   Korneliussen, T. *et al.* (2014) ANGSD: Analysis of Next Generation Sequencing Data.
409        *BMC Bioinformatics* 15, 356
410   14   Nielsen, R. *et al.* (2012) SNP calling, genotype calling, and sample allele frequency
411        estimation from new-generation sequencing data. *PLoS One* 7, e37558
412   15   Andrews, K.R. and Luikart, G. (2014) Recent novel approaches for population genomics
413        data analysis. *Mol. Ecol.* 23, 1661–1667
414   16   Casbon, J.A. *et al.* (2011) A method for counting PCR template molecules with
415        application to next-generation sequencing. *Nucleic Acids Res.* 39, e81
416   17   Gojobori, J. *et al.* (2007) Adaptive evolution in humans revealed by the negative
417        correlation between the polymorphism and fixation phases of evolution. *Proc. Natl. Acad.*
418        *Sci. U. S. A.* 104, 3907–12
419   18   Tine, M. *et al.* (2014) European sea bass genome and its variation provide insights into
420        adaptation to euryhalinity and speciation. *Nat. Commun.* 5, 5770
421   19   Lowry, D.B. *et al.* (2017) Breaking RAD: An evaluation of the utility of restriction site
422        associated DNA sequencing for genome scans of adaptation. *Mol. Ecol. Resour.* 17, 142–
423        152
424   20   Tiffin, P. and Ross-Ibarra, J. (2014) Advances and limits of using population genetics to
425        understand local adaptation. *Trends Ecol. Evol.* 29, 673–680
426   21   McKinney, G.J. *et al.* (2017) RADseq provides unprecedented insights into molecular
427        ecology and evolutionary genetics: comment on Breaking RAD by Lowry *et al* . (2016).
428        *Mol. Ecol. Resour.* 17, 356–361
429   22   Davies, R.W. *et al.* (2016) Rapid genotype imputation from sequence without reference
430        panels. *Nat. Genet.* 48, 965–969
431   23   Money, D. *et al.* (2015) LinkImpute: Fast and Accurate Genotype Imputation for
432        Nonmodel Organisms. *G3 Genes, Genomes, Genet.* 5, 2383–2390
433   24   Morin, R.D. *et al.* (2008) Profiling the HeLa S3 transcriptome using randomly primed
434        cDNA and massively parallel short-read sequencing. *Biotechniques* 45, 81–94
435   25   Meyer, E. *et al.* (2011) Profiling gene expression responses of coral larvae (Acropora
436        millepora) to elevated temperature and settlement inducers using a novel RNA-Seq
437        procedure. *Mol. Ecol.* 20, 3599–3616
438   26   Lohman, B.K. *et al.* (2016) Evaluation of TagSeq, a reliable low-cost alternative for
439        RNAseq. *Mol. Ecol. Resour.* 16, 1315–1321
440   27   Moll, P. *et al.* (2014) QuantSeq 3 ' mRNA sequencing for RNA quantification. *Nat.*
441        *Methods* 11, 25
442   28   Rostom, R. *et al.* (2017) Computational approaches for interpreting scRNA-seq data.
443        *FEBS Lett.* DOI: 10.1002/1873-3468.12684
444   29   Kenkel, C.D. and Matz, M. V (2016) Gene expression plasticity as a mechanism of coral
445        adaptation to a variable environment. *Nat. Ecol. Evol.* 1, 14
446   30   Dixon, G.B. *et al.* (2014) Bimodal signatures of germline methylation are linked with gene

447            expression plasticity in the coral Acropora millepora. *BMC Genomics* 15, 1109

448    31    Crowley, J.J. *et al.* (2015) Analyses of allele-specific gene expression in highly divergent
449            mouse crosses identifies pervasive allelic imbalance. *Nat. Genet.* 47, 353–60

450    32    The GTEx Consortium (2015) The Genotype-Tissue Expression (GTEx) pilot analysis:
451            Multitissue gene regulation in humans. *Science (80-. ).* 348, 648–660

452    33    Melé, M. *et al.* (2015) The human transcriptome across tissues and individuals. *Science*
453            *(80-. ).* 348, 660–665

454    34    Gamazon, E.R. *et al.* (2015) A gene-based association method for mapping traits using
455            reference transcriptome data. *Nat. Genet.* 47, 1091–1098

456    35    Wheeler, H.E. *et al.* (2016) Survey of the Heritability and Sparse Architecture of Gene
457            Expression Traits across Human Tissues. *PLOS Genet.* 12, e1006423

458    36    Lämke, J. and Bäurle, I. (2017) Epigenetic and chromatin-based mechanisms in
459            environmental stress adaptation and stress memory in plants. *Genome Biol.* 18, 124

460    37    Feng, S. *et al.* (2010) Conservation and divergence of methylation patterning in plants and
461            animals. *Proc. Natl. Acad. Sci. U. S. A.* 107, 8689–94

462    38    Dixon, G.B.. *et al.* (2016) Evolutionary Consequences of DNA Methylation in a Basal
463            Metazoan. *Mol. Biol. Evol.* 33, msw100

464    39    Sarda, S. *et al.* (2012) The Evolution of Invertebrate Gene Body Methylation. *Mol Biol*
465            *Evol* 29, 1907–1916

466    40    Uller, T. *et al.* (2013) Weak evidence for anticipatory parental effects in plants and
467            animals. *J. Evol. Biol.* 26, 2161–2170

468    41    Kitzman, J.O. (2016) Haplotypes drop by drop. *Nat. Biotechnol.* 34, 296–298

469    42    Gordon, D. *et al.* (2016) Long-read sequence assembly of the gorilla genome. *Science (80-*
470            *. ).* 352, aae0344

471    43    Michael, T.P. *et al.* (2017) High contiguity Arabidopsis thaliana genome assembly with a
472            single nanopore flow cell. *bioRxiv* at
473            <http://www.biorxiv.org/content/early/2017/06/14/149997>

474    44    Jansen, H.J. *et al.* (2017) Rapid de novo assembly of the European eel genome from
475            nanopore sequencing reads. *bioRxiv* at
476            <http://www.biorxiv.org/content/early/2017/01/20/101907>

477    45    Berlin, K. *et al.* (2015) Assembling large genomes with single-molecule sequencing and
478            locality-sensitive hashing. *Nat. Biotechnol.* 33, 623–630

479    46    Haas, B.J. *et al.* (2013) De novo transcript sequence reconstruction from RNA-seq using
480            the Trinity platform for reference generation and analysis. *Nat Protoc* 8, 1494–1512

481    47    Soumillon, M. *et al.* (2013) Cellular Source and Mechanisms of High Transcriptome
482            Complexity in the Mammalian Testis. *Cell Rep.* 3, 2179–2190

483    48    Cosart, T. *et al.* (2011) Exome-wide DNA capture and next generation sequencing in
484            domestic and wild species. *BMC Genomics* 12, 347

485    49    Bundock, P.C. *et al.* (2012) Enrichment of genomic DNA for polymorphism detection in a
486            non-model highly polyploid crop plant. *Plant Biotechnol. J.* 10, 657–667

487    50    De Wit, P. *et al.* (2015) SNP genotyping and population genomics from expressed
488            sequences - current advances and future possibilities. *Mol. Ecol.* 24, 2310–2323

489    51    Bay, R.A. and Palumbi, S.R. (2014) Multilocus adaptation associated with heat resistance
490            in reef-building corals. *Curr. Biol.* 24, 2952–6

491    52    Futschik, A. and Schlötterer, C. (2010) The Next Generation of Molecular Markers From
492            Massively Parallel Sequencing of Pooled DNA Samples. *Genetics* 186, 207–218

493    53    Kofler, R. *et al.* (2011) PoPoolation2: identifying differentiation between populations
494            using sequencing of pooled DNA samples (Pool-Seq). *Bioinformatics* 27, 3435–3436

495    54    Therkildsen, N.O. and Palumbi, S.R. (2017) Practical low-coverage genomewide
496            sequencing of hundreds of individually barcoded samples for population and evolutionary
497            genomics in nonmodel species. *Mol. Ecol. Resour.* 17, 194–208

498   55   Paradis, E. *et al.* (2004) APE: Analyses of phylogenetics and evolution in R language.
499        *Bioinformatics* 20, 289–290
500   56   Oksanen, J. *et al.* vegan: Community Ecology Package. , *R package ver. 2.4–3.* (2017)
501   57   Wright, R.M. *et al.* (2015) Gene expression associated with white syndromes in a reef
502        building coral, Acropora hyacinthus. *BMC Genomics* 16, 371
503   58   Strader, M.E. *et al.* (2016) Red fluorescence in coral larvae is associated with a diapause-
504        like state. *Mol. Ecol.* 25, 559–569
505   59   Jombart, T. (2008) adegenet: a R package for the multivariate analysis of genetic markers.
506        *Bioinformatics* 24, 1403–5
507   60   Langfelder, P. and Horvath, S. (2008) WGCNA: an R package for weighted correlation
508        network analysis. *BMC Bioinformatics* 9, 559
509   61   Bay, R.A. and Palumbi, S.R. (2017) Transcriptome predictors of coral survival and growth
510        in a highly variable environment. *Ecol. Evol.* 7, 4794–4803
511   62   Rose, N.H. *et al.* (2016) Gene Networks in the Wild: Identifying Transcriptional Modules
512        that Mediate Coral Resistance to Experimental Heat Stress. *Genome Biol. Evol.* 8, 243–
513        252
514   63   Lister, R. *et al.* (2009) Human DNA methylomes at base resolution show widespread
515        epigenomic differences. *Nature* 462, 315–22
516   64   Meissner, A. *et al.* (2008) Genome-scale DNA methylation maps of pluripotent and
517        differentiated cells. *Nature* 454, 766–770
518   65   Serre, D. *et al.* (2010) MBD-isolated Genome Sequencing provides a high-throughput and
519        comprehensive survey of DNA methylation in the human genome. *Nucleic Acids Res.* 38,
520        391–9
521   66   Jacinto, F. V. *et al.* Methyl-DNA immunoprecipitation (MeDIP): Hunting down the DNA
522        methylome. , *BioTechniques*, 44. (2008) , 35–43
523   67   Wang, S. *et al.* (2015) MethylRAD: a simple and scalable method for genome-wide DNA
524        methylation profiling using methylation-dependent restriction enzymes. *Open Biol.* 5,
525        150130
526   68   Flusberg, B.A. *et al.* (2010) Direct detection of DNA methylation during single-molecule,
527        real-time sequencing. *Nat. Methods* 7, 461–5
528   69   Feng, Z. *et al.* (2013) Detecting DNA Modifications from SMRT Sequencing Data by
529        Modeling Sequence Context Dependence of Polymerase Kinetic. *PLoS Comput. Biol.* 9,
530        e1002935
531   70   Stoiber, M.H. *et al.* (2017) De novo Identification of DNA Modifications Enabled by
532        Genome-Guided Nanopore Signal Processing. *bioRxiv* at
533        <http://www.biorxiv.org/content/early/2017/04/10/094672>
534   71   Koren, S. *et al.* (2017) Canu: scalable and accurate long-read assembly via adaptive k-mer
535        weighting and repeat separation. *Genome Res.* 27, 722–736
536   72   Chin, C.-S. *et al.* (2016) Phased diploid genome assembly with single-molecule real-time
537        sequencing. *Nat. Methods* 13, 1050–1054
538   73   Li, H. (2016) Minimap and miniasm: Fast mapping and de novo assembly for noisy long
539        sequences. *Bioinformatics* 32, 2103–2110
540   74   Vaser, R. *et al.* (2017) Fast and accurate de novo genome assembly from long uncorrected
541        reads. *Genome Res.* 27, 737–746
542   75   Walker, B.J. *et al.* (2014) Pilon: An integrated tool for comprehensive microbial variant
543        detection and genome assembly improvement. *PLoS One* 9, e112963
544
545

546

---

**Box 1: AFS models.**

In the world of OMOs we are usually dealing with samples from many populations, which would be hard or impossible to model simultaneously; moreover, there are usually many populations left unsampled. To infer meaningful demographic parameters in a sparsely sampled system of many populations, a practical solution is to perform two-dimensional AFS analysis of all population pairs [10]. Typical hypotheses and corresponding tests are:

- Are the two populations demographically separate?
    - compare model with split to model without split, under which the two compared populations are regarded as independent samples from the same population.
- If yes, is there still gene flow between them?
    - compare split models with and without migration.
- If yes, is the gene flow symmetric or asymmetric?
    - compare split model with two potentially different migration rates to a split model with a single symmetrical migration rate.
- Was population size stable or went through changes in the past?
    - compare single-population model involving population size change in the past to a standard neutral model.

Simple command-line scripts for AFS plotting and running basic pairwise models in *moments* can be found here: https://github.com/z0on/AFS-analysis-with-moments. To access the full potential of *moments*, however, the user is expected to compose python scripts of their own.

---

547    **Table 1. Genotyping approaches for genome scanning.**

| Approach | Features | Pros | Cons |
|---|---|---|---|
| Exome-seq [48,49] | Isolates and sequences only the protein-coding portion of genome. | Dense coverage of genes guarantees that coding variants and variants linked to cis-regulatory mutations are discovered. | Other (arguably less important) types of variation are not profiled (e.g., distant enhancers). |
| RNA-seq [50,51] | Sequences RNA. | Same as exome sequencing. | Genotyping quality of a gene depends on expression level. Allele-specific expression affects accuracy of heterozygote calls. |
| Pool-seq [52,53] | Sequences pooled DNA from multiple individuals from each population. | Dense whole-genome coverage with confident determination of allele frequencies in populations. | No possibility for individual–based analysis (such as STRUCTURE) or validation based on genotype-phenotype association across individuals. Must be confident in *a priori* population designations. |
| Low-coverage whole-genome sequencing (WGS) [54] | Sequences individual genomes at ~1-4x coverage. | Dense whole genome coverage at individual level. | Per-site genotypes are unreliable because of missing data; must use uncertainty-aware analysis such as *ANGSD*. |
| Ultra-low coverage WGS with imputation [22] | Sequences individual genomes at <2x coverage, imputes missing genotypes and corrects false homozygote calls | Dense whole genome coverage at individual level, phased data enables haplotype-based analysis | Rare alleles (minor allele frequency<0.05) are missed. Requires large sample sizes (depending on LD, hundreds or thousands of individuals). Accuracy of imputation must be experimentally validated for every new OMO. |

548
549

14

550 **Table 2. Gene expression analyses not relying on accurate gene annotations**

| Analysis | What does it do | Software | Applications |
|---|---|---|---|
| Principal coordinate analysis based on Manhattan distances (sum of all log-fold changes across genes) | Characterizes overall transcriptome differences across experimental groups. Measures fraction of variation attributable to each experimental factor. | R: package *ape*, function *pcoa* [55] package *vegan,* function *adonis* [56] | [57,58] |
| Differential analysis of principal components (DAPC) | Quantifies transcriptome differences between samples with respect to specified multivariate axis. Good for quantifying overall gene expression plasticity. | R: package *adegenet* [59] | [29] |
| Weighted gene co-expression network analysis (WGCNA) | Identifies co-regulated groups of genes, which are linked to experimental factors and traits *post hoc*. Method of choice for complex experimental designs (>20 samples) with many quantitative traits measured. | R: package WGCNA [60] | [29,61,62] |
| Rank-based functional summaries of KOG (euKaryotic Orthologous Groups) classes | Reveals broad functional trends in gene expression. Particularly useful for OMOs since it tolerates sparse and inaccurate annotations. Its main use is for statistical comparison of highly diverse datasets, even from different species. | R: package *KOGMWU* [11]. | [11,58] |

551

552

553 **Table 3. Methods for interrogating DNA methylation**

| Method | Features | Pros | Cons |
|---|---|---|---|
| Whole-Genome Bisulfite Sequencing (WGBS) [63] | Sequences complete genome after bisulfite conversion | Complete characterization of 5me-cythosine methylation at single-base resolution | High coverage is required to obtain quantitative data. In non-vertebrate OMOs, much sequencing effort is wasted since most of genome is not methylated. |
| RRBS-seq [64] | Bisulfite sequencing of genome fragments adjacent to all (methylated and un-methylated) CCGG sites | Saves costs dramatically compared to WGBS. | Only a fraction of all CpG sites is interrogated. Complicated library preparation protocol. Sequencing effort is wasted on non-methylated sites. |
| MBD-seq [65], meDIP [66] | Pull-down and sequencing of methylated DNA. | Optimizes sequencing effort by focusing on methylated DNA. | Complicated library preparation protocol. Resolution equals the length of pulled-down fragments (~300-500b). Pull-down procedure is not absolutely efficient, many reads still correspond to un-methylated genome regions. |
| methylRAD [67] | Direct sequencing of genomic fragments adjacent only to the methylated CCGG and CCWGG sites. | Very simple library prep protocol. Highly cost-efficient due to focus on methylated sites only. | Only a fraction of all CpG sites is interrogated. New method, requires further benchmarking. |
| PacBio [68,69] | Direct detection of modified DNA bases during normal SMRT sequencing, based on polymerase lags. | Robust detection of 4-methylcytosine, 8-oxoguanine, and N6-methyladenine. Single-base resolution. | Same as WGBS. 5-methylcytosine, the most common methylation mark in animals, is not reliably detected. |
| ONT [70] | Direct detection of modified DNA bases during normal nanopore sequencing, based on conductivity changes. | Detects all marks, including 5-methylcytosine. Single-base resolution. | Same as WGBS. |

554
555

556

557 **Table 4. Assembly pipelines for PacBio and ONT reads**

| Pipeline | Required coverage | Features | Pros | Cons |
|---|---|---|---|---|
| Canu +Quiver* [71] | >30x | Correct and trims reads before assembly. | Best accuracy at base, indel and assembly level. | Very computationally demanding for large genomes. Generates incomplete assemblies at low coverage. |
| Falcon + Quiver* [72] | >50x | Similar to Canu. | Standard for PacBio. | Very computationally demanding for large genomes. High reliance on reads >20kb. Highly incomplete assemblies at low coverage. |
| minimap + miniasm + racon [73,74] | <30x | Raw reads are assembled, correction is done post-assembly | Very fast even for large genomes. Works with lower coverage, shorter reads than Canu and Falcon. | The resulting accuracy is lower than with Canu + Quiver. |
| pilon [75] | NA (error correction method) | Performs additional correction post-assembly. | Boosts accuracy for any assembly. | Requires high-quality Illumina reads. |

558 *Quiver is a consensus polishing software that is now replaced by Arrow to handle PacBio
559 Sequel data (https://github.com/PacificBiosciences/GenomicConsensus ). Racon [74] can be used
560 instead of Quiver/Arrow [43].
561