



Peer Community In Evolutionary Biology

RESEARCH ARTICLE

 Open Access

 Open Peer-Review

 Open Data

 Open Code

1 A novel workflow to improve 2 multi-locus amplicon genotyping of 3 wildlife species: an experimental 4 set-up with a known model system

Cite as: Posted: XXX

Recommender:

François Rousset

Reviewers:

Thomas Bigot, Helena Westerdahl and
Sebastian Ernesto Ramos-Onsins

Correspondence:

mark.gillingham@uni-ulm.de;
mail@psc-santos.com

5 Mark A.F. Gillingham¹, B. Karina Montero^{1,2}, Kerstin Wihelm¹,
6 Kara Grudzus¹, Simone Sommer¹ & Pablo S.C. Santos¹

7 ¹ Institute of Evolutionary Ecology and Conservation Genomics, Ulm Universität – Ulm,
8 Germany

9 ² Zoological Institute, Animal Ecology and Conservation, Biocenter Grindel, Universität Ham-
10 burg – Hamburg, Germany

11 This article has been peer-reviewed and recommended by
12 *Peer Community In Evolutionary Biology*

13 **Keywords:** open-source genotyping pipeline; ACACIA; high-throughput sequencing; am-
14 plicon genotyping; allele dropout; PCR amplification bias; sequencing bias; multigene
15 family; MHC

16 Abstract

17 Genotyping novel complex multigene families is particularly challenging in non-model organ-
 18 isms. Target primers frequently amplify simultaneously multiple loci leading to high PCR
 19 and sequencing artefacts such as chimeras and allele amplification bias. Most genotyping
 20 pipelines have been validated in non-model systems whereby the real genotype is unknown
 21 and the generation of artefacts may be highly repeatable. Further hindering accurate geno-
 22 typing, the relationship between artefacts and genotype complexity (i.e. number of alleles
 23 per genotype) within a PCR remains poorly described. Here we investigated the latter by
 24 experimentally combining multiple known major histocompatibility complex (MHC) haplo-
 25 types of a model organism (chicken, *Gallus gallus*, 43 artificial genotypes with 2-13 alleles
 26 per amplicon). In addition to well defined “optimal” primers, we simulated a non-model
 27 species situation by designing “cross-species” primers, with sequence data from closely re-
 28 lated Galliforme species. We applied a novel open-source genotyping pipeline (ACACIA; https://gitlab.com/psc_santos/ACACIA), and compared its performance with another, previously
 29 published pipeline (AmplisAS). Allele calling accuracy was higher when using ACACIA (98.5% vs
 30 97% and 77.8% vs 75.2% for the “optimal” and “cross-species” datasets respectively). System-
 31 atic allele dropout of three alleles owing to primer mismatch in the “cross-species” dataset
 32 explained high allele calling repeatability (100% when using ACACIA) despite low accuracy,
 33 demonstrating that repeatability can be misleading when evaluating genotyping workflows.
 34 Genotype complexity was positively associated with non-chimeric artefacts, chimeric arte-
 35 facts (nonlinearly by leveling when amplifying more than 4-6 alleles) and allele amplification
 36 bias. Our study exemplifies and demonstrates pitfalls researchers should avoid to reliably
 37 genotype complex multigene families.

39 Introduction

40 A key challenge for molecular ecologists is that they frequently work on systems with limited
 41 to no knowledge of their genomes. Multigene complexes, such as resistance genes (R-genes)
 42 and self-incompatibility genes (SI-genes) in plants, immunoglobulin superfamily and major
 43 histocompatibility genes (MHC) in vertebrates, and homeobox genes in animals, plants and
 44 fungi, among many others, are particularly challenging to genotype in non-model organisms.
 45 Whilst a large number of *de novo* genomes has been published using short-sequencing tech-
 46 nology, thus far traditional genome assembly of non model organisms has not been able to as-
 47 semble the highly repetitive genomic regions of multigene families [46]. Recent long-read sin-
 48 gle molecule sequencing technologies offers a very promising avenue to characterise multi-
 49 gene families in the future [15, 19, 32, 66], but the high sequencing errors and the high cost
 50 associated with required sequencing depth continues to constrain characterisation of non
 51 model organisms, particularly if genotyping of a large number of individuals with highly com-
 52 plex multigene systems is required. Therefore, the development of a genotyping approach

for specific multigene families (i.e. amplicon-based genotyping) typically continues to rely on information from closely related species available in genetic databases which may have very different gene duplication and deletion events. As a result of high sequence similarity from recent gene duplication events, polymerase chain reaction (PCR) primers will frequently bind across multiple loci leading to the amplification of multiple allelic variants [3, 5, 6, 35, 36, 57, 61]. Assessing and validating genotyping methods can be particularly challenging when the number of loci targeted is unknown.

Unspecific locus amplification may lead to several biases during PCR since 1) chimeric sequences (hereafter “chimeras”; which may arise because of incomplete extension of sequences during a PCR cycle which are subsequently completed with a different allele template) are likely to become more frequent as more loci are amplified within an amplicon simply because there will be more gene variants from which chimeras can be generated [34]; 2) amplification bias of some gene variants relative to others may occur because primers preferentially bind to some alleles/loci (hereafter referred to as “PCR competition”) [38, 61]. Creative solutions in primer design and in PCR conditions, such as using pooled primers instead of degenerate primers [38], reducing the number of cycles and modifying elongation steps of PCRs [25, 34, 60], can significantly reduce amplification bias. However, even after the application of such methods, PCR biases will nonetheless persist and may lead to genotyping errors because: 1) chimeras may be difficult to distinguish from valid recombinant gene variants (frequent in multigene complexes [11]), resulting in either PCR artefacts being falsely validated as a true allelic variants (type I errors, hereafter referred to as “false positives”) or in true allelic variants being falsely rejected as an artefact (type II errors, hereafter referred to as “allele dropout”) and 2) poorly amplified allelic variants may not be sequenced resulting in allele dropout, particularly when the number of sequences per amplicon (a set of sequences of a target region generated within a PCR) is low [5, 16, 35, 36, 61].

The rapid dissemination of high-throughput DNA sequencing (HTS) platforms has provided molecular ecologists with an exciting opportunity to tackle the parallelised genotyping of multiple markers in numerous species, since it has allowed the generation of thousands of sequences (termed “reads”) per amplicon, at a fraction of the cost and time needed by previous methods, which typically involved laboriously isolating individual sequences via a cloning vector followed by Sanger sequencing [3, 36, 61]. However, HTS platforms have their own limitations, the most relevant being the relatively high amount of sequencing errors generated in a typical sequencing run [18, 23, 39, 55, 61]. For instance, Illumina, currently the mainstream technology for HTS amplicon sequencing, report an error rate (primarily substitutions of base pairs) of $\leq 0.1\%$ per base for $\geq 75\text{--}85\%$ of bases (see Glenn [18] for details), although final error rates are likely to be much higher and can reach up to 6% [39]. Indeed, previous genotyping studies in multi-locus-systems (>10) reported average amplification and sequencing artefact rates of 1.5% to 2.5% per amplicon [49, 52, 58]. Therefore, PCR competition when

93 amplifying multiple loci per amplicon means that sequences from some genuine allelic vari-
94 ants occur at a similar frequency to PCR artefacts or sequencing errors [5, 16, 35, 61]. In this
95 scenario, poorly amplified alleles cannot be easily distinguished from artefacts during allele
96 validation, leading to further false positives and allele dropout during genotyping.

97

98 The need to distinguish PCR and sequencing artefacts from valid allelic variants has led to
99 the development of multiple bioinformatic workflows (i.e. a set of bioinformatic steps during
100 processing of sequencing data which eventually leads to genotyping, hereafter referred to as
101 a “genotyping pipeline”). While all genotyping pipelines rely to some degree on the assump-
102 tion that artefacts are less frequent than genuine allelic variants, they vary in the approach
103 used to discriminate poorly amplified allelic variants from artefacts. Genotyping pipelines for
104 complex gene families have been extensively reviewed in Biedrzycka *et al* [5]. Recently devel-
105 oped pipelines cluster artefacts to their putative parental sequences thereby increasing the
106 read depths of true variants [36, 48, 57, 63]. Currently, the most commonly used pipeline for
107 MHC studies is the AmpliSAS web server pipeline [57]. After chimera removal, AmpliSAS uses
108 a clustering algorithm to discriminate between artefacts and allelic variants, which take into
109 account the error rate of a particular HTS technology and the expected lengths of the ampli-
110 fied sequences. This is achieved in a stepwise manner, whereby it first clusters the most com-
111 mon variant (according to specified error rates) and then moves on to the next most common
112 variant, until no variant remains to be clustered. Microbiome studies, which typically amplify
113 hypervariable regions of the 16S rRNA gene from very diverse bacterial communities within
114 a single amplicon, have used a similar strategy to AmpliSAS, whereby potential artefactual
115 variants are clustered to suspected parental sequences using Shannon entropy (referred to
116 as “Oligotyping” [14]) or other similar clustering methods [2, 10].

117

118 Most of the amplicon genotyping pipelines for multigene families available to molecular
119 ecologists have only been tested on non-model organisms for which the real genotype is un-
120 known (but see Sebastian *et al* [57]). As a consequence, studies have frequently depended on
121 repeatability of duplicated samples to justify genotyping pipeline reliability [5, 16, 36, 52, 57,
122 61]. However for a given set of PCR primers and sequencing technology, PCR and sequenc-
123 ing bias, and thus in turn the rate of false positives and allele dropout, will be consistently
124 repeatable [5]. For instance, the high rate of Illumina substitution errors are known to be
125 not random (see references within Sebastian *et al* [57]) and therefore variants which result
126 from substitution errors are highly repeatable between amplicons [5]. Furthermore, while
127 the generation of PCR and sequencing artefacts is well known, the precise relationship be-
128 tween artefacts and the number of alleles amplified within an amplicon for a given set of
129 primers and sequencing technology has never been described. Yet, having a clear indication
130 of this relationship is an important step in predicting what are the optimal pipelines settings
131 (e.g. predicting error rates) for a given number of loci amplified within an amplicon. The lat-
132 ter can only be achieved by experimentally manipulating the number of loci of *a priori* known

133 genotypes before PCR amplification and HTS sequencing.

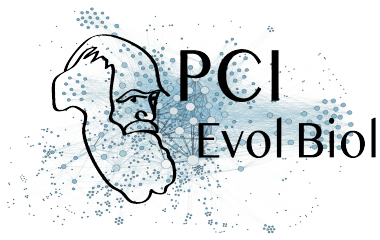
134

135 In this study, we artificially generated genotypes of known combinations of MHC alleles of
136 a model organism (the chicken, *Gallus gallus*) by mixing DNA samples from 7 haplotypes as
137 an example of a target multigene region of interest to molecular ecologists and to assess the
138 accuracy of amplicon-based genotyping. While we focus on the MHC hereafter, all methods
139 and results are applicable to any multigene family. Like many multigene complexes, MHC
140 genes are subject to multiple gene conversion, duplication and deletion [44, 45, 47] and MHC
141 gene copies vary considerably across and even within a species (reviewed in [30]). There-
142 fore, the number of MHC loci present in a non-model study system often remains unknown.
143 For instance, MHC class IIB copy number variation (CNV) was found to be as high as 21 in
144 some passerine species, resulting in up to 42 allelic variants amplified within an amplicon
145 and strong CNV between individuals [5]. In contrast, the chicken MHC B complex is unusu-
146 ally simple, leading it to be coined as a “minimal essential” system, with only two MHC class
147 I loci and two MHC class IIB loci [27–29]. The latter is therefore an ideal system to validate
148 MHC genotyping pipelines for the following reasons: 1.) the structure of the B complex is
149 well known with well-defined primers in conserved regions; 2.) the well characterised B com-
150 plex haplotype lineages can be used so that the expected MHC genotyping results are known
151 prior to sequencing and genotyping and 3.) The number of alleles amplified within an am-
152 plicon can be experimentally engineered by combining DNA samples from multiple MHC B
153 complex haplotypes.

154

155 To perform the genotyping of known chicken MHC haplotypes and extract data concerning
156 PCR and sequencing artefacts at each step of the genotyping workflow, we developed and cali-
157 brated our own genotyping pipeline (named ACACIA for **A**llele **C**alling pro**C**edure for **I**llumina
158 **A**mplicon sequencing data). ACACIA is written in Python and it takes advantage of several
159 previously published software dedicated to genomics (detailed in the methods), as well as
160 the widely used Biopython library [12] to handle genomic data. We experimentally gener-
161 ated a MHC dataset with a range of CNVs by combining DNA samples from multiple chicken
162 MHC B complex haplotypes. Since MHC B complex in chickens is well characterised, optimal
163 primers to amplify the entire exons which code for the antigen binding regions have been
164 developed [20, 59]. However in most wildlife species, such extensive genomic information
165 around the region of interest is unavailable. To avoid the problems associated with overfit-
166 ting ACACIA to one specific dataset and to replicate the challenge of designing primers for
167 a non-model species, we additionally designed primers within the exons coding for antigen-
168 binding regions using sequence data from closely related Galliforme species that were not
169 chickens (hereafter referred to as “cross-species” primers). The latter enabled us to gain in-
170 sight into the relative amount of artefacts generated by an intentionally sub-optimal set of
171 primers, for which we expected allele dropout.

172



173 Specifically, this study aimed to:

- 174 1. validate ACACIA using experimentally manipulated genotypes with different CNV that
175 are known *a priori*;
- 176 2. to investigate the relationship between multigene complexity (i.e. number of alleles am-
177 plified within an amplicon) and artefacts generated by PCR and sequencing (i.e. chimeras
178 and insertions/deletions)

179

180 **Materials and Methods**

181 **Samples and DNA extraction**

182 Chicken blood samples originated from experimental inbred lines kept at the Institute for An-
183 imal Health at Compton UK (lines 72, C, WL and N) and the Basel Institute for Immunology in
184 Basel Switzerland (lines H.B15 and H.B19+), as detailed in Jacob *et al* [24], Shaw *et al* [59] and
185 Wallny *et al* [65]. These lines carry seven common B haplotypes: B2 (line 72), B4 and B12 (line
186 C), B14 (line WL, sometimes referred as W), B15 (H.B15), B19 (H.B19) and B21 (line N). All the
187 lines used in this study are homozygotes (NCBI accession numbers: AJ248572 to AJ248586).
188 In each haplotype are two class IIB loci: BLB1 (previously known as BLBI or BLBminor) and
189 BLB2 (BLBII or BLBmajor), with alleles now designated as BLB1*02 and BLB2*02 from the B2
190 haplotype, etc. All alleles have different nucleotide sequences, except BLB1*12 and BLB1*19.
191 DNA was isolated from blood cells by a salting out procedure [42].

192

193 **Generating 43 artificial MHC genotypes**

194 We artificially generated 43 genotypes of varying CNV by combining equimolar amounts of
195 DNA samples from the seven MHC haplotypes mentioned above (Table 1; created genotypes
196 listed in Supplementary Table S1).

197

Table 1. In this study we generated 43 genotypes that were amplified twice (duplicated). The number of alleles per genotype and the number of genotypes with that number of all alleles are shown. The list of haplotypes used to artificially create the genotypes are listed in the Supplementary Table S1.

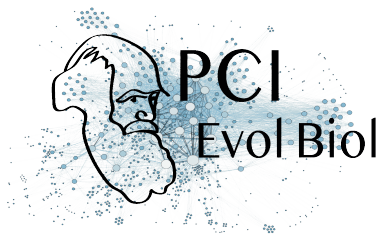
Number of alleles per genotype	Number of genotypes
2	7
4	7
6	7
8	7
10	7
11	5
12	2
13	1
Total	43

198 Optimal primers for chicken MHC class IIB

199 We targeted 241 bp of the 270 bp exon 2 of MHC class IIB, the polymorphic region known
 200 to code for antigen binding sites, using the primers OL284BL (5'-GTGCCCCGACGTTCTTC-3')
 201 and RV280BL (5'-TCCTCTGCACCGTGAAGG-3') [20]. The primers are not locus specific and bind
 202 to both loci of the chicken B complex.

204 Cross-species primer design for chicken MHC class IIB

205 To replicate designing primers without any *a priori* knowledge of the species MHC Class IIB
 206 structure or sequences, we downloaded 61 exon 2 MHC class IIB sequences from seven Gal-
 207 liform species (*Coturnix japonica*, *Crossoptilon crossoptilon*, *Meleagris gallopavo*, *Numida me-*
 208 *leagris*, *Pavo cristatus*, *Perdix perdix* and *Phasianus colchicus*, all accession numbers are listed
 209 in the Supplementary Table S2) from GenBank (<https://www.ncbi.nlm.nih.gov/genbank/>). We
 210 then used Primer3 [56, 64] to design the forward primer GagaF1 (5'-WTCTACAACCGGCAGCAGT-
 211 3') and the reverse primer GagaR2 (5'- TCCTCTGCACCGTGAWGGAC-3') aiming at amplifying
 212 151 bp of exon 2. No species were given more weight than others during primer design, and
 213 all default parameters of Primer3 (concerning melting temperatures and structural settings)
 214 were kept. The only exception is that we allowed up to two degenerate positions in the primer
 215 sequence.



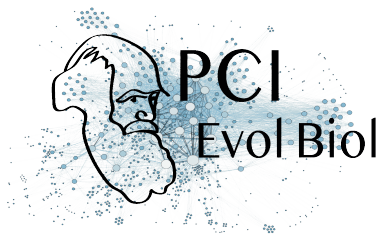
217 **PCR amplification, library preparation, and high-throughput sequencing**

218 For both datasets we replicated all individuals to estimate repeatability ($n_{\text{individuals}} = 43$ and
219 $n_{\text{amplicons}} = 86$). Individual PCR reactions were tagged with a 10-base pair identifier, using
220 a standardised Fluidigm protocol (Access Array™ System for Illumina Sequencing Systems,
221 ©Fluidigm Corporation). We first performed a target specific PCR with the CS1 adapter and
222 the CS2 adapter appended. To enrich base pair diversity of our libraries during sequencing,
223 we added four random bases to our forward primer. The CS1 and CS2 adapters were then
224 used in a second PCR to add a 10bp barcode sequence and the adapter sequences used by
225 the Illumina instrument during sequencing.

226
227 The first PCR consisted of 3–5 ng of extracted DNA, 0.5 units FastStart Taq DNA Polymerase
228 (Roche Applied Science, Mannheim, Germany), 1x PCR buffer, 4.5 mM MgCl_2 , 250 μM of each
229 dNTP, 0.5 μM primers, and 5% dimethylsulfoxide (DMSO). The PCR was carried out with an
230 initial denaturation step at 95°C for 4 min followed by 30 cycles at 95°C for 30 s, 60°C for 30 s,
231 72°C for 45 s, and a final extension step at 72°C for 10 min. The second PCR contained 2 μl of
232 the product generated by the initial PCR, 80 nM per barcode primer, 0.5 units FastStart Taq
233 DNA Polymerase, 1x PCR buffer, 4.5 mM MgCl_2 , 250 μM of each dNTP, and 5% dimethylsul-
234 foxide (DMSO) in a final volume of 20 μl . Cycling conditions were the same as those outlined
235 above but the number of cycles was reduced to ten. Reducing the number of PCR cycles, the
236 elongation time within PCR cycles and omitting the final extension step is recommended to
237 reduce the number of chimeras when co-amplifying multiple loci, because most incomplete
238 primer extensions which generate chimeras are thought to be formed in the final cycles of
239 PCRs and during the final extension step (see discussion) [25, 34, 60]. However we chose to
240 process samples using conventional PCR conditions, because a high number of cycles may be
241 necessary in some study systems and we wanted to replicate conditions used in most MHC
242 wildlife studies. Thus, we purposefully wanted to evaluate the robustness of our pipeline in
243 the more challenging setting where a high number of artefacts might be generated due to
244 sub-optimal PCR conditions.

245
246 PCR products were purified using an Agilent AMPure XP (Beckman Coulter) bead cleanup
247 kit. The fragment size and DNA concentration of the cleaned PCR products were estimated
248 with the QIAxcel Advanced System (Qiagen) and by UV/VIS spectroscopy on an Xpose instru-
249 ment (Trinean, Gentbrugge, Belgium). Samples were then pooled to equimolar amounts of
250 DNA. The library was prepared as recommended by Illumina (Miseq System Denature and
251 Dilute Libraries Guide 15039740 v05) and was loaded at 7.5 pM on a MiSeq flow cell with a
252 10% PhiX spike. Paired-end sequencing was performed over 2×251 cycles.

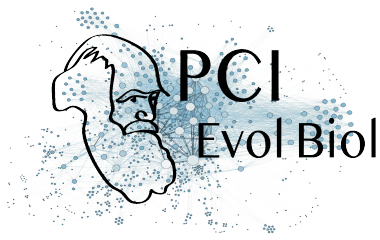
253



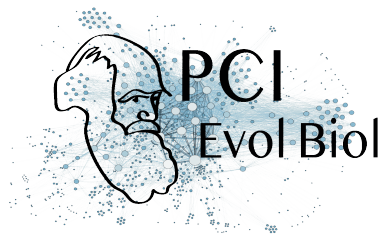
254 Data analysis with the ACACIA pipeline

255 ACACIA consists of 11 consecutive steps of data processing. The software requires two non-
256 standard python libraries (Pandas [40] and Biopython [12]) as well as six third-party soft-
257 ware (FastQC (www.bioinformatics.babraham.ac.uk/projects/fastqc/), FLASH [37], VSEARCH
258 [54], BLAST [1], MAFFT [26] and Oligotyping [14]), which can all be installed with one com-
259 mand. The input files are any number of FASTq files, which are the current canonical output
260 of the Illumina platform. The step-by-step workflow is described below:

- 261 1. **Generating Quality Reports.** Sequencing quality is assessed for each FASTq file yielded
262 by the sequencing platform, with the FastQC tool. Reports for each file are produced
263 in HTML format for visual inspection.
- 264 2. **Trimming low quality ends of forward and reverse reads (optional).** The informa-
265 tion generated in step #1 is crucial for an informed decision about how many (if any)
266 bases should be trimmed out of each read. If trimming is performed here, step #1 is
267 repeated. Shorter FASTq files are generated as output of this step.
- 268 3. **Merging paired-end reads (optional).** This concerns projects with paired-end sequenc-
269 ing only and should be skipped if using data from single-end sequencing (note: the
270 names of the paired forward and reverse FASTq files should be identical prior to the first
271 "_" character, e.g.: ID1S1L001_R1_001.fastq and ID1S1L001_R2_001.fastq). The reads of
272 file pairs are merged using FLASH [37]. The minimum and maximum lengths of over-
273 lap during merging can be adjusted by the user to improve performance (defaults are
274 zero and read length, respectively). New FASTq files with merged sequences are gen-
275 erated as output, as well as a series of .log files which allow users to monitor merging
276 performance.
- 277 4. **Trimming primers.** After prompting users to enter the sequences of the primers used
278 for target amplification, ACACIA trims primer sequences from both ends of the merged
279 sequences (IUPAC nucleotide ambiguity codes are allowed). Primerless sequences are
280 written into FASTq files which are the output of this step. The Python functions for
281 trimming primers and low-quality ends (step #2) are part of the core ACACIA pipeline.
282 External tools were avoided here to decrease dependency on further software.
- 283 5. **Quality-control.** Users are then prompted to enter the values of two parameters (q
284 and p) to filter sequences based on their mean phred-scores. First, q stands for qual-
285 ity and denotes a phred-score threshold that can take values from 0 to 40. Second,
286 p stands for percentage and denotes the proportion of bases, in any given sequence,
287 that have to achieve at least the quality threshold q for that sequence to pass the qual-
288 ity filter. ACACIA uses the default values $q = 30$ and $p = 90$ if users do not explicitly
289 change them. In practical terms, these thresholds correspond to an error probability
290 lower than 10^{-3} in at least 90% of bases for each sequence. All information on quality



- 291 data of sequences passing this filter is then removed and FASTA files with high-quality
292 sequences are given as the output of this step.
- 293 6. **Removing singletons.** A large proportion of sequences contain random errors inher-
294 ent to the sequencing technology [50]. To decrease file sizes without risking loss of
295 relevant allele information, ACACIA removes all singletons (sequences that appear one
296 single time) in an individual amplicon.
- 297 7. **Removing chimeras.** The chimera identification tool VSEARCH [54] is employed here,
298 with slightly altered settings (alignwidth = 0 and mindiffs = 1) aiming at increasing sen-
299 sitivity to chimeras that diverge very little from one of the “parent” sequences. FASTA
300 files with non-chimeric sequences, along with log files for each individual amplicon, are
301 given as output.
- 302 8. **Removing unrelated sequences.** All remaining sequences are then compared with
303 a set of reference sequences chosen by users. This step aims at removing sequences
304 that passed all filters so far but are products of unspecific priming during PCR. Typically,
305 sequences phylogenetically related to those being analyzed can be downloaded from
306 the GenBank (www.ncbi.nlm.nih.gov/genbank/). Users are prompted to provide one
307 FASTA file with reference sequences, which is converted by ACACIA to a local BLAST
308 database [1] and used for BLAST. Only sequences yielding high-scoring hits to the local
309 database ($E \leq 10^{-10}$) are written into new FASTA files as an output of this step, which is
310 the workflow’s last filtering procedure.
- 311 9. **Aligning.** The MAFFT aligner [26] is used to perform global alignments of sequences
312 that have passed filters. Since all sequences are pooled into one single alignment out-
313 put file, the individual IDs are now transferred from file names into the FASTA sequence
314 headers. We have successfully aligned up to 603,513 sequences in a desktop computer
315 with four CPUs and 32GB of RAM. Users with a significantly higher number of sequences
316 might find it useful to increase the computational parallelization of the aligner as de-
317 scribed recently [43].
- 318 10. **Calling candidate alleles.** The Oligotyping tool [14] is used to call candidate alleles.
319 Although originally conceived as a tool for identifying variants from microbiome 16S
320 rRNA amplicon sequencing projects, we recognised Oligotyping as ideal for other forms
321 of highly variable amplicon sequencing projects. This step consists of concatenating
322 high-information nucleotide positions (defined by entropy analysis of the alignment
323 produced in the previous step) and subsequently using entropy information to cluster
324 divergent variants, while grouping redundant information and filtering out artefacts.
325 Although Oligotyping was conceived as a supervised tool, we automated the selection
326 of parameter values aiming at high tolerance. This has the advantage of running an
327 unsupervised instance of Oligotype as a pipeline step, at the cost of keeping potential



328 false positives among the results. Report files with a list of candidate alleles grouped
329 by individual amplicons are the output of this step.

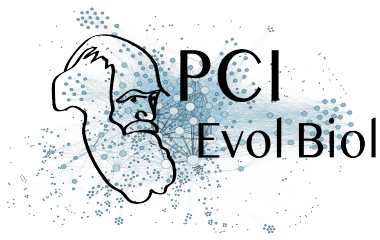
330 11. **Allele calling and final reporting.** A Python script is used to perform the final allele
331 calling by filtering out Oligotyping results according to the following criteria:

- 332 • Removal of unique allele variants (Y/N). Setting Y (yes) removes all alleles identified
333 in one single individual amplicon.
- 334 • Absolute number of reads (abs_nor): minimum number of sequences that need to
335 support an allele, otherwise the allele is considered an artefact. Ranges between
336 0 and 1000, with default = 10.
- 337 • Lowest proportion of reads (low_por): to be called in an individual amplicon, an
338 allele needs to be supported by at least the proportion of reads, within that in-
339 dividual amplicon, that is declared here. Ranges between 0 and 1, with default
340 = 0, while a value greater than 0 is recommended for data sets with ultra deep
341 sequencing depth, which can suffer more from false positives [5].

342 Subsequently, putative alleles with very low frequency (both at the individual and popu-
343 lation level) are scrutinised again. If the proportion of reads of a putative allele within an
344 individual amplicon is less than 10 times lower than the next higher ranking allele, and if it is
345 very similar (one single different base) to another, more frequent allele present in the same
346 individual amplicon, that putative allele is considered an artefact and removed. Finally, if an
347 individual amplicon has fewer than 50 sequences following all of the allele calling validation
348 steps, it is eliminated. Users are able to change all parameter values, but ACACIA recom-
349 mends settings based on our benchmarking. The output of this step consists of four files:

- 350 • **allelelereport.csv**: a brief allele report listing genotypes of all individual amplicons as
351 well as frequencies and abundances of all alleles found in the run;
- 352 • **allelelereport_XL.csv**: a detailed allele report including the number of reads supporting
353 each allele both within individuals and in the population;
- 354 • **pipelinereport.csv**: a pipeline report quantifying read counts and sequences failing or
355 passing each pipeline step described above;
- 356 • **alleles.fasta**: a FASTA sequence file of all alleles identified in the run.

357 To evaluate the pipeline, we calculated both allele calling accuracy and allele calling re-
358 peatability. Allele calling accuracy was calculated as the percentage of alleles that have been
359 correctly called across replicates. This was done by comparing the predicted genotype to
360 the genotype generated by ACACIA. All alleles that were dropped out or false positives were
361 marked as inaccurately called alleles. Allele calling repeatability on the other hand was cal-
362 culated as the percentage of alleles called in both replicates (including false positives). Note



here that allele calling accuracy and repeatability are not necessarily correlated if allele calling errors are highly repeatable, i.e. either false positives or allele dropout are consistent across replicates. However, allele calling accuracy can only be calculated if the genotype is known *a priori* which will not be available to most wildlife studies. We have therefore calculated both measures to investigate the pitfalls of relying on allele repeatability to validate a genotyping pipeline.

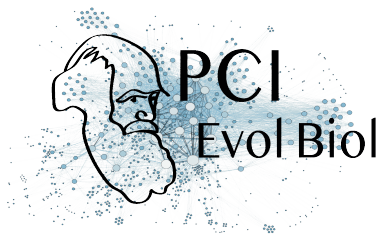
We investigated the best `abs_nor` and `low_por` settings for our datasets by first looking at the allele calling accuracy and repeatability at varying `abs_nor` values (range: 0-40, with `low_por` set at 0) first, and at varying `low_por` values (range: 0-0.02, with the optimal `abs_nor`, in our case 10) second. The latter is how we recommend users to find their optimal settings, although the range of `abs_nor` and `low_por` values to be investigated may vary across different datasets, depending on where the “peak” optimal setting lies.

The pipeline is supervised by a configuration text file (`config.ini`) which is appended every time users enter one of the settings mentioned above. Users can avoid running ACACIA interactively (and run the whole workflow in a “hands-free” mode) by providing a complete `config.ini` file at the beginning of the workflow. A template of a `config.ini` file is given in ACACIA's repository (https://gitlab.com/psc_santos/ACACIA/blob/master/config.ini).

Data analysis with the AmpliSAS pipeline

To compare how ACACIA performed relative to an existing relevant pipeline, we applied the web server AmpliSAS pipeline to our chicken datasets [57]. The default AmpliSAS parameters of a substitution error rate of 1% and an indel error rate of 0.001% for Illumina data was used. We then tested for the optimal ‘minimum dominant frequency’ clustering threshold for a given filtering threshold (i.e. 0.5% for the ‘minimum amplicon frequency’), by testing a set of thresholds of 10%, 15%, 20% and 25%. All clustering parameters tested gave an allele calling accuracy of 97%, but we chose the 25% clustering threshold because it was the only parameter which resulted in no false positives.

Subsequently, AmpliSAS filters for clusters that are likely to be artefacts, including chimeras and other low frequency artefacts that have filtered through the clustering step [57]. The default setting for the filtering of low frequency variants (i.e. ‘minimum amplicon frequency’) is 3%. However this value was far too high for our datasets, and we tested a range of filtering threshold between 0% and 1% at 0.1% intervals (i.e. 0%, 0.1%, 0.2% etc.). We assessed the optimal filtering threshold using both allele calling accuracy and repeatability.



400 Statistical analyses

401 To analyse the relationship between the number of alleles amplified and amplification bias/artefacts,
 402 we used generalized additive mixed models using the "mgcv" package [68] in R version 3.6.3
 403 [51]. The three response variables that were explored using a binomial error distributions
 404 corrected for over-dispersion (aka quasibinomial) were: proportion of reads assigned to an
 405 allele, proportion of reads that were non-chimeric artifacts, proportion of reads that were
 406 chimeras. The response variables number of chimeric variant and number of parental vari-
 407 ants that were generating chimeric variants were analysed using a Poisson error distribution,
 408 with the latter corrected for over-dispersion (aka quasipoisson). The fixed term number of
 409 alleles amplified was entered as a smoother and was limited to 6 estimated degrees of free-
 410 dom. To control for pseudo-replication, sample ID was entered as a random factor.

411 Results

412 Sequencing depth for each dataset and proportion of artefacts detected using ACACIA

413 A total of 530,101 paired-end reads were generated for the optimal primers dataset, which
 414 amounted to an average of 6,164 reads per amplicon ($n = 86$). For the cross-species primers
 415 dataset, 994,338 paired-end reads were generated, amounting to an average of 11,562 reads
 416 per amplicon ($n = 86$). The proportion of artefacts identified at each step of the ACACIA
 417 pipeline for the chicken datasets combined is illustrated in Figure 1. Workflow filtering re-
 418 moved the highest proportion of reads when filtering for singletons (13.6%) and chimeras
 419 (14.2%). After all filters, 66.4% of the original raw reads were used for allele calling.

421 Optimal settings of different workflows

422 We compared allele calling repeatability across a range of different `abs_nor` and `low_por` set-
 423 tings when using the ACACIA workflow to identify the optimal settings according to genotyp-
 424 ing accuracy for our datasets. We first fixed the `abs_nor` setting at 10 and tested different
 425 `low_por` values and found that the optimal setting (i.e. the highest accuracy values) was 0
 426 across both datasets (Figure 2a.). Setting higher `low_por` values resulted in a higher allele
 427 dropout rate, which led to lower accuracy and repeatability scores. We then tested the opti-
 428 mal `abs_nor` setting for a fixed `low_por` value of 0 and found that the optimal setting was 10
 429 across both datasets (Figure 2b.). An `abs_nor` value of 0 increased the rate of false positives,
 430 whilst a value above 10 increased the rate of allele dropout.

432 For the AmpliSAS workflow, we investigated the optimal filtering threshold and found dif-
 433 fering optimal values between datasets. For the optimal primer dataset, we found that the
 434 optimal filtering threshold was 0.3, whilst 0.5 was found to be optimal for the cross-species

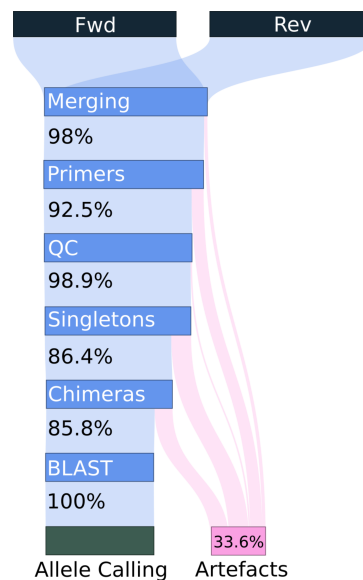


Figure 1. Flow diagram of reads and sequences from the two Illumina runs (a first run for the optimal primer dataset and the second for the cross-species dataset) analysed with ACACIA. Black bars denote the number of initial, raw reads (i.e. 100% of the reads generated by the Illumina runs). Blue bars correspond to filtering steps, and the percentages given correspond to the proportion of sequences from the previous step that were kept for the next stage of the workflow. The percentage given at the bottom right (Artefacts) refers to the total percentage of reads that were filtered from the total initial reads generated by Illumina, prior to any filtering steps. (Fwd & Rev) raw forward and reverse reads; (Merging) paired-end read merger, which includes a first quality filter; (Primers) primer trimming step, which also removes sequences lacking full primers; (QC) quality control; (Singletons) Singleton removal; (Chimeras) chimera removal; (BLAST) BLAST filter.

435 primer dataset (Figure 2c).

436

437 **AmplisAS vs ACACIA: optimal primers dataset**

438 When using the optimal settings of the ACACIA workflow, comparison of results with expected
 439 genotypes revealed that nine alleles dropped out, no false positives were found (Table 2) and
 440 as a result allele calling accuracy was 98.5% (Figure 2a. and b.). All instances of allele dropout
 441 derived from the B21 haplotype. For two genotypes, both BLB1*21 and BLB2*21 dropped out.
 442 For four genotypes, only BLB2*21 dropped out and for one genotype only BLB1*21 dropped
 443 out (Table 2). Allele calling repeatability was 97.7%.

444

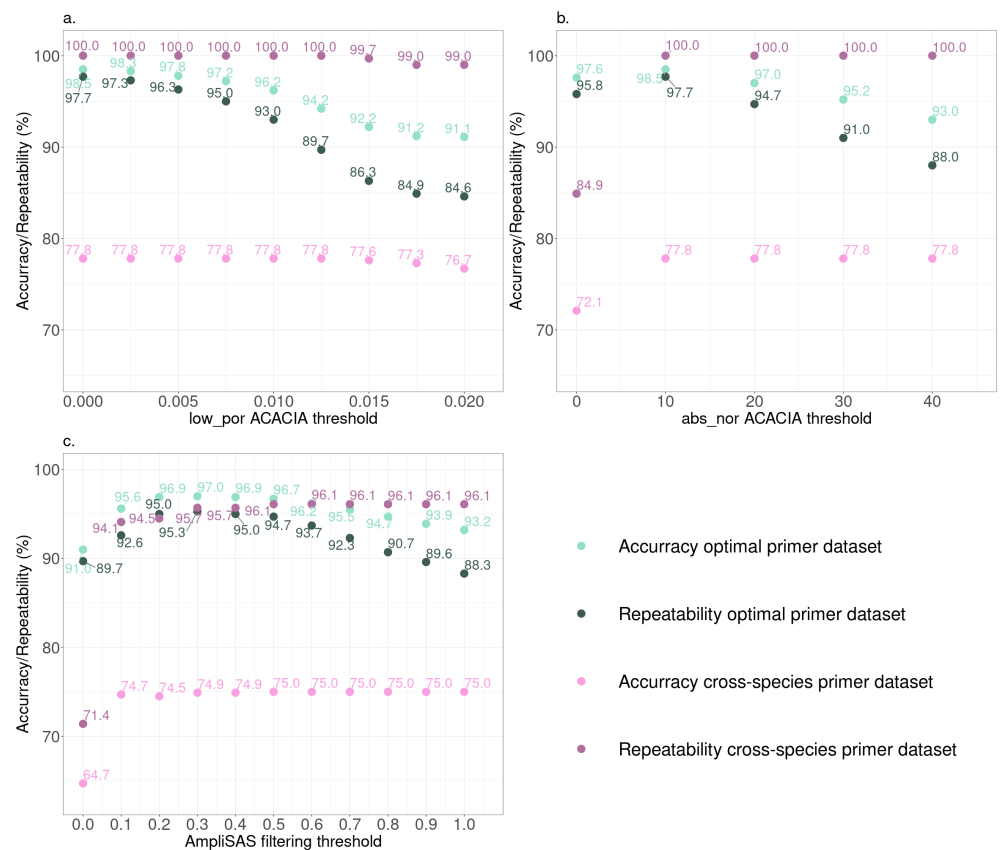


Figure 2. Allele calling accuracy and repeatability for the two datasets of this study (optimal primers or cross-species primers) at different low_por threshold settings with abs_nor set at 0 within the ACACIA pipeline (a.); at different low_por threshold settings with abs_nor set at 0 within the ACACIA pipeline (b.); and, at different filtering thresholds (i.e. ‘minimum amplicon frequency’) within the AmpliSAS pipeline (c.).

Table 2. Specific genotypes within replicates which had a genotyping error using ACACIA and AmpliSAS genotyping workflows (excluding allele dropout due to primer mismatch in the cross-species primers dataset). Genotypes, replicates (Rep.), predicted number of alleles (# Pred.All.), allele dropouts (Dropout) and false positives (F.P.) using ACACIA and AmpliSAS are shown.

Genotype	Rep.	# Pred.All.	Dropout ACACIA	Dropout AmpliSAS	F.P. AmpliSAS
a. Optimal primers dataset (BLB MHC class IIB)					
B2-B4-B12-B14-B19-B21	1	11	BLB2*21	BLB2*21	
				BLB1*21	
B4-B14-B15-B19-B21	1	10	BLB2*21	BLB1*21	
	2	10	BLB2*21	BLB2*21	
B4-B15-B19-B21	1	8	BLB2*21	BLB1*21	
B2-B4-B12-B14-B15-B19-B21	1	13	BLB2*21	BLB2*21	
			BLB1*21	BLB1*21	
B2-B4-B12-B14-B15-B21	1	12	BLB2*21	BLB2*21	
			BLB1*21	BLB1*21	
B2-B12-B14-B15-B19-B21	1	11	BLB1*21		
B2-B4-B12-B15-B19-B21	1	11		BLB1*21	
B2-B4-B12-B15-B21	1	10		BLB1*21	
B2-B4-B14-B15-B19-B21	1	12		BLB1*21	
B2-B4-B14-B15-B21	1	10		BLB1*21	
B2-B4-B15-B19-B21	1	10		BLB1*21	
	2	10		BLB1*21	
B4-B12-B21	1	6		BLB1*04	1
B4-B14-B15-B19-B21	2	10		BLB1*21	
b. Cross-species primers dataset (BLB MHC class IIB)					
B12-B14-B15-B21	1	5		BLB1*12 or *19	
	2	5		BLB1*12 or *19	
B14-B15-B19-B21	1	8		BLB1*12 or *19	
B2-B12-B14-B15	1	6		BLB1*12 or *19	
	2	6		BLB1*12 or *19	
B2-B12-B14-B15-B19-B21	1	11		BLB2*14	
B2-B14-B15-B19-B21	1	10		BLB1*12 or *19	
B2-B4-B12-B14-B15	1	10		BLB1*12 or *19	
	2	10		BLB1*12 or *19	
B2-B4-B12-B14-B15-B19	1	11		BLB2*14	
B2-B4-B12-B14-B15-B19-B21	1	13		BLB2*14	
B2-B4-B12-B14-B15-B21	1	12		BLB1*12 or *19	
B2-B4-B12-B14-B19-B21	1	11		BLB2*14	
B2-B4-B14-B15-B19-B21	1	12		BLB1*12 or *19	
B4-B12-B14-B15	1	8		BLB1*12 or *19	

	2	8		BLB1*12 or *19
B4-B14-B15-B19-B21	1	10		BLB1*12 or *19

Using the optimal settings in AmpliSAS, across 86 genotypes, a total of 17 alleles dropped out, one false positive was found (Table 2) which resulted in an allele calling accuracy of 97% (Figure 2c.). As with ACACIA, most allele dropouts (16 of 17) derived from the B21 haplotype. For three genotypes, both BLB1*21 and BLB2*21 dropped out. For nine genotypes, only BLB1*21 alleles dropped out and for one genotype only BLB2*21 allele dropped out. Finally for one genotype the allele dropout was BLB1*04 and the same genotype had a false positive allele (Table 2). Allele calling repeatability was 95.3%. Therefore, the ACACIA workflow resulted in higher allele calling accuracy and repeatability than the AmpliSAS workflow.

AmpliSAS vs ACACIA: chicken cross-species primers dataset

Using the optimal settings of ACACIA, we found a total of 134 allele dropouts across the 86 genotypes and allele calling accuracy was 77.8% (Figure 2a. and b.). However, all dropouts were from the alleles BLB1*04, BLB1*15 or BLB1*21 which were never called in the genotypes they were predicted to occur. Further comparison between the allelic reads and the primers revealed two mismatches at the 1st bp and 16th bp within the forward primer. Across the whole dataset, only 13 (0.001%), 114 (0.01%) and 11 (0.001%) reads prior to applying any downstream quality filtering steps after merging corresponded to BLB1*04, BLB1*15 or BLB1*21 respectively. By comparison, the range for all other alleles was between 25,812 (2.79%) and 115,489 (12.49%) and the range for all artifact reads were between one (0.0001%) and 5,535 (0.60%). Therefore all allele dropouts in the cross-species dataset when using the ACACIA workflow are explained by primer mismatch leading to very poor amplification and sequencing of these alleles which were well within the lower range of artefact reads. Since BLB1*04, BLB1*15 and BLB1*21 dropped out in all genotypes, allele calling repeatability between both replicates was 100% when using the ACACIA workflow, which highlights that relying on allele calling repeatability when validating a genotyping workflow can be misleading.

Using the optimal settings of AmpliSAS, we found 152 allele dropouts across all genotypes and allele calling accuracy was 75.2% (Figure 2c.). As above, 134 dropouts were due to a mismatch with the forward primer. The remaining 17 alleles that dropped out were BLB1*12 or *19 (13 alleles) and BLB2*14 (4 alleles) (Table 2). Allele calling repeatability between both replicates was 96.1%. Therefore, as with the optimal primer dataset the ACACIA workflow resulted in higher allele calling accuracy and repeatability than the AmpliSAS workflow.

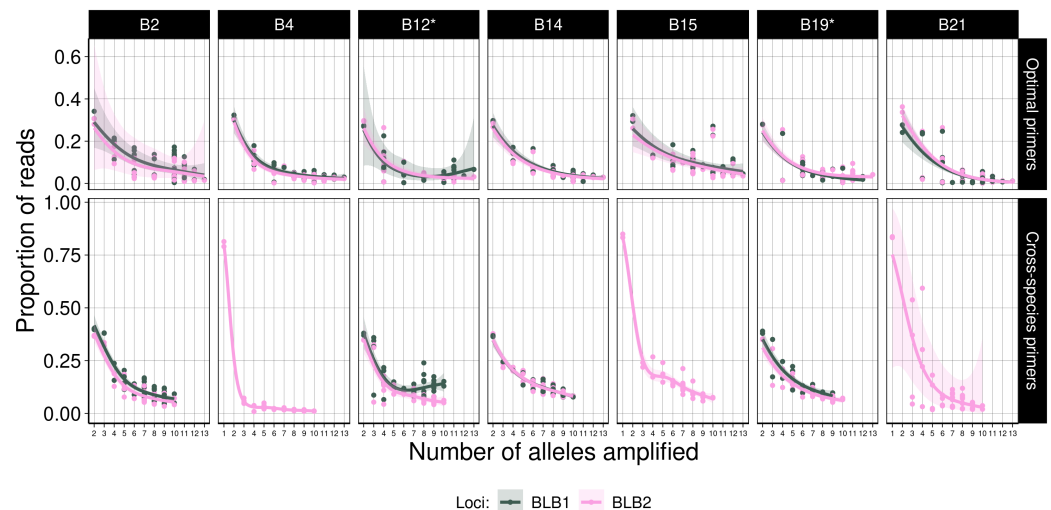


Figure 3. The relationship between the number of alleles amplified within a PCR and the proportion of reads assigned to alleles for each haplotype and locus. *Note that haplotype B12 and B19 have the same BLB01 allele, therefore for presentation purposes only, when both haplotypes were within the same genotype, BLB01*12 or *19 was assigned to haplotype B12 to avoid pseudo-replication. BLB1*04, BLB1*15 and BLB1*21 failed to amplify due to primer mismatch

478 Relationship between number of alleles amplified and artefacts

479 In the optimal primer dataset, when amplifying within haplotype, all alleles amplified and
 480 the proportion of reads assigned to alleles ranged from 0.24 to 0.36 (Figure 3). The latter
 481 confirms the suitability of the primer set design for this model system. In contrast, in the
 482 cross-species dataset, primer mismatch and systematic allele dropout for the alleles BLB1*04,
 483 BLB1*15 or BLB1*21 meant that three haplotypes had a single allele instead of two (Figure
 484 3). In both datasets, the contribution of allelic variants to the proportion of reads decreased
 485 sharply with increasing number of alleles when amplifying less than 4-6 alleles, but starts to
 486 level when amplifying more than 4-6 alleles (Figure 3, optimal dataset GAMM: $F_{3,477, 542}=237.3$;
 487 p-value <0.001; cross-species dataset GAMM: $F_{4,779, 420}=99.73$; p-value <0.001). Amplification
 488 efficiency was significantly different between alleles in both datasets (optimal dataset GAMM:
 489 $F_{12, 542}=10.63$; p-value <0.001; cross-species dataset GAMM: $F_{9, 420}=35.53$; p-value <0.001). Both
 490 alleles from the B4 and B21 haplotypes in the optimal dataset and the BLB2*04 allele in the
 491 cross-species primers dataset consistently amplified poorly when co-amplifying with alleles
 492 from other haplotypes (Figure 3; see Supplementary Figure S1 for multiple-comparison post-
 493 hoc of allele amplification). In the optimal primer dataset, the low amplification efficiency of
 494 the B21 haplotype when co-amplifying with other haplotypes explains the high allele dropout
 495 of alleles from this haplotype in more complex genotypes (i.e. when co-amplifying 10 or more

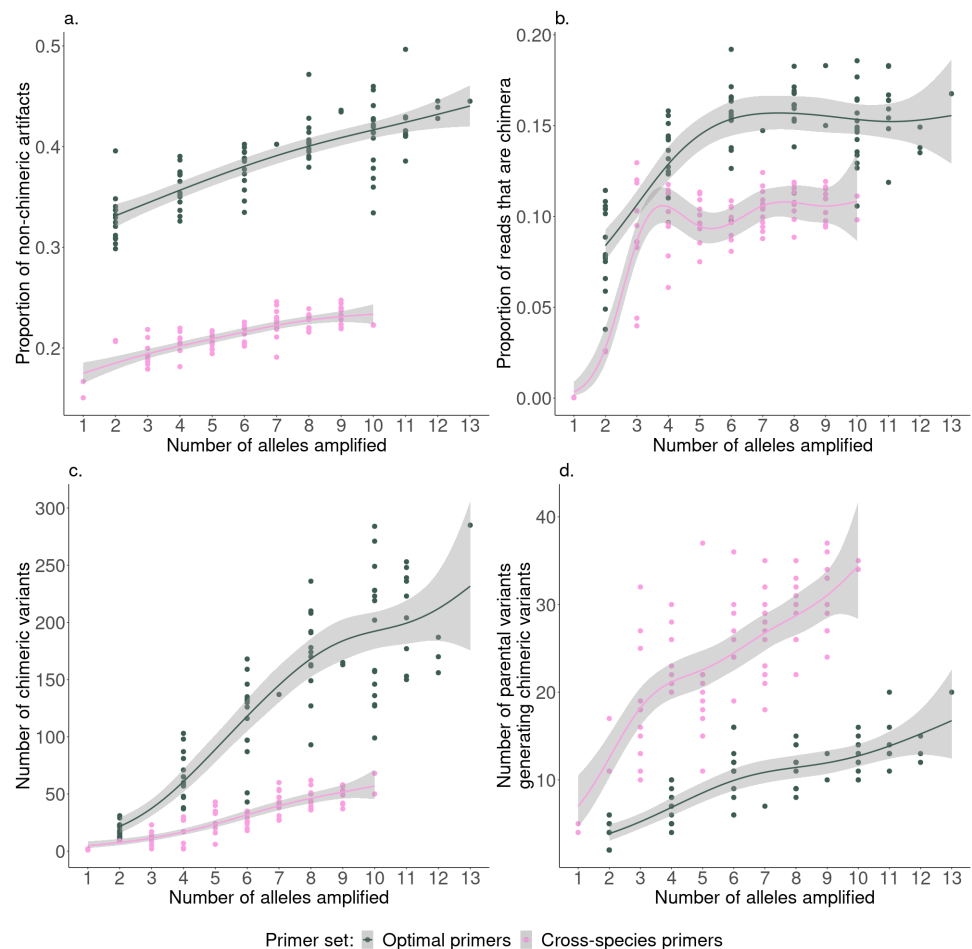


Figure 4. The relationship between the number of alleles amplified and: the proportion of reads that are non-chimeric artifacts (a.); the proportion of chimeric reads (b.); the absolute number of chimeric variants (c.); and, the absolute number of parental variants generating chimeric reads (d.).

alleles) (Figure 3). In contrast, the higher sequencing depth of the cross-species dataset meant that BLB2*04 allele did not dropout. However, we identified a primer mismatch between BLB2*04 allele and the second base pair of the reverse primer, explaining the lower amplification efficiency of this allele when co-amplified with other alleles.

The proportion of sequences classified as artefacts was much higher for PCRs using the optimal primer set than when using the cross-species primer set (Figure 4a. and 4b.; non-chimeric artefacts GMM: $F_{1, 74}=2669.1$; p-value <0.001; chimera: $F_{1, 74}=180.4$; p-value <0.001), which is likely due to the fact that the fragment length of the optimal primer dataset was

longer relative to the cross-species primers dataset (241 bp vs 151 bp, respectively; see discussion). For both datasets in this study, when considering non-chimeric artefacts, there was a positive relationship between the proportion of artefacts and the number of alleles amplified (Figure 4a.; GAMM: $F_{1, 74}=207.3$; p-value <0.001). There is a logarithmic relationship between the proportion of chimeric artefacts and the number of alleles amplified whereby the proportion of chimeric reads no longer increased with number of alleles amplified when amplifying more than 4-6 alleles (Figure 4b.; GAMM: $F_{4.857, 74}=35.77$; p-value <0.001). The total number of unique chimeric reads also tended to follow a logarithmic relationship, whereby the number of unique chimeric variants seemed to no longer increase with the number of alleles amplified when amplifying more than 10 alleles (Figure 3c.; GAMM: $F_{4.06, 74}=117.5$; p-value <0.001). The relationship between the total number of parental variants generating chimeras and the number of alleles amplified also levelled when amplifying more than six alleles (Figure 4d.; GAMM: $F_{4.06, 74}=117.5$; p-value <0.001).

518

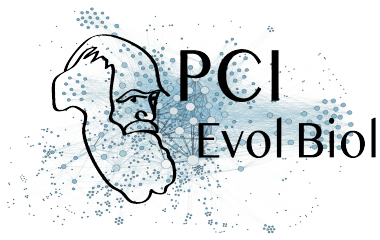
519 Discussion

Using known MHC genotypes for two datasets (chicken MHC Class IIB complex), we achieved higher allele calling accuracy ($\geq 98.5\%$) and repeatability ($\geq 97.7\%$) using ACACIA for the optimal primer dataset. With fewer allele dropouts and false positives, the ACACIA pipeline performed better than AmpliSAS. We demonstrated the “costs” of designing primers within MHC exon 2 in terms of allele dropout, with three common alleles failing to amplify when using primers designed from sequences of related Galliforme species. We also explored the relationship between artefacts and the number of alleles amplified per amplicon, and, as expected, found heterogeneous amplification efficiency of allelic variants when amplifying multiple loci within a PCR. Surprisingly, the relationship between the proportion of chimeric artefacts and number of alleles amplified was not linear but rather leveled when amplifying more than 4-6 alleles. However, non-chimeric artefacts did increase linearly with increasing number of alleles amplified. Below we discuss in further detail ACACIA, AmpliSAS and other genotyping pipelines, primer design for non-model organisms, the relationship between the number of alleles amplified and artefacts, the effect of chimera formation on genotyping pipelines and, finally, we conclude by advising users on important points to consider when genotyping complex multigene families in non-model organisms.

536

537 AmpliSAS vs ACACIA

Experimentally generating CNV of known chicken MHC class IIB genotypes allowed us to validate our ACACIA pipeline to genotype systems with high CNV complexity at high accuracy and repeatability across replicates in the optimal primer dataset. While we achieved higher allele



541 calling accuracy and repeatability using ACACIA than the AmpliSAS web server pipeline, we
542 do not claim that ACACIA will necessarily perform better than AmpliSAS with all datasets. To
543 demonstrate the latter we would need to test both pipelines on a larger number of datasets
544 and/or on simulated datasets. In addition, while our pipeline should suit data generated with
545 any high-throughput sequencing technologies, we have only tested ACACIA with paired-end
546 Illumina sequencing technology.

547

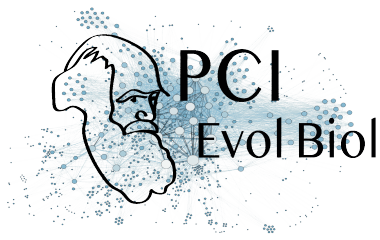
548 The most apparent benefit of using the AmpliSAS web server is that it is relatively easy to
549 use for users with limited knowledge of scripting languages (such as PYTHON, PERL, C++ or
550 R). However, we have noticed that a number of studies report results using default settings
551 when applying the AmpliSAS pipeline to their dataset. We find this concerning since, as our
552 study demonstrates, the default clustering and filtering parameters are unlikely to be opti-
553 mal for most datasets. Indeed, allele calling accuracy was much lower when using the default
554 settings (81.8%) as compared to the optimal settings (97%) in the optimal primer dataset in
555 our study, due to high allele dropout when using the default settings. We therefore strongly
556 discourage users from using default settings and advise to permutate between different fil-
557 tering and clustering parameters to find the best settings for their dataset when using the
558 AmpliSAS pipeline. As most wildlife studies cannot assess allele calling accuracies, duplicat-
559 ing samples and relying on repeatability is the only feasible method for most research to
560 optimise their amplicon-based genotyping workflow. However, authors should bear in mind
561 that due to the high recurrence of amplification and sequencing errors, high repeatability in
562 allele calling between replicates does not necessarily entail an error free workflow and may
563 be misleading. Therefore, careful design of primers and PCR conditions to reduce artefacts
564 during amplification are crucial to maximise amplicon-based genotyping accuracy regardless
565 of the bioinformatic tools used (see further discussion on this below).

566

567 An important disadvantage of the AmpliSAS web server is that at the time of writing, se-
568 quencing depth per amplicon was limited to 5,000 reads. The latter is particularly problematic
569 when wishing to genotype systems with complex CNV, which require high sequencing depth
570 to genotype with high repeatability [5]. For datasets with sequencing depth above 5000 reads,
571 AmpliSAS can be run locally but we found that, unlike the web server, the local version of Am-
572 pliSAS had limited documentation and troubleshooting was time consuming.

573

574 Once installed, ACACIA does not require users to have experience with scripting languages,
575 allows genotyping with virtually unlimited sequencing depth and provides output data report-
576 ing the number of reads kept at each step of the pipeline. The latter should aid users when
577 deciding upon optimal parameters and thresholds. As for the AmpliSAS pipeline, we advise
578 to not use default parameters of ACACIA without critically assessing different parameters for
579 each dataset. In particular, we urge users to permutate between different settings of `abs_nor`
580 and `low_por` parameters. We advise to first search for the optimal `abs_nor` setting with a fixed



low_por parameter of 0 because it is likely that it is only necessary to change the low_por parameter setting from 0 in datasets with ultra deep sequencing depth. If it is subsequently found that the optimal low_por setting is greater than 0, users should repeat the permuting step of abs_nor until the optimal settings are found. Of course finding optimal settings requires the inclusion of replicates for at least a subset of the dataset. We therefore recommend that a sufficient number of replicates are always included in genotyping runs to obtain sufficiently accurate repeatability values.

Comparing ACACIA to other pipelines

Prior to the development of AmpliSAS and ACACIA, researchers who wished to genotype complex multigene families generally relied on either earlier software such as SESAME [41] or jMHC [62] or their own customised scripts (e.g. [31, 70]). However while both SESAME and jMHC aided allele calling workflows by allowing users to demultiplex sequences and to generate tables which contains sequence variants and the number of reads, they do not allow users to apply an automated workflow to distinguish artefacts from real allelic variants.

Genotyping pipelines have evolved and matured in the last decade, however all genotyping pipelines rely to some degree on the assumption that artefacts are in general less frequent than genuine allelic variants. Genotyping pipelines vary in the methods used to discriminate poorly amplified allelic variants from artefacts. An early pipeline suggested by Radwan *et al* [52], which expanded from initial pipelines suggested by Kloch *et al* [31] and Zagalska-Neubauer *et al* [70], set a threshold below which all variants are considered artefacts (e.g. <1.5% per amplicon in Radwan *et al* [52]). This threshold is set by comparing rare variants to more common variants within an amplicon to determine whether the rare variant can be explained as an artefact (i.e. 1 to 2 bp mismatch compared to a common variant within an amplicon or a PCR chimera from two common parental variants within an amplicon). The weakness of this genotyping pipeline is that it relies on a single threshold below which all variants are considered artefacts, potentially making it particularly vulnerable to allele dropout [61]. A second method was suggested by Sommer, Courtiol, & Mazzoni [61], which relied on comparisons between duplicated amplicons and a series of decision making trees to discriminate between allelic variants and artefacts. While the pipeline of Sommer, Courtiol, & Mazzoni [61] also assumes that artefacts are less frequent than most allelic variants, it does not rely on a single threshold below which all sequences are considered artefacts. However, one potential weakness of this method is that it may be more vulnerable to repeatable artefacts and thus to false positives, particularly in systems highly diverse in terms of high copy number variation (CNV>10 [5]).

A further disadvantage of all the above early genotyping pipelines is that much of the se-

quencing depth data is wasted by simply discarding low threshold sequences. To maximise the available sequencing depth, recent genotyping methods have clustered artefactual (non-chimeric) sequences to their suspected parental variant to increase genotyping confidence. This trend has been particularly strong in the 16S rRNA microbiome community, which have traditionally clustered sequence variants to so called operational taxonomic units (OTUs) using a fixed similarity threshold (usually 97% similarity). More recent 16S rRNA clustering methods such as the entropy based Oligotyping tool used within ACACIA [14], as well as model based methods such as DADA2 [10] and Deblur [2], have used alternative and more sophisticated statistical methods to simple similarity thresholds to distinguish sequence variants that differ by as little as one base pair. The clear benefit of clustering is that it significantly reduces the number of reads with low abundances, while increasing the read counts from poorly amplified allelic variants. However even the most sophisticated clustering methods will retain some artefacts within datasets [2, 9, 14], hence the need for additional filtering steps following clustering. Downstream filtering strategies can also resemble the pre-clustering pipelines strategies mentioned above as was applied by Biedrzycka *et al* [5] using AmpliSAS in a highly complex system (19 to 42 allelic variants per amplicon). Biedrzycka *et al* [5] found a high agreement between genotyping methods as long as sequencing depth was sufficiently high. This will also likely be the case when applying ACACIA instead of AmpliSAS to such datasets.

Biedrzycka *et al* [5] had a <90% allele calling repeatability when coverage <5,000 sequences regardless of the genotyping workflow used and reached 99% with a sequencing depth of 20,000. While our study does not allow to extensively assess the relationship between CNV and sequencing depth using ACACIA, our results were consistent with [5] since allele calling repeatability was 97.7% for the optimal primer dataset and 100% for the cross-species primer dataset, which had an average sequencing depth of 6,164 and 11,562 reads per amplicon respectively. For the optimal primer dataset, regardless of the genotyping pipeline used, allele dropout occurred in genotypes with high CNV (for ACACIA 8 out of 9 and for AmpliSAS 12 out of 14 genotypes with allele dropouts had 10 alleles or more). Our optimal primers amplified all alleles at a similar efficiency when amplified within single haplotypes suggesting that the primers are indeed optimally designed. For all instances, allele dropout were alleles from the B4 and B21 haplotype which amplified poorly when coamplified with alleles from other haplotypes. Higher sequencing depth will reduce or even remove such allele dropout instances [5]. Indeed for the cross-species primer dataset, sequencing depth was nearly twice as high, and there were no instances of allele dropout due to the ACACIA pipeline (all allele dropouts were due to primer mismatch, see subsequent sub-section of the discussion) and allele calling repeatability was 100%. Therefore, in order to reach allele calling repeatability values <99%, we advise researchers to aim for a sequencing depth of at least 10,000 reads per amplicon when amplifying more than 4 alleles per amplicon and of 20,000 reads when amplifying more than 15 alleles regardless of the bioinformatic workflow used [5].

An important benefit of the Oligotyping tool in ACACIA is that unlike other clustering meth-

ods which use the entire sequence, it only uses the base pairs with the most discriminant information based on entropy analyses [14]. In the context of MHC genotyping in particular, such a strategy makes much intuitive sense, since most functional differences between MHC alleles will be within specific regions of the sequences which will contain the antigen-binding sites that are highly polymorphic as a result of strong positive selection.

The challenge of designing primers for non-model organisms

A common approach for primer design in complex genomic regions of non-model organisms includes aligning multiple sequences of phylogenetically related species. By building primers on consensus sequences, researchers assume that oligos will also amplify the target region in the species of interest. However, knowledge about related species is often limited to very few individuals. This means that primers can be designed in regions that are polymorphic in the target species. As a consequence, certain allelic variants are not amplified and homozygosity is overestimated. Indeed, this proved to be the case in our cross-species primers dataset, whereby two mismatches (1st bp and 16th bp) within the forward primer (19 bp long) were sufficient to prevent the amplification of three alleles (out of 13). Interestingly, a single base pair mismatch between the second base pair of the reverse primer and the BLB2*04 allele did not prevent the amplification of this allele, although it did suffer severely from low amplification efficiency when in competition with other alleles. However, high sequencing depth for the cross-species primer dataset prevented this allele from dropping out, regardless of the genotyping pipeline used. Our study therefore highlights the importance of carefully designed primers for amplicon based genotyping.

Two non-mutually exclusive strategies can be used to decrease allele dropout in non-model organism with no *a priori* information on the target region. First, designing multiple primers and combining them within a PCR reaction is known to reduce allele dropout due to primer mismatch and allele amplification bias [38]. In addition combining multiple primers within PCR reactions reduces the need to sequence multiple primer sets separately, considerably reducing the cost of using multiple primer sets to genotype a novel target region. Second, the recent development of long-read HTS of single molecules, such as Pacific Biosciences Single Molecules Real-Time sequencing or Oxford Nanopore sequencing, offers much promise in characterising the structure of novel target regions and consequently to enable more informed primer design [46]. For instance Fuselli *et al* [15] were recently able to develop an assembly pipeline that combined long-read HTS from a single sample with *de novo* short-read HTS assembly of six samples to characterise a 9Kb region of the MHC Class II (*DRB*) locus in Alpine chamois *Rupicapra rupicapra*. This approach allowed the authors to conclude with some degree of confidence that there is a single copy of the *DRD* locus within the six individuals sampled in their study and provides a reference template for future genotyping

697 strategies for a larger number of individuals [15]. More complex multi-gene families have
698 also been characterised in primate species [32], including the MHC [19, 66] using long-read
699 HTS. An important benefit of long-read HTS is that they may allow the design of locus specific
700 primer design, which would reduce allele dropout due primer mismatch and allow assigning
701 alleles to loci which could provide complete genotype information. Indeed, an important set-
702 back of solely relying on short-read HTS is that assigning alleles to loci in complex systems
703 is frequently impossible, even when using recent phasing algorithms [22], limiting the use
704 of many population genetic analyses for these loci. In addition, CNV is currently likely to be
705 underestimated in many species since recent gene duplication means that different genes of-
706 ten carry identical alleles [69]. An important set-back with long-read HTS technologies, is that
707 they have higher sequencing error rates than short read HTS (the reported error rates is 16%
708 per base compared to 0.1% for Illumina HTS) and lower sequencing output (5-20 Gb/run com-
709 pared to 200-600 for Illumina HTS) [4, 53]. Therefore, whilst long read HTS will undoubtedly
710 improve our understanding of MHC and other multigene complexes structure and conse-
711 quently non-model species primer design, population studies genotyping a large number of
712 individuals are likely to continue to rely on amplicon-based genotyping from short read HTS
713 for the foreseeable future.

714

715 **Relationship between number of alleles amplified and artefacts**

716 By knowing the exact alleles to expect for the chicken genotypes, we were able to quantify
717 chimeric artefacts precisely (Figure 1). There was a higher proportion of chimeric and non-
718 chimeric artefacts in the optimal primer dataset than in the cross-species primer dataset. The
719 most likely explanation for the latter is the shorter sequence for the cross-species primer
720 dataset (151 bp) compared to the optimal primer dataset (241 bp). A shorter fragment re-
721 duces the number of base pairs that can be erroneously substituted/deleted and the number
722 of breaking points for chimera formation. In addition, it is likely that the probability of incom-
723 plete elongation is inversely related to fragment length. Thus, fragment length appears to be
724 the dominant factor predicting the proportion of artefactual reads.

725

726 As expected, the proportion of reads that were non-chimeric artefacts increased linearly
727 as the number of alleles amplified with an amplicon increased, which can be explained simply
728 by the fact that there is an increasing number of possible artefacts that can be generated as
729 the number of initial template variants increases. Thus, reads that failed to be completely
730 elongated within the PCR cycles are more likely to be erroneously elongated during the final
731 extension step.

732

733 An unexpected result was that the proportions of chimeras did not increase with increas-
734 ing number of alleles amplified with an amplicon, when amplifying more than 4-6 alleles.

Similarly, when amplifying more than 10 alleles, the number of chimeric variants no longer increased with increasing number of alleles amplified within an amplicon. Such saturation in chimera generation beyond a threshold of alleles amplified is likely to be a by-product of allele PCR competition. Indeed, as demonstrated by our own data, there is amplification bias whereby some gene variants are amplified preferentially relative to others [38, 61]. Therefore, a few gene variants (~3-6 gene variants) are preferentially amplified and most chimeras originate from these dominantly amplified variants and few chimeras are generated from the poorly amplified variants. Indeed, we found that the number of parental variants generating chimeras in our dataset did not increase with increasing number of alleles amplified when amplifying more than 4-6 alleles. The non-linear relationship between chimera generation and number of alleles amplified have important implications when considering sequencing depth needed to accurately genotype complex multigene families, since it suggests that linearly increasing sequencing depth for increasing CNV is not necessarily the optimal strategy. The challenges of dealing with chimeras in genotyping pipelines is discussed below in detail.

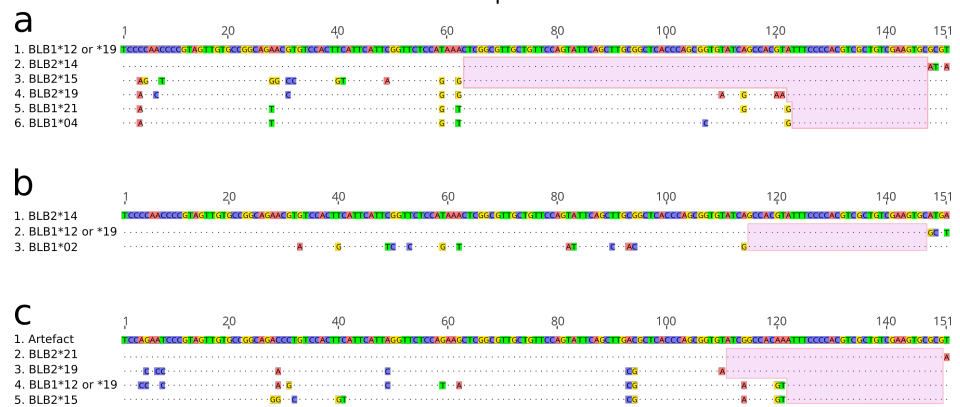
Chimeras in genotyping pipelines

The formation of artificial chimeras during amplification is an important source of artefacts in amplicon sequencing projects [34, 60], including those with newer sequencing technologies [33]. Chimeras are challenging to identify as artefacts because they resemble real alleles generated by recombination, particularly in multigene families under high rates of interlocus genetic exchange ("concerted evolution"), which is common in many MHC systems [7, 8, 13, 17, 21, 67]. Our results suggest that chimeras are more prevalent, harder to identify and potentially more reproducible across technical replicates than previously assumed. We expect the same to be true for similar projects with conserved, yet variable amplification targets such as the MHC.

One allele erroneously called as a real variant (i.e. a false positive) by the AmpliSAS pipeline in the optimal primer dataset was actually a chimera between the BLB1*21 and BLB2*21 alleles. Furthermore, when using the AmpliSAS pipeline, 15 allele dropouts in the cross-species primers dataset were due to erroneous assignment of real allelic variants as chimera artefacts. Indeed, the BLB1*12 or *19 allele was identical to potential chimeric artefact sequences between BLB2*14 (85 possible breakpoints) and any of the following alleles: BLB1*04, BLB2*15, BLB2*19, BLB2*21 or BLB1*21 (Figure 5a.). In addition, BLB2*14 dropped out because it is identical to a chimera formed between the BLB1*02 and BLB1*12 or *19 alleles (33 breakpoints; Figure 5b.).

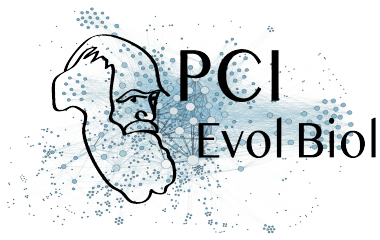
We have identified two factors which seemed to enhance chimera formation and challenge the distinction between artefact and real allelic variants. First, the combination of multiple

Figure 5. Three alignments with examples of sequences which can be classified as chimeras. The points denote identity to the first sequence in each alignment, while the differences to it are highlighted. The shaded areas indicate possible chimera-yielding breakpoints. (a) The allele BLB1*12 or *19 could be a chimera of BLB2*14 with any of the four other allele sequences depicted, in a case of multiple potential parent pairs. (b) BLB2*14 can be interpreted as a chimera between BLB1*12 or *19 minor and BLB1*02. (c) Actual chimera with multiple potential parents and a peripheral breakpoint, and therefore very similar to one of its parents.



real “parent” sequences can yield the same chimeras, as illustrated in our examples in Figure 5a. and Figure 5b., whereby any breakpoint in the shaded areas leads to the same chimeras. Second, peripheral breakpoints (Figure 5c.) can generate chimeras that differ to parental sequences by as little as a single base pair. For instance, a chimera could be a product of the allele BLB2*21 combined with any of the other alleles shown in the alignment, with a breakpoint within the shaded area (Figure 5c.). Since the potential breakpoints are at the very end of the sequence, the chimera is very similar to one of its parents (in this example, it is different from BLB2*21 by only one base). In an attempt to deal with this issue as much as possible, we changed the default settings of VSEARCH so that chimeras can be detected even if they differ from one parent by one single base. Both the “multiple parents” and the “peripheral breakpoints” issues are likely to contribute to making chimeras reproducible across replicates.

Our study highlights the challenges of chimeras for amplicon-based genotyping. In our study, we purposefully used conventional PCR conditions to replicate methods used by most wildlife MHC studies. However, the formation of most chimeras is known to occur during the final cycles of PCR amplification when dNTP and primer concentrations are low and when incompletely elongated sequences are high [25, 34, 60]. When target primers and dTNPs concentration are low during the latter stages of PCR cycles, incompletely elongated sequences act as primers and bind with the wrong sequences generating chimeras. Chimera forma-



tion during amplification can be simply reduced by adjusting the ratio of DNA template to dNTP and primer concentrations, reducing the number of cycles, increasing the extension step within PCR cycles and omitting the final extension step (which elongates high concentrations of incomplete chimeric sequences) [25, 34, 60]. Therefore, prior to any amplicon based genotyping study, we advise researchers to reduce artefacts, including chimeras, during the wet lab stage of their workflow by applying carefully designed and optimal PCR conditions. Practices during the wet lab that reduce artefacts generation in the first place is likely to be the most effective way of reducing genotyping errors regardless of the bioinformatic allele calling workflow used.

Conclusion

Genotyping accuracy and artefacts are intrinsically linked. We have demonstrated that the ACACIA genotyping pipeline provides high allele calling accuracy and repeatability. Regardless of the pipeline used, however, users should critically assess the optimal parameters to be used concerning both the wet lab and bioinformatic pipelines. We are convinced that universal default settings for optimal genotyping accuracy cannot be achieved, since optimal parameters will depend on dataset-specific generation of artefacts. The latter, in turn, varies according to species-specific CNV, DNA quality, and the conditions of PCR (e.g. extension time, number of cycles and the polymerase used) and sequencing (e.g. quality and depth). High sequencing depth allows detecting alleles that amplify poorly in complex (multigene) systems. Furthermore simple steps prior to sequencing can greatly reduce the number of artefacts generated and improve genotyping accuracy: designing more than one PCR primer pair, reducing the number of PCR cycles, increasing PCR in-cycle extension time, and omitting the final extension step. Reducing chimera formation during PCRs is particularly critical, because they are difficult to distinguish from real alleles generated by inter-locus recombination.

Data accessibility

Raw sequences of all datasets, example input files, suggested settings and the source code at the time of this publication are available at FigShare (<https://figshare.com/projects/ACACIA/66485> and doi.org/10.6084/m9.figshare.9952520). ACACIA is freely available on the GitLab at https://gitlab.com/psc_santos/ACACIA (this paper's code is available as a snapshot tagged as V1.0, https://gitlab.com/psc_santos/ACACIA/-/tags/V1.0), under an MIT license.

824 Author contributions

825 MG and PS conceived the study. PS wrote ACACIA. MG did the data analysis in R. MG, PS
826 and KM ran the allele calling workflows. KM did the AmpliSAS analysis. KW participated in
827 and supervised the lab work. KG did the lab work. SS instigated the study and heads the lab
828 where the work was carried out. MG, KM and PS wrote the first draft of the paper and all
829 authors commented and approved subsequent versions.

830 Supplementary material

Supplementary Table S1 The chicken MHC B complex haplotypes and combined haplotypes which formed experimental genotypes with varying number of alleles.

Combined haplotypes	Number of alleles
B2	2
B4	2
B12	2
B14	2
B15	2
B19	2
B21	2
B2-B4	4
B2-B12	4
B4-B12	4
B12-B14	4
B12-B21	4
B14-B15	4
B19-B21	4
B2-B4-B19	6
B2-B14-B19	6
B2-B15-B19	6
B4-B12-B21	6
B4-B14-B19	6
B12-B14-B21	6
B15-B19-B21	6
B2-B4-B12-B14	8
B2-B12-B14-B15	8
B2-B14-B19-B21	8
B4-B12-B14-B15	8
B4-B15-B19-B21	8

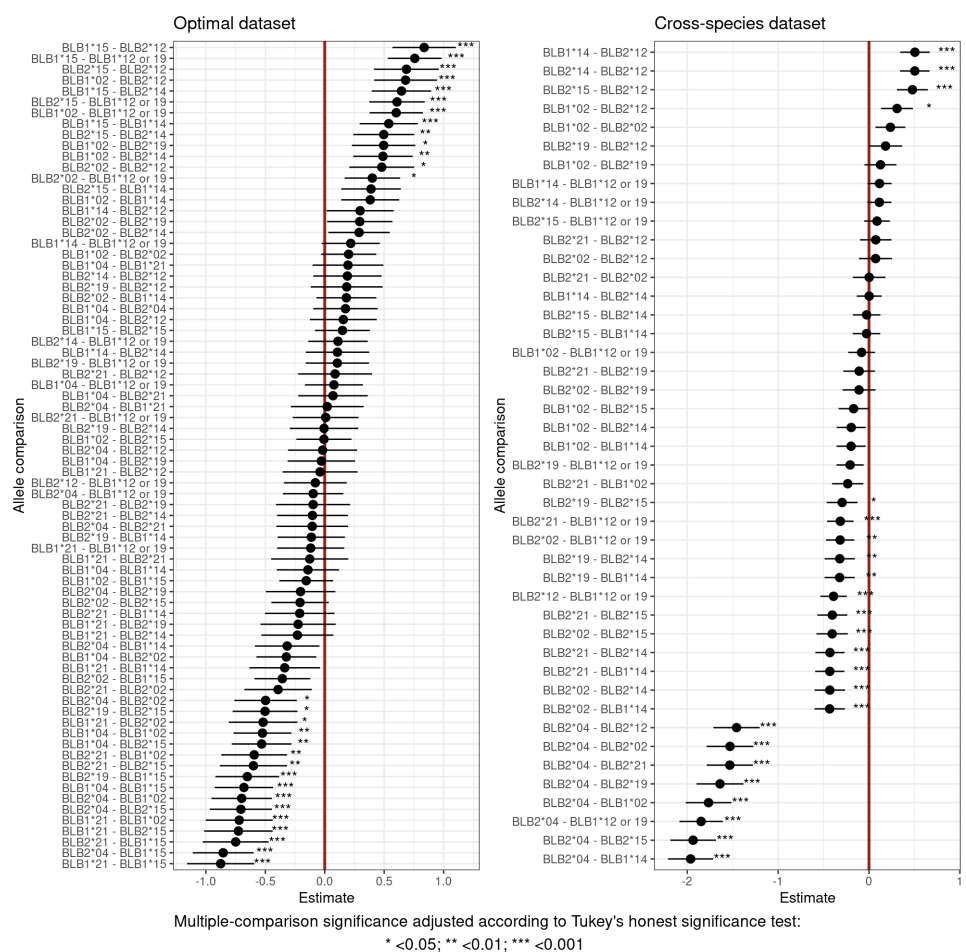
B12-B14-B15-B21	8
B14-B15-B19-B21	8
B2-B4-B12-B14-B15	10
B2-B4-B12-B14-B21	10
B2-B4-B12-B15-B21	10
B2-B4-B14-B15-B21	10
B2-B4-B15-B19-B21	10
B2-B14-B15-B19-B21	10
B4-B14-B15-B19-B21	10
B2-B4-B12-B14-B15-B19	11
B2-B4-B12-B14-B15-B21	12
B2-B4-B12-B14-B19-B21	11
B2-B4-B12-B15-B19-B21	11
B2-B4-B14-B15-B19-B21	12
B2-B12-B14-B15-B19-B21	11
B4-B12-B14-B15-B19-B21	11
B2-B4-B12-B14-B15-B19-B21	13

Supplementary Table S2 The list of NCBI accession numbers, along with species name, of the 62 MHC Class IIB exon 2 used to design the "cross-species" primers.

Species	Accession number
<i>Coturnix japonica</i>	AB110466
<i>Coturnix japonica</i>	AB110468
<i>Coturnix japonica</i>	AB110475
<i>Coturnix japonica</i>	AB110477
<i>Coturnix japonica</i>	AB110478
<i>Coturnix japonica</i>	AB110482
<i>Coturnix japonica</i>	AB181862
<i>Coturnix japonica</i>	AB181866
<i>Coturnix japonica</i>	AB181867
<i>Coturnix japonica</i>	AB181868
<i>Coturnix japonica</i>	AB181871
<i>Coturnix japonica</i>	AB181872
<i>Coturnix japonica</i>	AB181873
<i>Coturnix japonica</i>	AB181874
<i>Coturnix japonica</i>	AB181875
<i>Coturnix japonica</i>	AB181876
<i>Coturnix japonica</i>	AB181877
<i>Coturnix japonica</i>	AB264281

<i>Coturnix japonica</i>	AB264282
<i>Coturnix japonica</i>	AB282647
<i>Coturnix japonica</i>	AB282648
<i>Coturnix japonica</i>	AB282649
<i>Coturnix japonica</i>	AB282650
<i>Coturnix japonica</i>	AB282651
<i>Coturnix japonica</i>	XM_015878560
<i>Crossoptilon crossoptilon</i>	JQ001779
<i>Meleagris gallopavo</i>	AM233486
<i>Meleagris gallopavo</i>	FJ946995
<i>Meleagris gallopavo</i>	FJ946997
<i>Meleagris gallopavo</i>	GU189283
<i>Meleagris gallopavo</i>	GU189285
<i>Meleagris gallopavo</i>	GU189286
<i>Numida meleagris</i>	DQ885563
<i>Numida meleagris</i>	EF643464
<i>Numida meleagris</i>	EU030445
<i>Numida meleagris</i>	XM_021413450
<i>Numida meleagris</i>	XM_021413509
<i>Pavo cristatus</i>	AY928093
<i>Pavo cristatus</i>	AY928094
<i>Pavo cristatus</i>	AY928096
<i>Pavo cristatus</i>	AY928097
<i>Pavo cristatus</i>	AY928098
<i>Pavo cristatus</i>	AY928100
<i>Pavo cristatus</i>	AY928101
<i>Pavo cristatus</i>	JQ001780
<i>Perdix perdix</i>	KF007890
<i>Perdix perdix</i>	KF007892
<i>Perdix perdix</i>	KF007894
<i>Perdix perdix</i>	KF007895
<i>Perdix perdix</i>	KF007896
<i>Perdix perdix</i>	KF007897
<i>Perdix perdix</i>	KF007898
<i>Perdix perdix</i>	KF007900
<i>Perdix perdix</i>	KF007901
<i>Perdix perdix</i>	KY040298
<i>Perdix perdix</i>	KY040299
<i>Perdix perdix</i>	KY040300
<i>Perdix perdix</i>	KY040302

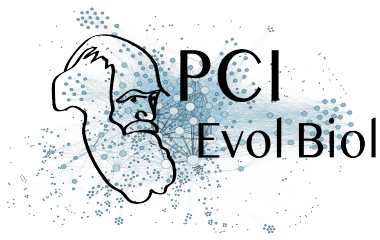
<i>Phasianus colchicus</i>	AJ224346
<i>Phasianus colchicus</i>	AJ224347
<i>Phasianus colchicus</i>	AJ224349



Supplementary Figure S1 Multiple-comparison post-hoc of allele amplification according to GAMM estimates for the optimal and cross-species datasets. Dots represent the coefficient estimates and the thin lines are 95% confidence intervals.

831 Acknowledgements

832 MG was supported by a DFG grant (DFG Gi 1065/2-1). We are very grateful to Jim Kaufman and
833 his lab members for providing the chicken DNA samples used in this study and for his com-
834 ments on a previous version of this work. Version 3 of this preprint has been peer-reviewed



and recommended by Peer Community In Evolutionary Biology (<https://doi.org/10.24072/pci.evolbiol.100092>). We thank the PCI reviewers Thomas Bigot, Helena Westerdahl and Sebastian Ernesto Ramos-Onsins, the PCI recommender François Rousset, and two anonymous reviewers for their comments which improved our manuscript.

Conflict of interest disclosure

The authors of this preprint declare that they have no financial conflict of interest with the content of this article. François Rousset is the recommender for PCI Evolutionary Biology.

References

- [1] Altschul SF, Gish W, Miller W, Myers EW, and Lipman DJ. Basic local alignment search tool. *Journal of Molecular Biology* 215 (Oct. 5, 1990), 403–410. issn: 0022-2836. doi: 10.1016/S0022-2836(05)80360-2.
- [2] Amir A, McDonald D, Navas-Molina JA, Kopylova E, Morton JT, Xu ZZ, Kightley EP, Thompson LR, Hyde ER, Gonzalez A, and Knight R. Deblur Rapidly Resolves Single-Nucleotide Community Sequence Patterns. *mSystems* 2 (Apr. 21, 2017), e00191–16. issn: 2379-5077. doi: 10.1128/mSystems.00191-16.
- [3] Babik W. Methods for MHC genotyping in non-model vertebrates. *Molecular Ecology Resources* 10 (2010), 237–251. issn: 1755-0998. doi: 10.1111/j.1755-0998.2009.02788.x.
- [4] Bansal V and Boucher C. Sequencing Technologies and Analyses: Where Have We Been and Where Are We Going? *iScience*. RECOMB-Seq 2019 18 (Aug. 30, 2019), 37–41. issn: 2589-0042. doi: 10.1016/j.isci.2019.06.035.
- [5] Biedrzycka A, Sebastian A, Migalska M, Westerdahl H, and Radwan J. Testing genotyping strategies for ultra-deep sequencing of a co-amplifying gene family: MHC class I in a passerine bird. *Molecular Ecology Resources* 17 (July 1, 2017), 642–655. issn: 1755-0998. doi: 10.1111/1755-0998.12612.
- [6] Burri R, Promerová M, Goebel J, and Fumagalli L. PCR-based isolation of multigene families: lessons from the avian MHC class II B. *Molecular Ecology Resources* 14 (2014), 778–788. issn: 1755-0998. doi: 10.1111/1755-0998.12234.
- [7] Burri R, Hirzel HN, Salamin N, Roulin A, and Fumagalli L. Evolutionary patterns of MHC class II B in owls and their implications for the understanding of avian MHC evolution. *Molecular Biology and Evolution* 25 (June 2008), 1180–91. issn: 1537-1719. doi: 10.1093/molbev/msn065.
- [8] Burri R, Salamin N, Studer RA, Roulin A, and Fumagalli L. Adaptive Divergence of Ancient Gene Duplicates in the Avian MHC Class II β . *Molecular Biology and Evolution* 27 (Oct. 1, 2010), 2360–2374. issn: 0737-4038, 1537-1719. doi: 10.1093/molbev/msq120.

- [9] Callahan BJ, McMurdie PJ, and Holmes SP. Exact sequence variants should replace operational taxonomic units in marker-gene data analysis. *The ISME Journal* 11 (Dec. 2017), 2639–2643. issn: 1751-7370. doi: 10.1038/ismej.2017.119.
- [10] Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, and Holmes SP. DADA2: High-resolution sample inference from Illumina amplicon data. *Nature Methods* 13 (July 2016), 581–583. issn: 1548-7105. doi: 10.1038/nmeth.3869.
- [11] Chen JM, Cooper DN, Chuzhanova N, Férec C, and Patrinos GP. Gene conversion: mechanisms, evolution and human disease. *Nature Reviews Genetics* 8 (Oct. 2007), 762–775. issn: 1471-0064. doi: 10.1038/nrg2193.
- [12] Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, Friedberg I, Hamelryck T, Kauff F, Wilczynski B, and Hoon MJL de. Biopython: freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics* 25 (June 1, 2009), 1422–1423. issn: 1367-4803. doi: 10.1093/bioinformatics/btp163.
- [13] Edwards S, Grahn M, and Potts W. Dynamics of Mhc evolution in birds and crocodilians: amplification of class II genes with degenerate primers. *Molecular Ecology* 4 (1995), 719–729.
- [14] Eren AM, Maignien L, Sul WJ, Murphy LG, Grim SL, Morrison HG, and Sogin ML. Oligotyping: differentiating between closely related microbial taxa using 16S rRNA gene data. *Methods in Ecology and Evolution* 4 (2013), 1111–1119. issn: 2041-210X. doi: 10.1111/2041-210X.12114.
- [15] Fuselli S, Baptista RP, Panziera A, Magi A, Guglielmi S, Tonin R, Benazzo A, Bauzer LG, Mazzoni CJ, and Bertorelle G. A new hybrid approach for MHC genotyping: high-throughput NGS and long read MinION nanopore sequencing, with application to the non-model vertebrate Alpine chamois (*Rupicapra rupicapra*). *Heredity* 121 (Oct. 2018). Number: 4 Publisher: Nature Publishing Group, 293–303. issn: 1365-2540. doi: 10.1038/s41437-018-0070-5.
- [16] Galan M, Guivier E, Caraux G, Charbonnel N, and Cosson JF. A 454 multiplex sequencing method for rapid and reliable genotyping of highly polymorphic genes in large-scale studies. *BMC Genomics* 11 (May 11, 2010), 296. issn: 1471-2164. doi: 10.1186/1471-2164-11-296.
- [17] Gillingham MaF, Courtiol A, Teixeira M, Galan M, Bechet A, and Cezilly F. Evidence of gene orthology and trans-species polymorphism, but not of parallel evolution, despite high levels of concerted evolution in the major histocompatibility complex of flamingo species. *Journal of Evolutionary Biology* 29 (2016), 438–454. issn: 1420-9101. doi: 10.1111/jeb.12798.
- [18] Glenn TC. Field guide to next-generation DNA sequencers. *Molecular Ecology Resources* 11 (2011), 759–769. issn: 1755-0998. doi: 10.1111/j.1755-0998.2011.03024.x.

- [19] Gordon D, Huddleston J, Chaisson MJP, Hill CM, Kronenberg ZN, Munson KM, Malig M, Raja A, Fiddes I, Hillier LW, Dunn C, Baker C, Armstrong J, Diekhans M, Paten B, Shendure J, Wilson RK, Haussler D, Chin CS, and Eichler EE. Long-read sequence assembly of the gorilla genome. *Science* 352 (Apr. 1, 2016). Publisher: American Association for the Advancement of Science Section: Research Article. issn: 0036-8075, 1095-9203. doi: 10.1126/science.aae0344.
- [20] Goto RM, Afanassieff M, Ha J, Iglesias GM, Ewald SJ, Briles WE, and Miller MM. Single-strand conformation polymorphism (SSCP) assays for major histocompatibility complex B genotyping in chickens. *Poultry Science* 81 (Dec. 1, 2002), 1832–1841. issn: 0032-5791. doi: 10.1093/ps/81.12.1832.
- [21] Hess C and Edwards S. The evolution of the major histocompatibility complex in birds. *Bioscience* 52 (2002), 423–431.
- [22] Huang K, Zhang P, Dunn DW, Wang T, Mi R, and Li B. Assigning alleles to different loci in amplifications of duplicated loci. *Molecular Ecology Resources* 19 (2019). _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/1755-0998.13036>, 1240–1253. issn: 1755-0998. doi: 10.1111/1755-0998.13036.
- [23] Huse SM, Huber JA, Morrison HG, Sogin ML, and Welch DM. Accuracy and quality of massively parallel DNA pyrosequencing. *Genome Biology* 8 (July 20, 2007), R143. issn: 1474-760X. doi: 10.1186/gb-2007-8-7-r143.
- [24] Jacob JP, Milne S, Beck S, and Kaufman J. The major and a minor class II β -chain (B-LB) gene flank the Tapasin gene in the B-F /B-L region of the chicken major histocompatibility complex. *Immunogenetics* 51 (Feb. 1, 2000), 138–147. issn: 0093-7711, 1432-1211. doi: 10.1007/s002510050022.
- [25] Judo MSB, Wedel AB, and Wilson C. Stimulation and suppression of PCR-mediated recombination. *Nucleic Acids Research* 26 (Apr. 1, 1998), 1819–1825. issn: 0305-1048. doi: 10.1093/nar/26.7.1819.
- [26] Katoh K and Standley DM. MAFFT Multiple Sequence Alignment Software Version 7: Improvements in Performance and Usability. *Molecular Biology and Evolution* 30 (Apr. 1, 2013), 772–780. issn: 0737-4038. doi: 10.1093/molbev/mst010.
- [27] Kaufman J, Jacob J, Shaw I, Walker B, Milne S, Beck S, and Salomonsen J. Gene organisation determines evolution of function in the chicken MHC. *Immunological reviews* 167 (Feb. 1999), 101–17. issn: 0105-2896.
- [28] Kaufman J, Milne S, Göbel TW, Walker Ba, Jacob JP, Auffray C, Zoorob R, and Beck S. The chicken B locus is a minimal essential major histocompatibility complex. *Nature* 401 (Oct. 1999), 923–5. issn: 0028-0836. doi: 10.1038/44856.
- [29] Kaufman J, Völk H, and Wallny HJ. A “Minimal Essential Mhc” and an “Unrecognized Mhc”: Two Extremes in Selection for Polymorphism. *Immunological Reviews* 143 (1995), 63–88. issn: 1600-065X. doi: 10.1111/j.1600-065X.1995.tb00670.x.

- 944 [30] Kelley J, Walter L, and Trowsdale J. Comparative genomics of major histocompatibility
945 complexes. *Immunogenetics* 56 (Jan. 1, 2005), 683–695. issn: 1432-1211. doi: 10.1007/
946 s00251-004-0717-7.
- 947 [31] Kloch A, Babik W, Bajer A, Siński E, and Radwan J. Effects of an MHC-DRB genotype and
948 allele number on the load of gut parasites in the bank vole *Myodes glareolus*. *Molecular*
949 *Ecology* 19 Suppl 1 (Mar. 2010), 255–65. issn: 1365-294X. doi: 10.1111/j.1365-294X.2009.
950 04476.x.
- 951 [32] Larsen PA, Heilman AM, and Yoder AD. The utility of PacBio circular consensus sequenc-
952 ing for characterizing complex gene families in non-model organisms. *BMC Genomics*
953 15 (Aug. 26, 2014), 720. issn: 1471-2164. doi: 10.1186/1471-2164-15-720.
- 954 [33] Laver TW, Caswell RC, Moore KA, Poschmann J, Johnson MB, Owens MM, Ellard S, Paszkiewicz
955 KH, and Weedon MN. Pitfalls of haplotype phasing from amplicon-based long-read se-
956 quencing. *Scientific Reports* 6 (Feb. 17, 2016), 21746. issn: 2045-2322. doi: 10.1038 /
957 srep21746.
- 958 [34] Lenz TL and Becker S. Simple approach to reduce PCR artefact formation leads to re-
959 liable genotyping of MHC and other highly polymorphic loci — Implications for evolu-
960 tionary analysis. *Gene* 427 (Dec. 31, 2008), 117–123. issn: 0378-1119. doi: 10.1016/j.
961 gene.2008.09.013.
- 962 [35] Lighten J, Oosterhout Cv, and Bentzen P. Critical review of NGS analyses for de novo
963 genotyping multigene families. *Molecular Ecology* 23 (2014), 3957–3972. issn: 1365-294X.
964 doi: 10.1111/mec.12843.
- 965 [36] Lighten J, Oosterhout C, Paterson IG, McMullan M, and Bentzen P. Ultra-deep Illumina
966 sequencing accurately identifies MHC class IIb alleles and provides evidence for copy
967 number variation in the guppy (*Poecilia reticulata*). *Molecular Ecology Resources* 14 (Jan. 9,
968 2014), 753–767. issn: 1755-098X. doi: 10.1111/1755-0998.12225.
- 969 [37] Magoč T and Salzberg SL. FLASH: Fast Length Adjustment of Short Reads to Improve
970 Genome Assemblies. *Bioinformatics* (Sept. 7, 2011), btr507. issn: 1367-4803, 1460-2059.
971 doi: 10.1093/bioinformatics/btr507.
- 972 [38] Marmesat E, Soriano L, Mazzoni CJ, Sommer S, and Godoy JA. PCR Strategies for Com-
973 plete Allele Calling in Multigene Families Using High-Throughput Sequencing Approaches.
974 *PLOS ONE* 11 (June 13, 2016), e0157402. issn: 1932-6203. doi: 10.1371/journal.pone.
975 0157402.
- 976 [39] McElroy KE, Luciani F, and Thomas T. GemSIM: general, error-model based simulator of
977 next-generation sequencing data. *BMC Genomics* 13 (Feb. 15, 2012), 74. issn: 1471-2164.
978 doi: 10.1186/1471-2164-13-74.
- 979 [40] McKinney W. Data structures for statistical computing in Python. *Proceedings of the 9th*
980 *Python in Science Conference* (2010), 51–56.

- 981 [41] Megléc E, Piry S, Desmarais E, Galan M, Gilles A, Guivier E, Pech N, and Martin JF.
982 SESAME (SEquence Sorter & AMplicon Explorer): genotyping based on high-throughput
983 multiplex amplicon sequencing. *Bioinformatics (Oxford, England)* 27 (Jan. 2011), 277–8.
984 issn: 1367-4811. doi: 10.1093/bioinformatics/btq641.
- 985 [42] Miller SA, Dykes DD, and Polesky HF. A simple salting out procedure for extracting DNA
986 from human nucleated cells. *Nucleic Acids Research* 16 (Feb. 11, 1988), 1215. issn: 0305-
987 1048.
- 988 [43] Nakamura T, Yamada KD, Tomii K, and Katoh K. Parallelization of MAFFT for large-scale
989 multiple sequence alignments. *Bioinformatics* 34 (July 15, 2018), 2490–2492. issn: 1367-
990 4803. doi: 10.1093/bioinformatics/bty121.
- 991 [44] Nei M, Gu X, and Sitnikova T. Evolution by the birth-and-death process in multigene fam-
992 ilies of the vertebrate immune system. *Proceedings of the National Academy of Sciences*
993 94 (July 22, 1997). 00575, 7799–7806. issn: 0027-8424, 1091-6490.
- 994 [45] Nei M and Rooney AP. Concerted and Birth-and-Death Evolution of Multigene Families.
995 *Annual review of genetics* 39 (2005), 121–152. issn: 0066-4197. doi: 10.1146/annurev.
996 genet.39.073003.112240.
- 997 [46] O'Connor EA, Westerdahl H, Burri R, and Edwards SV. Avian MHC Evolution in the Era of
998 Genomics: Phase 1.0. *Cells* 8 (Oct. 2019). Number: 10 Publisher: Multidisciplinary Digital
999 Publishing Institute, 1152. doi: 10.3390/cells8101152.
- 1000 [47] Parham P and Ohta T. Population Biology of Antigen Presentation by MHC Class I Molecules.
1001 *Science* 272 (Apr. 5, 1996), 67–74. issn: 0036-8075, 1095-9203. doi: 10.1126/science.272.
1002 5258.67.
- 1003 [48] Pavé SA, Sevellec M, Adam W, Normandeau E, Lamaze FC, Gagnaire PA, Filteau M,
1004 Hebert FO, Maaroufi H, and Bernatchez L. Nonparallelism in MHCII β diversity accom-
1005 panies nonparallelism in pathogen infection of lake whitefish (*Coregonus clupeaformis*)
1006 species pairs as revealed by next-generation sequencing. *Molecular Ecology* 22 (2013),
1007 3833–3849. issn: 1365-294X. doi: 10.1111/mec.12358.
- 1008 [49] Promerová M, Babik W, Bryja J, Albrecht T, Stuglik M, and Radwan J. Evaluation of two
1009 approaches to genotyping major histocompatibility complex class I in a passerine—CE-
1010 SSCP and 454 pyrosequencing. *Molecular Ecology Resources* 12 (2012), 285–292. issn:
1011 1755-0998. doi: 10.1111/j.1755-0998.2011.03082.x.
- 1012 [50] Quail MA, Smith M, Coupland P, Otto TD, Harris SR, Connor TR, Bertoni A, Swerdlow
1013 HP, and Gu Y. A tale of three next generation sequencing platforms: comparison of Ion
1014 Torrent, Pacific Biosciences and Illumina MiSeq sequencers. *BMC Genomics* 13 (July 24,
1015 2012), 341. issn: 1471-2164. doi: 10.1186/1471-2164-13-341.
- 1016 [51] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for
1017 Statistical Computing. Vienna, Austria, 2020.

- 1018 [52] Radwan J, Zagalska-Neubauer M, Cichoń M, Sendek J, Kulma K, Gustafsson L, and
1019 Babik W. MHC diversity, malaria and lifetime reproductive success in collared flycatch-
1020 ers. *Molecular Ecology* 21 (2012), 2469–2479. issn: 1365-294X. doi: 10.1111/j.1365-
1021 294X.2012.05547.x.
- 1022 [53] Reinert K, Langmead B, Weese D, and Evers DJ. Alignment of Next-Generation Sequenc-
1023 ing Reads. *Annual Review of Genomics and Human Genetics* 16 (2015). _eprint: <https://doi.org/10.1146/annurev-genom-090413-025358>, 133–151. doi: 10.1146/annurev-genom-090413-025358.
- 1024 [54] Rognes T, Flouri T, Nichols B, Quince C, and Mahé F. VSEARCH: a versatile open source
1025 tool for metagenomics. *PeerJ* 4 (Oct. 18, 2016), e2584. issn: 2167-8359. doi: 10.7717/
1026 peerj.2584.
- 1027 [55] Ross MG, Russ C, Costello M, Hollinger A, Lennon NJ, Hegarty R, Nusbaum C, and Jaffe
1028 DB. Characterizing and measuring bias in sequence data. *Genome Biology* 14 (May 29,
1029 2013), R51. issn: 1474-760X. doi: 10.1186/gb-2013-14-5-r51.
- 1030 [56] Rozen S and Skaletsky H. Primer3 on the WWW for General Users and for Biologist
1031 Programmers. In: *Bioinformatics Methods and Protocols*. Ed. by Misener S and Krawetz
1032 SA. Methods in Molecular Biology™. Totowa, NJ: Humana Press, 1999, pp. 365–386. isbn:
1033 978-1-59259-192-3. doi: 10.1385/1-59259-192-2:365.
- 1034 [57] Sebastian A, Herdegen M, Migalska M, and Radwan J. amplisas: a web server for mul-
1035 tilocus genotyping using next-generation amplicon sequencing data. *Molecular Ecology*
1036 *Resources* 16 (Mar. 1, 2016), 498–510. issn: 1755-0998. doi: 10.1111/1755-0998.12453.
- 1037 [58] Sepil I, Moghadam HK, Huchard E, and Sheldon BC. Characterization and 454 pyrosequencing of Major Histocompatibility Complex class I genes in the great tit reveal complexity in a passerine system. *BMC Evolutionary Biology* 12 (May 15, 2012), 68. issn: 1471-2148. doi: 10.1186/1471-2148-12-68.
- 1038 [59] Shaw I, Powell TJ, Marston DA, Baker K, Hateren A van, Riegert P, Wiles MV, Milne S, Beck S, and Kaufman J. Different evolutionary histories of the two classical class I genes BF1 and BF2 illustrate drift and selection within the stable MHC haplotypes of chickens. *The Journal of Immunology* 178 (2007), 5744–5752.
- 1039 [60] Smyth RP, Schlub TE, Grimm A, Venturi V, Chopra A, Mallal S, Davenport MP, and Mak J. Reducing chimera formation during PCR amplification to ensure accurate genotyping. *Gene* 469 (Dec. 1, 2010), 45–51. issn: 0378-1119. doi: 10.1016/j.gene.2010.08.009.
- 1040 [61] Sommer S, Courtiol A, and Mazzoni CJ. MHC genotyping of non-model organisms using next-generation sequencing: a new methodology to deal with artefacts and allelic dropout. *BMC Genomics* 14 (Aug. 9, 2013), 542. issn: 1471-2164. doi: 10.1186/1471-2164-14-542.
- 1041 [62] Stuglik MT, Radwan J, and Babik W. jMHC: software assistant for multilocus genotyping of gene families using next-generation amplicon sequencing. *Molecular Ecology Resources* 11 (2011), 739–742. issn: 1755-0998. doi: 10.1111/j.1755-0998.2011.02997.x.
- 1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055

- 1056 [63] Stutz WE and Bolnick DI. Stepwise Threshold Clustering: A New Method for Genotyping
1057 MHC Loci Using Next-Generation Sequencing Technology. *PLOS ONE* 9 (July 18, 2014),
1058 e100587. issn: 1932-6203. doi: 10.1371/journal.pone.0100587.
- 1059 [64] Untergasser A, Cutcutache I, Koressaar T, Ye J, Faircloth BC, Remm M, and Rozen SG.
1060 Primer3—new capabilities and interfaces. *Nucleic Acids Research* 40 (Aug. 1, 2012), e115–
1061 e115. issn: 0305-1048. doi: 10.1093/nar/gks596.
- 1062 [65] Wallny HJ, Avila D, Hunt LG, Powell TJ, Riegert P, Salomonsen J, Skjødt K, Vainio O, Vilbois
1063 F, Wiles MV, and Kaufman J. Peptide motifs of the single dominantly expressed class I
1064 molecule explain the striking MHC-determined response to Rous sarcoma virus in chick-
1065 ens. *Proceedings of the National Academy of Sciences of the United States of America* 103
1066 (Jan. 31, 2006), 1434–1439. issn: 0027-8424, 1091-6490. doi: 10.1073/pnas.0507386103.
- 1067 [66] Westbrook CJ, Karl JA, Wiseman RW, Mate S, Koroleva G, Garcia K, Sanchez-Lockhart M,
1068 O'Connor DH, and Palacios G. No assembly required: Full-length MHC class I allele dis-
1069 covery by PacBio circular consensus sequencing. *Human Immunology*. Single-Molecule
1070 DNA Sequencing 76 (Dec. 1, 2015), 891–896. issn: 0198-8859. doi: 10.1016/j.humimm.
1071 2015.03.022.
- 1072 [67] Wittzell H, Bernot A, Auffray C, and Zoorob R. Concerted evolution of two Mhc class II
1073 B loci in pheasants and domestic chickens. *Molecular Biology and Evolution* 16 (Apr. 1,
1074 1999), 479–490. issn: 0737-4038. doi: 10.1093/oxfordjournals.molbev.a026130.
- 1075 [68] Wood S. *Generalized Additive Models: An Introduction with R*. CRC Press, Feb. 27, 2006.
1076 412 pp. isbn: 978-1-58488-474-3.
- 1077 [69] Worley K, Gillingham M, Jensen P, Kennedy LJ, Pizzari T, Kaufman J, and Richardson DS.
1078 Single locus typing of MHC class I and class II B loci in a population of red jungle fowl.
1079 *Immunogenetics* 60 (May 2008), 233–47. issn: 0093-7711. doi: 10.1007/s00251-008-
1080 0288-0.
- 1081 [70] Zagalska-Neubauer M, Babik W, Stuglik M, Gustafsson L, Cichoń M, and Radwan J. 454
1082 sequencing reveals extreme complexity of the class II Major Histocompatibility Com-
1083 plex in the collared flycatcher. *BMC Evolutionary Biology* 10 (Dec. 31, 2010), 395. issn:
1084 1471-2148. doi: 10.1186/1471-2148-10-395.