

1 Inference of ploidy by leveraging read depth from amplicon sequencing

2

3 Thomas A. Delomas*¹, Stuart C. Willis², Andrea Schreier³, Shawn Narum²

4

5 ¹ Pacific States Marine Fisheries Commission / Idaho Department of Fish and Game, Eagle Fish

6 Genetics Laboratory, Eagle, ID, USA

7 ² Columbia River Inter-Tribal Fish Commission, Hagerman Genetics Lab, 3059F National Fish

8 Hatchery Rd, Hagerman ID, USA

9 ³ Genomic Variation Laboratory, Department of Animal Science, University of California Davis,

10 Davis, CA 95616 USA

11

12 *Corresponding author, thomas.delomas@idfg.idaho.gov

13

14

15

16

17 Abstract

18 Variation in ploidy occurs naturally in select plant and animal species. Ploidy variation
19 can also occur spontaneously or be induced during artificial propagation of fish and shellfish.
20 Studying species and systems that have variable ploidy requires techniques to infer ploidy of
21 individuals. Massively parallel sequencing of biallelic SNPs has been used to infer ploidy, but
22 existing techniques have several drawbacks. These include being limited to only comparing a
23 fixed number of ploidies (diploidy, triploidy, and tetraploidy) and requiring that heterozygous
24 genotypes in an individual be identified prior to ploidy inference. We describe a method of
25 inferring ploidy from sequencing of biallelic SNPs based on beta-binomial mixture models. This
26 method is generalized to apply to any ploidy and does not require prior identification of
27 heterozygous genotypes. We demonstrate efficacy of this method for comparing ancestral
28 octoploidy, decaploidy, and dodecaploidy (tetraploidy, pentaploidy, and hexaploidy for the
29 sequenced SNPs) in white sturgeon and diploidy and triploidy in Chinook salmon with amplicon
30 sequencing (GT-seq) data. Results indicated that ploidy could be reliably estimated for
31 individuals based on distinct distribution of log-likelihood ratios (LLR) for known ploidy
32 samples of both species that were tested. Confidence in ploidy estimates increased with
33 sequencing depth. We encourage users to explore the sequencing depths and LLR critical values
34 that provide reliable estimates of ploidy for a given organism and set of SNPs. We expect that the
35 R package provided will empower studies of genetic variation and inheritance in organisms that
36 vary in ploidy naturally or as a result of artificial propagation practices.

37

38

39

40 Introduction

41 The number of haploid copies of the genome in somatic cells of an individual, termed
42 ploidy, is naturally variable in some plant and animal species (Lamatsch & Stöck, 2009; Mock et
43 al., 2012; Yamashita, Jiang, Onozato, Nakanishi, & Nagahama, 1993; Zhang & Arai, 1999).
44 Ploidy variation in fish has been observed to occur spontaneously (Gold & Avise, 1976; Machado
45 et al., 2012; Utsunomia et al., 2014). Some artificial spawning and rearing practices may increase
46 the rate of spontaneous ploidy variation (Aegerter & Jalabert, 2004; Cherfas, Gomelsky, Ben-
47 Dom, & Hulata, 1995; Delomas & Dabrowski, 2016; Flajšhans, Kvasnička, & Ráb, 1993; Glover
48 et al., 2015; Thorgaard et al., 1982; Van Eenennaam et al., 2020) which could lead to more ploidy
49 variation in systems stocked with hatchery-origin fish. Alterations in ploidy can also be induced
50 in fish and shellfish to yield individuals with advantageous qualities for cultivation, such as
51 sterility (Benfey, 1999; Nell, 2002), or stocking in water bodies where reproduction of stocked
52 fish is undesirable (Cassinelli, Meyer, Koenig, Vu, & Campbell, 2018). Ploidy variation is linked
53 to reproductive isolation (Husband & Sabara, 2004; Husband, Schemske, Burton, & Goodwillie,
54 2002), speciation (Ptacek, Gerhardt, & Sage, 1994; Wood et al., 2009), changes in gamete ploidy
55 and reduced fertility (Delomas & Dabrowski, 2018; Feindel, Benfey, & Trippel, 2010; Liu et al.,
56 2001), and differences in metabolism (Hyndman, Kieffer, & Benfey, 2003; Leal, Clark, Van
57 Eenennaam, Schreier, & Todgham, 2018). Studies in systems and species with variable ploidy
58 that do not account for such variation therefore risk ignoring an important confounding factor.

59 Techniques to infer the ploidy of individual samples are required when the study system
60 and species have variable ploidy. Ploidy is commonly inferred using flow cytometry to directly
61 measure nuclear DNA content in cells from a blood or solid tissue sample (Delomas &
62 Dabrowski, 2018). A Coulter counter is also commonly used in animals when a fresh blood

63 sample can be easily obtained (Wattendorf, 1986). However, these techniques can be untenable
64 for ploidy determination when sampling is conducted in remote locations, far from a flow
65 cytometer, Coulter counter, or the reagents and consumables needed to fix samples for future
66 flow cytometry. A researcher may also wish to determine the ploidy of an archived tissue sample,
67 such as from a museum specimen. Because fresh or specially fixed tissues are required for
68 Coulter counter or flow cytometry analyses, another method is necessary to determine ploidy of
69 archived samples.

70 To address this shortcoming, methods have been developed to infer ploidy from massively
71 parallel sequencing data, namely read counts at biallelic single nucleotide polymorphisms
72 (SNPs). One graphical technique assisted by ploidyNGS (Augusto Corrêa dos Santos, Goldman,
73 & Riaño-Pachón, 2017) is to visually inspect histograms of allele depth ratios (Figure 1). The
74 number and location of peaks corresponding to heterozygous genotypes can be used to classify
75 samples. A drawback of this technique is that it requires visual, not statistical, evaluation of
76 histograms.

77 Several methods that allow more automated ploidy inference have been developed. The R
78 package `gbs2ploidy` (Gompert & Mock, 2017) uses read counts to estimate relative allele dosage
79 at heterozygous loci. These estimates and observed heterozygosity are then used to categorize
80 sample ploidy with clustering algorithms. One drawback of this method is that it requires prior
81 identification of heterozygous genotypes. At higher levels of ploidy (e.g. octoploidy), confidently
82 separating homozygous genotypes from genotypes with only one copy of the minor allele can
83 require a large number of reads. In some situations, the depth required may be difficult to achieve
84 across a sufficient number of loci due to low sample quality or cost constraints. Additionally, the
85 current implementation of `gbs2ploidy` is limited to only discriminating between diploidy,

86 triploidy, and tetraploidy.

87 The program nQuire (Weiß, Pais, Cano, Kamoun, & Burbano, 2018) models observed
88 ratios of allele depth at heterozygous SNPs with a Gaussian mixture model. The means of the
89 Gaussians correspond to the allele dosage expected with a given ploidy (e.g., 1/3 and 2/3 for
90 triploidy). The use of Gaussian distributions, as compared to binomials, allows higher levels of
91 dispersion in the data to be modeled. The authors additionally demonstrate that a uniform noise
92 component can be added to model spurious observations. However, a drawback of this approach
93 is that modelling the ratios of allele depth, and not the read counts, ignores the relationship
94 between variance and depth. Second, the method requires identification of loci with heterozygous
95 genotypes in each individual as homozygous genotypes are not modelled. As mentioned above,
96 confidently separating homozygous genotypes from genotypes with one copy of the minor allele
97 can require a prohibitive number of reads at higher levels of ploidy. Third, this approach is
98 currently only implemented in nQuire for discriminating between diploidy, triploidy, and
99 tetraploidy.

100 A method based on a likelihood ratio statistic was developed to address some of the
101 shortcomings of the above approaches, as well as account for variable sequencing error and
102 allelic bias between loci (Delomas, 2019). This method excludes homozygous loci using a
103 binomial test and then calculates a likelihood ratio comparing diploidy and triploidy. The
104 likelihoods are calculated by assuming the read counts are binomial random variables. One
105 drawback of this method is that it is limited to comparing diploidy and triploidy. Additionally,
106 while it does not require the user to identify heterozygous loci, it attempts to identify them using
107 a binomial test. A final drawback is that modelling the read counts as binomials does not allow
108 for overdispersion which can be present in some sequencing data. As demonstrated (Delomas,

109 2019), this method performs well for differentiating diploids and triploids with amplicon
110 sequencing data, but a similar strategy may not be suitable for differentiating ploidy levels with
111 more similar allele dosages.

112 Our goal was to develop a method of inferring ploidy from high throughput sequencing
113 data for biallelic SNPs that addressed the drawbacks of existing methods and was applicable to
114 any ploidy. Our motivation stems from the case of white sturgeon (*Acipenser transmontanus*).
115 White sturgeon are ancestral octoploids (8n) (Drauch Schreier, Gille, Mahardja, & May, 2011),
116 but spontaneous autopolyploidy has been observed in hatchery settings, producing dodecaploids
117 (12n) (Van Eenennaam et al., 2020). Crossing individuals with these two ploidy levels then yields
118 decaploids (10n). Although flow cytometry and Coulter counter analysis can be used to
119 accurately distinguish between white sturgeon of different ploidies (Fiske et al., 2019), these
120 techniques cannot be used for archived tissue samples. However, a panel of biallelic SNPs was
121 developed (Willis et al., 2020) that are detected in four copies in the genomes of the ancestrally
122 octoploid white sturgeon. We developed a method to efficiently distinguish between tetraploidy,
123 pentaploidy, and hexaploidy using these SNPs, corresponding to octoploidy, decaploidy, and
124 dodecaploidy on the ancestral scale, respectively. We here describe the method and validate it
125 both with white sturgeon of three ploidies and Chinook salmon of two ploidies.

126

127 Methods

128 *Beta-binomial mixture model*

129 For the purpose of this explanation, we use the term “reference allele” to refer arbitrarily
130 to one of the alleles of a biallelic SNP. A biallelic SNP in an individual with ploidy x has $x + 1$
131 possible states corresponding to 0, 1, ..., x copies of the reference allele. As such, we model

132 counts of the reference allele as a mixture with $x + 1$ components. The weights of the components
133 correspond to the proportion of SNPs in each state. A similar approach is implemented in nQuire
134 except that components corresponding to states 0 and x are not included (Weiß et al., 2018).

135 Observation of a read of the reference allele can be considered a Bernoulli random
136 variable with probability of success dependent upon the true state of the SNP, x , the rate of
137 sequencing error, and allelic bias. Gerard et al. (2018) derived an equation for calculating this
138 probability. Because individual reads can be considered Bernoulli random variables, it is natural
139 to model the read count as a binomial random variable. Considering a mixture of binomials, the
140 likelihood of observing c_i counts of the reference allele given n_i total reads and a ploidy of x is

141

$$L = \sum_{s=0}^x w_s \text{Bin}(c_i; n_i, p_{is}),$$

142 where w_s is the weight of the component for state s and p_{is} is the probability of observing a read
143 of the reference allele for locus i and state s calculated according to Gerard et al. (2018).

144 Overdispersion is commonly present in sequencing data. One solution is to model read
145 counts with a beta-binomial distribution and an overdispersion parameter, τ (Gerard et al., 2018).

146 We adopt this approach, and the likelihood now becomes

147

$$L = \sum_{s=0}^x w_s \text{BB}(c_i; n_i, \alpha_{is}, \beta_{is}),$$

148

149 with

$$\alpha_{is} = p_{is} \frac{1 - \tau_s}{\tau_s},$$

$$\beta_{is} = (1 - p_{is}) \frac{1 - \tau_s}{\tau_s}.$$

150 Gerard et al. (2018) developed a method to call genotypes in groups of individuals of known
151 ploidy and inferred values of τ that were specific to each locus. In the current application, we
152 wanted to avoid using information from multiple individuals to simplify application of this
153 method and to eliminate the need to either assume that τ is constant across ploidies or develop a
154 mechanism for estimating different values of τ for each possible ploidy using samples of
155 unknown ploidy. Therefore, a value of τ was defined for each state, τ_s , within a given individual
156 and ploidy. This decision implies that variance depends on the true state of a locus. nQuire
157 similarly fits the variance of Gaussian distributions separately between states within an individual
158 (Weiß et al., 2018). Assuming independence between loci, the total likelihood is the product of
159 the likelihoods of each locus.

160 Addition of a uniform noise component to the Gaussian mixture model used by nQuire
161 was demonstrated to be helpful when analyzing noisy data (Weiß et al., 2018). The same addition
162 is possible with a beta-binomial mixture model. The model then has $x + 2$ components and the
163 likelihood for one locus is

$$L = \sum_{s=0}^x w_s BB(c_i; n_i, \alpha_{is}, \beta_{is}) + w_{x+1} BB(c_i; n_i, 1, 1).$$

164 This model can be fit using an expectation-maximization (EM) algorithm. In the
165 expectation step, weights of each component are updated as typical for a mixture model

$$w_s = \frac{1}{N} \sum_{i=1}^N \frac{w_s BB(c_i; n_i, \alpha_{is}, \beta_{is})}{\sum_{j=0}^x w_j BB(c_i; n_i, \alpha_{ij}, \beta_{ij}) + w_{x+1} BB(c_i; n_i, 1, 1)},$$

166 with all elements of w being updated using the values of w from the previous iteration. The
167 maximization step updates the values of τ by maximizing the log-likelihood using the method of
168 Byrd et al. (1995) informed by the analytical gradient.

169 Functions to fit this model as well as the intermediary models described (binomial
170 mixture model and beta-binomial without uniform noise) are implemented in an R package,
171 `tripsAndDipR` v0.2.0, available at www.github.com/delomast/tripsAndDipR.

172

173 *Inferring ploidy*

174 To infer ploidy of an individual, models with different ploidies need to be compared. The
175 maximum likelihood estimate (MLE) is the ploidy corresponding to the model with the highest
176 likelihood. Models can be more quantitatively compared by calculating a log-likelihood ratio
177 between two models. If more than two ploidies are possible, a series of log-likelihood ratios
178 (LLR) comparing the most likely model to all others can be calculated. Rejection of less likely
179 models depends on the distribution of the log-likelihood ratios. These distributions are unknown
180 but can be empirically approximated by assessing samples of known ploidy. If a reliable method
181 of simulating data given ploidy is available for a particular sequencing protocol, a Monte Carlo
182 method for approximating these distributions could also be used.

183

184 *Assessing the method*

185 To assess the efficacy of this method, we applied it to three groups of samples. All
186 samples were fin clips. The first group, subsequently referred to as “known ploidy sturgeon”,

187 consisted of 19 octoploid and 23 dodecaploid white sturgeon. These sturgeon were from a Central
188 California caviar farm (original broodstock source: Sacramento River). Ploidy was confirmed
189 using a Coulter counter at the University of California Davis. These samples were genotyped at
190 325 SNPs according to Willis et al. (2020). Samples were initially sequenced targeting a depth
191 ten times higher than Willis et al. (2020) recommended to achieve high genotyping success, and
192 reads were then randomly down-sampled at levels of 10%, 30%, and 50% to evaluate the effect
193 of sequencing depth on ploidy inference.

194 The second group, subsequently referred to as “presumed decaploid sturgeon”, consisted
195 of 17 full-sibling white sturgeon produced at a caviar farm in Central California (broodstock
196 source: Sacramento River). Individual parents had been identified as octoploid and dodecaploid
197 through the method described here, and so the method was applied to their offspring who were
198 presumed to be decaploid. This group of sturgeon was genotyped according to Willis et al.
199 (2020). In both groups of sturgeon, tetraploid, pentaploid, and hexaploid models were considered.

200 The third group of samples, subsequently referred to as “known ploidy salmon”, consisted
201 of 93 triploid Chinook salmon *Oncorhynchus tshawytscha* from Idaho Department of Fish and
202 Game’s Nampa Fish Hatchery whose ploidy was confirmed by flow cytometry and 80 diploid
203 Chinook salmon from Idaho Department of Fish and Game’s Rapid River Fish Hatchery
204 (diploidy confirmed by successful reproduction). Fin clips were taken from these fish and
205 genotyped according to the GT-seq method of amplicon sequencing (Campbell, Harmon, &
206 Narum, 2015) with a panel of 342 SNPs. Diploid and triploid models were considered for these
207 samples.

208 We assumed no allelic bias and sequencing error rates of 0.01 (1%) for all loci in all
209 analyses.

210

211 Results

212 In the known ploidy sturgeon, mean depth of SNPs within individuals had mean \pm SD of
213 528 ± 249 , 1585 ± 746 , 2642 ± 1243 , 5283 ± 2486 reads per locus at subsampling levels of 10,
214 30, 50, and 100%, respectively. Mean depth of SNPs within individuals of the presumed
215 decaploid sturgeon and known ploidy salmon was 752 ± 284 and 300 ± 92 , respectively. The
216 MLEs for the known ploidy sturgeon in all subsampling levels and the known ploidy salmon
217 were correct. The MLEs for the presumed decaploid sturgeon were all decaploid, fitting
218 expectations based on the inferred ploidy of their parents. The LLRs comparing the true ploidy
219 with alternative ploidies were centered away from zero (Figures 2, 3, and 4) and the distance
220 increased with increasing depth (Figure 2).

221 Fitting the described beta-binomial model and comparing ploidies through LLR
222 accurately separated individuals according to true ploidy. The magnitude of separation of
223 different ploidies with a given set of SNPs was dependent on sequencing depth (Figure 2). The
224 lowest down-sampling level (10%), which corresponds to targeting the depth recommended for
225 genotyping by Willis et al. (2020), gave accurate MLEs for ploidy. While not demonstrated in
226 these analyses, the statistical model implied that the magnitude of separation was also influenced
227 by the variability of the SNPs in the sequenced individuals. Genotypes with relative allele
228 dosages that were shared between ploidies did not contribute information about ploidy. This
229 included homozygous genotypes and, in comparisons of tetraploidy and hexaploidy, genotypes
230 with equal numbers of both alleles (relative dosage of 0.5).

231 We found that the distribution of LLRs for white sturgeon samples of known and
232 presumed ploidy could be used to set critical values for rejecting less likely models. With the

233 panel of SNPs and depth targeted in this study, a critical value of 10 was appropriate for rejecting
234 alternative ploidy models. Very few individuals had an LLR less than 10, and this critical value
235 did not result in any false classifications of the known ploidy samples (Figures 2 and 3).

236

237 Discussion

238 We describe here a new statistical model for inferring ploidy from sequencing data and
239 demonstrate its efficacy using amplicon sequencing for two species of varying ploidy levels.
240 Increased sequencing depth increased the likelihood of the correct ploidy (Figure 2). The
241 relationship between mean depth and accuracy of inferred ploidy depends on the ploidies being
242 assessed, the number and variability of loci, and the desired level of confidence. Additionally, the
243 observed variance in the read counts over what would be expected from a binomial random
244 variable (overdispersion) impacts the depth required. As such, we recommend users evaluate
245 minimum mean depth requirements for their panel and species.

246 Unlike previous methods (Delomas, 2019; Gompert & Mock, 2017; Weiß et al., 2018),
247 the current method is generalized to assess any ploidy and does not require identification of
248 heterozygous genotypes in an individual prior to ploidy inference. However, as noted by Gompert
249 and Mock (2017), inferring ploidy from sequencing data cannot separate individuals of lower and
250 higher ploidy when the higher ploidy is formed solely by duplicating a lower ploidy genome. An
251 example is when a tetraploid is formed by suppression of the first mitotic division in an embryo.
252 This is because the allelic ratios for the higher ploidy are identical to those expected in the lower
253 ploidy. Polyploidy of this kind is relatively rare, and so the method described here is expected to
254 apply in most circumstances.

255 We suggested a critical value of 10 for the white sturgeon panel based on visual

256 evaluation of the LLR distributions for the known and presumed ploidy white sturgeon. With
257 larger sample sizes, a more quantitative choice of critical values is possible by using those
258 samples to estimate false positive and false negative rates for a given critical value and
259 comparison.

260 When differentiating between ploidies that are multiples of each other (e.g. diploid and
261 tetraploid), the set of all possible models of the higher ploidy contains all possible models of the
262 lower ploidy. nQuire addresses this for the case of diploids and tetraploids by fixing all three
263 component weights of the tetraploid mixture model at 1/3, effectively assuming that the
264 heterozygous genotype states occur in fixed proportions (Weiß et al., 2018). While they
265 demonstrate the efficacy of this approach, it is unclear whether ploidy inferences would still be
266 accurate when the true genotype proportions have large deviations from those assumed. Gompert
267 and Mock (2017) did not restrict proportions of genotype states. They relied on tetraploids being
268 sufficiently separated from diploids by posterior allele dosages for a clustering algorithm to
269 separate the two categories. The current method also does not restrict proportions of genotype
270 states. As such, when comparing ploidies that are multiples of each other the larger ploidy will
271 always have a likelihood higher than or equal to that of the smaller ploidy (apart from deviations
272 due to the threshold at which convergence is assumed). The larger ploidy can have a higher
273 likelihood due to over-fitting. Ploidy can still be inferred, however, as the distribution of LLR
274 should be approximately bimodal: samples with a true smaller ploidy will have smaller LLR
275 (distributed close to zero), and those with the larger ploidy will have larger LLR (distributed
276 further away from zero). Additionally, the current method estimates the proportion of loci in each
277 genotype state (the component weights) and these can be compared with expectations based on
278 the species' biology. For example, when comparing diploidy and tetraploidy for a sample and

279 fitting the tetraploid model, if the proportions of genotypes in states 1 and 3 (genotypes of ABBB
280 and AAAB) are estimated to be close to zero, then it may be reasonable to categorize this sample
281 as diploid. This logic is similar to that of restricting the proportions of genotype states.

282 When integrated into pipelines utilizing amplicon sequencing data, e.g. GT-seq (Campbell
283 et al., 2015), the routine presented herein provides a straightforward and effective method by
284 which samples can be simultaneously genotyped and ploidy inferred from archived as well as
285 fresh tissue samples of diverse types. We provide a convenient R package by which this can be
286 accomplished (tripsAndDipR v 0.2.0 available at www.github.com/delomast/tripsAndDipR).
287 While we encourage users to explore the sequencing depths, heterozygosity, and LLR critical
288 values that provide reliable and robust estimates of ploidy in each particular organism, we expect
289 that this package will empower studies of genetic variation and inheritance in organisms that vary
290 in ploidy naturally or as a result of artificial propagation practices.

291

292 Acknowledgments

293 We would like to thank staff at Idaho Department of Fish and Game's Nampa Fish Hatchery,
294 Rapid River Fish Hatchery, and Fish Health Laboratory for providing tissue samples of Chinook
295 salmon and the staff of the Idaho Department of Fish and Game's Eagle Fish Genetics Laboratory
296 for assistance in genotyping the Chinook salmon samples. We would like to thank Ken Lepla and
297 Idaho Power Company for assistance in acquiring samples. For white sturgeon, Lori Maxwell of
298 CRITFC assisted with collection of laboratory data. Joel Van Eenennaam collected fin clips from
299 known octoploid and dodecaploid white sturgeon. Funding was provided by Bonneville Power
300 Administration project 2008-907-00.

301

302 Conflict of Interest

303 No conflict of interest to report.

304

305 Author Contributions

306 TAD derived the statistical model, wrote the R package, assessed the model, and drafted the
307 manuscript. SN and SCW sequenced the white sturgeon samples. AS provided the sturgeon
308 samples and collected the presumed decaploid sturgeon samples. All authors participated in
309 editing and revising the manuscript.

310

311

312 References

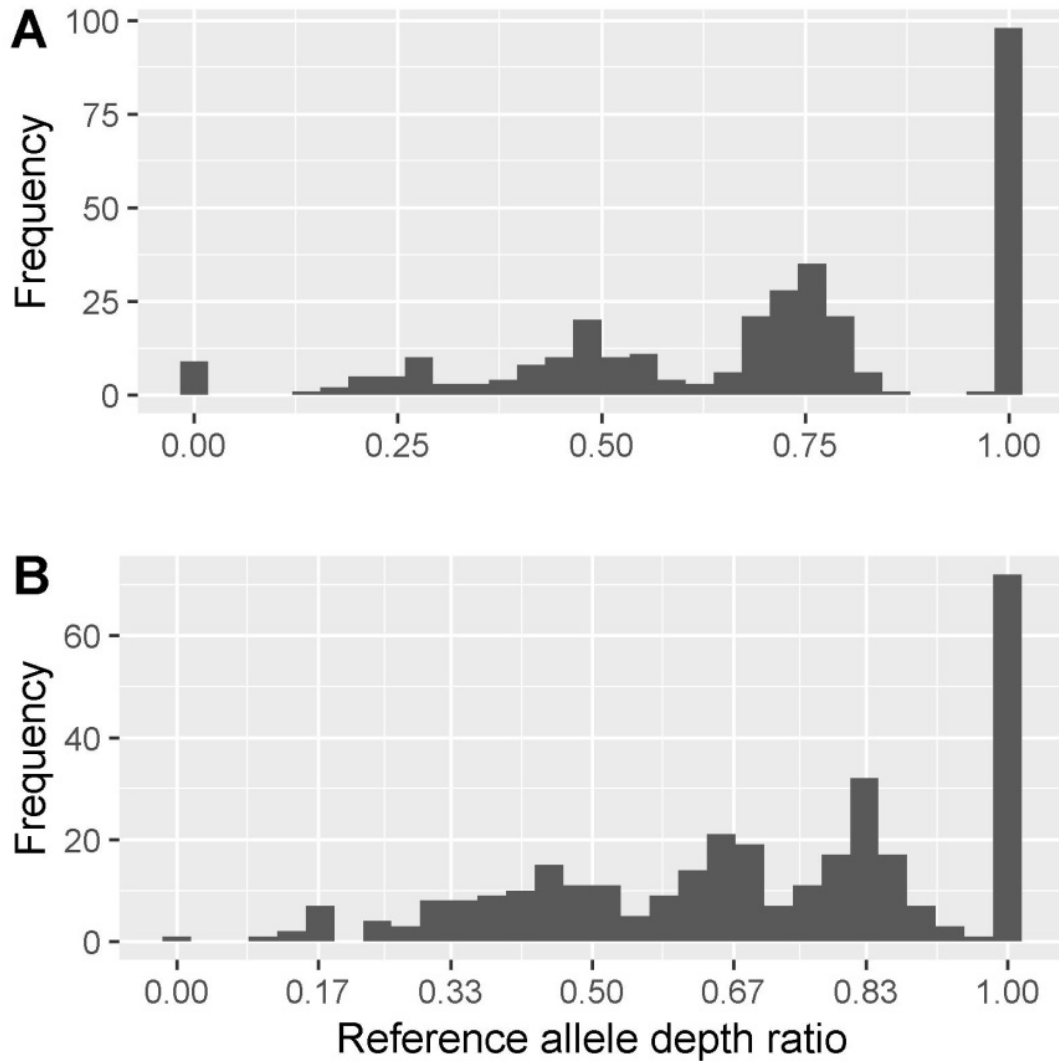
- 313 Aegerter, S., & Jalabert, B. (2004). Effects of post-ovulatory oocyte ageing and temperature on
314 egg quality and on the occurrence of triploid fry in rainbow trout, *Oncorhynchus mykiss*.
315 *Aquaculture*, 231(1–4), 59–71. doi: 10.1016/j.aquaculture.2003.08.019
- 316 Augusto Corrêa dos Santos, R., Goldman, G. H., & Riaño-Pachón, D. M. (2017). ploidyNGS:
317 visually exploring ploidy with next generation sequencing data. *Bioinformatics*, 33(16),
318 2575–2576. doi: 10.1093/bioinformatics/btx204
- 319 Benfey, T. J. (1999). The physiology and behavior of triploid fishes. *Reviews in Fisheries*
320 *Science*, 7(1), 39–67. doi: 10.1080/10641269991319162
- 321 Byrd, R. H., Lu, P., Nocedal, J., & Zhu, C. (1995). A limited memory algorithm for bound
322 constrained optimization. *SIAM Journal on Scientific Computing*, 16(5), 1190–1208. doi:
323 10.1137/0916069
- 324 Campbell, N. R., Harmon, S. A., & Narum, S. R. (2015). Genotyping-in-Thousands by
325 sequencing (GT-seq): A cost effective SNP genotyping method based on custom amplicon
326 sequencing. *Molecular Ecology Resources*, 15(4), 855–867. doi: 10.1111/1755-0998.12357
- 327 Cassinelli, J. D., Meyer, K. A., Koenig, M. K., Vu, N. V., & Campbell, M. R. (2018).
328 Performance of diploid and triploid westslope cutthroat trout fry stocked into Idaho alpine
329 lakes. *North American Journal of Fishery Management*, 39(1), 112–123. doi:
330 10.1002/nafm.10254
- 331 Cherfas, N., Gomelsky, B., Ben-Dom, N., & Hulata, G. (1995). Evidence for the heritable nature
332 of spontaneous diploidization in common carp, *Cyprinus carpio* L., eggs. *Aquaculture*
333 *Research*, 26(4), 289–292. doi: 10.1111/j.1365-2109.1995.tb00914.x
- 334 Delomas, T. A. (2019). Differentiating diploid and triploid individuals using single nucleotide
335 polymorphisms genotyped by amplicon sequencing. *Molecular Ecology Resources*, 19(6),
336 1545–1551. doi: 10.1111/1755-0998.13073
- 337 Delomas, T. A., & Dabrowski, K. (2016). Zebrafish embryonic development is induced by carp
338 sperm. *Biology Letters*, 12(11), 20160628.
- 339 Delomas, T. A., & Dabrowski, K. (2018). Why are triploid zebrafish all male? *Molecular*
340 *Reproduction and Development*, 85(7), 612–621. doi: 10.1002/mrd.22998
- 341 Drauch Schreier, A., Gille, D., Mahardja, B., & May, B. (2011). Neutral markers confirm the
342 octoploid origin and reveal spontaneous autoploidy in white sturgeon, *Acipenser*
343 *transmontanus*. *Journal of Applied Ichthyology*, 27(SUPPL. 2), 24–33. doi: 10.1111/j.1439-
344 0426.2011.01873.x
- 345 Feindel, N. J., Benfey, T. J., & Trippel, E. A. (2010). Competitive spawning success and fertility
346 of triploid male Atlantic cod *Gadus morhua*. *Aquaculture Environment Interactions*, 1, 47–
347 55. doi: 10.2307/24864017
- 348 Fiske, J. A., Van Eenennaam, J. P., Todgham, A. E., Young, S. P., Holem-Bell, C. E., Goodbla,
349 A. M., & Schreier, A. D. (2019). A comparison of methods for determining ploidy in white
350 sturgeon (*Acipenser transmontanus*). *Aquaculture*, 507, 435–442. doi:
351 10.1016/j.aquaculture.2019.03.009
- 352 Flajšhans, M., Kvasnička, P., & Ráb, P. (1993). Genetic studies in tench (*Tinca tinca* L.): high
353 incidence of spontaneous triploidy. *Aquaculture*, 110(3–4), 243–248. doi: 10.1016/0044-
354 8486(93)90372-6
- 355 Gerard, D., Ferrão, L. F. V., Garcia, A. A. F., & Stephens, M. (2018). Genotyping polyploids
356 from messy sequencing data. *Genetics*, 210(3), 789–807. doi: 10.1534/genetics.118.301468

- 357 Glover, K. A., Madhun, A. S., Dahle, G., Sørvik, A. G. E., Wennevik, V., Skaala, Ø., ... Fjellidal,
358 P. G. (2015). The frequency of spontaneous triploidy in farmed Atlantic salmon produced in
359 Norway during the period 2007–2014. *BMC Genetics*, *16*(1), 37. doi: 10.1186/s12863-015-
360 0193-0
- 361 Gold, J. R., & Avise, J. C. (1976). Spontaneous triploidy in the California roach *Hesperoleucus*
362 *symmetricus* (Pisces: Cyprinidae). *Cytogenetic and Genome Research*, *17*(3), 144–149. doi:
363 10.1159/000130706
- 364 Gompert, Z., & Mock, K. E. (2017). Detection of individual ploidy levels with genotyping-by-
365 sequencing (GBS) analysis. *Molecular Ecology Resources*, *17*(6), 1156–1167. doi:
366 10.1111/1755-0998.12657
- 367 Husband, B. C., & Sabara, H. A. (2004, March 1). Reproductive isolation between
368 autotetraploids and their diploid progenitors in fireweed, *Chamerion angustifolium*
369 (Onagraceae). *New Phytologist*, Vol. 161, pp. 703–713. doi: 10.1046/j.1469-
370 8137.2004.00998.x
- 371 Husband, B. C., Schemske, D. W., Burton, T. L., & Goodwillie, C. (2002). Pollen competition as
372 a unilateral reproductive barrier between sympatric diploid and tetraploid *Chamerion*
373 *angustifolium*. *Proceedings of the Royal Society of London. Series B: Biological Sciences*,
374 *269*(1509), 2565–2571. doi: 10.1098/rspb.2002.2196
- 375 Hyndman, C. A., Kieffer, J. D., & Benfey, T. J. (2003). Physiology and survival of triploid brook
376 trout following exhaustive exercise in warm water. *Aquaculture*, *221*(1–4), 629–643. doi:
377 10.1016/S0044-8486(03)00119-4
- 378 Lamatsch, D. K., & Stöck, M. (2009). Sperm-dependent parthenogenesis and hybridogenesis in
379 teleost fishes. In I. Schön, K. Martens, & P. Dijk (Eds.), *Lost Sex* (pp. 399–432). doi:
380 10.1007/978-90-481-2770-2_19
- 381 Leal, M. J., Clark, B. E., Van Eenennaam, J. P., Schreier, A. D., & Todgham, A. E. (2018). The
382 effects of warm temperature acclimation on constitutive stress, immunity, and metabolism in
383 white sturgeon (*Acipenser transmontanus*) of different ploidies. *Comparative Biochemistry*
384 *and Physiology -Part A*: *Molecular and Integrative Physiology*, *224*, 23–34. doi:
385 10.1016/j.cbpa.2018.05.021
- 386 Liu, S., Liu, Y., Zhou, G., Zhang, X., Luo, C., Feng, H., ... Yang, H. (2001). The formation of
387 tetraploid stocks of red crucian carp × common carp hybrids as an effect of interspecific
388 hybridization. *Aquaculture*, *192*(2), 171–186. doi: 10.1016/S0044-8486(00)00451-8
- 389 Machado, S. N., Neto, M. F., Bakkali, M., Vicari, M. R., Artoni, R. F., Oliveira, C. de, & Foresti,
390 F. (2012). Natural triploidy and B chromosomes in *Astyanax scabripinnis* (Characiformes,
391 Characidae): a new occurrence. *Caryologia*, *65*(1), 40–46. doi:
392 10.1080/00087114.2012.678086
- 393 Mock, K. E., Callahan, C. M., Islam-Faridi, M. N., Shaw, J. D., Rai, H. S., Sanderson, S. C., ...
394 Wolf, P. G. (2012). Widespread triploidy in western North American aspen (*Populus*
395 *tremuloides*). *PLoS ONE*, *7*(10), e48406. doi: 10.1371/journal.pone.0048406
- 396 Nell, J. A. (2002). Farming triploid oysters. *Aquaculture*, *210*(1–4), 69–88. doi: 10.1016/S0044-
397 8486(01)00861-4
- 398 Ptacek, M. B., Gerhardt, H. C., & Sage, R. D. (1994). Speciation by polyploidy in treefrogs:
399 Multiple origins of the tetraploid, *Hyla versicolor*. *Evolution*, *48*(3), 898–908. doi:
400 10.1111/j.1558-5646.1994.tb01370.x
- 401 Thorgaard, G. H., Rabinovitch, P. S., Shen, M. W., Gall, G. A. E., Propp, J., & Utter, F. M.

- 402 (1982). Triploid rainbow trout identified by flow cytometry. *Aquaculture*, 29(3–4), 305–
403 309. doi: 10.1016/0044-8486(82)90144-2
- 404 Utsunomia, R., Pansonato Alves, J. C., Paiva, L. R. S., Costa Silva, G. J., Oliveira, C., Bertollo,
405 L. A. C., & Foresti, F. (2014). Genetic differentiation among distinct karyomorphs of the
406 wolf fish *Hoplias malabaricus* species complex (Characiformes, Erythrinidae) and report of
407 unusual hybridization with natural triploidy. *Journal of Fish Biology*, 85(5), 1682–1692.
408 doi: 10.1111/jfb.12526
- 409 Van Eenennaam, J. P., Fiske, A. J., Leal, M. J., Cooley-Rieders, C., Todgham, A. E., Conte, F.
410 S., & Schreier, A. D. (2020). Mechanical shock during egg de-adhesion and post-ovulatory
411 ageing contribute to spontaneous autoploidy in white sturgeon culture (*Acipenser*
412 *transmontanus*). *Aquaculture*, 515, 734530. doi: 10.1016/j.aquaculture.2019.734530
- 413 Wattendorf, R. J. (1986). Rapid identification of triploid grass carp with a Coulter counter and
414 channelyzer. *The Progressive Fish Culturist*, 48(2), 125–132. doi: 10.1577/1548-
415 8640(1986)48<125:RIOTGC>2.0.CO;2
- 416 Weiß, C. L., Pais, M., Cano, L. M., Kamoun, S., & Burbano, H. A. (2018). nQuire: a statistical
417 framework for ploidy estimation using next generation sequencing. *BMC Bioinformatics*,
418 19(1), 122. doi: 10.1186/s12859-018-2128-z
- 419 Willis, S. C., Delomas, T. A., Parker, B., Miller, D., Anders, P., & Narum, S. (2020). Single
420 nucleotide polymorphism genotypes and ploidy estimates for ploidy variable species
421 generated with massively parallel amplicon sequencing. *Preprint/Submitted*, 0.
- 422 Wood, T. E., Takebayashi, N., Barker, M. S., Mayrose, I., Greenspoon, P. B., & Rieseberg, L. H.
423 (2009). The frequency of polyploid speciation in vascular plants. *Proceedings of the*
424 *National Academy of Sciences of the United States of America*, 106(33), 13875–13879. doi:
425 10.1073/pnas.0811575106
- 426 Yamashita, M., Jiang, J., Onozato, H., Nakanishi, T., & Nagahama, Y. (1993). A tripolar spindle
427 formed at meiosis I assures the retention of the original ploidy in the gynogenetic triploid
428 crucian carp, ginbuna *Carassius auratus langsdorfii*. *Develop. Growth & Differ*, 35(6), 631–
429 636. doi: 10.1111/j.1440-169X.1993.00631.x
- 430 Zhang, Q., & Arai, K. (1999). Distribution and reproductive capacity of natural triploid
431 individuals and occurrence of unreduced eggs as a cause of polyploidization in the loach,
432 *Misgurnus anguillicaudatus*. *Ichthyological Research*, 46(2), 153–161. doi:
433 10.1007/BF02675433
- 434
- 435

437 Figures

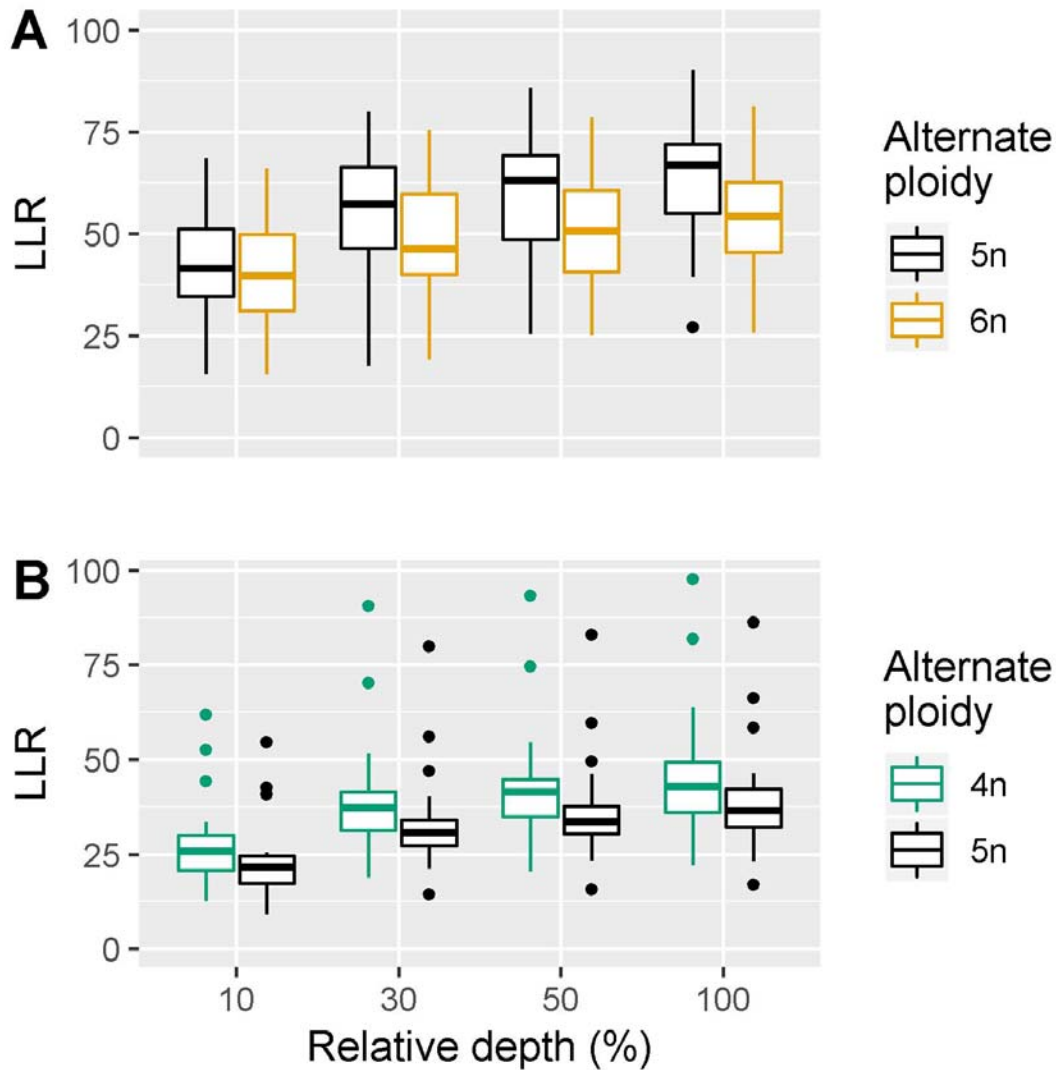
438 Figure 1. Distribution of the reference allele depth ratio (depth of reference allele / total depth)
439 for two white sturgeon of known ploidy. These graphs were generated similarly to the method
440 used by ploidyNGS (Augusto Corrêa dos Santos et al., 2017) A) True ancestral octoploid
441 (tetraploid for the genotyped SNPs) B) True ancestral dodecaploid (hexaploid for the genotyped
442 SNPs)



443

444

445 Figure 2. Boxplots of LLR for known ploidy white sturgeon at various depths. Values of LLR are
446 comparing the true ploidy with the alternate ploidy. 4n, 5n, 6n represent tetraploid, pentaploid,
447 and hexaploid, respectively A) True ancestral octoploids (tetraploid for the genotyped SNPs) B)
448 True ancestral dodecaploids (hexaploid for the genotyped SNPs)



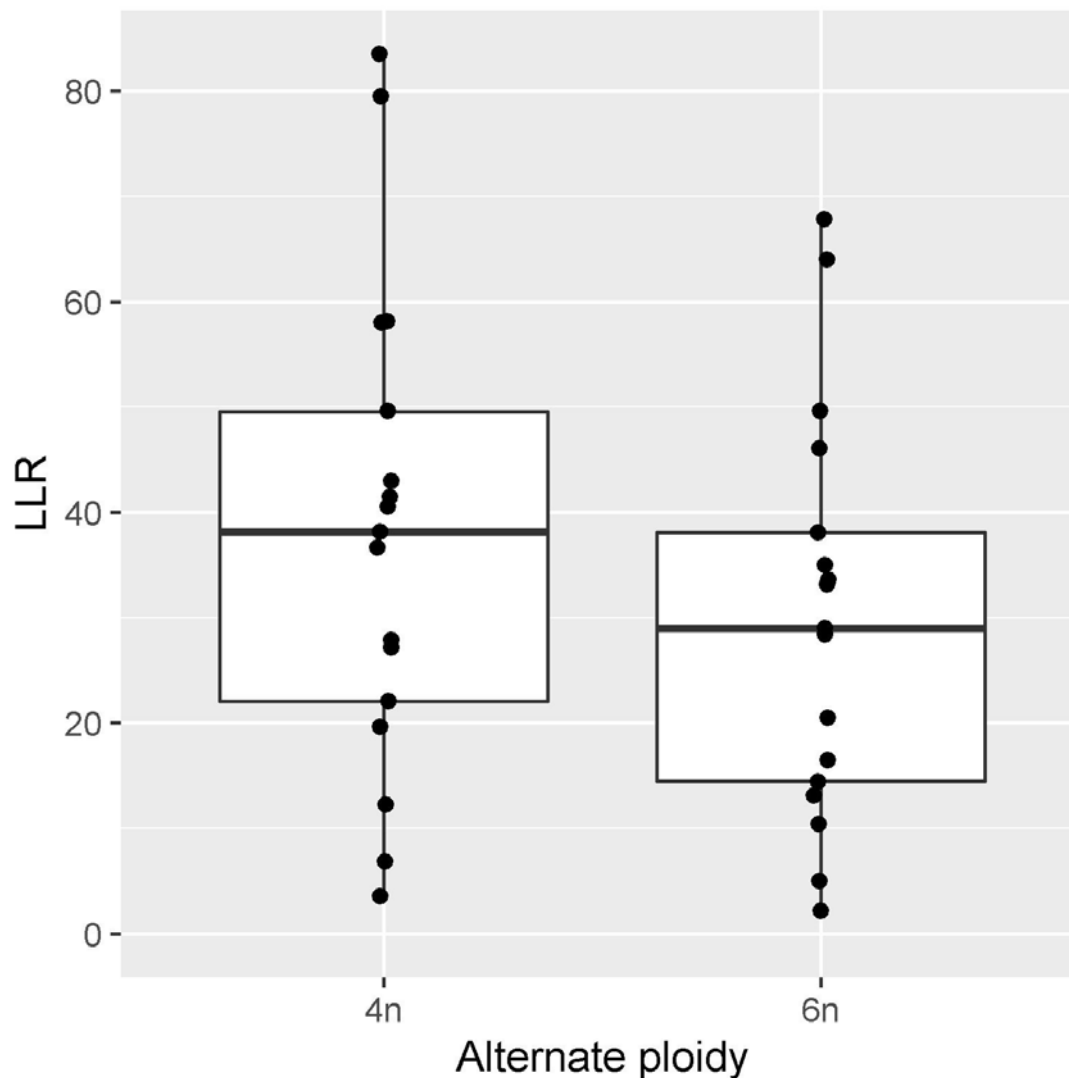
449

450

451

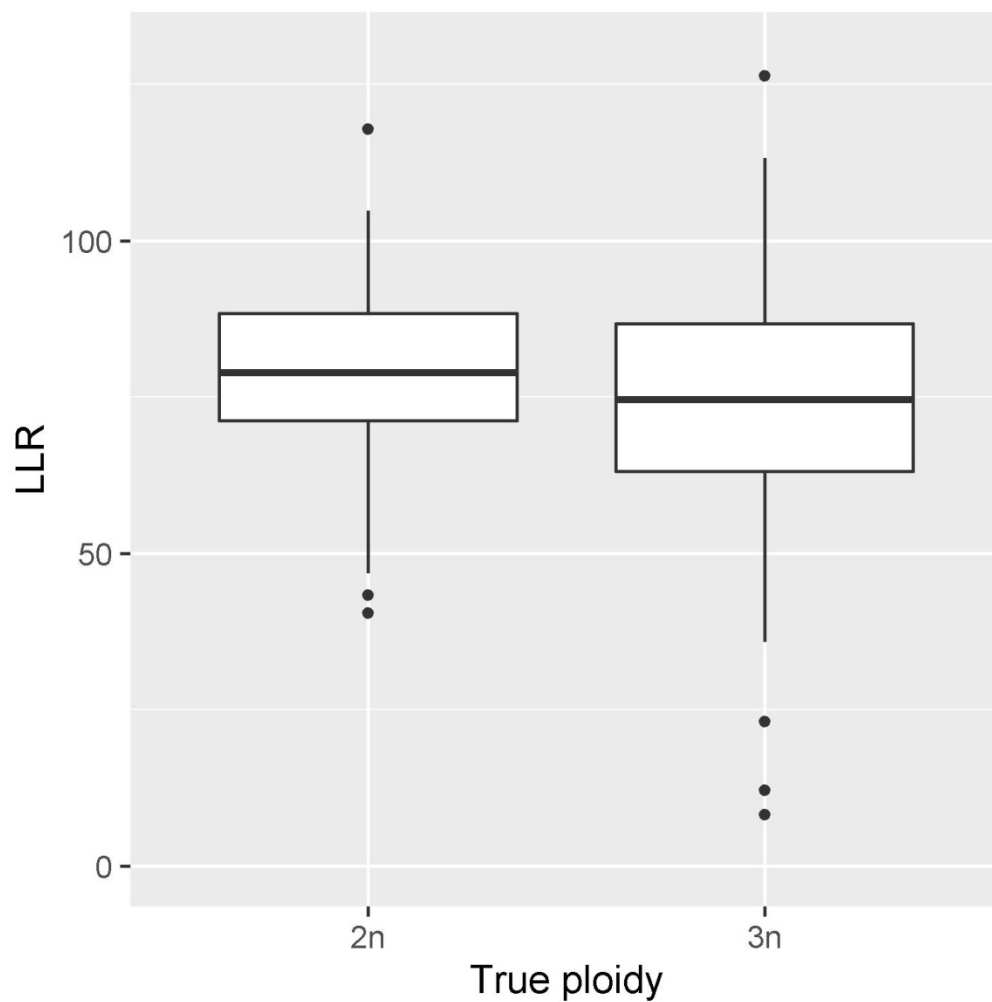
452

453 Figure 3. Boxplots of LLR for ancestral decaploid (pentaploid for the genotyped SNPs) white
454 sturgeon. Values of LLR are comparing pentaploidy with the alternate ploidy. All data points are
455 plotted overlying the boxplots. 4n and 6n represent tetraploid and hexaploid, respectively.



456
457
458

459 Figure 4. Boxplots of LLR for known ploidy Chinook salmon. Values of LLR are comparing the
460 true ploidy (x-axis) with the opposing ploidy. 2n and 3n represent diploid and triploid,
461 respectively



462

463

464