1    **A simple approach to multiplexed PCR amplicon sequencing of plumage-**

2    **associated loci in *Vermivora* warblers**

3

4    David P. L. Toews and Catharine E. Besch

5

6    Department of Biology, The Pennsylvania State University, 619 Mueller

7    Laboratory, University Park, State College, PA 16802. Corresponding Author:

8    toews@psu.edu

9

10    **Abstract:**

11    The dynamics of hybridization between golden-winged (*Vermivora chrysoptera*) and

12    blue-winged warblers (*V. cyanoptera*) has been of interest for over a century. Whole

13    genome analysis found only a small number of genomic regions that differed between

14    the species. We previously developed a restriction enzyme-based RFLP approach to

15    genotype large numbers of individuals at each of these loci. Here we extend this

16    approach to an amplicon sequencing method to genotype individuals at these six

17    plumage-associated regions. We demonstrate the efficacy using preliminary data from 4

18    golden-winged and 4 blue-winged warblers as well as provide the data and scripts

19    necessary to analyze these data for other interested in replicating this approach. Our

20    hope is that these data are useful for other researchers interested in genotyping

21    *Vermivora* warblers.

22

23    **Additional methods: https://github.com/david-toews/vermivora_amplicon**
24

**Introduction:**

The dynamics of hybridization between golden-winged (*Vermivora chrysoptera*) and blue-winged warblers (*V. cyanoptera*) has been of interest for over a century (Faxon 1913). Our previous worked initially sequencing whole genomes of the two species identified only six regions that were divergent between them within an otherwise homogeneous genomic background. Within nearly all of these regions housed genes involved in different aspects of plumage pigmentation—including melanogenesis and carotenoid-based traits—presumably underlying the phenotypic differences between the taxa (Toews et al. 2016, Baiz et al. 2020, Baiz et al. 2021).

We previously developed a restriction enzyme-based RFLP approach to genotype large numbers of individuals at each of these loci. However, the time involved in running—and sometimes re-running—agarose gels to obtain individual genotypes for hundreds of birds across the six loci was significant. Given the advances and accessibility of genomic analysis and amplicon sequencing, we therefore endeavored to optimize a method to genotype hundreds of individuals at the same time using amplicon sequencing. Given there are a number of research labs interested in the hybridization dynamics in *Vermivora*—and have been using the RFLP method—our goal here is to provide a simple outline of our methods such that others can replicate the same amplicon sequencing approach, if desired.

**Methods:**

*Sample and library preparation*

For a proof of concept, we obtained blood samples from 8 *Vermivora* warblers, 4 phenotypic golden-winged warblers and 4 phenotypic blue-winged warblers. From each bird, we collected blood samples from the brachial vein and stored them in Queen's lysis buffer. We conducted DNA extractions using the Qiagen DNeasy Blood and Tissue kit following the nucleated blood spin-column protocol. We used six primers (Table S1 from Toews et al. 2016) and prepared a single primer mix containing all forward and reverse primers each at a final concentration of 0.7uM. We ordered the six previously developed primers with P5 [TCGTCGGCAGCGTCAGCTGTGTATAAGAGACAG] and P7 [GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAG] adapter overhangs (note there is an error in Table S1 in Toews et al. [2016] regarding the restriction enzymes cut sites, which should show that locus "24-563" cuts GW – C with BW – A, and locus "25-653" cuts BW – T, with GW – C).

For each sample, we conducted individual 15uL PCR reactions containing: 3 uL 5X Platinum Buffer (Invitrogen 14966005), 0.3 uL 10mM dNTP mix (Promega U151A), 1.5 uL of the pooled primer mix, 0.12 uL of Platunim Taq HS II (Invitrogen 14966005), 8.58 uL of nuclease-free water, and 1.5 uL of the DNA samples (Qubit concentrations: 22.6-62.4 ng/uL). All primers were previously designed to have an annealing temperature of ~60ºC.

The thermocycling protocol consisted of:

        a. 3 min at 94ºC
        b. 35 cycles of:
            i. 30 sec at 94ºC
            ii. 30 sec at 63ºC

71           iii.    30 sec at 72ºC

72        c.   10 min at 72ºC

73        d.   Hold at 4ºC

74

75 For the indexing PCR, we initially diluted PCR products by 10X. We then added

76 combinatorial indexes based on Faircloth and Glenn (2012). The 30 uL indexing

77 reaction contained: 15 uL KAPA HiFi Hotstart ReadyMix, 6 uL of combined i5 and i7

78 indexes (3 uL of each index, if they are being added individually), and 9 uL of PCR1.

79      Reaction conditions for the indexing PCR consisted of:

80        a.   45 sec at 98ºC

81        b.   7 cycles of:

82           i.    15 sec at 98ºC

83           ii.    15 sec at 60ºC

84           iii.    15 sec at 72ºC

85        c.   1 min at 72ºC

86        d.   Hold at 12ºC

87

88 We then pooled equal volumes (5 uL) of each sample into a pool (to a total volume of

89 40 uL). We then cleaned the reaction using 1.7X SeraPure SPRI-beads (i.e. 68uL of

90 bead solution), washing with ethanol, and eluting in a final volume of 40 uL of water. We

91 checked the pool by TapeStation (D5000; Figure 1) and qPCR and sequenced it on a

92 MiSeq NANO (500 cycle, i.e. 250 x 250 bp reads) at the Penn State Huck Genomics

93 Core facility.

94

95 *Bioinformatics*

96      We created a single scaffold as a reference that consisted of six contigs,

97 covering the six plumage loci, each separated by 50 ambiguous nucleotides "N's" to aid

98 alignment. We aligned the resulting reads to this scaffold using *Bowtie2* with the "--very-

99    sensitive-local" set of presets. We then used *Samtools* to convert .SAM to .BAM files,

100   sort, and index each alignment.

101       We called SNPs on each individual alignment using GATK's (v 4.6.1.0)

102   HaplotypeCaller, with the ERC flag set to "GVCF" output and output mode to

103   "EMIT_ALL_CONFIDENT_SITES". We forced GATK to call SNPs at only our focal

104   SNPs using the -L flag and supplying it with the locations of the SNPs in a .bed file

105   (available at https://github.com/david-toews/vermivora_amplicon).

106       Typically we would use GATK to called genotypes jointly across all "g.vcf" files,

107   however even at the stage of combining genotypes GATK had issues with inserting

108   missing data when the raw genotypes had high coverage. The extremely high coverage

109   of the amplicon data may have produced issues at this stage, thus we simply extracted

110   the raw genotypes from the individual .vcf files using a custom script (available at the

111   GitHub repository).

112

113   **Results and Discussion:**

114       The pooled library resulted in three peaks on the TapeStation, one of which (at

115   168bp) was likely adapter dimer (Figure 1). The peaks above 300bp represented a size

116   distribution that matched the predicted sizes of the six amplicons.

117       Sequencing resulted in approximately 140,000 reads per sample, which is much

118   higher than is be necessary for calling genotypes. A target coverage is likely be between

119   10-50,000 reads per individual for future applications based on previous experience with

120   meta-barcoding datasets, though even lower may be possible.

121        Genotypes across all six plumage associated loci match expectations based on

122   RFLP cut sites and known phenotypes. In particular, the locus on Chromosome 20 is

123   perfectly associated with the black throat color phenotype observed in all four golden-

124   wings in this small sample (where a single heterozygous genotype would result in a

125   plain, white throat). This locus includes the *ASIP* gene involved in the melanogenesis

126   pathway (Toews et al. 2016, Baiz et al. 2020). Heterozygous bases were also seemingly

127   correctly called by the pipeline—based on manual inspection of the raw alignments—

128   such as the heterozygous site (Table 1) also clearly identified in Figure 3 in sample GW-

129   23-155.

130        Multiplexing many more individuals is feasible and we have plans to include 384

131   individuals in a single run using a NextSeq P1 chemistry (600 cycles i.e. 300 x 300 bp),

132   likely resulting in >250,000 reads per individual, as well as adding additional loci.

133   Therefore, this appears as an efficient and cost effective method for unambiguously

134   genotyping *Vermivora* at an important set of SNP loci across a large number of

135   individuals simultaneously.

136

137   **Acknowledgements:**

144    Pennsylvania State University, including funds from the Eberly College of Science and

145    the Huck Institutes of the Life Sciences.

146

147    **References:**

148    Baiz, M.D., Kramer, G.R., Streby, H.M., Taylor, S.A., Lovette, I.J. and Toews, D.P.L.,

149         2020. Genomic and plumage variation in Vermivora hybrids. *The Auk*. 137,

150         ukaa027.

151

152    Baiz, M.D., Wood, A.W. and Toews, D.P.L., 2021. Rare hybrid solves "genetic problem"

153         of linked plumage traits. *Ecology*, 102, e03424.

154

155    Faxon, W. 1913. Brewster's warbler (*Helminthophila leucobronchlais*) a hybrid between

156         the golden-winged warbler (*H. chrysoptera*) and the blue-winged warbler (*H.*

157         *pinus*). *Memoirs of the Museum of Comparative Zoology*. 40, 311–316.

158

159    Toews, D.P.L., Taylor, S.A., Vallender, R., Brelsford, A., Butcher, B.G., Messer, P.W. and

160         Lovette, I.J., 2016. Plumage genes and little else distinguish the genomes of

161         hybridizing warblers. *Current Biology*. 26, 2313-2318.

162

163

164

165

166    **Table 1.** Genotype results from amplicon sequencing approach across four blue-winged

167    warblers ("BW") and four golden-winged warblers ("GW"). The sites are relative to three

168    references: the scaffolds described in the original Toews et al. (2016) paper, the "v2"

169    reference from Baiz et al. (2021), and the position in the artificial reference scaffold

170    generated for this analysis.

| mywa_v1.1 scaffold | mywa_v2 chromosome | Pos. in Amplicon Reference Scaffold | BW-23-116 | BW-23-119 | BW-23-125 | BW-23-126 | GW-23-026 | GW-23-118 | GW-23-154 | GW-23-155 |
|---|---|---|---|---|---|---|---|---|---|---|
| 120 | chr4a | 228 | 1/1 | 1/1 | 1/1 | 0/0 | 0/0 | 0/0 | 0/0 | 0/0 |
| 299 | chr20 | 611 | 0/0 | 0/0 | 0/0 | 0/0 | 1/1 | 1/1 | 1/1 | 1/1 |
| 563 | chr24 | 1085 | 0/0 | 0/1 | 0/0 | 0/0 | 0/0 | 0/1 | 1/1 | 1/1 |
| 653 | chr25 | 1555 | 1/1 | 1/1 | 1/1 | 1/1 | 0/0 | 0/0 | 0/1 | 0/0 |
| 24 | chrZ | 1802 | 0/0 | 0/0 | 0/0 | 0/0 | 1/1 | 1/1 | 1/1 | 0/1 |
| 38 | chrZ | 2264 | 1/1 | 1/1 | 1/1 | 0/1 | 0/0 | 0/0 | 0/0 | 0/0 |

171

172

173

174

175 **Figure 1.** TapeStation results from pooled amplicon sequencing of 8 *Vermivora*

176 warblers. The 521bp peak represents the main 4 amplicons, the small peak to the right

177 is the 580 amplicon, and the smaller peak at ~330 is the 299bp amplicon. The 168 peak
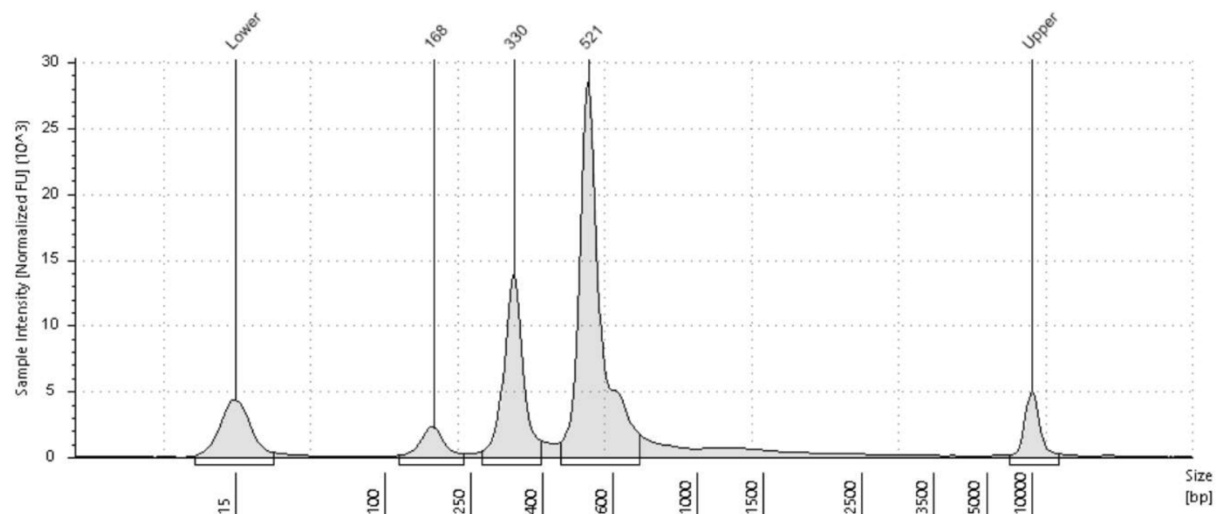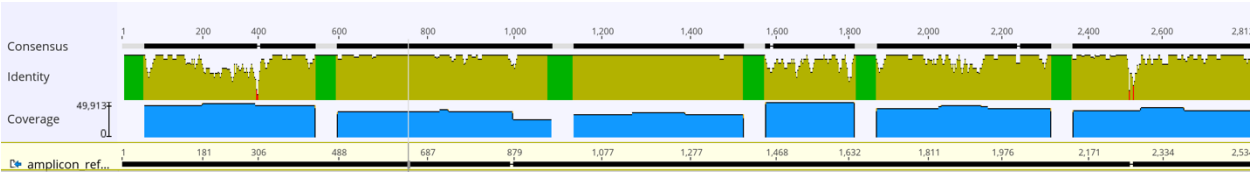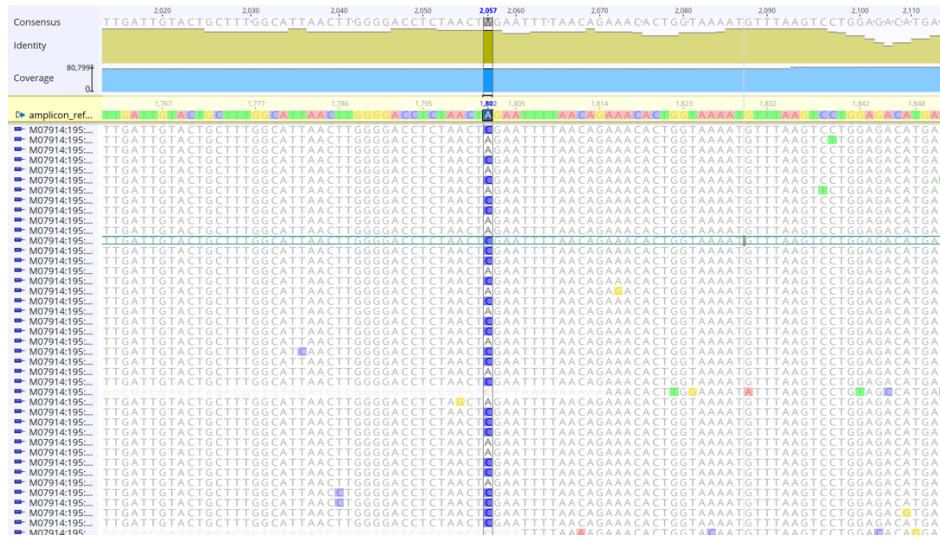
178 is likely adapter dimer.



179

180

181 **Figure 2.** Coverage output from Geneious of sample BW 23-116. The green areas of

182 the "identity" plot represent the ambiguous nucleotides seperating the 6 amplicon

183 sequences.



184

190 **Figure 3.** A heterozygous at site 1,802 bp for sample GW 23-155.



191