1    **Patterns of transcriptome divergence in the male accessory gland of two closely**

2    **related species of field crickets.**

3

4

5    ANDRÉS, J.A.[1], LARSON, E.L.[2], BOGDANOWICZ, S.M.[2], AND HARRISON, R.G.[2]

6    [1]*Department of Biology, University of Saskatchewan, Saskatoon, S.K., Canada*

7    [2]*Department of Ecology and Evolutionary Biology, Cornell University, Ithaca, N.Y., USA*

8

9    **Running Title:** Transcriptome differentiation

10

11    **Key Words:** SNP identification, reproductive isolation, allopatric differentiation, *Gryllus*.

12

13    Corresponding author:

14    Jose A. Andrés

15    Department of Biology

16    University of Saskatchewan

17    Saskatoon, SK, S7N-1E5. Canada.

18    E-mail: jose.andres@usask.ca

19

**Abstract**

20

21    One of the central questions in evolutionary genetics is how much of the genome is

22    involved in the early stages of divergence between populations, causing them to be

23    reproductively isolated. In this paper we investigate genomic differentiation in a pair of

24    closely related field crickets (*Gryllus firmus* and *G. pennsylvanicus*). These two species

25    are the result of allopatric divergence and now interact along an extensive hybrid zone in

26    eastern North America. Genes encoding seminal fluid proteins (SFPs) are often divergent

27    between species and it has been hypothesized that these proteins may play a key role in

28    the origin and maintenance of reproductive isolation between diverging lineages. Hence,

29    we chose to scan the accessory gland transcriptome to enable direct comparisons of

30    differentiation for genes known to encode SFPs with differentiation at a much larger set

31    of genes expressed in the same tissue. We have characterized differences in allele

32    frequency between two populations for >6,000 SNPs and >26,000 contigs. About 10% of

33    all SNPs showed nearly fixed differences between the two species. Genes encoding SFPs

34    did not have significantly elevated numbers of fixed SNPs per contig, nor do they seem to

35    show larger differences than expected in their average allele frequencies. The distribution

36    of allele frequency differences across the transcriptome is distinctly bimodal, but the

37    relatively high proportion of fixed SNPs does not necessarily imply "ancient" divergence

38    between these two lineages. Further studies of linkage disequilibrium and introgression

39    across the hybrid zone are needed to direct our attention to those genome regions that are

40    important for reproductive isolation.

41

## Introduction

The study of speciation, defined as the origin of intrinsic barriers to gene exchange (Mayr 1942; Harrison 1998; Coyne and Orr 2004), relies on comparisons of phenotypes and genotypes among diverging populations, strains, subspecies or closely related species. In recently diverged taxa observed differences in genotypes or phenotypes are likely to be associated with the origin of reproductive barriers and less likely to be differences that have accumulated subsequent to initial divergence. As Templeton (1981) emphasized, our ultimate goal is to understand the genetics of speciation, not simply the genetics of species differences.

It is now widely recognized that the amount of divergence between populations or species will vary across the genome, due to selective and random lineage sorting from polymorphic ancestral populations and differential introgression when diverging taxa hybridize where their distributions overlap (Harrison 1991; Nosil *et al*. 2009; Turner *et al*. 2005; Wu 2001). Barton and Hewitt (1981) explicitly argued that gene exchange between hybridizing taxa will depend on genome region. Differential introgression has been widely discussed and documented in the hybrid zone literature (Harrison 1990; Payseur 2010; Rieseberg *et al*. 1999), where species boundaries have been described as semi-permeable. Chromosome regions that harbor genes that contribute to reproductive isolation or local adaptation will have reduced levels of gene flow.

The notion that genomes should be viewed as mosaics of different evolutionary histories also emerges from observations of discordance among individual gene genealogies for

65    closely related groups of species (Beltran *et al*. 2002; Carneiro *et al*. 2010; Dopman *et al*.

66    2005; Geraldes *et al*. 2008; Machado and Hey 2003; Putnam *et al*. 2007; White *et al*.

67    2009). Some loci reveal species to be reciprocally monophyletic or exclusive; at other

68    loci haplotypes do not sort by species and may even be shared across species. Genome

69    scans of allele frequencies for microsatellite loci, AFLPs, or SNPs also suggest

70    substantial heterogeneity in amounts of differentiation, with "$F_{ST}$ outliers" marking

71    regions that have become (or have remained) differentiated (Emelianov *et al*. 2004;

72    Grahame *et al*. 2006; Via and West 2008; Wood *et al*. 2008). These observations have

73    given rise to a diversity of terms that refer to the fact that genome divergence is

74    heterogeneous and that we can exploit this heterogeneity to identify gene regions that

75    either contribute to reproductive isolation and/or have experienced a recent selective

76    sweep. Increasingly sophisticated molecular tools and the ease with which we can

77    generate massive amounts of sequence data make it far easier to scan the genome (or

78    parts of the genome with reduced complexity) and search for regions that exhibit fixed

79    differences or major shifts in allele frequencies between recently diverged taxa.

80

81    An alternative to the genome scan approach is to identify candidate genes/proteins that

82    might account for phenotypic differences responsible for reproductive barriers.

83    Considerable attention has recently focused on the evolution of seminal fluid proteins

84    (SFPs) in a wide variety of taxa (Clark *et al*. 2006; Dorus *et al*. 2004; Walters and

85    Harrison 2010; Walters and Harrison 2011). In insects, male accessory glands are the site

86    of synthesis and secretion of SFPs that are transferred from male to female during

87    copulation (Gillott 2003; Wolfner 1997). Evolutionary genetic analyses have revealed

88  that, although most of these proteins are subject to selective constraints, many are rapidly

89  evolving, partly as the result of differential selection pressures (Andrés *et al*. 2006; Dean

90  *et al*. 2009; Dean *et al*. 2008; Ramm *et al*. 2009; Walters and Harrison 2010, 2011).

91  Although the functional and evolutionary consequences of this rapid divergence are not

92  fully understood, experimental work suggests that SFPs may play a key role in

93  reproductive isolation between diverging lineages (Andrés and Arnqvist 2001; Marshall

94  *et al*. 2011; Turner and Hoekstra 2008). Thus, *a priori*, we might expect genes encoding

95  SFPs to show elevated rates of molecular evolution and greater divergence between

96  closely related species.

97

98  Here we combine the candidate gene and genome scan approaches, using high throughput

99  sequencing to survey the male accessory gland transcriptomes of two closely related

100  species of field crickets that interact in a well-characterized hybrid zone in North

101  America. The two cricket species (*Gryllus firmus* and *G*. *pennsylvanicus*) are estimated to

102  have diverged about 200,000 years ago (Broughton and Harrison 2003; Maroja *et al*.

103  2009). Attempts to identify fixed differences between the species have met with only

104  limited success. Allozyme surveys and sequencing of mtDNA and nuclear gene introns

105  failed to identify the two species as exclusive groups (Harrison and Arnold 1982; Willett

106  *et al*. 1997; Broughton and Harrison 2003). However, analysis of anonymous nuclear

107  RFLPs did uncover four apparently diagnostic loci (Harrison and Bogdanowicz 1997).

108  These data suggest that much of the field cricket genome has remained undifferentiated

109  following the origin of reproductive barriers. In contrast, recent proteomic analysis of

110  spermatophore contents identified two SFP genes that exhibit nearly fixed differences

5

111 and strong evidence that positive selection has been responsible for patterns of

112 differentiation (Andres et al. 2008, Maroja et al. 2009).

113

114 In this paper, we used both Sanger and 454 sequencing to assemble and characterize the

115 transcriptome of the male cricket accessory gland. To detect SNPs, we then aligned

116 millions of pooled Illumina reads from allopatric populations of each species to the

117 Sanger/454 reference transcriptome. We characterized differences in allele frequency

118 between the two populations for >6,000 SNPs and >26,000 contigs and identified a

119 subset of highly differentiated SNPs and contigs showing strong allele frequency

120 differences. Using Sanger sequencing in a larger sample of crickets from the same

121 allopatric populations, we confirmed that a sample of divergent contigs identified from

122 Illumina reads indeed represents sequences that are highly divergent between the two

123 cricket populations. Finally, we compared the patterns of transcriptome differentiation for

124 SFP genes with genes expressed in the male accessory gland that are not SFPs.

125

126 **Materials and Methods**

127 *Cricket samples*

128 All crickets used in this study came from allopatric populations of the two species, *G.*

129 *firmus* from Guilford, CT and *G. pennsylvanicus* from Ithaca, NY. Guilford is close to the

130 hybrid zone and may show limited introgression of *G. pennsylvanicus* alleles. Ithaca is

131 more distant from the hybrid zone and is essentially "pure" *G. pennsylvanicus*.

132 Independent samples from these populations were used for constructing each of the

133 libraries described below (Sanger, 454, Illumina) and for subsequent SNP validation.

134

135 *Normalized Sanger library*

136 Accessory glands were dissected from 10 anesthetized (chilled) adult male *G. firmus*

137 from Guilford, CT. Total RNA was extracted in TRIZOL (Invitrogen). A single pooled

138 RNA sample was constructed using equimolar amounts of total RNA from each male.

139 First-strand cDNA was prepared using the Creator SMART cDNA Synthesis Kit

140 (Clontech). Briefly, complementary DNA was synthesized from the RNA pool, amplified

141 by 11-13 PCR-cycles using a 5' PCR primer (5'-AAGCAGTGGTATCAACGCAGAGT-

142 3') and normalized with a TRIMMER kit (AXXORA). Normalized cDNA was digested

143 with the restriction enzyme *SfiI*, and ligated to pDNR-Lib for directional cloning. Ligated

144 cDNA was used to transform Electromax DH5-α E cells (Invitrogen). Transformations

145 were spread onto Luria-Bertani (LB) plates containing 30 mg/mL chloramphenicol.

146 Colonies were randomly picked into 384-well plates containing 50 $\mu$L 0.5x AE buffer

147 (Qiagen) per well. Plates were heated at 96°C for five minutes, and one $\mu$L of supernatant

148 was used as template for PCR with Platinum *Taq* polymerase (Invitrogen) and M13

149 primers. PCR products were treated with Exonuclease I (New England Biolabs) and

150 Shrimp Alkaline Phosphatase (GE Healthcare) and sequenced with an M13 primer and

151 BigDye v 3.1 terminators (Applied Biosystems).

152

153 *Normalized Roche/454 library*

154 To further characterize the accessory gland transcriptome, total RNA from a single

155 Guilford *G. firmus* male accessory gland was extracted as described above. The

156 concentration and quality of the total RNA was determined using an Agilent Bioanalyzer

157    2100. One $\mu$g of total RNA was combined with 12 pmol SMART 3' oligo dT primer (5'-

158    AAGCAGTGGTATCAACGCAGAGTATTTTCTTTTTTCTTTTTTVN-3"), 12 pmol

159    SMART™ 5' rG primer (5-AAGCAGTGGTATCAACGCAGAGTACGCrGrGrG-3')

160    and RNase-free water in a 10 $\mu$L volume. The sample was heated at 65°C for five

161    minutes, placed on ice, and the following was added to the reaction: two $\mu$L 10x RT

162    buffer, four $\mu$L 25 mM $MgCl_2$, two $\mu$L 0.1 M DTT, one $\mu$L RNaseOUT™, and one $\mu$L

163    SuperScript reverse transcriptase (Invitrogen). The reaction was incubated at 45°C for

164    one hour, followed by 85°C for five minutes. The sample was then diluted 40-fold, and

165    one $\mu$L was PCR–amplified with Platinum *Taq* polymerase (Invitrogen) and a SMART™

166    IIA PCR primer (5'-AAGCAGTGGTATCAACGCAGAGTA-3') for 14 cycles. We used

167    the TRIMMER cDNA normalization kit to normalize the cDNA pool, following the

168    manufacturer's instructions.

169

170    Normalized cDNA was nebulized, and fragments were subjected to end-

171    repair/adenylation by incubating the cDNA with ATP, dNTPs, T4 polymerase, T4

172    Polynucleotide Kinase, and *Taq* polymerase (New England Biolabs). Samples were

173    incubated at room temperature for 20 minutes (end repair), and 72°C for 20 minutes

174    (adenylation of 3' ends by *Taq* polymerase). Normalized cDNA (500 ng) was ligated to

175    Titanium adapters A and B. Ligation products were diluted 100-fold and amplified with

176    the appropriate PCR primers (see File S1). PCR products were checked on agarose gels,

177    pooled, purified with a Qiaquick™ PCR purification kit (Qiagen), and submitted to the

178    Genomics Facility of the Life Sciences Core Laboratory Center at Cornell for

179    quantification, bead titration, and 454 sequencing.

180

*Illumina tag libraries*

Total RNA was extracted from the accessory gland of individual crickets as described

above, and quantified on a Nanodrop ND-1000 spectrophotometer. Equal amounts of

total RNA from each cricket were combined in two separate pools, representing 15 *G.*

*firmus* adult males from Guilford, CT and 15 *G. pennsylvanicus* adult males from Ithaca,

NY. First strand synthesis, PCR amplification, and normalization of cDNA for each pool

was carried out as described above for 454 sequencing. The two pools of normalized

cDNA were then submitted to the Genomics Facility at the Cornell Life Sciences Core

Laboratories Center for nebulization, end repair, and construction/sequencing of Illumina

paired-end fragments (2 x 86 bp). Each species pool of cDNA was run on a single

channel on a Solexa Genome Analyzer IIx.

*Transcriptome assembly and annotation*

Initial quality check of the 454 sequences was performed using Newbler (Margulies *et al*.

2005). Sanger and 454 reads were combined in a *de novo* assembly with NGen v2

(Lasergene 8.1.1, DNASTAR). Reads were end-trimmed (window length five nucleotides

with a minimum quality of 20) and scanned/trimmed for the plasmid pDNR-Lib and 454

adapters (mer length = 9, minimum matches = 3, trim length = 20, trim to end = 25). We

used assembly parameters that are similar to those reported in other transcriptome

assemblies (e.g. Renaut *et al*. 2010; Vera *et al*. 2008), including an estimate of 100x fixed

coverage, an estimated transcriptome length of 20 megabases, gap penalty = 25, match

size = 19, mismatch penalty = 25, and minimum match percent = 85. Both unassembled

203    (single-read) and multi-read contigs (*i.e.* transcripts) were saved to a SeqmanPro file. The

204    resulting fasta file was then used as a reference transcriptome for alignment of Illumina

205    reads generated from pools of *G. firmus* and *G. pennsylvanicus* adult male accessory

206    gland cDNAs. Functional annotation was performed using BLAST2GO using BLASTX

207    and the default parameters.

208

209    *SNP discovery*

210    We screened the accessory gland transcriptome for SNPs using NextGene v1.99

211    (Softgenetics) in a series of sequential steps. First, we converted raw Illumina tags to

212    fasta files (median score threshold ≥20, maximum number of uncalled bases ≤ 3, called

213    base number for each read ≥ 25, trim or reject read when ≥ 3 bases with score ≤ 16).

214    Second, the converted reads were sequentially trimmed by the following sequences:

215    SMART$^{TM}$ 5' rG primer, SMART$^{TM}$ PCR primer IIA, and the 3' half of SMART$^{TM}$ PCR

216    primer IIA. Third, the trimmed reads of the *G. firmus* and *G. pennsylvanicus* pools were

217    independently aligned to the reference transcriptome using one round of condensation

218    and one round of alignment (unambiguous mapping, matching requirement ≥ 12 bases, ≥

219    90% identity, mutation filter ≤ 5, SNP allele > 1 count, coverage > 20, forward/reverse

220    balance ≤ 0.05, and read library size range of 50-300 bases). Fourth, resulting alignments

221    were compared and screened for SNPs using the variant comparison tool. In our SNP

222    analyses we only included single-base substitutions. We excluded deletions/insertions

223    and multiple base substitutions.

224

225    *Transcriptome scan*

10

226    We defined interspecific SNPs as those homologous sites that show base frequency

227    differences between *G. firmus* and *G. pennsylvanicus*. Therefore, the ability to correctly

228    identify and quantify interspecific differences critically depends on the quality and

229    coverage of the SNPs. Here we used a highly stringent screen strategy aimed at finding

230    reliable interspecific SNPs. First, we only considered SNPs with a high base quality score

231    ($\geq$ 12) and a high total coverage ($\geq$ 40X; $\geq$ 20X in each species). If the coverage was <

232    100x, we only recognized a SNP if the rarer nucleotide variant was observed at least three

233    times. This allowed us to further reduce spurious SNP identification due to sequencing

234    errors. Otherwise, we considered SNPs with a minimum minor allele frequency (MAF) of

235    1%. This strategy allowed us to identify high quality SNPs for which accurate estimates

236    of allele frequency differences between the two species could be obtained. For each of

237    the identified SNPs we then defined the interspecific differentiation index (*D*) as:

238
$$D = \left| P_{Gf} - P_{Gp} \right|$$

239    That is, *D* is the absolute value of the relative frequency difference of alleles between the

240    two species (see Renaut *et al.* 2010). As opposed to other differentiation statistics (e.g.

241    $F_{ST}$) this estimator is not sensitive to unequal sample sizes (*i.e.* unequal coverage) in the

242    two species (Renaut *et al.* 2010). Moreover, $F_{ST}$ estimators depend on both within- and

243    between-population variation and thus the precise cause of $F_{ST}$ outliers can be difficult to

244    infer. Therefore, absolute allele frequency differences may be a better indicator of recent

245    selection (Strasburg *et al.* 2012). To identify candidate genes that may contribute to local

246    adaptation and reproductive isolation between the two species we screened for those

247    transcripts (*i.e.* contigs) that show the largest shifts in their average allele frequencies

248     between the two species. We calculated the total number of fixed SNPs for each contig,

249     as well as the mean interspecific differentiation index:

250     $$\overline{D} = \frac{1}{n}\sum_{i=1}^{n} D_i,$$

251     where $n$ is the total number of SNPs in a contig. Since $\overline{D}$ is not a very informative

252     estimator of interspecific divergence for contigs containing a small number of SNPs, our

253     analyses were limited to those contigs containing three or more SNPs. Furthermore,

254     because $\overline{D}$ depends on allele frequency differences (as does $F_{ST}$) and not on amount of

255     sequence divergence between alleles, we also estimated the number of fixed differences

256     per site and used this as a metric of sequence divergence.

257     The estimated SNP frequency differences are likely to depend on the set of

258     bioinformatics parameters we used to analyze our data; coverage and base quality scores

259     are of particular importance. Therefore, we did test for the robustness of our results by

260     varying parameter combinations and comparing the resulting data using permutation

261     analyses ($2 \times 10^3$ simulations) as implemented in R v2.11.1(R Development Core Team

262     2010). Specifically we tested if either increasing the quality base threshold ($\geq 15$ $vs. \geq 12$)

263     or reducing the coverage ($\geq 30X$ $vs. \geq 40X$) had a significant effect on our results.

264     For each contig, we also estimated the number of amino acid replacement SNPs per non-

265     synonymous site ($pN$) relative to the number of silent SNPs per synonymous site ($pS$).

266     This index is equivalent to $\omega$ ($dN/dS$) ratios and, therefore, provides insight into the

267     evolutionary forces driving molecular divergence between closely related lineages. We

268     first generated all possible open reading frames (ORFs, minimum length 200 nucleotides)

269     using Getorf (European Molecular Biology Open Software) and kept the longest ORF of

270     each contig as the most probable coding region of the gene. Then, we used a maximum

271    likelihood method to estimate *pN/pS* using PAML 4.2 (runmode= 0, CodonFreq= 2,

272    model= 2; Yang 2007). Putative mitochondrial and nuclear contigs were run separately

273    using icode=4 and 0 respectively. All analyses were carried out using R v2.11.1 and

274    dnds,R, a specific code kindly provided by Sébastien Renaut (see Renaut *et al*. 2010).

275    Mean *pN/pS* values were estimated by resampling (Bustamante et al. 2002) excluding

276    those contigs with *pS* = 0 and infinite *pN/pS*.

277    *Candidate gene approach*

278    In animals with internal fertilization, a subset of genes encoding SFPs are rapidly

279    evolving and often positively selected; they represent potential candidate barrier genes.

280    We might then expect elevated allele frequency differences in genes encoding seminal

281    fluid proteins. To test this hypothesis we identified high quality SNPs by aligning the

282    Illumina tags to a reference set of 70 previously described SFPs (Andrés *et al*. 2006;

283    Andrés *et al*. 2008) keeping the same parameters used in our transcriptome scan. Using

284    permutation analyses (R v2.11.1; R Development Core Team 2010), we first compared

285    $\overline{D}$ between SFPs and a subset of genes that, based on BLAST results, do not encode

286    SFPs. However, because SFPs are often rapidly evolving and their functions not

287    necessarily determined, it is possible that some fraction of the accessory gland contigs

288    represent SFPs even if they are not currently annotated. To minimize this potential bias

289    we extended our analysis to compare contigs with or without predicted signal peptides (as

290    a proxy for putative SFPs and non-SFPs respectively). The significance of all

291    permutation analyses was assessed using $2x10^3$ simulations.

292

293    *Intraspecific polymorphism*

294    To estimate the levels of intraspecific variation we screened the transcriptome of each

295    species for homologous sites exhibiting sequence variation. As above, we applied a

296    stringent SNP definition and only considered those variable sites with a quality score ≥

297    12 and MAF of 1%. Then, for each contig we estimated the average number of nucleotide

298    differences ($\pi$). Correlational analyses between polymorphism and divergence were

299    carried out using Spearman's $\rho$. Significance of permutation analyses was assessed as

300    above. All statistical analyses were performed using R v2.11.1 (R Development Core

301    Team 2010).

302

303    *SNP validation and gene genealogies*

304    To validate our transcriptome assembly and SNP identification methods we used Sanger

305    sequencing to characterize variation and divergence for a subset of contigs (n=10) that

306    showed at least three fixed SNPs, and high interspecific allelic divergence values ($\overline{D}$

307    ranging from 0.85-1). We used a new panel of 32 crickets, 16 each from our two focal

308    populations (*G. firmus*: Guilford, CT and *G. pennsylvanicus*: Ithaca, NY). Collectively,

309    the contigs we targeted contained a total of 60 putative SNPs. Selecting these highly

310    divergent contigs allowed us to validate putative regions of genomic differentiation

311    between the two crickets. In addition, we generated gene genealogies for this subset of

312    genes and compared them with the genealogies of two highly differentiated SFP genes

313    (*AG-0005F* and *AG-0334P*) that show almost exclusive relationships between the two

314    species (Andrés *et al.* 2008; Maroja *et al.* 2009).

315

316    Predicted SNPs for the 10 loci were validated through PCR amplification and Sanger

317    sequencing (primer sequences and conditions available upon request). Resulting

318    amplicons were sequenced on a 3130xl DNA analyzer (Applied Biosystems) using

319    BigDye v3.1 terminators. Haplotypes were reconstructed using PHASE (Stephens *et al*.

320    2001) implemented in ARLEQUIN v3.5 (Schneider *et al*. 2000). For each contig the

321    optimal substitution model was determined using hierarchical likelihood ratio test

322    searches implemented in JMODELTEST v0.1.1 (Posada 2008). Gene genealogies were

323    reconstructed using the neighbor-joining algorithm in PAUP v4.0 (Swofford 2003). We

324    calculated nodal support using 1,000 heuristic non-parametric bootstrap replicates.

325

326    **Results**

327    *Transcriptome assembly and annotation*

328    We have taken advantage of both Sanger sequencing and 454 sequencing technologies to

329    characterize the transcriptome of the male accessory glands in two *Gryllus* species. After

330    quality control, the resulting library contains $5.1 \times 10^5$ sequences, $9.2 \times 10^6$ nucleotides,

331    26,565 unique contigs (mean length 434 bp), and an average coverage of 4x. Although

332    our transcriptome assembly may contain information on alternatively spliced variants

333    (contig isoforms) we did not include this information in our assembly. Therefore, all

334    contigs represent the longest isoforms. Illumina reads mapped to a subset of ~15,000

335    contigs (average coverage: *G. firmus* = 40x, *G. pennsylvanicus* = 41x); half of them

336    (n=7,931) showing similarities with either functionally annotated genes or other insect

337    genomes and ESTs (TBLASTX, e≤ $10^{-5}$, Table S1). A significant fraction of these contigs

338    (~30%, 4,635/15,000) have a predicted signal peptide, and as expected, some of them (n=

15

339    126) represent 42 previously described *Gryllus* SFPs. Almost 60% of the annotated genes

340    (4,669/7,931) show strong similarities to other known genes and do not seem to be

341    components of the seminal fluid (TBLASTX, e$\leq 10^{-10}$).

342

343    *Frequency differences between G. firmus and G. pennsylvanicus*

344    To identify SNPs we mapped 7.6 million and 7.3 million high quality Illumina reads for

345    *G. firmus* and *G. pennsylvanicus* respectively onto the reference transcriptome. If a

346    substantial number of the predicted SNPs were the result of random sequencing and

347    assembly errors a transition:transversion ratio of 1:2 is expected. However, the observed

348    transition:transversion ratio for the our dataset is 1.55:1, suggesting that most of the SNPs

349    are not false positives. A total of 9,731 SNPs met our criteria for inferring allele

350    frequencies. The distribution of the allele frequency divergence values for these SNPs is

351    shown in Figure 1A. Many SNPs show low allelic divergence values ($D < 0.2$), but

352    11.6% (n = 1,133) of the inferred SNPs are highly differentiated ($D \geq 0.98$) between the

353    two species. Thus, the distribution of allele frequency differences is distinctly bimodal.

354    The distribution of the highly differentiated SNPs also varies among loci, with few

355    contigs showing an accumulation of differentiated sites (Figure 1B).

356    The mean allele frequency divergence value ($\bar{D}$) for the 1175 contigs that have $\geq 3$ SNPs

357    (see Methods) ranged from 0.041 to 1, and approximately 4% of these contigs showed an

358    average allele frequency difference ($\bar{D}$) of at least 0.98 (Figure 2, Table S2). Several of

359    these highly divergent contigs show significant similarities with currently annotated

360    genes (Table 1), including some genes encoded in the mitochondrial genome.

361

362    Increasing the SNP calling stringency from 92 to 95 had a drastic effect on SNP

363    discovery, reducing the total number of predicted SNPs (9,731 *vs*. 5,207, p < 0.0001) and

364    our estimates of SNPs/site (0.0063 *vs*. 0.0034, p< 0.0001) by almost half. However, this

365    only had a limited impact on the distribution of SNP frequency differences (File S2).

366    Similarly, decreasing the coverage stringency by 10X significantly increased the number

367    of predicted SNPs (p < 0.0001) but this difference also seems to have had little impact on

368    our divergence estimates.

369    *Variation in selective constraints across the accessory gland transcriptome*

370    Across the transcriptome, we found an average of 1.7 non-synonymous SNPs per 1000

371    non-synonymous sites SNPs and 5 synonymous SNPs per 1000 synonymous sites. The

372    estimates of *pN/pS* range from 0 (only synonymous polymorphic sites present) to infinite

373    (only non-synonymous polymorphic sites present). For the subset of contigs showing at

374    least one synonymous SNP, the overall *pN/pS* obtained by resampling contigs is 0.105

375    (95% confidence interval [CI]: 0.076-0.135), a value similar to the one obtained for the

376    subset of highly divergent genes (those with $\overline{D} \geq 0.98$) (permutation test p= 0.51).

377    However, the inferred proportion of highly divergent contigs showing only non-

378    synomymous variation is higher for the subset of highly divergent contigs than for the

379    rest of the transcriptome (permutation test p= 0.046).

380

381    *Candidate gene approach*

382    We have previously identified *Gryllus* SFPs by proteomic analysis (Andrés *et al*. 2006;

383    unpublished data). Therefore, we were able to compare average allele frequency

384    differences ( $\overline{D}$) between SFPs (*n*= 28) and $\overline{D}$ between "housekeeping" genes (*n*= 1,621)

385　that meet our functional annotation, quality and coverage criteria ($n$= 1,621, (see

386　methods). Although some of the SFP genes showed fixed differences between the two

387　species, on average, these genes do not seem to show larger shifts in their average allele

388　frequencies ( $\overline{D}_{\text{SFPs}} = 0.336$ , $\overline{D}_{\text{housekeeping}} = 0.416, P = 0.069$). Similar results were obtain

389　by comparing either SFP encoding genes to the subset of "housekeeping" genes with a

390　predicted signal peptide ($n$= 231, $\overline{D}_{\text{SFPs}} = 0.336$ , $\overline{D}_{\text{housekeeping\_signal}} = 0.402, P = 0.158$), or

391　annotated contigs that lack a predicted signal peptide ($n$= 1,130) with those that have it

392　($n$=491, $\overline{D}_{\text{signal}} = 0.381$ , $\overline{D}_{\text{no signal}} = 0.386, P = 0.754$). Table 2 summarizes values of $D$

393　and numbers of fixed SNPs for SFP genes that we have characterized. Only two of these

394　genes have $\overline{D} > 0.95$.

395

396　*Intraspecific polymorphism*

397　We observed 5,996 and 5,085 polymorphic nucleotides in *G. firmus* and *G.*

398　*pennsylvanicus* respectively. Polymorphism levels are similar between the two species

399　($\pi_{Gf}$= 1.38 10-5, $\pi_{Gp}$ = 1.29, $P = 0.761$). Indeed, estimates of $\pi$ for many of the contigs are

400　roughly equal in both species ($\rho$ =0.48, $P < 0.0001$). There is considerable variance in

401　polymorphism and divergence across contigs. Figure 3 shows for each species the

402　negative correlation between interspecific divergence estimated as $\overline{D}$ and and

403　intraspecific polymorphism (*Gf*: $\rho$ =- 0.23, $P < 0.0001$; *Gp*: $\rho$ =- 0.34, $P < 0.0001$).

404　Several contigs show high levels of divergence and low levels of intraspecific

405　polymorphism, a pattern consistent with recent selective sweeps.

406

407　*SNP validation and gene genealogies*

408    Thirty-two individual crickets (16 each from the two allopatric populations) were Sanger

409    sequenced for a subset of 10 highly differentiated contigs containing a total of 60 putative

410    SNPs (see Methods). Of these, 6 SNPs could not be typed because of primer design

411    constraints. All of the remaining predicted SNPs ($n= 54$) were validated by the

412    amplification and sequencing of PCR products of individual crickets. Allele frequency

413    estimates based on the pooled samples (Illumina tagging) were very similar to those

414    obtained from individual sequencing of the new panel of crickets (Table 3, File S3).

415    Accordingly, the gene genealogies (Figure 4) for these contigs show that observed

416    variation is indeed partitioned among populations (*i.e*. species). However, contigs 5368,

417    6023 and 1774 each have one haplotype shared between the two species (Figure 4).

418    Overall, NJ trees for the targeted contigs reveal greater differences than do similar trees

419    for *AG-0005F* and *AG-0334P*, the SFP genes that, from previous studies, were found to

420    be the most differentiated between the cricket species.

421

422    **Discussion**

423

424    *Genetic mosaics, transcriptome scans, and patterns of differentiation*

425    During the process of speciation, the degree of differentiation between diverging lineages

426    will vary across the genome, which is therefore a mosaic of different evolutionary

427    histories (Harrison 1991; Nosil *et al*. 2009; Rieseberg *et al*. 1999; Turner *et al*. 2005; Wu

428    2001). During the early stages of differentiation in allopatry, chromosome regions

429    harboring genes that contribute to local adaptation will diverge most rapidly.  In

430    secondary contact (or in cases of divergence with gene flow), chromosome regions that

19

431    contain genes that contribute to reproductive isolation will have reduced levels of gene

432    flow. These insights provide the foundation for a growing list of population genomics

433    studies that aim to identify genomic regions contributing to reproductive isolation (e.g.

434    Apple *et al*. 2010; Galindo *et al*. 2010; Manel *et al*. 2009; Nosil *et al*. 2008; Renaut *et al*.

435    2010; Schwarz *et al*. 2009; Michel *et al*. 2010; Fan *et al.* 2012, Nadeu *et al.* 2012).

436

437    Next generation sequencing has made it possible to effectively scan the genome for

438    specific genes (or gene regions) that exhibit low gene flow (i.e. fixed differences or major

439    shifts in allele frequencies) between recently diverged taxa. We chose to scan the

440    accessory gland transcriptome of hybridizing field crickets to enable direct comparisons

441    of differentiation for genes known to encode SFPs with differentiation at a much larger

442    set of genes expressed in the same tissue. It has been hypothesized that, in crickets, genes

443    encoding SFPs are likely to contribute to post-mating reproductive isolation between

444    closely related species, and we have previously shown (using comparisons of $d_N$ and $d_S$)

445    that some SFPs exhibit rapid evolution and evidence of positive selection (Andrés *et al*.

446    2006). Results of our transcriptome scan suggest that SNPs fixed between species

447    constitute about 10% of all identified SNPs, and that many contigs contain multiple fixed

448    SNPs. These contigs are as differentiated (or more differentiated) than are *AG-0005F* and

449    *AG-0334P*, the two highly divergent cricket SFP-encoding genes characterized

450    previously. In the process of scanning the transcriptome we have revealed evidence of

451    highly divergent SNPs between the two populations in four other SFP genes (some of

452    these with multiple fixed or nearly fixed SNPs; see Table 2). However, on average, SFPs

453    did not seem to have higher rates of divergence than other genes expressed in the

454  accessory gland, a result that may reflect the heterogeneity in evolutionary rates

455  previously observed in SFPs of field crickets (Andrés *et al*. 2006). Over a decade of

456  research on the evolution of SFPs has emphasized that a subset of SFP genes are among

457  the most rapidly evolving genes and that these divergent genes/proteins contribute to

458  reproductive isolation. However, it must be recognized that relatively few SFP genes

459  accumulate fixed differences and that a significant fraction of SFP genes show evidence

460  of evolutionary constraint (Andrés *et al*. 2006; Dean *et al*. 2009; Findlay *et al*. 2008;

461  Walters and Harrison 2011). Thus, an overall increase of evolutionary rate in SFPs genes

462  should not necessarily be expected.

463

464  Most of the divergent contigs have no identified homologues or known function, so it is

465  not yet possible to speculate about the ultimate causes or consequences of observed

466  divergence. However, the estimated proportion of loci showing *pN/pS* ratios consistent

467  with divergent (directional) selection is significantly higher for this subset of contigs than

468  for the rest of the transcriptome, supporting the hypothesis that many of highly divergent

469  loci are likely to be involved in local adaptation and perhaps in reproductive isolation.

470  Consistently, we found a negative association between intraspecific variation and

471  divergence between species, a pattern similar to that found in *Anopheles* mosquitoes,

472  where SNPs with $F_{ST} > 0.6$ have significantly reduced polymorphism (Neafsey *et al*.

473  2010). This negative correlation is a pattern that might be expected if the accessory gland

474  transcriptome differences between *G. firmus* and *G. pennsylvanicus* have mostly been

475  driven by directional selection. However, this correlation should be interpreted with

476  caution. The two species of field crickets exhibit large amounts of shared ancestral

477    polymorphism (Broughton and Harrison 2003). Therefore, the effects of variation in the

478    rate of recombination across the genome might explain the negative correlation, if there is

479    a reduction in $N_e$ in low recombining regions due to background selection.

480    Does the discovery of substantial differentiation between *G. firmus* and *G.*

481    *pennsylvanicus* imply that our previous assessment of "recent" divergence (estimated at

482    about 200,000 years) is wrong? The observed distribution of allele frequency differences

483    is distinctly bimodal, and many of the highly differentiated SNPs represent nearly fixed

484    differences between these two species (Figure 1). Unfortunately, it is difficult to compare

485    our results with those from other recent genome scans of strains, races, or closely related

486    species. Most of these studies report $F_{ST}$ values and identify "$F_{ST}$ outliers," but do not

487    provide information on fixed (or nearly fixed) SNPs. There are a few exceptions. Host

488    races of the budmoth show no markers completely fixed for alternative AFLP genotypes

489    (Emelianov *et al*. 2004). In contrast, genes with fixed amino acid substitutions between

490    forms occur "throughout the genome" in comparisons of the M and S forms of *Anopheles*

491    *gambiae* (Lawniczak et al. 2010). Because we have only sampled single populations of

492    the two cricket species, it is likely that a fraction of the highly divergent SNPs found in

493    our study represent frequency differences between populations rather than frequency

494    differences between species and that our divergence estimate is therefore elevated.

495    However, it is clear that the hybridizing field crickets are not as recently diverged as

496    many insect host races (e.g., budmoth, apple maggot, pea aphid), many of which have

497    been cited as exemplars of ecological speciation and/or sympatric speciation (Emelianov

498    *et al*. 2004; Michel *et al*. 2010; Via and West 2008). The observed pattern of

499    trancriptome divergence in crickets is reminiscent of the summary figures showing

500   divergence for allozyme loci in subspecies or semispecies in the *Drosophila willistoni*

501   group (Avise 1976; Ayala 1975; Ayala *et al*. 1974) and in *Lepomis* sunfish (Avise 1994).

502   The *D. willistoni* group and the genus *Lepomis* both provided early model systems for

503   studying genetic differentiation during the process of geographic speciation. Although

504   conspecific populations exhibited allele frequency differences at some loci, only in

505   recognized subspecies or semispecies was there a small proportion of loci with fixed or

506   nearly fixed differences. These loci, it was suggested, were those important for local

507   adaptation (Avise 1976; Ayala 1975). The proportion of loci with fixed differences

508   increased dramatically (to >30% of all loci) when sibling species were compared.

509   Allozyme studies reveal differences in the frequencies of charge-changing amino acid

510   substitutions, a presumably small subset of the differences that we can identify in

511   transcriptome scans. Consistent with this interpretation is the earlier observation that

512   there are no allozyme loci that exhibit fixed differences in allele frequency between *G.*

513   *firmus* and *G. pennsylvanicus*. Although by no means conclusive, these observations are

514   consistent with recent divergence of the two cricket species, at least relative to other

515   model systems for geographic speciation.

516

517   *Ascertainment Bias, Mapping Bias, and Sampling Error*

518   Our use of pooled DNA samples for SNP discovery and transcriptome wide scans of

519   allele frequencies could raise questions about ascertainment bias, mapping bias and

520   sampling error. Because of the relatively high sequencing error associated with high

521   throughput sequencing, SNP detection has focused on minimizing the false-positive rate

522   by considering only SNPs occurring more than a predefined number of times (e.g.

523    Galindo *et al.* 2010; Renaut *et al.* 2010), a SNP-calling criterion that generates a

524    systematic bias by excluding many rare alleles from the data. This, in turn, may lead to

525    biased estimates of several population genetic parameters, potentially compromising the

526    ability to identify outlier loci (see Helyar *et al.* 2011). Mapping bias can arise from the

527    assembly of tags from one lineage to a reference transcriptome from a different lineage.

528    This bias is likely to be more severe in highly differentiated regions of the genome and in

529    comparisons involving distantly related lineages. Sampling error in pooled samples has

530    two different sources. First, the number of individuals included in the pool and second,

531    the unequal representation of individual alleles. This second error source arises because

532    of variation in RNA amounts among individuals contributing to the pool, and because

533    some alleles are sequenced repeatedly whereas other alleles may not be sequenced at all.

534

535    In this paper we have attempted to minimize the concerns raised above. First, to reduce

536    ascertainment bias and sampling error we have identified putative SNPs using a relatively

537    large panel of alleles ($2n= 60$), and we have considered only those SNPs with high

538    coverage ($\geq 20$x). Two recent studies suggest that variation associated with heterogeneity

539    in the probe material (RNA) is not a serious problem and can be kept small by combining

540    relatively large pools ($2n >100$) with relatively deep (10-60x) sequence coverage

541    (Futschik and Schlotterer 2010; Galindo *et al.* 2010).  Our results strongly suggest that

542    relatively modest coverage (20x) and smaller pools still result in reliable identification of

543    SNPs. In fact, our validation experiment verified 90% of the predicted SNPs, a fraction

544    similar to results from other organisms without a reference genome (e.g. Williams *et al.*

545    2010; You *et al.* 2011). Likewise, we found a strong correspondence between the

546    predicted allele frequencies based on the pooled samples and those obtained from Sanger

547    sequencing of a different sample of crickets. This result is similar to those reported in

548    other SNP discovery experiments with comparable coverage (Van Tassell *et al*. 2008;

549    Wiedmann *et al*. 2008).

550

551    Second, instead of defining candidate loci by generating an expected neutral distribution

552    of differentiation values and identifying outlier loci (see Butlin *et al*. 2008), we have

553    defined candidate loci as those that show a high proportion of fixed (or almost) fixed

554    SNPs between species. This approach is similar to that of studies in which candidate

555    genes are defined as those that reveal closely related taxa to be reciprocally monophyletic

556    or exclusive groups (e.g. Andrés *et al*. 2008; Dopman *et al*. 2005). By using D and $\overline{D}$ we

557    avoid any potential biases associated with the estimation of "neutral" distributions.

558

559    However, it is also important to recognize that $\overline{D}$ (the average divergence across a

560    contig) may not be a reliable indicator of functional differences. Some contigs have

561    several fixed differences, but also many sites that are segregating within species-specific

562    allelic classes. In these cases, $\overline{D}$ can be low, but haplotypes in the two species may be

563    functionally distinct. Both *AG-0005F* and *AG-0334P* might fall into this category. In

564    addition, some fixed differences detected by traditional Sanger sequencing do not show

565    up as fixed SNPs in the Illumina reads, because the relevant sites fall below our

566    thresholds for coverage or sequence quality. Thus the numbers of fixed SNPs for *AG-*

567    *0005F* and *AG-0334P* reported in Table 3 are less than the numbers we know to be

568    present from earlier Sanger sequencing (Andrés *et al*. 2008). Moreover, the significance

569   of fixed SNPs is still uncertain because fixation may be a consequence of linkage to a

570   different causative locus. Follow-up studies and a detailed linkage map are therefore

571   critical to establish the possible link between functional divergence and elevated $\overline{D}$

572   values.

573

574   Finally, mapping bias does not seem to be important in our study. Although *G.*

575   *pennsylvanicus* shows lower levels of intraspecific polymorphisms as expected if there

576   was reduced ability to map *G. pennsylvanicus* sequences onto a *G. firmus* reference, this

577   result is also consistent with the $\theta$ values previously estimated using nuclear introns

578   (Broughton and Harrison 2003). Moreover, the total number of *G. firmus* reads mapped

579   onto the reference is only 4% higher than the number of mapped *G. pennsylvanicus* tags,

580   suggesting only a small bias, if any.

581

582   *The importance of fixed SNPs*

583   In the study of speciation, a focus on recently diverged taxa is important. This partly

584   explains the current attention devoted to recently diverged (still diverging) sympatric

585   populations or ecotypes, in which rapid adaptive divergence occurs in the face of gene

586   flow. In this paper we examine genomic divergence between a pair of species that are the

587   result of a more "conventional" model of allopatric divergence, a model that may

588   represent a majority of speciation events across all animal taxa. As discussed above, the

589   discovery of many fixed SNPs does not necessarily imply "ancient" divergence, and the

590   hybrid zone between *G. firmus* and *G. pennsylvanicus* remains an important model to

591   study the origins of reproductive isolation. Hybrid zones that result from allopatric

26

592    divergence and secondary contact (a majority of hybrid zone systems; see Barton and

593    Hewitt 1985) provide unique insights into the mechanistic and genetic basis of

594    reproductive isolation. These zones represent many generations of hybridization and

595    recombination between differentiated populations, and therefore patterns of introgression

596    across hybrid zones and patterns of linkage disequilibrium within hybrid zones direct our

597    attention to genome regions that are important for reproductive isolation or regions that

598    have recently experienced selection. The fixed SNPs we have discovered will allow

599    careful dissection of patterns of introgression and linkage disequilibrium within the field

600    cricket hybrid zone (see Gompert and Buerkle 2009; Payseur 2010; Teeter *et al*. 2008;

601    Teeter *et al*. 2010). This will bring us a step closer to our ultimate goal, to identify the

602    differences in genotypes or phenotypes that are more likely associated with the origin of

603    reproductive barriers and less likely to have accumulated subsequent to initial divergence.

604

**Table 1**. Annotation of most divergent contigs between *G. firmus* and *G. pennsylvanicus*

(*i.e.* those showing interspecific differentiation index ($\overline{D}$) greater than 0.98. Contigs in

bold correspond to mitochondrial loci. ns= non-significant (*E*-value $> 10^{-3}$) similarity.

| Contig | SNPs | | *pN/pS* | TBLASTX similarity |
| | Total | Per site | | |
|---|---|---|---|---|
| 70 | 10 | 0.0056 | ∞ | **Cytochrome b** |
| 310 | 9 | 0.0082 | 0.348 | **NADH dehydrogenase subunit 2 (*ND2*)** |
| 454 | 6 | 0.0073 | - | ***Teleogryllus emma* mitochondrion** |
| 755 | 6 | 0.0026 | 0.612 | ns |
| 618 | 5 | 0.0037 | ∞ | Conserved protein (similar to Cyclin-D1-binding protein 1) |
| 1341 | 5 | 0.0044 | 0.076 | Citrate lyase beta-like protein |
| 1699 | 5 | 0.0018 | - | Similar to *Tribolium castaneum* ADP ribosylation factor |
| 1774 | 5 | 0.0040 | 0 | ns |
| 1903 | 5 | 0.0026 | 0.090 | ns |
| 1978 | 5 | 0.0075 | ∞ | Similar to conserved hypothetical protein |
| 5368 | 5 | 0.0068 | 0.411 | ns |
| 1309 | 4 | 0.0023 | 0 | Similar to *Tribolium castaneum* B52 CG10851-PA |
| 1412 | 4 | 0.0041 | - | Insect conserved protein |
| 1721 | 4 | 0.0061 | 0 | Similar to *Gryllus bimaculatus* mRNA, GBcontig28218 |
| 5711 | 4 | 0.0053 | 0.878 | GalNAc transferase 6-like |
| 7164 | 4 | 0.0065 | - | ns |
| 14713 | 4 | 0.0092 | - | ns |
| 87 | 3 | 0.0021 | 0 | Similar to *Gryllus bimaculatus* mRNA, GBcontig31800 |
| 580 | 3 | 0.0037 | - | Similar to *Nasonia vitripennis* p15-2a protein |
| 937 | 3 | 0.0021 | ∞ | Dynactin subunit 4 (*Dctn4*) |
| 963 | 3 | 0.0038 | - | Similar to growth hormone-inducible soluble protein |
| 1101 | 3 | 0.0017 | 0.130 | ns |
| 1275 | 3 | 0.0021 | ∞ | Protease regulatory subunit S10B |
| 1306 | 3 | 0.0045 | - | Similar to translocase of outer membrane 7 |
| 1374 | 3 | 0.0033 | ∞ | Conserved protein: unknown |
| 1415 | 3 | 0.0023 | - | Myosin essential light chain |
| 1513 | 3 | 0.0026 | ∞ | UBX domain-containing protein |
| 1667 | 3 | 0.0033 | - | Similar to *Gryllus bimaculatus* mRNA, GBcontig12028 |
| 2182 | 3 | 0.0023 | - | Histone h2a |
| 2658 | 3 | 0.0022 | 0.562 | ns |

| | | | | |
|---|---|---|---|---|
| 3084 | 3 | 0.0019 | - | Similar to *Glossina morsitans* mRNA |
| 3432 | 3 | 0.0053 | ∞ | Similar to *DnaJ* (*Hsp40*) |
| 3566 | 3 | 0.0036 | - | ns |
| 3758 | 3 | 0.0026 | ∞ | **NADH dehydrogenase. Mitochondrial** |
| 3843 | 3 | 0.0017 | 0 | Translation initiation factor 4 gamma |
| 4655 | 3 | 0.0037 | 0.148 | Conserved protein: unknown |
| 5777 | 3 | 0.0015 | ∞ | Similar to transport and Golgi organization 1 (*Tango1*) |
| 6030 | 3 | 0.0035 | 0.147 | Ethanolaminephosphotransferase |
| 8373 | 3 | 0.0044 | ∞ | Asparagine synthetase |
| 9851 | 3 | 0.0050 | 0 | ns |
| 14741 | 3 | 0.0115 | 0.570 | Similar to eritrophin-like protein 1 |
| 6271 | 3 | 0.0114 | - | ns |
| 6026 | 3 | 0.0035 | - | Omega-amidase (*NIT2-B*) |
| 4450 | 3 | 0.0035 | - | Similar to *Gryllus bimaculatus* mRNA, GBcontig24459 |
| 861 | 4 | 0.0042 | 0.143 | Ribulose-5-phosphate-3-epimerase mRNA |

613

614

615 **Table 2**. Mean interspecific differentiation index ($\overline{D}$) for the subset of identified genes

616 encoding seminal fluid proteins. $N_t$ = Total number of SNPs in each gene. $N_{fix}$ = Number

617 of SNPs showing allele frequency differences (D) > 0.9 between *G. firmus* and *G.*

618 *pennsylvanicus*.

619

| SFP-Gene | Functional homology | $\overline{D}$ | $N_t$ | $N_{fix}$ |
|---|---|---|---|---|
| AG-0202F | Lectin similar | 0.965 | 2 | 2 |
| AG-0383F | Chaperonin | 0.772 | 4 | 3 |
| AG-0501F | Proteasome | 0.689 | 6 | 4 |
| AG-0509F | Proteasome | 0.528 | 1 | 0 |
| AG-0005F | Unknown | 0.447 | 23 | 1 |
| AG-0010F | Serine Protease | 0.360 | 4 | 1 |
| AG-0085F | Unknown | 0.277 | 4 | 0 |
| AG-0334P | Unknown | 0.264 | 16 | 3 |
| AG-0115F | Unknown | 0.237 | 48 | 0 |
| AG-0076F | Unknown | 0.214 | 28 | 0 |
| AG-0159F | Serine Protease | 0.210 | 32 | 0 |
| AG-0312F | Unknown | 0.203 | 12 | 0 |
| AG-0090F | Unknown | 0.199 | 14 | 0 |
| AG-0517F | Lectin similar | 0.188 | 35 | 0 |
| AG-0001F | Unknown | 0.177 | 9 | 0 |
| AG-0188F | Carboxipeptidase | 0.167 | 1 | 0 |
| AG-0254F | Chemiosensory protein | 0.164 | 3 | 0 |
| AG-0273F | Chymotrypsin | 0.159 | 12 | 0 |
| AG-0315F | Unknown | 0.157 | 40 | 0 |
| AG-0025F | Serine Protease | 0.153 | 5 | 0 |
| AG-0055F | Unknown | 0.151 | 13 | 0 |
| AG-0056F | Unknown | 0.137 | 1 | 0 |
| AG-0099F | Unknown | 0.119 | 47 | 0 |
| AG-0042F | Unknown | 0.112 | 8 | 0 |
| AG-0313F | Unknown | 0.107 | 16 | 0 |
| AG-0197P | Unknown | 0.100 | 37 | 0 |
| AG-0020F | Unknown | 0.097 | 2 | 0 |

620

621

622   **Table 3**. Comparison of the allele frequencies differences, estimated as $\overline{D}$, between *G.*

623   *firmus* and *G. pennsylvanicus* for a subset of highly differentiated contigs using pooled

624   Illumina tags and individual Sanger sequencing (see Methods). For each experiment we

625   independently sampled the same two allopatric populations (Guilford, CT and Ithaca, NY

626   respectively). $N_{ind}$ = Total number of individuals sequenced in each experiment. $N_{SNPs}$=

627   Number of SNPs typed in each contig.

628

| Contig | Illumina ($N_{ind}$=30) | | Sanger ($N_{ind}$=32) | | |
|---|---|---|---|---|---|
| | $N_{SNPs}$ | $\overline{D}$ | $N_{SNPs}$ | $\overline{D}$ | *pN/pS* |
| 5214 | 7 | 0.865 | 5 | 0.829 | ∞ |
| 5368 | 5 | 1 | 5 | 0.969 | 0.411 |
| 1002 | 9 | 0.855 | 6 | 0.911 | 0.141 |
| 6023 | 7 | 0.867 | 7 | 0.848 | - |
| 142 | 9 | 0.897 | 7 | 0.853 | 0.237 |
| 7153 | 9 | 0.899 | 8 | 0.823 | 0.096 |
| 14741 | 3 | 1 | 3 | 0.979 | 0.570 |
| 4655 | 3 | 1 | 3 | 0.990 | 0.148 |
| 1774 | 5 | 1 | 5 | 0.969 | 0 |
| 1231 | 5 | 0.816 | 4 | 0.969 | 0.917 |

629

630

631

632    **Figure 1**. (A) Frequency distribution of the interspecific differentiation index (D) for

633    each of the 6,761 predicted SNPs in *G. firmus* and *G. pennsylvanicus*. For any given SNP,

634    *D* represents allele frequency differences between the two species (see Methods). (B)

635    Frequency distribution of the number of highly differentiated SNPs (D≥ 0.98) per contig.

636

637    **Figure 2**. Ranked distribution of the mean interspecific differentiation index ($\overline{D}$)

638    between *G. firmus* and *G. pennsylvanicus* for each of the 1,157 contigs that showed high

639    coverage (≥ 20x) and at least 3 SNPs (see Methods). Dashed vertical lines represent the

640    standard error.

641

642    **Figure 3**. Correlation between polymorphism within species ($\pi$) and  divergence between

643    *G. firmus* and *G. pennsylvanicus*.

644

645    **Figure 4**. DNA gene genealogies for a subset of 10 highly differentiated contigs and two

646    seminal fluid protein genes (*AG-0005F* and *AG-0334P*). *Gryllus firmus* is represented by

647    white circles and *G. pennsylvanicus* by black circles. Size of symbols is proportional to

648    the frequency of the haplotype. Numbers on the branches represent bootstrap support

649    values over 75%.

650
651

652 **References**

653 Andrés, J. A., and G. Arnqvist, 2001 Genetic divergence of the seminal signal-receptor

654     system in houseflies: the footprints of sexually antagonistic coevolution?

655     Proceedings of the Royal Society of London, Series B: Biological Sciences **268:**

656     399-405.

657 Andrés, J. A., L. S. Maroja, S. M. Bogdanowicz, W. J. Swanson and R. G. Harrison,

658     2006 Molecular evolution of seminal proteins in field crickets. Molecular Biology

659     and Evolution **23:** 1574-1584.

660 Andrés, J. A., L. S. Maroja and R. G. Harrison, 2008 Searching for candidate speciation

661     genes using a proteomic approach: seminal proteins in field crickets. Proceedings

662     of the Royal Society of London, Series B: Biological Sciences **275:** 1975-1983.

663 Apple, J. L., T. Grace, A. Joern, P. S. Amand and S. M. Wisely, 2010 Comparative

664     genome scan detects host-related divergent selection in the grasshopper

665     *Hesperotettix viridis*. Molecular Ecology **19:** 4012-4028.

666 Avise, J., 1976 Genetic differentiation during speciation, pp. 106- 122 in *Molecular*

667     *Evolution*, edited by F. J. Ayala. Sinauer Associates, Sunderland, MA.

668 Avise, J., 1994 *Molecular Markers, Natural History, and Evolution*. Chapman and Hall,

669     New York, NY.

670 Ayala, F. J., 1975 Genetic differentiation during the speciation process, pp. 1-78 in

671     *Evolutionary Biology*, edited by T. Dobzhansky, M. K. Hecht and W. C. Steere.

672     Plenum Press, New York, NY.

673 Ayala, F. J., M. L. Tracey, D. Hedgecock and R. C. Richmond, 1974 Genetic

674     differentiation during speciation process in *Drosophila*. Evolution **28:** 576-592.

675 Barton, N. H., and G. M. Hewitt, 1981 A chromosomal cline in the grasshopper *Podisma*

676     *pedestris*. Evolution **35:** 1008-1018.

677 Barton, N. H., and G. M. Hewitt, 1985 Analysis of hybrid zones. Annual Review Of

678     Ecology And Systematics **16:** 113-148.

679 Bazykin, A. D., 1969 Hypothetical mechanism of speciaton. Evolution **23:** 685-687.

680 Beltran, M., C. D. Jiggins, V. Bull, M. Linares, J. Mallet *et al*., 2002 Phylogenetic

681     discordance at the species boundary: Comparative gene genealogies among

682       rapidly radiating *Heliconius* butterflies. Molecular Biology And Evolution **19:**

683       2176-2190.

684   Broughton, R. E., and R. G. Harrison, 2003 Nuclear gene genealogies reveal historical,

685       demographic and selective factors associated with speciation in field crickets.

686       Genetics **163:** 1389-1401.

687   Bustamante, C. D., R. Nielsen and D. L. Hartl, 2002. A maximum likelihood method for

688       analyzing pseudogene evolution: implications for silent site evolution in humans

689       and rodents. Molecular Biolology and Evolution **19:**110–7.

690   Butlin, R. K., J. Galindo and J. W. Grahame, 2008 Sympatric, parapatric or allopatric: the

691       most important way to classify speciation? Proceedings of the Royal Society of

692       London, Series B: Biological Sciences **363:** 2997-3007.

693   Carneiro, M., J. A. Blanco-Aguiar, R. Villafuerte, N. Ferrand and M. W. Nachman, 2010

694       Speciation in the European rabbit (Oryctogalus cuniculus):islands of

695       differentiation on the X chromsome and autosomes. Evolution **64:** 3443-3460.

696   Clark, N. L., J. E. Aagaard and W. J. Swanson, 2006 Evolution of reproductive proteins

697       from animals and plants. Reproduction **131:** 11-22.

698   Coyne, J. A., and H. A. Orr, 2004 *Speciation*. Sinauer Associates, Sunderland, MA.

699   Dean, M. D., N. L. Clark, G. D. Findlay, R. C. Karn, X. Yi *et al*., 2009 Proteomics and

700       comparative genomic investigations reveal heterogeneity in evolutionary rate of

701       male reproductive proteins in mice (*Mus domesticus*). Molecular Biology and

702       Evolution **26:** 1733-1743.

703   Dean, M. D., J. M. Good and M. W. Nachman, 2008 Adaptive evolution of proteins

704       secreted during sperm maturation: an analysis of the mouse epididymal

705       transcriptome. Molecular Biology and Evolution **25:** 383-392.

706   Dopman, E. B., L. Perez, S. M. Bogdanowicz and R. G. Harrison, 2005 Consequences of

707       reproductive barriers for genealogical discordance in the European corn borer.

708       Proceedings of the National Academy of Sciences, USA **102:** 14706-14711.

709   Dorus, S., P. D. Evans, G. J. Wyckoff, S. S. Choi and B. T. Lahn, 2004 Rate of molecular

710       evolution of the seminal protein gene SEMG2 correlates with levels of female

711       promiscuity. Nature Genetics **36:** 1326-1329.

712   Emelianov, I., F. Marec and J. Mallet, 2004 Genomic evidence for divergence with gene
713        flow in host races of the larch budmoth. Proceedings of the Royal Society of
714        London, Series B: Biological Sciences **271:** 97-105.
715   Fan, S., K. R. Elmer, and A. Meyer, 2012 Genomics of adaptation and speciation in
716        cichlid fishes: recent advances and analyses in African and Neotropical lineages.
717        Philosophical Transactions of the Royal Society **367:** 385-394.
718   Findlay, G., X. Yi, M. Maccoss and W. J. Swanson, 2008 Proteomics reveals novel
719        Drosophila seminal fluid proteins transferred at mating. PLoS Biology **6:** 1417-
720        1426.
721   Futschik, A., and C. Schlotterer, 2010 The next generation of molecular markers from
722        massively parallel sequencing of pooled DNA samples. Genetics **186:** 207-218.
723   Galindo, J., J. W. Grahame and R. K. Butlin, 2010 An EST-based genome scan using 454
724        sequencing in the marine snail *Littorina saxatilis*. Journal of Evolutionary
725        Biology **23:** 2004-2016.
726   Geraldes, A., P. Basset, B. Gibson, K. L. Smith, B. Harr *et al.*, 2008 Inferring the history
727        of speciation in house mice from autosomal, X-linked, Y-linked and
728        mitochondrial genes. Molecular Ecology **17:** 5349-5363.
729   Gillott, C., 2003 Male accessory gland secretions: modulators of female reproductive
730        physiology and behavior. Annual Review of Entomology **48:** 163-184.
731   Gompert, Z., and C. A. Buerkle, 2009 A powerful regression-based method for admixture
732        mapping of isolation across the genome of hybrids. Molecular Ecology **18:** 1207-
733        1224.
734   Grahame, J. W., C. S. Wilding and R. K. Butlin, 2006 Adaptation to a steep
735        environmental gradient and an associated barrier to gene exchange in *Littorina*
736        *saxatilis*. Evolution **60:** 268-278.
737   Harrison, R., and S. J. Arnold, 1982 A narrow hybrid zone between closely related
738        cricket species. Evolution **36:** 535-552.
739   Harrison, R. G., 1990 Hybrid zones: windows on evolutionary process, pp. 69-128 in
740        *Oxford Surveys in Evolutionary Biology*, edited by D. Futuyma and J. Antonovics.
741        Oxford University Press, New York, NY.

742    Harrison, R. G., 1991 Molecular changes at speciation. Annual Review of Ecology and

743          Systematics **22:** 281-308.

744    Harrison, R. G., 1998 Linking evolutionary pattern and process: the relevance of species

745          concepts for the study of speciation, pp. 19-31 in *Endless forms: species and*

746          *speciation*, edited by D. J. Howard and S. H. Berlocher. Oxford University Press,

747          New York, NY.

748    Harrison, R. G., and S. M. Bogdanowicz, 1997 Patterns of variation and linkage

749          disequilibrium in a field cricket hybrid zone. Evolution **51:** 493-505.

750    Helyar, S. J., J. Hemmer-Hansen, D. Bekkevold, M. I. Taylor, R. Ogden *et al*., 2011

751          Application of SNPs for population genetics of nonmodel organisms: new

752          opportunities and challenges. Molecular Ecology Resources **11:** 123-136.

753    Lawniczak, M.K. N., S. J.Emrich, A. K. Holloway, A. P. Regier, M. Olson, *et al*., 2010.

754          Widespread divergence between incipient *Anopheles gambiae* species revealed by

755          whole genome sequences. Science **330:** 512-514.

756    Machado, C. A., and J. Hey, 2003 The causes of phylogenetic conflict in a classic

757          *Drosophila* species group. Proceedings of the Royal Society of London, Series B:

758          Biological Sciences **270:** 1193-1202.

759    Manel, S., C. Conord and L. Despres, 2009 Genome scan to assess the respective role of

760          host-plant and environmental constraints on the adaptation of a widespread insect.

761          BMC Evolutionary Biology **9:** 288.

762    Margulies, M., M. Egholm, W. E. Altman, S. Attiya, J. S. Bader *et al*., 2005 Genome

763          sequencing in microfabricated high-density picolitre reactors. Nature **437:** 376-

764          380.

765    Maroja, L. S., J. A. Andrés and R. G. Harrison, 2009 Genealogical discordance and

766          patterns of introgression and selection across a cricket hybrid zone. Evolution **63:**

767          2999-3015.

768    Marshall, J. L., D. L. Huestis, C. Garcia, Y. Hiromasa, S. Wheeler *et al*., 2011

769          Comparative proteomics uncovers the signature of natural selection acting on the

770          ejaculate proteomes of two cricket species isolated by postmating, prezygotic

771          phenotypes. Molecular Biology and Evolution **28:** 423-435.

772    Mayr, E., 1942 Systematics and the Origins of Species. Columbia University Press, New
773        York.
774    Michel, A. P., S. Sim, T. H. Q. Powell, M. S. Taylor, P. Nosil *et al*., 2010 Widespread
775        genomic divergence during sympatric speciation. Proceedings of the National
776        Academy of Sciences, USA **107:** 9724-9729.
777    Nadeau, N. J., A. Whibley, R. T. Jones, J. W. Davey, K. K. Dasmahapatra, *et al*., 2012
778        Genomic islands of divergence in hybridizing *Heliconius* butterflies identified by
779        large-scale targeted sequencing. Philosophical Transactions of the Royal Society
780        **367:** 343-353.
781    Neafsey, D. E., M. K. N. Lawniczak, D. J. Park, S. N. Redmond, M. B. Coulibaly *et al*.,
782        2010 SNP genotyping defines complex gene-flow boundaries among African
783        malaria vector mosquitoes. Science **330:** 514-517.
784    Nosil, P., S. P. Egan and D. J. Funk, 2008 Heterogeneous genomic differentiation
785        between walking-stick ecotypes: "Isolation by adaptation" and multiple roles for
786        divergent selection. Evolution **62:** 316-336.
787    Nosil, P., D. J. Funk and D. Ortiz-Barrientos, 2009 Divergent selection and
788        heterogeneous genomic divergence. Molecular Ecology **18:** 375-402.
789    Payseur, B. A., 2010 Using differential introgression in hybrid zones to identify genomic
790        regions involved in speciation. Molecular Ecology Resources **10:** 806-820.
791    Posada, D., 2008 jModelTest: phylogenetic model averaging. Molecular Biology and
792        Evolution **25:** 1253-1256.
793    Putnam, A. S., J. M. Scriber and P. Andolfatto, 2007 Discordant divergence times among
794        Z-chromosome regions between two ecologically distinct swallowtail butterfly
795        species. Evolution **61:** 912-927.
796    Ramm, S. A., L. Mcdonald, J. L. Hurst, R. J. Beynon and P. Stockley, 2009 Comparative
797        proteomics reveals evidence for evolutionary diversification of rodent seminal
798        fluid and its functional significance in sperm competition. Molecular Biology and
799        Evolution **26:** 189-198.
800    Renaut, S., A. W. Nolte and L. Bernatchez, 2010 Mining transcriptome sequences
801        towards identifying adaptive single nucleotide polymorphisms in lake whitefish
802        species pairs (*Coregonus* spp. Salmonidae). Molecular Ecology **19:** 115-131.

803    Renaut, S. E. Maillet, E. Normandeau, C. Sauvage, N. Derome, *et al*., 2012 Genome-
804         wide patterns of divergence during speciation: the lake whitefish case study.
805         Philosophical Transactions of the Royal Society **367:** 354-363.

806    Rieseberg, L. H., J. Whitton and K. Gardner, 1999 Hybrid zones and the genetic
807         architecture of a barrier to gene flow between two sunflower species. Genetics
808         **152:** 713-727.

809    Schneider, S., D. Roessli and L. Excoffier, 2000 Arlequin, Version 2.0, a software for
810         population genetics data analysis, pp., University of Geneva.

811    Schwarz, D., H. M. Robertson, J. L. Feder, K. Varala, M. E. Hudson *et al*., 2009
812         Sympatric ecological speciation meets pyrosequencing: sampling the
813         transcriptome of the apple maggot *Rhagoletis pomonella*. BMC Genomics **10:**
814         633.

815    Stephens, M., N. J. Smith and P. Donnelly, 2001 A new statistical method for haplotype
816         reconstruction from population data. American Journal of Human Genetics **68:**
817         978-989.

818    Strasburg, J. L., N. A. Sherman, K. M. Wright, L. C. Moyle, J. H. Willis, et al., 2012
819         What can patterns of differentiation across plant genomes tell us about adaptation
820         and speciation? Philosophical. Transactions of the Royal Society **367:** 364-373.

821    Swofford, D. L., 2003 PAUP, pp. Sinauer, Sunderland, MA.

822    Teeter, K. C., B. A. Payseur, L. W. Harris, M. A. Bakewell, L. M. Thibodeau *et al*., 2008
823         Genome-wide patterns of gene flow across a house mouse hybrid zone. Genome
824         Research **18:** 67-76.

825    Teeter, K. C., L. M. Thibodeau, Z. Gompert, C. A. Buerkle, M. W. Nachman *et al*., 2010
826         The variable genomic architecture of isolation between hybridizing species of
827         house mice. Evolution **64:** 472-485.

828    Templeton, A. R., 1981 Mechanisms of speciation: a population genetic approach.
829         Annual Review of Ecology and Systematics **12:** 23-48.

830    Turner, L. M., and H. E. Hoekstra, 2008 Causes and consequences of the evolution of
831         reproductive proteins. International Journal of Developmental Biology **52:** 769-
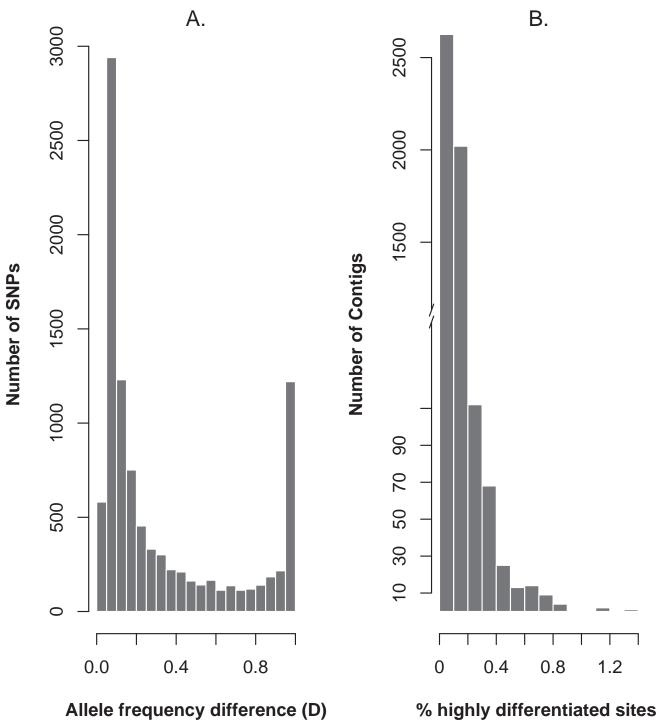832         780.

833 Turner, T. L., M. W. Hahn and S. V. Nuzhdin, 2005 Genomic islands of speciation in
834   *Anopheles gambiae*. PLoS Biology **3:** 1572-1578.

835 Van Tassell, C. P., T. P. L. Smith, L. K. Matukumalli, J. F. Taylor, R. D. Schnabel *et al*.,
836   2008 SNP discovery and allele frequency estimation by deep sequencing of
837   reduced representation libraries. Nature Methods **5:** 247-252.

838 Vera, J. C., C. W. Wheat, H. W. Fescemyer, M. J. Frilander, D. L. Crawford *et al*., 2008
839   Rapid transcriptome characterization for a nonmodel organism using 454
840   pyrosequencing. Molecular Ecology **17:** 1636-1647.

841 Via, S., and J. West, 2008 The genetic mosaic suggests a new role for hitchhiking in
842   ecological speciation. Molecular Ecology **17:** 4334-4345.

843 Walters, J. R., and R. G. Harrison, 2010 Combined EST and proteomic analysis identifies
844   rapidly evolving seminal fluid proteins in *Heliconius* butterflies. Molecular
845   Biology and Evolution **27:** 2000-2013.

846 Walters, J. R., and R. G. Harrison, 2011 Decoupling of rapid and adaptive evolution
847   among seminal fluid proteins in *Heliconius* butterflies with divergent mating
848   systems. Evolution **65:** 2855-2871.

849 White, M. A., C. Ané, C. N. Dewey, B. R. Larget and B. A. Payseur, 2009 Fine-scale
850   phylogenetic discordance across the house mouse genome. PLoS Genetics **5:**
851   e1000729.

852 Wiedmann, R. T., T. P. L. Smith and D. J. Nonneman, 2008 SNP discovery in swine by
853   reduced representation and high throughput pyrosequencing. BMC Genetics **9:** 81.

854 Willett, C. S., M. J. Ford and R. G. Harrison, 1997 Inferences about the origin of a field
855   cricket hybrid zone from a mitochondrial DNA phylogeny. Heredity **79:** 484-494.

856 Williams, L. M., X. Ma, A. R. Boyko, C. D. Bustamante and M. F. Oleksiak, 2010 SNP
857   identification, verification, and utility for population genetics in a non-model
858   genus. BMC Genetics **11:** 32.

859 Wolfner, M. F., 1997 Tokens of love: functions and regulation of *Drosophila* male
860   accessory gland products. Insect Biochemistry and Molecular Biology **27:** 179-
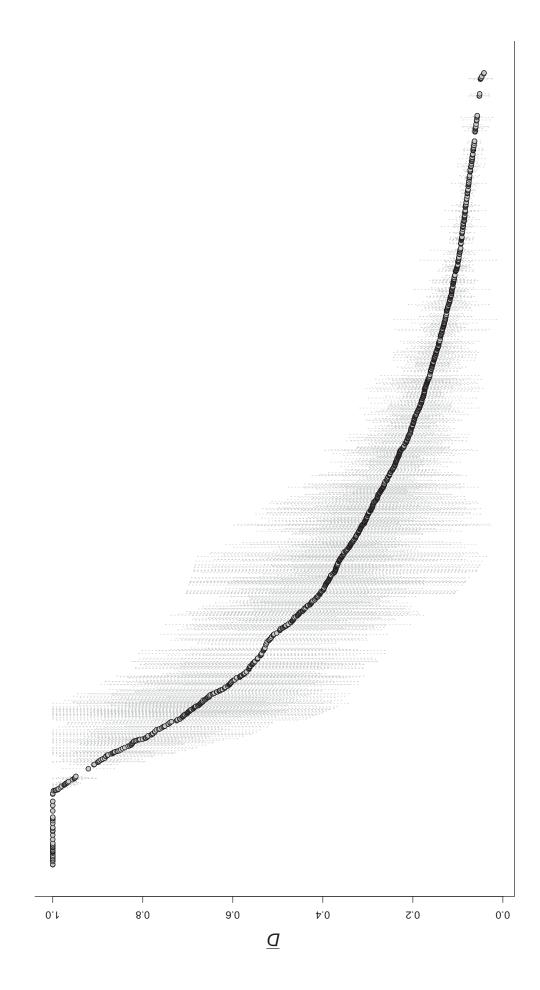861   192.

862    Wood, H. M., J. W. Grahame, S. Humphray, J. Rogers and R. K. Butlin, 2008 Sequence
863        differentiation in regions identified by a genome scan for local adaptation.
864        Molecular Ecology **17:** 3123-3135.

865    Wu, C. I., 2001 The genic view of the process of speciation. Journal of Evolutionary
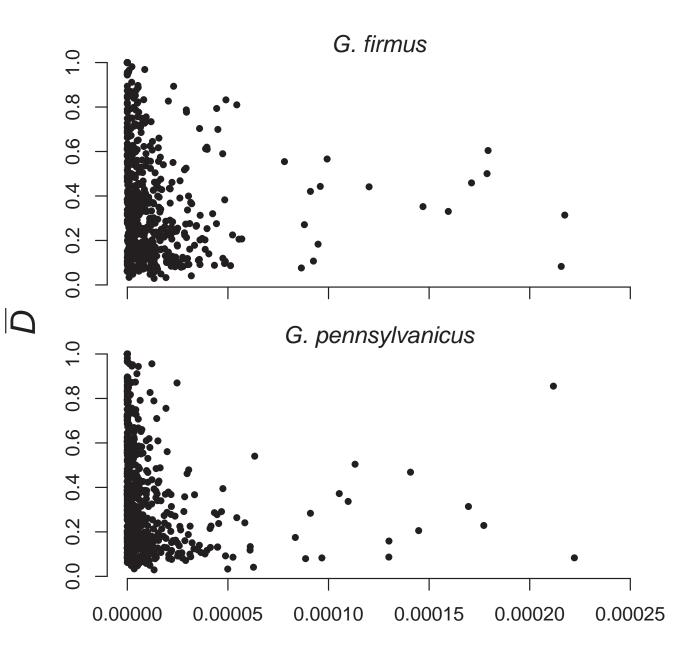866        Biology **14:** 851-865.

867    Yang, Z. H., 2007 Paml 4: a program package for phylogenetic analysis by maximum
868        likelihood. Molecular Biology and Evolution **24:** 1586–1591.

869    You, F. M., N. X. Huo, K. R. Deal, Y. Q. Gu, M. C. Luo *et al*., 2011 Annotation-based
870        genome-wide SNP discovery in the large and complex *Aegilops tauschii* genome
871        using next-generation sequencing without a reference genome sequence. BMC
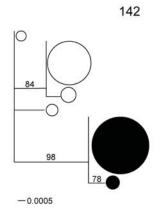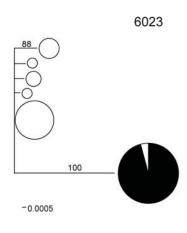872        Genomics **12:** 59.

873

874

860     Wood, H. M., J. W. Grahame, S. Humphray, J. Rogers and R. K. Butlin, 2008 Sequence
861          differentiation in regions identified by a genome scan for local adaptation.
862          Molecular Ecology **17:** 3123-3135.

863     Wu, C. I., 2001 The genic view of the process of speciation. Journal of Evolutionary
864          Biology **14:** 851-865.

865     Yang, Z. H., 2007 Paml 4: a program package for phylogenetic analysis by maximum
866          likelihood. Molecular Biology and Evolution **24:** 1586–1591.

867     You, F. M., N. X. Huo, K. R. Deal, Y. Q. Gu, M. C. Luo *et al.*, 2011 Annotation-based
868          genome-wide SNP discovery in the large and complex *Aegilops tauschii* genome
869          using next-generation sequencing without a reference genome sequence. BMC
870          Genomics **12:** 59.

871

872

41

A.

B.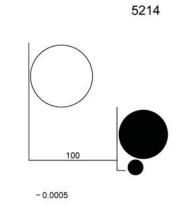