*Mouse single nucleotide polymorphic targets for cross hybridization in rodents*

1
2
3
4  **Utility of a high-resolution mouse single nucleotide polymorphism microarray assessed for**
5  **rodent comparative genomics**
6
7
8
9  **Short title***: Mouse single nucleotide polymorphic targets for cross hybridization in rodents*
10
11
12
13  Rachel D. Kelly, Maja Milojevic, Freda Qi, Kathleen A. Hill*
14
15
16
17  [1]Department of Biology, *The* University *of* Western Ontario, London, Ontario, Canada
18
19
20  * Corresponding author
21
22  Email: khill22@uwo.ca (KH)
23
24
25
26
27

1

## 28    Abstract

29    In the study of genetic diversity in non-model species there is a notable lack of the low-cost, high

30    resolution tools that are readily available for model organisms. Genotyping microarray

31    technology for model organisms is well-developed, affordable, and potentially adaptable for

32    cross-species hybridization. The Mouse Diversity Genotyping Array (MDGA), a single

33    nucleotide polymorphism (SNP) genotyping tool designed for *Mus musculus*, was tested as a tool

34    to survey genomic diversity of wild species for inter-order, inter-genus, and intra-genus

35    comparisons. Application of the MDGA cross-species provides genetic distance information that

36    reflects known taxonomic relationships reported previously between non-model species, but

37    there is an underestimation of genetic diversity for non-Mus samples, indicated by a plateau in

38    loci genotyped beginning 10-15 millions of years divergence from the house mouse. The number

39    and types of samples included in datasets genotyped together must be considered in cross-species

40    hybridization studies. The number of loci with heterozygous genotypes mapped to published

41    genome sequences indicates potential for cross-species MDGA utility. A case study of seven

42    deer mice yielded 159,797 loci (32% of loci queried by the MDGA) that were genotyped in these

43    rodents. For one species, *Peromyscus maniculatus*, 6,075 potential polymorphic loci were

44    identified. Cross-species utility of the MDGA provides needed genetic information for non-

45    model species that are lacking genomic resources. Genotyping arrays are widely available,

46    developed tools that are capable of capturing large amounts of genetic information in a single

47    application, and represent a unique opportunity to identify genomic variation in closely related

48    species that currently have a paucity of genomic information available. A candidate list of

49    MDGA loci that can be utilized in cross-species hybridization studies was identified and may

50    prove to be informative for rodent species that are known as environmental sentinels. Future

51    studies may evaluate the utility of candidate SNP loci in populations of non-model rodents.

52

53

## 54    Author Summary

55    There is a need for a tool that can assay DNA sequence differences in species for which there is

56    little or no DNA information available. One method of analyzing differences in DNA sequences

57    in species with well-understood genomes is through a genotyping microarray, which has

58    demonstrated utility cross-species. The Mouse Diversity Genotyping Array (MDGA) is a tool

59    designed to examine known differences across the genome of the house mouse, *Mus musculus*.

60    Given that related organisms share genetic similarity, the MDGA was tested for utility in

61    identifying genome variation in other wild mice and rodents. Variation identified from distantly

62    related species that were not of the same genus as the house mouse was an underestimate of the

63    true amount of variation present in the genomes of wild species. Utility of the MDGA for wild

64    species is best suited to mice from the same genus as the house mouse, and candidate variation

65    identified can be tested in rodent populations in future studies. Identifying changes in genetic

66    variation within populations of wild rodents can help researchers understand the links between

67    specific genome changes and the ability to adapt to pressures in the environment, as well as

68    better understand the evolution of rodents.

69 **Introduction**

70 The study and characterization of genomic diversity of non-model organisms is complicated by

71 limitations in knowledge and genomic resources available [1]. By contrast, researchers studying

72 model organisms benefit from the advantage of working with species that have sequenced and

73 annotated genomes, and high throughput platforms to survey genetic diversity at low cost. There

74 is a lack of genomic sequence information available for non-model species, and a deficit of tools

75 to assay genomic diversity in understudied organisms [2–4]. There is a need for custom tools to

76 survey genomic diversity in non-model organisms, but the creation of these tools can be time

77 consuming and expensive. There is an opportunity to explore existing technologies designed for

78 model organisms and test the applicability of these tools in non-model species.

79

80 Genotyping arrays are convenient tools that obtain large amounts of genetic diversity

81 information in a single assay at low cost [5]. Genotyping arrays are designed to capture a large

82 swath of diversity within a species, but the technology is typically tailored to the model species

83 of interest. Hybridization of microarray oligos targeted to unique locations in test DNA of the

84 organism of interest provides a picture of the genomic landscape of that sample [6].  Single

85 nucleotide polymorphisms (SNPs) are single base pair genome variations found in at least one

86 percent of individuals in a population, and are an informative type of genomic diversity that is

87 captured by genotyping arrays [6,7]. SNPs are found in abundance throughout the genome, and

88 this variation can be used as a metric of genomic diversity when comparing different individuals

89 in a population, or different species of interest [8].

90

91    There is a precedent for exploring the possibility of applying existing genotyping array

92    technologies to related, non-model species. The majority of research examining the applicability

93    of existing mammalian genotyping arrays in cross-species analyses focus on applying array

94    technologies designed for agricultural and domestic breeding purposes to related species [2–4,9–

95    14]. Researchers using a bovine genotyping array were able to identify a panel of over 100

96    candidate SNPs conserved within two species of wild oryx, despite a 23 million year divergence

97    time between oryx and modern cows [2]. Other researchers have applied domestic arrays to non-

98    model organisms that diverged from the model species millions of years ago to identify SNPs

99    associated with an ideal physical trait that would inform breeding strategies [4], or to identify

100    sexually selected traits that are associated with the fitness of a non-model organism [11].

101

102    Looking at the research performed in the field of cross-species genotyping array use, we identify

103    three metrics of success for the application of existing genotyping arrays to non-model species.

104    The first metric of success for applicability of genotyping arrays cross-species is the

105    identification of a panel of candidate SNPs that may be conserved between the model and non-

106    model organisms. This panel of SNPs represents variation that can be successfully genotyped in

107    the non-model organism of interest. While one metric of success for genotyping array use is the

108    number of loci or positions in the genome that can be accurately genotyped, the ability to detect

109    heterozygous loci is the second metric. Heterozygous loci, or positions in the genome in which

110    both the major and minor alleles in a population can be genotyped, are key when surveying

111    diversity in populations [2,3,15,16]. The third and final metric of success we have identified is

112    the ability to validate the candidate panel of SNPs and heterozygous loci either through *in silico*

113    methods for non-model species with some sequence information available, or by testing for the

114    candidate SNPs in populations using alternative experimental methods.

115

116    Genotyping arrays have demonstrated utility in identifying polymorphic SNPs, or sites of

117    variability within non-model organisms, which is an important goal for conservation studies of

118    endangered species, and molecular ecology [2,3,17]. In one particular study, researchers

119    Hoffman et al. (2013) applied a Canine HD Beadchip genotyping array to a population of

120    Antarctic fur seals, despite a 44 million year divergence time between the species of seal and

121    dogs [3]. Using the Canine HD Beadchip which queries over 173,000 SNP loci in dogs, the

122    researchers were able to identify a panel of 173 polymorphic SNP loci that were conserved

123    between the Antarctic fur seals and dogs [3]. A subset of the loci genotyped were validated *in*

124    *silico* using available transcriptomic data. Gene ontology analysis of shared loci between dogs

125    and seals showed that the panel of loci were involved in energy metabolism, suggesting the

126    genomic markers conserved between dogs and seals were a part of a highly conserved functional

127    pathway.

128

129    The identification of SNPs in non-model species can be used as markers of rapid evolution

130    between populations [18], and a genotyping array would allow researchers to identify large lists

131    of candidate SNP loci in a single application. The characterization of SNPs across the genomes

132    of wild organisms is of keen interest to population geneticists as molecular markers for

133    comparative studies [19]. Cross-species genotyping can provide information regarding variants

134    that are involved in sexual selection [11], and variants tied to a phenotype of interest, which can

135    inform breeding strategies [4]. A study by More et al. (2019) tested the utility of a Bovine SNP50

136    array on the alpaca species *Vicugno pacos* which is bred domestically for its hair fiber that is

137    economically valued [4]. The cross-species application of the Bovine SNP50 array allowed

138    researchers to identify a panel of 400 polymorphic SNPs in the alpaca, and they were able to map

139    209 SNPs to alpaca gene sequence information that was available [4]. This study helped identify

140    a number of SNP markers with utility cross-species that is currently needed to help guide

141    breeding practices for the alpaca species that will maximize high-quality hair fiber yield in the

142    future. This study also highlights the need for the development of genomic tools capable of

143    genotyping non-model species of interest.

144

145    There are a number of different genotyping arrays that have been applied cross-species to non-

146    model organisms, but there has been little research focusing on cross-species genotyping in mice

147    and other rodents. Mice are a peculiarity in that most genomic tools are designed for classical,

148    inbred mice used in research, but mice and related rodents can be found all across the world.

149    There is a need for a tool that can survey diversity in non-model mice and other rodents.

150    Wild rodents represent unique research opportunities because of the unique selective pressures

151    that are placed on them through human influence, and their ability to rapidly adapt to changing

152    human environments [18,20]. For instance, deer mice from the genus Peromyscus make

153    interesting candidates for non-model research as they can be found across North America and

154    despite lacking fully sequenced and annotated reference genomes, they have been previously

155    used as sentinels of environmental contaminants [21], and are becoming key organisms for

156    evolutionary studies and molecular genetics [18,22,23]. While some genetic resources are

157    available for deer mice and other rodents of interest, there remains a paucity of genomic

158    information for these understudied species and few low-cost tools to assess genomic variation in

159    a high-throughput manner.

160

161    The Mouse Diversity Genotyping Array (MDGA) is a tool designed to survey hundreds of

162    thousands of SNP loci across the genome of the house mouse and was specifically created to

163    maximize the amount of SNP diversity that can be identified within laboratory mouse strains and

164    crosses [24]. After testing and the removal of poorly performing SNP probes, the MDGA was

165    found to genotype 493,290 SNP loci within the genome of the house mouse [25]. The aim of our

166    study is to explore the use of the MDGA for its utility as a cross-species genotyping tool. The

167    MDGA was tested on 44 samples ranging in relatedness to the model house mouse, *Mus*

168    *musculus*, that span different Genus, Family, and Orders of taxonomic classification (Table 1, S1

169    Table). The goal was to identify the three metrics of success that define MDGA cross-species

170    utility in related organisms. This study represents an advance in the field of mammalian cross-

171    species genotyping that will add to the paucity of genomic sequence and SNP information

172    available for non-model mice and rodents (Fig 1). It was hypothesized that application of the

173    MDGA to wild rodent DNA samples will help elucidate potential polymorphic loci, or the

174    number of loci that can detect both the A and B allele in a population, and that can be used cross-

175    species in non-model organisms.

176

177

178

179

180

181 **Table 1** Genotyping sets of study
182
183
184
185
186
187

| Genotyping Test Sets | | | Common Name | Scientific Name |
|---|---|---|---|---|
| | | | House Mouse | *Mus musculus* |
| | | | South-Eastern House Mouse | *Mus musculus castaneus* |
| | | | Earth-Colored Mouse | *Mus dunni/Mus terricolor* |
| | | | Servant Mouse/Bonhote's Mouse | *Mus famulus* |
| | | | Sheath-Tailed Mouse | *Mus fragilicauda* |
| | | | Ryukyu Mouse | *Mus caroli* |
| | | | Fawn-Colored Mouse | *Mus cervicolor* |
| Inter-Order | Inter-Genus | Intra-Genus | Cook's Mouse | *Mus cookii* |
| | | | Flat-Haired Mouse | *Mus platythrix* |
| | | | Rock-Loving Mouse | *Mus saxicola* |
| | | | Gairdner's Shrewmouse | *Mus pahari* |
| | | | African Pygmy Mouse | *Mus (nannomys) minutoides* |
| | | | Orange Pygmy Mouse | *Mus (nannomys) orangiae* |
| | | | Matthey's Mouse | *Mus (nannomys) mattheyi* |
| | | | Wood Mouse | *Apodemus sylvaticus* |
| | | | Sprague Dawley Rat | *Rattus norvegicus* |
| | | | Wistar Rat | *Rattus norvegicus* |
| | | | Aztec Mouse | *Peromyscus aztecus* |
| | | | California Mouse | *Peromyscus californicus* |
| | | | North American Deer Mouse | *Peromyscus maniculatus* |
| | | | Sonoran Deer Mouse | *Peromyscus maniculatus* |
| | | | Plateau Deer Mouse | *Peromyscus melanophrys* |
| | | | Oldfield Mouse/Beach Mouse | *Peromyscus polionotus* |
| | | | White-Footed Mouse | *Peromyscus leucopus* |
| | | | Squirrel | Sciuridae |
| | | | Naked Mole Rat | *Heterocephalus glaber* |
| | | | African Black Rhino | *Diceros bicornis* |
| | | | Mountain Tapir/Wooly Tapir | *Tapirus pinchaque* |

Genotyping sets organized in descending order according to bounds of taxonomic classification and differences in maximum genetic divergence of a test set from the reference C57BL/6J (*Mus musculus*) organism

188 **RESULTS**

189 **Cross-species test sets exceed maximum genetic diversity of the training set**

190 A training set of DNA samples from 114 classical, inbred laboratory mice was used in training

191 the genotyping algorithm employed by Affymetrix Power Tools to provide accurate genotypes

192 (S2 Table). Genetic distances reflect the relatedness between samples and were obtained from

193 calculations of SNP distances derived from raw genotyping results. The maximum genetic

194 distance of the training set is approximately 0.225 with respect to the reference C57BL/6J house

195 mouse (Fig 2). The intra-genus test set of 27 species from the genus Mus has a maximum genetic

196 distance value of 0.836 and is over three times larger than the maximum genetic distance of the

197 reference set of 114 classical inbred mice (Fig 2). A case study of seven Peromyscus samples

198 genotyped together has a maximum genetic distance of 0.941 from the house mouse, and far

199 exceeds the diversity of the training set. Also, the maximum genetic distance of the inter-order

200 test set (n=44, 96 MYD) is 0.938, and is over four times larger than the maximum genetic

201 diversity represented in the training set (Fig 2). The training set used does not encapsulate the

202 genetic diversity of the test sets.

203

204 The samples of the inter-order test set are significantly different in genotypic composition and

205 allelic frequency (P<0.0001; Fisher's exact test with Monte Carlo simulation). The samples of

206 the intra-genus test set of only Mus samples are also significantly different in genotypic and

207 allelic frequency (P<0.0001). Two *R. norvegicus* samples were compared to one another as a

208 control and the genotypic composition is not significantly different (p=0.0934). Differences in

209 allelic composition between *R. norvegicus* samples are also not significant (p = 0.2232). The four

11

210   *H. glaber* (naked mole rat) samples genotyped together are significantly different in the genotype

211   composition (p<0.0001), but not allelic composition (p=0.0038).

212

213   **Underestimation of genetic diversity occurs when genotyping across multiple genera**

214   For the inter-order genotyping set (n=44), a general decrease is observed in the percentage of loci

215   genotyped as divergence time increases from *M. musculus* (r = -0.57; p-value<0.0001; Fig 3A).

216   As divergence time increases from *M. musculus*, the number of 'no calls', or inability to

217   determine a genotype at a locus, increases. A plateau in the percentage of loci genotyped is

218   observed between 10-15 MYD for non-Mus samples from the inter-genus test set. Loci with

219   heterozygous genotypes were of particular interest, as those loci have the potential to identify

220   both the major and minor alleles in a population (polymorphic loci). The percentage of loci that

221   had a heterozygous genotype increases as divergence time from the house mouse increases (Fig

222   3B). There is a positive correlation between increasing percent heterozygosity and the known

223   divergence times from the house mouse (r = 0.67; p-value<0.0001). Similar to the percentage of

224   loci genotyped, a plateau in percent heterozygosity is also observed to begin between 10-15

225   million years divergence from *M. musculus* (Fig 3B).

226

227   **MDGA captures the genetic diversity of wild samples from the genus Mus**

228   As seen in the inter-order test set, there is a general decrease in the percentage of loci that were

229   genotyped in samples of the intra-genus test set (Fig 4A). There is a negative correlation between

230   the percentage of loci genotyped and the known divergence times from *M. musculus* (r = -0.76;

231   p<0.0001). In the intra-genus test set, heterozygosity increases as divergence time increases (Fig

232   4B). The increase in percent heterozygosity of Mus samples is positively correlated with an

233    increase in divergence times (r = 0.93; p-value<0.0001). There is no plateau or obvious

234    underestimate of genetic diversity for samples in the intra-genus test set.

235

236    A tree of relatedness derived from SNP-based genetic distance values differentiates Mus samples

237    of the intra-genus test set from one another at a species level (Fig 5). Enough genetic diversity is

238    captured using the MDGA to reflect the known taxonomic relationships between the intra-genus

239    samples at a species level. At 9.5 MYD, the pygmy mouse subspecies *M. n. minutoides* is

240    grouped with the subspecies *M. n. orangiae* and not the replicate data file of the same species.

241

242    **Peromyscus case study**

243    Seven Peromyscus species were genotyped together as a case study to determine if the MDGA

244    could provide useful results that reflect known biological diversity for a number of species of a

245    different genus from Mus. Of the Peromyscus samples queried, approximately 52% of loci

246    queried by the array produce a genotype (Table 2). There are 159,797 loci genotyped across all

247    seven samples (32% of loci queried by the array) despite a 32.7 million-year divergence time

248    from *M. musculus*. SNP-based genetic distances of Peromyscus species were utilized to produce

249    trees of genetic relatedness that reflect the known divergence times of these species (Fig 6). Top

250    KEGG pathway annotations of the genotyped loci in Peromyscus samples are associated with

251    neurological signaling (Table 3).

252

**Table 2** Percentage of loci genotyped and percent heterozygosity in a Peromyscus case study (n=7)

| MDGA Data (CEL) File | Sample Scientific Name | Loci Genotyped (%) | Heterozygosity (%) |
|---|---|---|---|
| SNP_mDIV_B2-660_102109.CEL | *P. melanophrys* | 51.31 | 34.83 |
| SNP_mDIV_B1-659_102109.CEL | *P. aztecus* | 52.03 | 36.02 |
| SNP_mDIV_B3-661_102109.CEL | *P. californicus* | 52.13 | 36.27 |
| SNP_mDIV_B4-662_102109.CEL | *P. m. sonoriensis* | 52.26 | 35.95 |
| SNP_mDIV_B5-663_102109.CEL | *P. m. bairdii* | 52.27 | 36.71 |
| SNP_mDIV_B6-664_102109.CEL | *P. polionotus* | 52.57 | 37.02 |
| SNP_mDIV_B8-666_102109.CEL | *P. leucopus* | 52.62 | 36.55 |

253

254

255

**Table 3** Top KEGG[1] (Kyoto Encyclopedia of Genes and Genomes) pathways determined using the DAVID functional annotation tool

| KEGG Pathway associated with SNP loci genotyped in Peromyscus species (p<0.001) |
|---|
| Glutamatergic synapse |
| Circadian entrainment |
| Axon guidance |
| Retrograde endocannabinoid signaling |
| Dopaminergic synapse |
| Morphine addiction |
| Long-term depression |
| Hippo signaling pathway |
| cAMP signaling pathway |
| Cholinergic synapse |
| Rap1 signaling pathway |
| Long-term potentiation |
| GABAergic synapse |

256   [1]Enriched KEGG (Kyoto Encyclopedia of Genes and Genomes) pathways determined using
257   DAVID (Database for Annotation, Visualization, and Integrated Discovery) functional
258   annotation tool
259

260

14

261    ***In silico* cross-validation of potential polymorphic loci**

262    *P. maniculatus* was examined given that there is a partial genome sequence available online for

263    *in silico* search of unique and perfect 25 nt MDGA probe target sequence matches. There are

264    226,265 loci on the MDGA genotyped (~52%) for both *P. maniculatus bairdii* and *P.*

265    *maniculatus sonoriensis* within this study. Of the loci that were genotyped, there are 143,971 loci

266    that were genotyped as heterozygous in both *P. maniculatus* samples (Table 4). Heterozygous

267    loci represent potential polymorphic loci that can query both the common and uncommon allele

268    in a population. There are 6,075 MDGA probe sequences that perfectly match a unique position

269    within the *P. maniculatus* genome, and 481 of the *in silico* sequence matches are associated with

270    heterozygous loci.

271
272    **Table 4** In silico validation of potential polymorphic loci conserved cross-species

| Common Name | Scientific Name | MYD from *M. musculus* | Unique in silico matches of alleles queried | Number of heterozygous loci in all samples | Number of candidate polymorphic loci |
|---|---|---|---|---|---|
| Ryukyu Mouse | *Mus caroli* | 7.41 | 303,680 | 147,452 | 9,413 |
| Gairdner's Shrewmouse | *Mus pahari* | 8.29 | 152,971 | 251,902 | 9.341 |
| Rat | *Rattus norvegicus* | 20.9 | 61,372 | 85,926 | 1,019 |
| Deer Mouse | *Peromyscus maniculatus* | 32.7 | 6,075 | 143,971 | 481 |
| Naked Mole Rat | *Heterocephalus glaber* | 73 | 1,179 | 91,324 | 52 |

273    MYD=Millions of Years Divergence
274    Candidate Polymorphic Loci are the number of loci identified that had a heterozygous genotype
275    call for the samples using the Mouse Diversity Genotyping Array Data that could also be mapped
276    to the available genomic sequences for these organisms.

277

278    An average of 382,968 loci were genotyped between three available *M. caroli* CEL files using

279    the MDGA, and there are 303,680 unique theoretical matches to the *M. caroli* genome

280    determined through an *in silico* search using E-MEM (Table 4). A shrew mouse (*M. pahari*)

281    applied to the array has 411,514 loci that were genotyped experimentally using the MDGA.

282    Theoretically, there are 152,971 unique sequences from the MDGA that are present in the shrew

283    mouse only once (Table 4). The pathways associated with genotyped loci in *M. musculus*, *M.*

284    *caroli*, and *M. pahari* that are shared between these three species are primarily signaling

285    pathways and pathways involved in maintaining the structural integrity of a cell, such as focal

286    adhesion and adherens junction (Table 5). The Sprague Dawley rat (*R. norvegicus*) has a fully

287    sequenced and annotated genome available online. There are 170,156 loci that were genotyped

288    experimentally in both *R. norvegicus* samples using the MDGA. Using the E-MEM *in silico*

289    program, 61,372 sequences were determined to be theoretically present within the genome

290    (Table 4).

291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310

311 **Table 5** Top KEGG pathways enriched for house mouse gene annotations with genotype
312 assignments across wild Mus species

| KEGG pathways[1] significant (p<0.001) in reference house mouse (build 38) and wild Mus test samples[2] |
| :---: |
| Focal adhesion |
| Rap1 signaling pathway |
| Adherens junction |
| cAMP signaling pathway |
| ErbB signaling pathway |
| cGMP-PKG signaling pathway |
| Neuroactive ligand-receptor interaction |
| Platelet activation |
| Calcium signaling pathway |
| Purine metabolism |
| Phosphatidylinositol signaling system |
| Amoebiasis |
| Regulation of actin cytoskeleton |
| PI3K-Akt signaling pathway |
| Oxytocin signaling pathway |

313 [1]Enriched KEGG (Kyoto Encyclopedia of Genes and Genomes) pathways determined using
314 DAVID (Database for Annotation, Visualization, and Integrated Discovery) functional
315 annotation tool
316 [2]Mus test samples are *M. pahari* and *M. caroli* species
317 KEGG pathways are shared between the reference *M. musculus*, *M. pahari*, and *M. caroli* species
318

319

320 Special attention was given to potential polymorphic loci that were genotyped as heterozygous in

321 samples using the MDGA and could be cross-validated as being present in the genome using an

322 *in-silico* search of publicly available genome sequences. There is a trend of there being more

323 heterozygous loci genotyped using the MDGA than the number of those loci that can be cross

324 validated as present in the publicly available genome sequence (Table 4). There are 147,452

325 heterozygous loci genotyped in all three *M. caroli* samples, and 9,413 of these loci were

326 validated as present in the publicly available genome sequence (Table 4). There are 9,341 of the

327 147,452 heterozygous loci genotyped in a *M. pahari* sample that were cross validated as

328     potential polymorphic SNP loci (Table 4). In two *R. norvegicus* samples, there are 85,926 loci

329     that were genotyped empirically using the MDGA, and 1,019 loci that were cross-validated using

330     an *in-silico* genome sequence search (Table 4).

331

332

333

334

335  **Discussion**

336  Specialized genotyping arrays have been successfully applied cross-species to closely related

337  organisms in previous research [2–4,9–13,16,32–34]. Here we present evidence that the MDGA

338  can be applied to wild rodents to produce SNP genotyping results that reflect the known

339  taxonomic relationships between test samples and the reference house mouse. The identification

340  of polymorphic SNPs within non-model organisms is of great interest, as these genetic markers

341  can be used to assay diversity in wild populations in studies of population genetics [2–

342  4,9,16,18,34,35]. Panels of candidate polymorphic SNPs have been identified for wild species of

343  the genus Mus and Peromyscus. This study is a first step in contributing where there is a paucity

344  of information available for non-model rodent species.

345

346  Outside of the genus Mus, the plateau in SNP loci genotyped and the percentage of heterozygous

347  loci is attributed to off-target mutations that hinder DNA hybridization to array probe sequences.

348  When DNA hybridizes to a probe on the MDGA, the hybridization does not have to be a perfect

349  25 nt match, where incomplete hybridization of the sample DNA to the probe is enough to result

350  in a genotype assignment [36]. Determination of the divergence time from *M. musculus* at which

351  genetic diversity is underestimated is limited by the samples available for use in this study. A

352  greater number of species genotyped using the MDGA that have a divergence time between 10-

353  15 MYD from the house mouse would be beneficial in identifying situations where

354  underestimations of genetic diversity occur. Miller et al. (2012), found previously that applying

355  the Bovine, Ovine, and Equine SNP50 Beadchip arrays cross-species resulted in a linear decrease

356  in genotyped loci as the millions of years of divergence from the model species increased [10].

357  Previous studies that have examined the utility of the cross-species application of commercially

358    available genotyping array technology have identified trends of decreasing ability to genotype

359    loci as divergence time from the model organism increases as well [8,9]. This study is unique as

360    it tests the array technology on a wide range of species spanning multiple millions of years

361    divergence from the reference house mouse.

362

363    Previous research has determined potentially conserved sequences between model organisms and

364    the wild species of interest through application of commercial arrays to test samples [2,4]. This

365    study of the MDGA cross-validates genotyped loci in rodent samples with an *in silico* analysis of

366    available genomic sequences for wild species. The heterozygous SNP variation in rodent samples

367    of this study cross-validated through *in silico* analyses represents candidate polymorphic SNPs

368    that can be tested for conservation in populations of wild species of Mus and Peromyscus. To be

369    truly considered a polymorphic SNP conserved cross-species, the variation must be validated in

370    wild populations with the alternate, or minor allele present in at least 1% of the population.

371

372    A major difficulty in cross-species genotyping using array data is the assembly of appropriate

373    test sets that would allow for accurate genotyping. Previous research has demonstrated that the

374    genotyping algorithm recommended by Affymetrix, BRLMM-P, is sensitive to the composition

375    of the samples included in a test set [37,38]. Samples in a test set that are more similar to one

376    another genetically will produce fewer false genotyping results [38]. The number of loci

377    genotyped can become inflated if the samples in the test set are too genetically different, as was

378    seen when samples of different orders of classification were genotyped together in the inter-order

379    test set. The greater genetic homogeneity of only Mus samples in the intra-genus test set

380    produced genotyping results that matched what was expected of the species based on divergence

381    times. The linear decrease in loci genotyped in Mus samples as divergence time increased

382    reflected previous cross-species findings [10]. Recommendations for the construction of a test set

383    of samples for an experiment utilizing the MDGA cross-species would be dependent upon the

384    hypothesis tested. A large number of samples are needed to establish whether the minor allele of

385    a SNP is present in populations of non-model species for at least 1% of the population [7,39].

386    Technical replicates should be included to assess the quality of DNA hybridization to array

387    probes for a particular species. Optimization of hybridization conditions should be made to

388    reduce differences in array hybridization intensities and the resulting differences in genotype

389    assignments between technical replicates.

390

391    The use of a training set that has sufficient genetic diversity to encompass that of the

392    experimental test sets can assist in producing accurate genotyping of samples [40,41]. The

393    training set of 114 classical inbred strains of mice used in this study does not encompass the high

394    relative genetic diversity of the sample sets of this cross-species study. A training set optimized

395    for cross-species genotyping would be composed of members of the same species as the test set

396    and would be validated using another method such as sequencing. Inclusion of male and female

397    samples would ensure more accurate genotype assignments on the X chromosome, as

398    hemizygous males are assigned a diploid homozygous genotype [42]. Analyzing SNPs on the X

399    chromosome separately from autosomal SNPs and separating male and female samples would

400    aid in fewer false genotype assignments.

401

402    In comparing the research knowledge gained through this study using the deer mouse (*P.*

403    *maniculatus*) to the knowledge obtained from the study of Antarctic fur seals by Hoffman et al.

404    (2013), similar metrics of utility were obtained through cross-species genotyping (Table 6).

405    Given that the mouse array targets over two times the number of positions than the canine array

406    targets, there is a much larger number of loci that can be genotyped in the deer mouse than the

407    Antarctic fur seal. Future studies will focus on validating a panel of SNPs that are polymorphic

408    in deer mouse populations. Pathway analyses are limited by the information assayed by each

409    technology and are with respect to the annotations of the model organisms. As new sequence

410    information and genome annotations become available for the deer mouse, it will be interesting

411    to see which SNP markers associated with conserved pathways will be found to be shared

412    between the house mouse and the deer mouse. The deer mouse is an intriguing sentinel of

413    environmental effects and a model for population studies that has a surprising lack of genomic

414    information available [18,43]. Cross-species array use may be one technique to identify SNP

415    diversity in these relevant species until genome sequencing prices become more affordable for

416    non-model species. The use of a rat genotyping array in the future may be of use, as the deer

417    mouse and rat share greater genetic synteny than with the mouse [44].

418    **Table 6. Comparison of the *Hoffman et al. (2013)* model study with the current study**

| Hoffman *et al.* | Comparison | Kelly et al. |
|---|---|---|
| Antarctic fur seal | Non-model species | Deer Mouse |
| CanFam2.0 | Reference genome for array | *Mus musculus* |
| 44 | MYD from model species | 32.7 |
| 173,662 | Loci queried by array | 493,290 |
| 33,324 | Loci genotyped | ~226,000 |
| 2 of 5 | Loci validated in silico | 3,195 |
| 173 | Polymorphic loci | 481 |
| 2 | Loci validated in a population | Future |
| Energy metabolism | Pathways shared between model and non-model species | Neurological signaling |

419    MYD = Million years divergence

420

*Mouse single nucleotide polymorphic targets for cross hybridization in rodents*

421    There is a great potential for cross-species MDGA utility for wild Mus species in providing

422    genomic markers for research in mouse population genetics and studies of rodent evolution.

423    Genotype data generated from application of the MDGA captured enough genetic diversity to

424    differentiate Mus samples at a species level. Further testing is required to determine if the

425    MDGA can capture enough diversity to differentiate between subspecies. As in the case of the

426    deer mouse, wild Mus species represent an untapped wealth of genomic information that would

427    benefit researchers of environmental mutagens, evolution, and population genetics. With newer

428    mouse array technologies becoming available that have greater capacity for high-throughput

429    analysis, novel polymorphic SNPs in non-model rodents can be identified through a low-cost and

430    efficient manner.

431

432    Utilizing the Mouse Diversity Genotyping Array for cross-species genotyping represents a first

433    step towards development of a tool that can rapidly identify SNP variation in wild rodent species.

434    A panel of candidate SNPs on the MDGA have been identified for use with wild mouse species

435    and was cross validated using an *in silico* genome search. Future work may address the

436    validation of this candidate cross-species panel in wild populations. This research highlights the

437    need for greater genomic resources for wild rodents and demonstrates the potential of the MDGA

438    as a high-throughput genotyping tool for non-model organisms. The development of novel tools

439    specialized for non-model species opens up previously inaccessible avenues of research. Next-

440    generation sequencing technologies are often not accessible and too costly for a majority of

441    researchers with population-based research questions that require rapid, high-throughput genome

442    wide analysis of variation. Until the price of sequencing and the complexity of assembling new

443    reference genome assemblies is reduced, the adaptation of existing genomic tools for use with

444    closely related species is one method researchers can use to combat the genomic disparity

445    between studying model and non-model species.

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467    **MATERIALS AND METHODS**

468    **Cross-species samples**

469    Forty publicly available MDGA raw data (CEL) files were downloaded from the Center for

470    Genome Dynamics at the Jackson Laboratory (2012, The Jackson Laboratory;

471    ftp.jax.org/petrs/MDA/). Four MDGA CEL files of *H. glaber* DNA cross-hybridization to the

472    MDGA were generated in-house. The forty-four samples consist of twenty-seven Mus CEL files,

473    two Rattus CEL files, seven Peromyscus CEL files, one Apodemus CEL file, and CEL files

474    representing more highly diverged species including a squirrel, four naked mole rats, a tapir, and

475    an African Black Rhino (S3 Table). CEL file raw array intensity images were analyzed for

476    quality and abnormalities in array images were noted. Two CEL files (S1 Fig) were noted for

477    having an abnormal spot with uneven DNA hybridization to the array. Due to the redundancy of

478    probes on the MDGA, it was determined that abnormal CEL file images still had sufficient

479    genomic coverage to be used for further analysis and were not removed from the study.

480

481    **SNP genotyping**

482    Samples were genotyped using the protocol outlined by Locke et al. [25].  Affymetrix Power

483    Tools was used to generate genotype calls of AA, AB, BB, or No Call (numerical representations

484    0, 1, 2, -1, respectively) using the BRLMM-P algorithm for 493,290 SNPs [25] (Affymetrix

485    Power Tools (APT) Release 1.16.0). A training set of 114 classical laboratory mouse CEL files

486    obtained from a set of 351 mice utilized by Didion et al. (2012) was used in conjunction with

487    BRLMM-P to train the algorithm in accurate assignment of genotypes [26].  The samples were

488    organized into three test sets that were genotyped separately from one another. The first

489    genotyping set (known as the inter-order test set) consists of all 44 CEL files representing species

490    spanning different orders of classification and a maximum divergence time of 96 million years of

491    divergence (MYD) from the reference house mouse, *Mus musculus* (Table 1). The second test set

492    (the intra-genus test set) is composed of the 27 samples from the genus Mus and has a maximum

493    divergence of 9.5 MYD from the house mouse (Table 1). The third test set (Peromyscus case

494    study test set) was composed of seven deer mouse species from the genus Peromyscus that have

495    32.7 MYD from the house mouse (Table 1). The genotyping results obtained were analyzed and

496    compared to reference genotyping data from *Mus musculus*. The reference *Mus musculus* data

497    was obtained by averaging the genotyping results from 8 *Mus musculus* samples (percentage of

498    loci genotyped > 99%).

499

500    **Estimation of divergence times**

501    The estimated divergence time of each species from the reference house mouse was obtained

502    using an evolutionary timetree of life (http://www.timetree.org/) [27] with a few exceptions. The

503    estimated divergence of the subspecies *M. m. castaneus* was determined through previous work

504    by Geraldes et al. (2012) [28], and the evolutionary divergence time of the pygmy mouse species

505    from the house mouse was determined by Kouassi et al. (2008) [29].

506

507    **Statistical analyses**

508    A Fisher's exact test was utilized to assess the extent of genetic differences between samples

509    genotyped together. A nonparametric, unordered, Fisher-Freeman-Halton exact test (Monte

510    Carlo simulation) was performed using the StatXact statistical analysis software package

511    (CYTEL Software, Cambridge, MA). Pearson's r was used in tests of significance of correlations

512    between the genotyping results of the test set samples using Graphpad Prism 8 software.

513

**Genetic distance calculations**

515 Pairwise comparison of SNP genotypes between species in the inter-order test set was utilized to

516 create SNP-based distance matrices using R. The distance matrix values used to create

517 phenograms (SNP trees) were generated using an in-house R script courtesy of Marjorie E.

518 Osbourne Locke. The in-house script utilized the 'bionj' R package to create a tree of genetic

519 relatedness using the neighbour-joining method [30]. The resulting trees were modified using

520 Figtree (v1.4.3) software. Pairwise genetic distances were computed by dividing the total number

521 of genotypic differences between two samples by the total number of loci queried by the MDGA,

522 where 493,290 total loci were used in this study [25]. The values in the distance matrix are a

523 numerical representation of the amount of genetic diversity between test species analyzed and the

524 reference house mouse. A genetic distance value of zero indicates the species are genetically the

525 same at the loci queried, and a value of one indicates the species compared are completely

526 genetically dissimilar from one another at the loci queried. The estimated evolutionary

527 relationships seen in the SNP trees generated were compared to the divergence times of test

528 samples from the reference house mouse provided in literature and the Timetree database [27–

529 29].

530

**In silico validation of MDGA loci genotyped cross-species and pathway analysis**

532 *In silico* validation of loci genotyped from MDGA data was performed using the program E-

533 MEM (efficient computation of maximal exact matches for very large genomes) designed by

534 Khiste and Ilie (2015) [31]. The publicly available genomes of rodents were searched for the

535 unique presence of MDGA probe sequences. E-MEM was employed to search a publicly

536    available genome of wild rodents available on NCBI (S4 Table) for perfect 25 nt MDGA SNP

537    probe target sequences that have only one genomic match (ftp.ncbi.nlm.nih.gov/genomes/).

538    Unique MDGA matches discovered via E-MEM were identified and then compared with the list

539    of heterozygous loci genotyped using the MDGA. Ensembl gene IDs associated with candidate

540    loci genotyped were analyzed using the Database for Annotation, Visualization, and Integrated

541    Discovery (DAVID).

542

## Acknowledgements

## Figure Legends

**Fig 1. Summary of published research on mammalian cross-species genotyping using SNP genotyping microarrays**
(A) Published research is organized in increasing order of genetic divergence in millions of years divergence (MYD) of non-model test samples from the model reference organism. Authors, publication year, genotyping microarray technology, and approximate number of loci queried (in thousands) are listed for each publication. (B) The sample of publications on mammalian cross-species array studies with the 13th representing the contributions of this thesis to the cross-species genotyping array field.

**Fig 2. Genetic diversity of test sets exceeds maximum genetic diversity of training set**
Boxplots representing the minimum, first quartile, median, third quartile, and maximum genetic distances for the training set (n=114), intra-genus test set (n=27), case study of Peromyscus (n=7), and inter-order test set (n=44). All genetic distances are with respect to the reference house mouse *Mus musculus*.

**Fig 3. Underestimation of genetic diversity for highly diverged species in cross-species genotyping**
(A) The percentage of loci genotyped from the inter-order test set (n=44). (B) The percentage of loci from the inter-order test set with a heterozygous genotype call. MYD = Millions of years divergence, with respect to the reference *Mus musculus*.

577 **Fig 4. Genetic diversity of wild Mus species**
578 (A) The percentage of loci genotyped from the intra-genus test set (n=27). (B) The percentage of
579 loci from the intra-genus test set with a heterozygous genotype call. MYD = Millions of years
580 divergence, with respect to the reference *Mus musculus*.
581
582 **Fig 5. SNP distance-based tree of genetic relatedness reflects known taxonomic**
583 **relationships between Mus species**
584 SNP distance-based tree of genetic relatedness of samples from the intra-genus test set (n = 27).
585 At 9.5 MYD a pygmy mouse subspecies *M. n. orangiae* has SNP-based genetic distances that
586 reflect greater genetic similarity to another pygmy mouse subspecies *M. n. minutoides* than the
587 replicate MDGA data file of the same *M. n. orangiae* sample. MYD = Millions of years
588 divergence, with respect to the reference *Mus musculus*.
589
590 **Fig 6. SNP distance-based tree of genetic relatedness reflects known taxonomic**
591 **relationships between Peromyscus species**
592 Pairwise SNP distance-based tree of genetic relatedness of samples from the intra-genus test set
593 of Peromyscus species (n=7).
594
595
596
597 # References

598 1.    Russell JJ, Theriot JA, Sood P, Marshall WF, Landweber LF, Fritz-Laylin L, et al. Non-
599       model model organisms. BMC Biol. 2017;15: 55. doi:10.1186/s12915-017-0391-5

600 2.    Ogden R, Baird J, Senn H, McEwing R. The use of cross-species genome-wide arrays to
601       discover SNP markers for conservation genetics: A case study from Arabian and scimitar-
602       horned oryx. Conserv Genet Resour. 2012;4: 471–473. doi:10.1007/s12686-011-9577-2

603 3.    Hoffman JI, Thorne MAS, McEwing R, Forcada J, Ogden R. Cross-Amplification and
604       Validation of SNPs Conserved over 44 Million Years between Seals and Dogs. PLoS One.
605       2013;8: 1–10. doi:10.1371/journal.pone.0068365

606 4.    More M, Gutiérrez G, Rothschild M, Bertolini F, Ponce de León FA. Evaluation of SNP
607       Genotyping in Alpacas Using the Bovine HD Genotyping Beadchip. Front Genet.
608       2019;10: 361. doi:10.3389/fgene.2019.00361

609 5.    Maresso K, Broeckel U. Genotyping Platforms for Mass-Throughput Genotyping with
610       SNPs, Including Human Genome-Wide Scans. Adv Genet. 2008;60: 107–139.
611       doi:10.1016/S0065-2660(07)00405-1

612 6.    LaFramboise T. Single nucleotide polymorphism arrays: a decade of biological,
613       computational and technological advances. Nucleic Acids Res. 2009;37: 4181–93.
614       doi:10.1093/nar/gkp552

615 7.    Wang DG, Fan JB, Siao CJ, Berno A, Young P, Sapolsky R, et al. Large-scale
616       identification, mapping, and genotyping of single-nucleotide polymorphisms in the human

617      genome. Science. 1998;280: 1077–82. doi:10.1126/science.280.5366.1077

618   8.   Xia W, Luo T, Zhang W, Mason AS, Huang D, Huang X, et al. Development of High-
619      Density SNP Markers and Their Application in Evaluating Genetic Diversity and
620      Population Structure in Elaeis guineensis. Front Plant Sci. 2019;10: 130.
621      doi:10.3389/fpls.2019.00130

622   9.   Kharzinova VR, Sermyagin AA, Gladyr EA, Okhlopkov IM, Brem G, Zinovieva NA. A
623      study of applicability of SNP chips developed for bovine and ovine species to whole-
624      genome analysis of reindeer rangifer tarandus. J Hered. 2015;106: 758–761.
625      doi:10.1093/jhered/esv081

626   10.   Miller JM, Kijas JW, Heaton MP, McEwan JC, Coltman DW. Consistent divergence times
627      and allele sharing measured from cross-species application of SNP chips developed for
628      three domestic species. Mol Ecol Resour. 2012;12: 1145–1150. doi:10.1111/1755-
629      0998.12017

630   11.   Miller JM, Festa-Bianchet M, Coltman DW. Genomic analysis of morphometric traits in
631      bighorn sheep using the Ovine Infinium ® HD SNP BeadChip. PeerJ. 2018;6: e4364.
632      doi:10.7717/peerj.4364

633   12.   Moravcikova N, Kirchner R, Sidlova V, Kasarda R, Trakovicka A. Estimation of genomic
634      variation in cervids using cross-species application of SNP arrays.
635      Poljoprivreda/Agriculture. 2015;21: 33–36. doi:10.18047/poljo.21.1.sup.6

636   13.   Pertoldi C, Wójcik JM, Tokarska M, Kawałko A, Kristensen TN, Loeschcke V, et al.
637      Genome variability in European and American bison detected using the BovineSNP50
638      BeadChip. Conserv Genet. 2010;11: 627–634. doi:10.1007/s10592-009-9977-y

639   14.   vonHoldt BM, Pollinger JP, Lohmueller KE, Han E, Parker HG, Quignon P, et al.
640      Genome-wide SNP and haplotype analyses reveal a rich history underlying dog
641      domestication. Nature. 2010;464: 898–902. doi:10.1038/nature08837

642   15.   Aslam ML, Bastiaansen JW, Elferink MG, Megens H-J, Crooijmans RP, Blomberg L, et
643      al. Whole genome SNP discovery and analysis of genetic diversity in Turkey (Meleagris
644      gallopavo). BMC Genomics. 2012;13: 391. doi:10.1186/1471-2164-13-391

645   16.   McCue ME, Bannasch DL, Petersen JL, Gurr J, Bailey E, Binns MM, et al. A high density
646      SNP array for the domestic horse and extant Perissodactyla: Utility for association
647      mapping, genetic diversity, and phylogeny studies. PLoS Genet. 2012;8.
648      doi:10.1371/journal.pgen.1002451

649   17.   Minias P, Dunn PO, Whittingham LA, Johnson JA, Oyler-McCance SJ. Evaluation of a
650      Chicken 600K SNP genotyping array in non-model species of grouse. Sci Rep. 2019;9: 1–
651      10. doi:10.1038/s41598-019-42885-5

652   18.   Harris SE, Munshi-South J, Obergfell C, Neill O. Signatures of Rapid Evolution in Urban
653      and Rural Transcriptomes of White-Footed Mice (Peromyscus leucopus) in the New York

654           Metropolitan Area. PLoS One. 2013;8: 74938. doi:10.1371/journal.pone.0074938

655    19.    Williams LM, Ma X, Boyko AR, Bustamante CD, Oleksiak MF. SNP identification,
656           verification, and utility for population genetics in a non-model genus. BMC Genet.
657           2010;11: 32. doi:10.1186/1471-2156-11-32

658    20.    Hulme-Beaman A, Dobney K, Cucchi T, Searle JB. An Ecological and Evolutionary
659           Framework for Commensalism in Anthropogenic Environments. Trends Ecol Evol.
660           2016;31: 633–645. doi:10.1016/j.tree.2016.05.001

661    21.    Rodríguez-Estival J, Smits JEG. Small mammals as sentinels of oil sands related
662           contaminants and health effects in northeastern Alberta, Canada. Ecotoxicol Environ Saf.
663           2016;124: 285–295. doi:10.1016/j.ecoenv.2015.11.001

664    22.    Bedford NL, Hoekstra HE. Peromyscus mice as a model for studying natural variation.
665           Elife. 2015;4: 1–13. doi:10.7554/eLife.06813

666    23.    Vrana PB, Shorter KR, Szalai G, Felder MR, Crossland JP, Veres M, et al. *Peromyscus*
667           (deer mice) as developmental models. Wiley Interdiscip Rev Dev Biol. 2014;3: 211–230.
668           doi:10.1002/wdev.132

669    24.    Yang H, Ding Y, Hutchins LN, Szatkiewicz J, Bell TA, Paigen BJ, et al. A customized
670           and versatile high-density genotyping array for the mouse. Nat Methods. 2009;6: 663–
671           666. doi:10.1038/nmeth.1359

672    25.    Locke MEO, Milojevic M, Eitutis ST, Patel N, Wishart AE, Daley M, et al. Genomic copy
673           number variation in Mus musculus. Additional file 1. BMC Genomics. 2015;16: 497.
674           doi:10.1186/s12864-015-1713-z

675    26.    Didion JP, Yang H, Sheppard K, Fu C-P, McMillan L, de Villena F, et al. Discovery of
676           novel variants in genotyping arrays improves genotype retention and reduces
677           ascertainment bias. BMC Genomics. 2012;13: 34. doi:10.1186/1471-2164-13-34

678    27.    Hedges SB, Marin J, Suleski M, Paymer M, Kumar S. Tree of life reveals clock-like
679           speciation and diversification. Mol Biol Evol. 2015;32: 835–845.
680           doi:10.1093/molbev/msv037

681    28.    Geraldes A, Basset P, Smith KL, Nachman MW. Higher differentiation among subspecies
682           of the house mouse (Mus musculus) in genomic regions with low recombination. Mol
683           Ecol. 2012;20: 4722–4736. doi:10.1111/j.1365-294X.2011.05285.x

684    29.    Kouassi SK, Nicolas V, Aniskine V, Lalis A, Cruaud C, Couloux A, et al. Taxonomy and
685           biogeography of the African Pygmy mice, Subgenus Nannomys (Rodentia, Murinae, Mus)
686           in Ivory Coast and Guinea (West Africa). Mammalia. 2008;72: 237–252.
687           doi:10.1515/MAMM.2008.035

688    30.    Gascuel O. BIONJ: an improved version of the NJ algorithm based on a simple model of
689           sequence data. Mol Biol Evol. 1997;14: 685–695.

690          doi:10.1093/oxfordjournals.molbev.a025808

691    31.   Khiste N, Ilie L. E-MEM: efficient computation of maximal exact matches for very large
692          genomes. Bioinformatics. 2015;31: 509–514. doi:10.1093/bioinformatics/btu687

693    32.   Vonholdt BM, Pollinger JP, Lohmueller KE, Han E, Parker HG, Quignon P, et al.
694          Genome-wide SNP and haplotype analyses reveal a rich history underlying dog
695          domestication. Nature. 2010;464: 898–902. doi:10.1038/nature08837

696    33.   Haynes GD, Latch EK. Identification of Novel Single Nucleotide Polymorphisms (SNPs)
697          in Deer (Odocoileus spp.) Using the BovineSNP50 BeadChip. Breuker C, editor. PLoS
698          One. 2012;7: e36536. doi:10.1371/journal.pone.0036536

699    34.   Michelizzi VN, Wu X, Dodson M V, Michal JJ, Zambrano-Varon J, McLean DJ, et al. A
700          global view of 54,001 single nucleotide polymorphisms (SNPs) on the Illumina
701          BovineSNP50 BeadChip and their transferability to water buffalo. Int J Biol Sci. 2010;7:
702          18–27. doi:10.7150/ijbs.7.18

703    35.   Williams LM, Ma X, Boyko AR, Bustamante CD, Oleksiak MF. SNP identification,
704          verification, and utility for population genetics in a non-model genus. BMC Genet.
705          2010;11: 32. doi:10.1186/1471-2156-11-32

706    36.   Binder H, Preibisch S. Specific and nonspecific hybridization of oligonucleotide probes on
707          microarrays. Biophys J. 2005;89: 337–52. doi:10.1529/biophysj.104.055343

708    37.   Miclaus K, Wolfinger R, Vega S, Chierici M, Furlanello C, Lambert C, et al. Batch effects
709          in the BRLMM genotype calling algorithm influence GWAS results for the Affymetrix
710          500K array. Pharmacogenomics J. 2010;10: 336–46. doi:10.1038/tpj.2010.36

711    38.   Hong H, Su Z, Ge W, Shi L, Perkins R, Fang H, et al. Assessing batch effects of genotype
712          calling algorithm BRLMM for the Affymetrix GeneChip Human Mapping 500 K array set
713          using 270 HapMap samples. 2008; doi:10.1186/1471-2105-9-S9-S17

714    39.   Akey JM. The Effect of Single Nucleotide Polymorphism Identification Strategies on
715          Estimates of Linkage Disequilibrium. Mol Biol Evol. 2003;20: 232–242.
716          doi:10.1093/molbev/msg032

717    40.   Zhang P, Zhan X, Rosenberg NA, Zöllner S. Genotype Imputation Reference Panel
718          Selection Using Maximal Phylogenetic Diversity. 2013; doi:10.1534/genetics.113.154591

719    41.   Huang L, Jakobsson M, Pemberton TJ, Ibrahim M, Nyambo T, Omar S, et al. Haplotype
720          variation and genotype imputation in African populations. Genet Epidemiol. 2011;35:
721          766–80. doi:10.1002/gepi.20626

722    42.   Zhao S, Jing W, Samuels DC, Sheng Q, Shyr Y, Guo Y. Strategies for processing and
723          quality control of Illumina genotyping arrays. Brief Bioinform. 2018;19: 765–775.
724          doi:10.1093/bib/bbx012

725  43.  Rodríguez-Estival J, Smits JEG. Small mammals as sentinels of oil sands related
726        contaminants and health effects in northeastern Alberta, Canada. Ecotoxicol Environ Saf.
727        2016;124: 285–295. doi:10.1016/J.ECOENV.2015.11.001

728  44.  Ramsdell CM, Lewandowski AA, Glenn J, Vrana PB, O'Neill RJ, Dewey MJ.
729        Comparative genome mapping of the deer mouse (Peromyscus maniculatus) reveals
730        greater similarity to rat (Rattus norvegicus) than to the lab mouse (Mus musculus). BMC
731        Evol Biol. 2008;8: 65. doi:10.1186/1471-2148-8-65

732

733

## Supporting Information

735 **S1 Fig. Abnormalities in two MDGA raw intensity CEL file images.** CEL file raw array
736 intensity images were analyzed for quality control purposes and abnormalities in array images
737 were noted for two CEL files. The two samples were not removed from analysis.
738
739 **S1 Table. Forty-four MDGA data (CEL) files of the present study.** [1]MDGA data (CEL) files
740 were downloaded from the Center for Genome Dynamics at the Jackson Laboratory, with the
741 exception of the four *H. glaber* CEL files (generated in-house). [2]Divergence time is given in
742 millions of years from the reference house mouse, *M. musculus* (timetree.org). [3]"redo" files are a
743 technical replicate of the CEL file with the same sample identifier code. Ex: SNP_mDIV_D3-
744 639_101509-redo is a technical replicate of SNP_mDIV_D3-639_91809, where D3-639 is the
745 sample identifier. [4]Only family level information available for CEL file SNP_mDIV_B9-
746 667_102109; Genus and species of sample are unknown.
747
748 **S2 Table. Training set of samples for genotyping algorithm (n=114).** CEL files of 114
749 classically inbred laboratory mouse strains were downloaded from the Jackson Laboratory
750 Center for Genome Dynamics for genotyping algorithm training.
751
752 **S3 Table. Genotype summary of 44 study samples genotyped at 493,290 single nucleotide**
753 **polymorphic loci located across the genome of *Mus musculus*.** Genotyping summary results
754 for all 44 Mouse Diversity Genotyping Array data files.
755
756 **S4 Table. Study species evaluated with publicly available nuclear genome sequence**
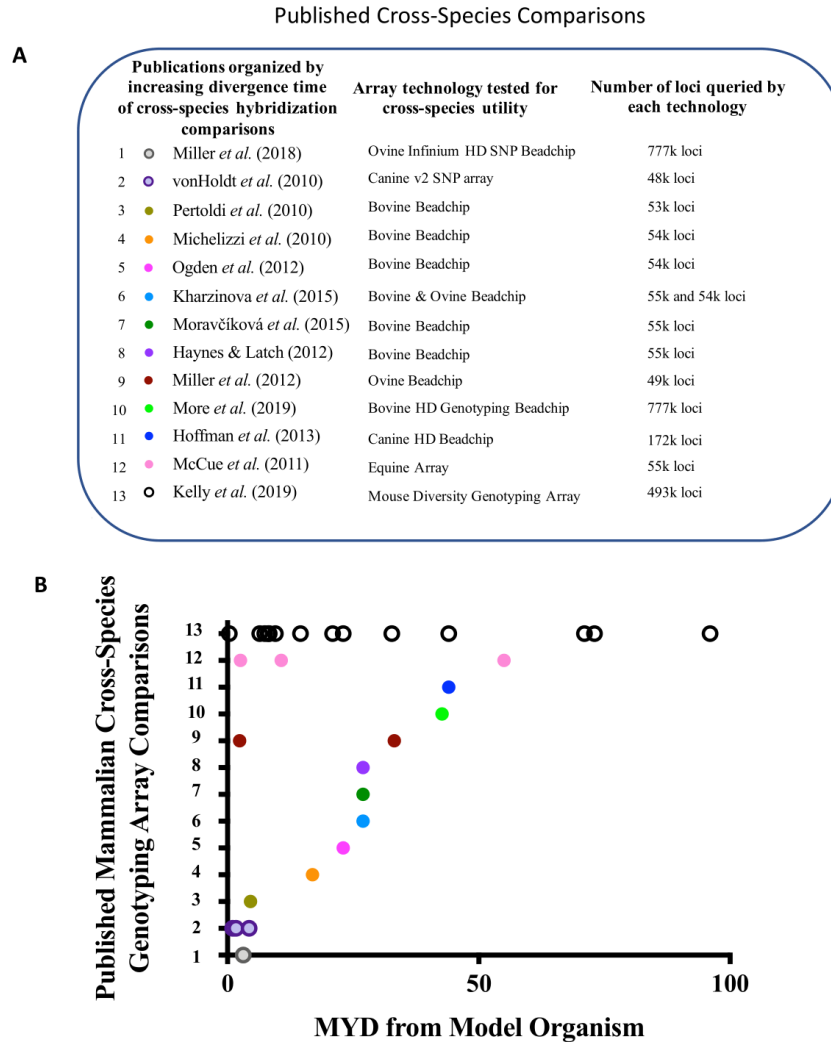757 **information.** Genomes accessed through the NCBI Genomes FTP site of samples under study
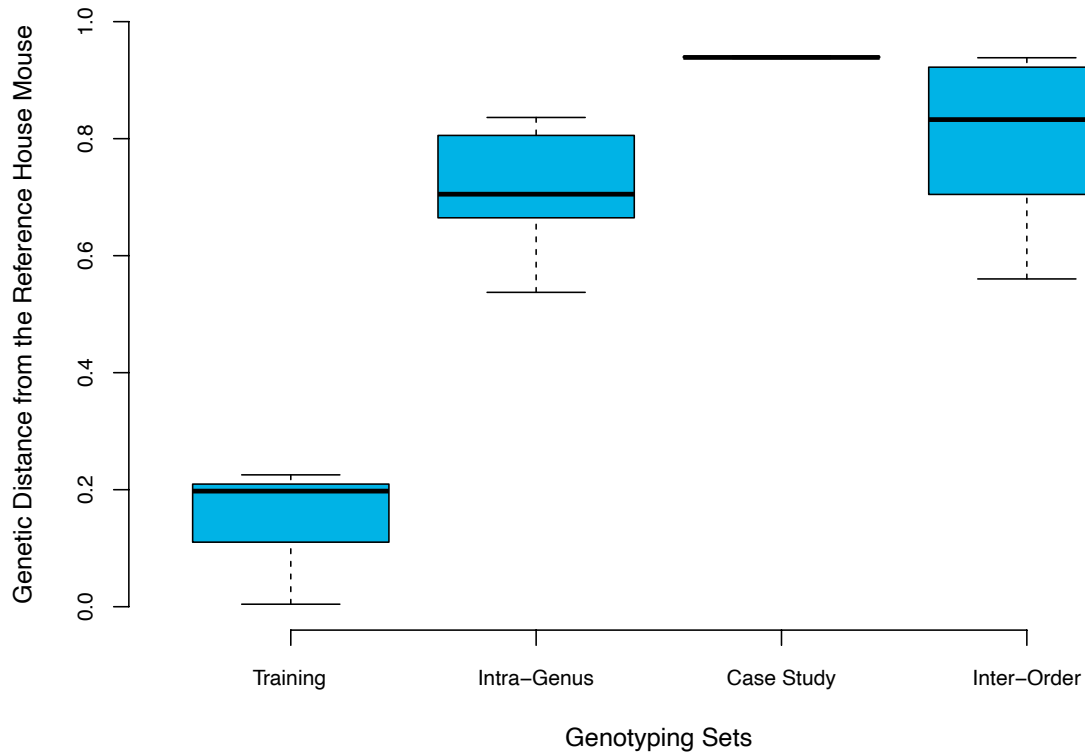758 (ftp.ncbi.nlm.nih.gov/genomes/).

759

760

761

762

763

764 **Figures**



765

**Fig 1. Summary of published research on mammalian cross-species genotyping using SNP genotyping microarrays**
(A) Published research is organized in increasing order of genetic divergence in millions of years divergence (MYD) of non-model test samples from the model reference organism. Authors, publication year, genotyping microarray technology, and approximate number of loci queried (in thousands) are listed for each publication. (B) The sample of publications on mammalian cross-species array studies with the 13th representing the contributions of this thesis to the cross-species genotyping array field.

774

**Fig 2. Genetic diversity of test sets exceeds maximum genetic diversity of training set**
Boxplots representing the minimum, first quartile, median, third quartile, and maximum genetic distances for the training set (n=114), intra-genus test set (n=27), case study of Peromyscus (n=7), and inter-order test set (n=44). All genetic distances are with respect to the reference house mouse *Mus musculus*.
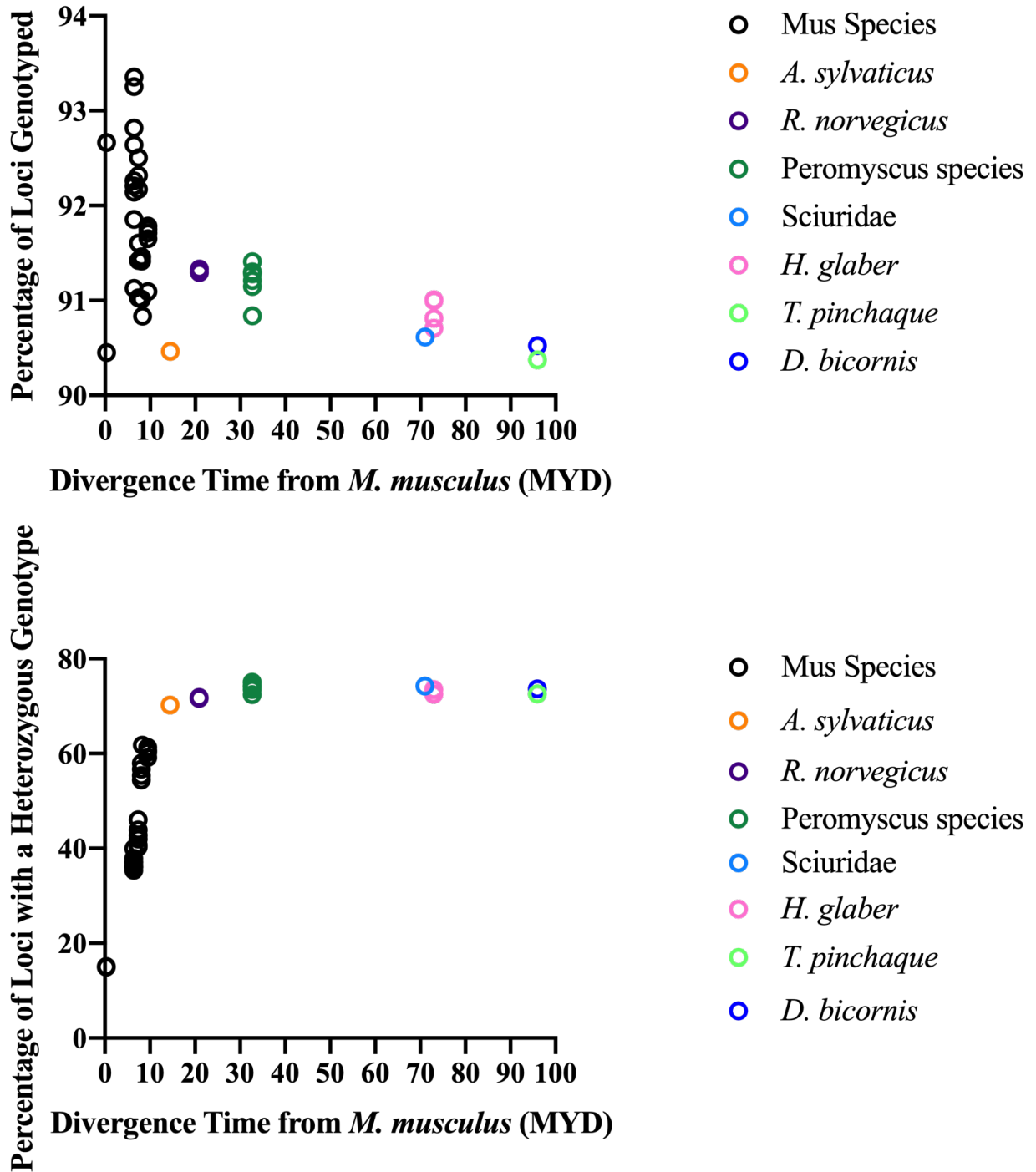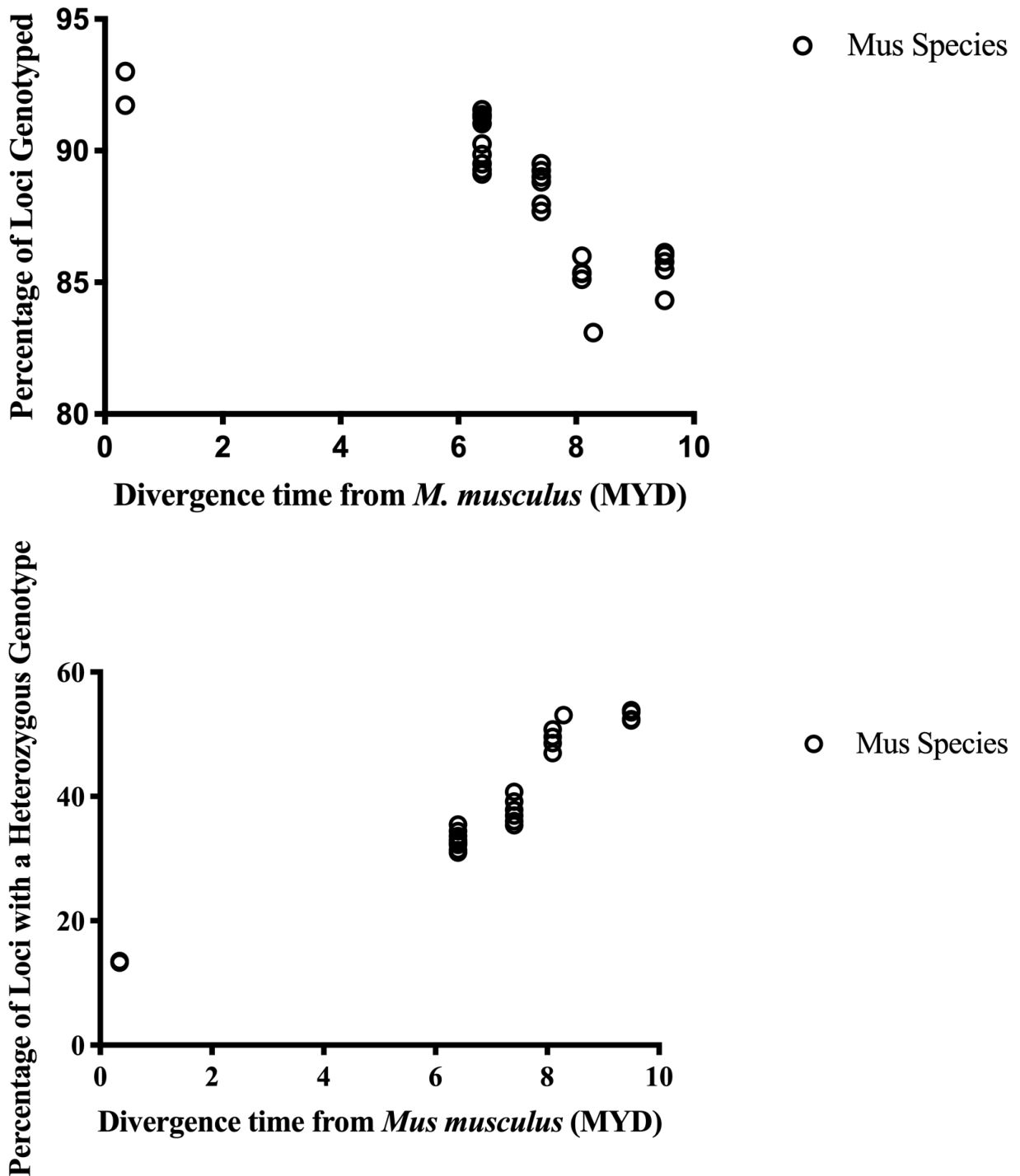
784



785
786
787 **Fig 3. Underestimation of genetic diversity for highly diverged species in cross-species**
788 **genotyping**
789 (A) The percentage of loci genotyped from the inter-order test set (n=44). (B) The percentage of
790 loci from the inter-order test set with a heterozygous genotype call. MYD = Millions of years
791 divergence, with respect to the reference *Mus musculus*.
792

793



794
795

**Fig 4. Genetic diversity of wild Mus species**

796   (A) The percentage of loci genotyped from the intra-genus test set (n=27). (B) The percentage of
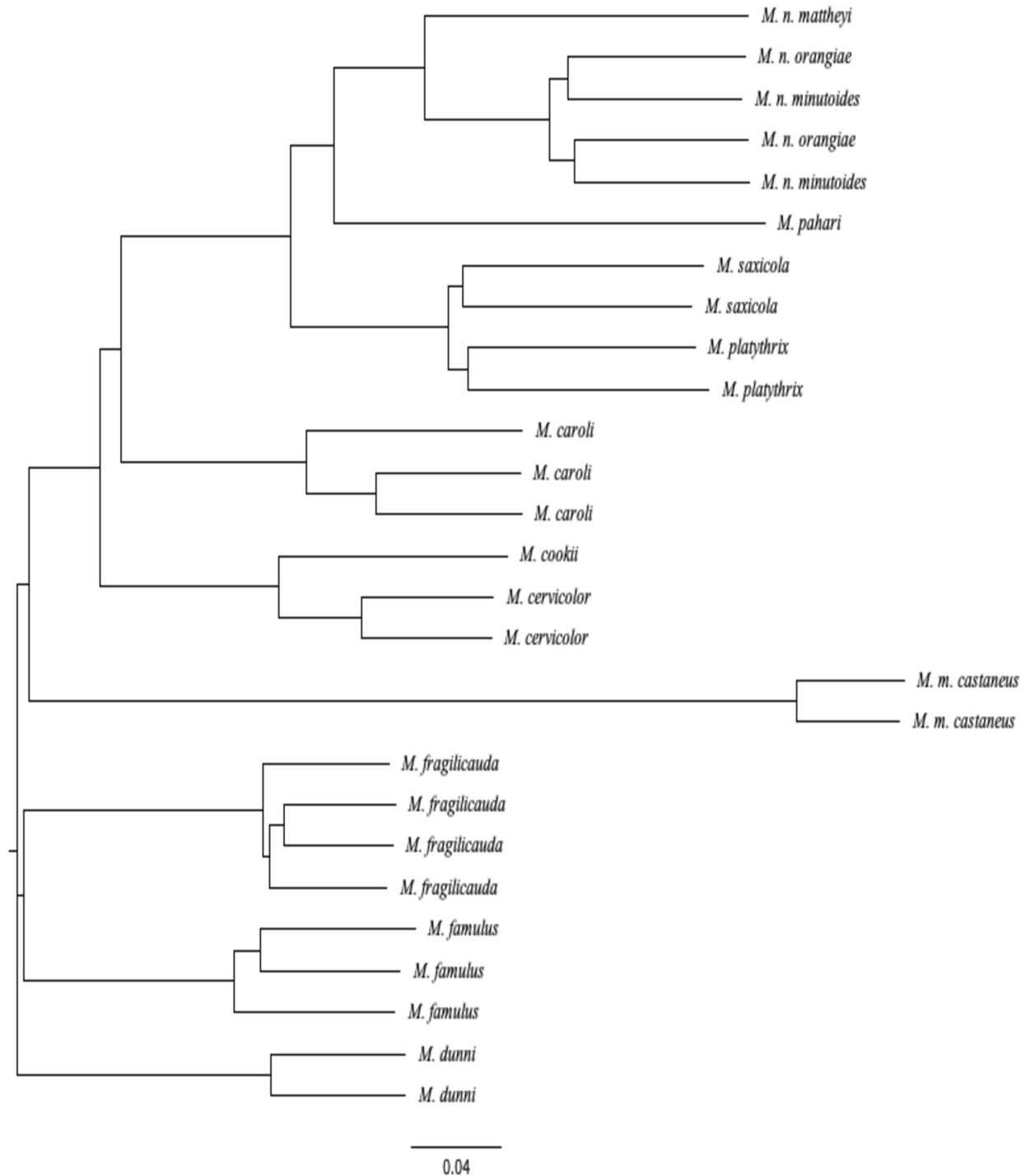797   loci from the intra-genus test set with a heterozygous genotype call. MYD = Millions of years
798   divergence, with respect to the reference *Mus musculus*.
799
800
801

38

802

**Fig 5. SNP distance-based tree of genetic relatedness reflects known taxonomic relationships between Mus species**

SNP distance-based tree of genetic relatedness of samples from the intra-genus test set (n = 27). At 9.5 MYD a pygmy mouse subspecies *M. n. orangiae* has SNP-based genetic distances that reflect greater genetic similarity to another pygmy mouse subspecies *M. n. minutoides* than the replicate MDGA data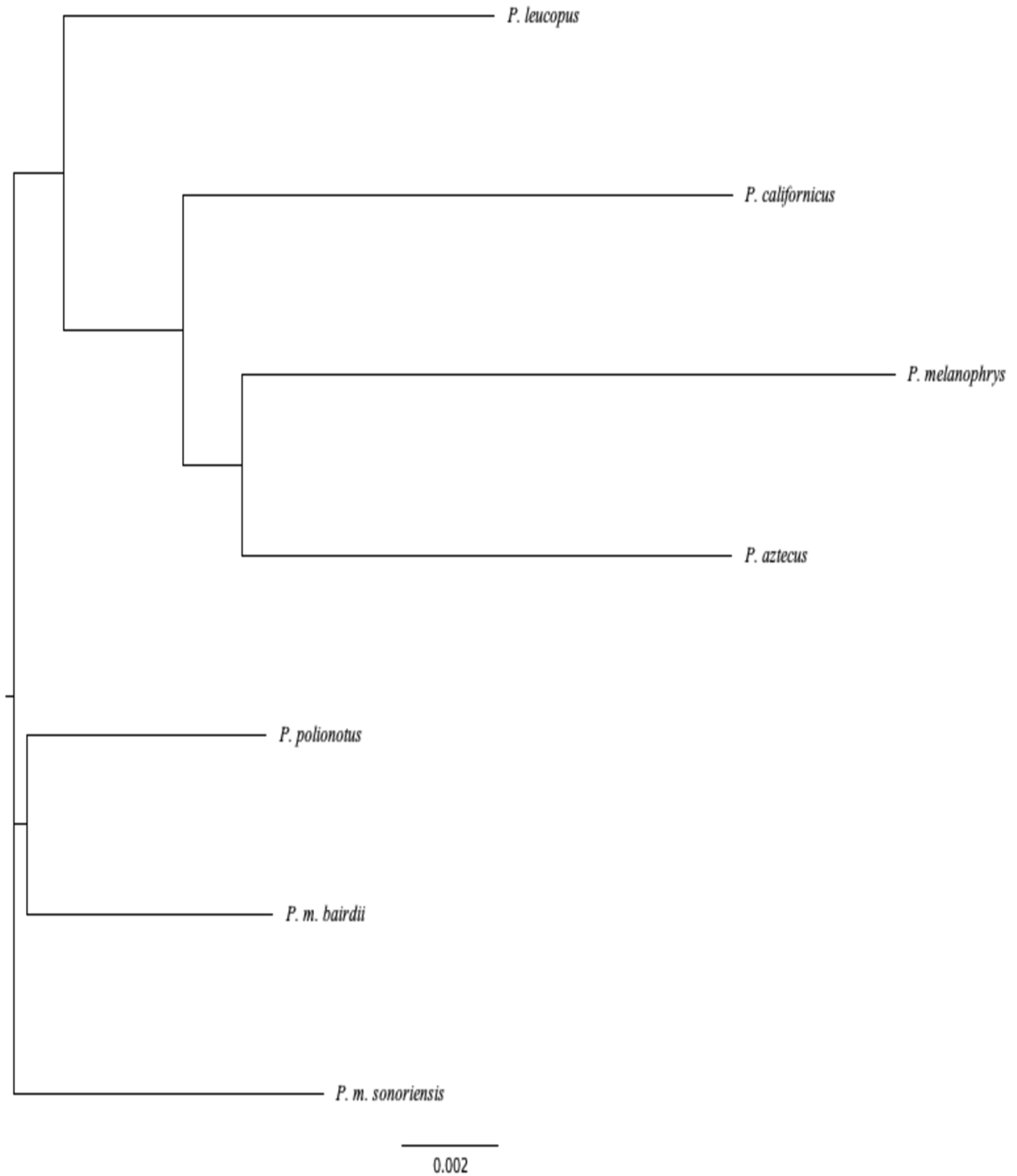 file of the same *M. n. orangiae* sample. MYD = Millions of years divergence, with respect to the reference *Mus musculus*.

810
**Fig 6. SNP distance-based tree of genetic relatedness reflects known taxonomic**
**relationships between Peromyscus species**
Pairwise SNP distance-based tree of genetic relatedness of samples from the intra-genus test set
of Peromyscus species (n=7).

815

816

817

818    **Supporting Information**
819



820
821    **S1 Fig. Abnormalities in two MDGA raw intensity CEL file images.** CEL file raw array
822    intensity images were analyzed for quality control purposes and abnormalities in array images
823    were noted for two CEL files. The two samples were not removed from analysis.
824
825    **S1 Table. Forty-four MDGA data (CEL) files of the present study.**

| CEL File[1] | Sex of Organism | Common Name | Scientific Name | Divergence Time[2] from *Mus musculus* (MYD) |
|---|---|---|---|---|
| SNP_mDIV_A7-7_081308.CEL | Male | House Mouse Reference | *Mus musculus* | 0 |
| SNP_mDIV_D3-639_101509-redo[3] | Female | Southeastern Asian House Mouse | *Mus musculus castaneus* | 0.35 |
| SNP_mDIV_D3-639_91809 | Female | Southeastern Asian House Mouse | *Mus musculus castaneus* | 0.35 |
| SNP_mDIV_D9-647_101509-redo | Male | Earth-Colored Mouse | *Mus dunni* | 6.4 |
| SNP_mDIV_D9-647_91809 | Male | Earth-Colored Mouse | *Mus dunni* | 6.4 |
| SNP_mDIV_D4-640_101509-redo | Male | Servant Mouse | *Mus famulus* | 6.4 |
| SNP_mDIV_D4-640_91809 | Male | Servant Mouse | *Mus famulus* | 6.4 |

*Mouse single nucleotide polymorphic targets for cross hybridization in rodents*

| | | | | |
|---|---|---|---|---|
| SNP_mDIV_D8-474_012209 | Male | Servant Mouse | *Mus famulus* | 6.4 |
| SNP_mDIV_D5-642_101509-redo | Male | Sheath-Tailed Mouse | *Mus fragilicauda* | 6.4 |
| SNP_mDIV_D5-642_91809 | Male | Sheath-Tailed Mouse | *Mus fragilicauda* | 6.4 |
| SNP_mDIV_D6-643_101509-redo | Male | Sheath-Tailed Mouse | *Mus fragilicauda* | 6.4 |
| SNP_mDIV_D6-643_91809 | Male | Sheath-Tailed Mouse | *Mus fragilicauda* | 6.4 |
| SNP_mDIV_D7-644_101509-redo | Male | Ryukyu Mouse | *Mus caroli* | 7.41 |
| SNP_mDIV_D7-644_91809 | Male | Ryukyu Mouse | *Mus caroli* | 7.41 |
| SNP_mDIV_D6-472_012209 | Male | Ryukyu Mouse | *Mus caroli* | 7.41 |
| SNP_mDIV_D8-646_101509-redo | Male | Fawn-Coloured Mouse | *Mus cervicolor* | 7.41 |
| SNP_mDIV_D8-646_91809 | Male | Fawn-Coloured Mouse | *Mus cervicolor* | 7.41 |
| SNP_mDIV_A2-645_102109 | Male | Cook's Mouse | *Mus cookii* | 7.41 |
| SNP_mDIV_A3-648_102109 | Male | Flat-Haired Mouse | *Mus platythrix* | 8.1 |
| SNP_mDIV_A4-649_102109 | Male | Flat-Haired Mouse | *Mus platythrix* | 8.1 |
| SNP_mDIV_A5-650_102109 | Male | Rock-Loving Mouse | *Mus saxicola* | 8.1 |
| SNP_mDIV_A6-651_102109 | Male | Rock-Loving Mouse | *Mus saxicola* | 8.1 |
| SNP_mDIV_D7-473_012209 | Male | Shrew Mouse | *Mus pahari* | 8.29 |
| SNP_mDIV_D11-653_101509-redo | Male | African Pygmy Mouse | *Mus nannomys minutoides* | 9.5 |
| SNP_mDIV_D11-653_91809 | Male | African Pygmy Mouse | *Mus nannomys minutoides* | 9.5 |

### Mouse single nucleotide polymorphic targets for cross hybridization in rodents

| | | | | |
|---|---|---|---|---|
| SNP_mDIV_D10-652_101509-redo | Male | Orange Mouse | *Mus nannomys orangiae* | 9.5 |
| SNP_mDIV_D10-652_91809 | Male | Orange Mouse | *Mus nannomys orangiae* | 9.5 |
| SNP_mDIV_A7-654_102109 | Male | Matthey's Mouse | *Mus nannomys mattheyi* | 9.5 |
| SNP_mDIV_B8-1190_082410 | Male | Wood Mouse | *Apodemus sylvaticus* | 14.5 |
| SNP_mDIV_A9-656_102109 | Male | Sprague Dawley rat | *Rattus norvegicus* | 20.9 |
| SNP_mDIV_A10-657_102109 | Male | Outbred Wistar rat | *Rattus norvegicus* | 20.9 |
| SNP_mDIV_B1-659_102109 | Male | Aztec Mouse | *Peromyscus aztecus* | 32.7 |
| SNP_mDIV_B3-661_102109 | Male | California Mouse | *Peromyscus californicus* | 32.7 |
| SNP_mDIV_B5-663_102109 | Male | North American Deer Mouse | *Peromyscus maniculatus bairdii* | 32.7 |
| SNP_mDIV_B4-662_102109 | Male | Sonoran Deer Mouse | *Peromyscus maniculatus sonoriensis* | 32.7 |
| SNP_mDIV_B2-660_102109 | Male | Plateau Deer Mouse | *Peromyscus melanophrys* | 32.7 |
| SNP_mDIV_B6-664_102109 | Male | Oldfield Mouse | *Peromyscus polionotus* | 32.7 |
| SNP_mDIV_B8-666_102109 | Male | White-Footed Mouse | *Peromyscus leucopus* | 32.7 |
| SNP_mDIV_B9-667_102109 | Male | Squirrel | Sciuridae[4] | 71 |
| DNA3337 | Female | Naked Mole Rat | *Heterocephalus glaber* | 73 |
| DNA3338 | Female | Naked Mole Rat | *Heterocephalus glaber* | 73 |
| DNA3339 | Male | Naked Mole Rat | *Heterocephalus glaber* | 73 |

| | | | | |
|---|---|---|---|---|
| DNA3340 | Male | Naked Mole Rat | *Heterocephalus glaber* | 73 |
| SNP_A2-GES11_4907_AGT-JLP-120115-24-35517 | Male | African Black Rhino | *Diceros bicornis* | 96 |
| SNP_A1-GES11_4902_AGT-JLP-120115-24-35517 | Male | Mountain Tapir | *Tapirus pinchaque* | 96 |

826  [1]MDGA data (CEL) files were downloaded from the Center for Genome Dynamics at the
827  Jackson Laboratory, with the exception of the four *H. glaber* CEL files (generated in-house).
828  [2]Divergence time is given in millions of years from the reference house mouse, *M. musculus*
829  (timetree.org). [3]"redo" files are a technical replicate of the CEL file with the same sample
830  identifier code. Ex: SNP_mDIV_D3-639_101509-redo is a technical replicate of
831  SNP_mDIV_D3-639_91809, where D3-639 is the sample identifier. [4]Only family level
832  information available for CEL file SNP_mDIV_B9-667_102109; Genus and species of sample
833  are unknown.
834
835  **S2 Table. Training set of samples for genotyping algorithm (n=114).**

| 114 training set CEL file names | Sample name | Sex |
|---|---|---|
| SNP_mDIV_B3-387_022709.CEL | 129P1/ReJ | M |
| SNP_mDIV_B4-388_012709.CEL | 129P3/J | M |
| SNP_mDIV_A1-1_081308.CEL | 129S1/SvImJ | M |
| SNP_mDIV_A8-199_091708.CEL | 129S6 | M |
| SNP_mDIV_B5-389_012709.CEL | 129T2/SvEmsJ | M |
| SNP_mDIV_D6-254_111308.CEL | 129X1/SvJ | M |
| SNP_mDIV_A2-2_081308.CEL | A/J | M |
| SNP_mDIV_B7-391_012709.CEL | A/WySnJ | M |
| SNP_mDIV_B4-118_091708.CEL | AEJ/GnLeJ | M |
| SNP_mDIV_B8-392_012709.CEL | AEJ/GnRk | M |
| SNP_mDIV_A3-3_081308.CEL | AKR/J | M |
| SNP_mDIV_A6-119_090908.CEL | ALR/LtJ | M |
| SNP_mDIV_C9-120_090908.CEL | ALS/LtJ | M |
| SNP_mDIV_A4-4_081308.CEL | BALB/cByJ | M |
| SNP_mDIV_D5-253_111308.CEL | BALB/cJ | M |
| SNP_mDIV_B9-393_012709.CEL | BDP/J | M |
| SNP_mDIV_B5-123_091708.CEL | BPH/2J | M |
| SNP_mDIV_B3-316_120908.CEL | BPL/1J | M |
| SNP_mDIV_B6-124_091708.CEL | BPN/3J | M |
| SNP_mDIV_A5-5_081308.CEL | BTBR T<+> Itpr3<tf>-Fbxl3<Ovtm>/J | M |

### Mouse single nucleotide polymorphic targets for cross hybridization in rodents

| | | |
|---|---|---|
| SNP_mDIV_C11-125_090908.CEL | BUB/BnJ | M |
| SNP_mDIV_B10-394_012709.CEL | BXSB/MpJ | M |
| SNP_mDIV_A6-6_081308.CEL | C3H/HeJ | M |
| SNP_mDIV_D1-126_090908.CEL | C3HeB/FeJ | M |
| SNP_mDIV_B8-85_090908.CEL | C57BL/10J | M |
| SNP_mDIV_A5-378_121608.CEL | C57BL/6J | M |
| SNP_mDIV_A1-SNP08_001_103008.CEL | C57BL/6J | F |
| SNP_mDIV_A2-SNP08_001_103008.CEL | C57BL/6J | F |
| SNP_mDIV_A3-SNP08_001_103008.CEL | C57BL/6J | F |
| SNP_mDIV_A4-SNP08_002_103008.CEL | C57BL/6J | M |
| SNP_mDIV_A5-SNP08_002_103008.CEL | C57BL/6J | M |
| SNP_mDIV_A6-SNP08_002_103008.CEL | C57BL/6J | M |
| SNP_mDIV_A7-7_081308.CEL | C57BL/6J | M |
| SNP_mDIV_A9-382_012709.CEL | C57BL/6NCI | M |
| SNP_mDIV_B1-385_012709.CEL | C57BL/6NCI | M |
| SNP_mDIV_A8-381_012709.CEL | C57BL/6Crl | M |
| SNP_mDIV_A10-SNP08_004_103008.CEL | C57BL/6NJ | M |
| SNP_mDIV_A11-SNP08_004_103008.CEL | C57BL/6NJ | M |
| SNP_mDIV_A7-SNP08_003_103008.CEL | C57BL/6NJ | F |
| SNP_mDIV_A8-SNP08_003_103008.CEL | C57BL/6NJ | F |
| SNP_mDIV_A9-SNP08_003_103008.CEL | C57BL/6NJ | F |
| SNP_mDIV_B1-SNP08_004_103008_4.CEL | C57BL/6NJ | M |
| SNP_mDIV_A10-383_012709.CEL | C57BL/6Tc | M |
| SNP_mDIV_A11-384_012709.CEL | C57BL/6Tc | M |
| SNP_mDIV_B9-86_090908.CEL | Wrong sample name (not C57BLKS/J, close to C57L/J) | M |
| SNP_mDIV_D2-SNP09_024_022709.CEL | C57BLKS/J | M |
| SNP_mDIV_B10-87_090908.CEL | C57BR/cdJ | M |
| SNP_mDIV_B11-88_090908.CEL | C57L/J | M |
| SNP_mDIV_C1-89_090908.CEL | C58/J | M |
| SNP_mDIV_B4-15_081308.CEL | CBA/CaJ | M |

45

### *Mouse single nucleotide polymorphic targets for cross hybridization in rodents*

| | | |
|---|---|---|
| SNP_mDIV_D8-256_111308.CEL | CBA/J | M |
| SNP_mDIV_D2-128_090908.CEL | CE/J | M |
| SNP_mDIV_D3-129_090908.CEL | CHMU/LeJ | M |
| SNP_mDIV_B7-18_081308.CEL | DBA/1J | M |
| SNP_mDIV_C3-398_012709.CEL | DBA/1LacJ | M |
| SNP_mDIV_C4-399_012709.CEL | DBA/2DeJ | F |
| SNP_mDIV_C5-400_012709.CEL | DBA/2HaSmnJ | M |
| SNP_mDIV_B8-19_081308.CEL | DBA/2J | M |
| SNP_mDIV_A8-56_082108.CEL | DDK/Pas | F |
| SNP_mDIV_B9-20_081308.CEL | DDY/JclSidSeyFrkJ | M |
| SNP_mDIV_D4-130_090908.CEL | DLS/LeJ | M |
| SNP_mDIV_D5-131_090908.CEL | EL/SuzSeyFrkJ | M |
| SNP_mDIV_B10-21_081308.CEL | FVB/NJ | M |
| SNP_mDIV_B8-132_091708.CEL | HPG/BmJ | M |
| SNP_mDIV_B2-90_091708.CEL | I/LnJ | M |
| SNP_mDIV_A11-431_022709.CEL | WSP2 | F |
| SNP_mDIV_A9-429_022709.CEL | WSR2 | M |
| SNP_mDIV_A8-427_022709.CEL | COLD2 | M |
| SNP_mDIV_A6-424_022709.CEL | HOT1 | M |
| SNP_mDIV_A7-425_022709.CEL | HOT2 | M |
| SNP_mDIV_B2-433_022709.CEL | ILS/IbgTejJ | M |
| SNP_mDIV_B1-432_022709.CEL | ISS/IbgTejJ | M |
| SNP_mDIV_C2-91_090908.CEL | JE/LeJ | M |
| SNP_mDIV_B11-22_081308.CEL | KK/HlJ | M |
| SNP_mDIV_C3-92_090908.CEL | LG/J | M |
| SNP_mDIV_C4-93_090908.CEL | LP/J | M |
| SNP_mDIV_C6-401_012709.CEL | LT/SvEiJ | M |
| SNP_mDIV_D7-134_090908.CEL | MRL/MpJ | M |
| SNP_mDIV_C6-30_081308.CEL | NOD/ShiLtJ | M |
| SNP_mDIV_C9-404_012709.CEL | NOD/ShiLtJ | M |
| SNP_mDIV_A2-48_082108.CEL | NON/ShiLtJ | M |
| SNP_mDIV_C11-406_012709.CEL | NONcNZO5/LtJ | M |
| SNP_mDIV_A3-49_082108.CEL | NOR/LtJ | M |
| SNP_mDIV_D9-136_090908.CEL | NU/J | M |
| SNP_mDIV_A1-50_091708.CEL | NZB/BlNJ | M |
| SNP_mDIV_C5-94_090908.CEL | NZL/LtJ | M |
| SNP_mDIV_D10-137_090908.CEL | NZM2410/J | M |
| SNP_mDIV_C7-31_081308.CEL | NZO/HlLtJ | M |
| SNP_mDIV_C8-32_081308.CEL | NZW/LacJ | M |
| SNP_mDIV_B9-138_091708.CEL | P/J | M |

| | | |
|---|---|---|
| SNP_mDIV_C6-95_090908.CEL | PL/J | M |
| SNP_mDIV_A1-147_111308.CEL | SI/Col Tyrp1 Dnahc11/J | M |
| SNP_mDIV_D11-139_090908.CEL | PN/nBSwUmabJ | M |
| SNP_mDIV_B11-141_091708.CEL | RF/J | M |
| SNP_mDIV_B9-142_103008_3.CEL | RHJ/LeJ | M |
| SNP_mDIV_C8-97_090908.CEL | RIIIS/J | M |
| SNP_mDIV_B10-143_103008_3.CEL | RSV/LeJ | M |
| SNP_mDIV_D9-144_103008_3.CEL | SB/LeJ | M |
| SNP_mDIV_D10-145_103008_3.CEL | SEA/GnJ | M |
| SNP_mDIV_D2-408_012709.CEL | SEC/1GnLeJ | M |
| SNP_mDIV_D3-409_012709.CEL | SEC/1ReJ | M |
| SNP_mDIV_D11-146_103008_3.CEL | SH1/LeJ | M |
| SNP_mDIV_D4-410_012709.CEL | SJL/Bm | M |
| SNP_mDIV_D1-36_081308.CEL | SJL/J | M |
| SNP_mDIV_A2-148_111308.CEL | SM/J | M |
| SNP_mDIV_A4-150_111308_2.CEL | SSL/LeJ | M |
| SNP_mDIV_D5-411_012709.CEL | ST/bJ | M |
| SNP_mDIV_D6-412_012709.CEL | STX/Le | M |
| SNP_mDIV_A5-151_111308.CEL | SWR/J | M |
| SNP_mDIV_A6-152_111308.CEL | TALLYHO/JngJ | M |
| SNP_mDIV_A7-153_111308.CEL | TKDU/DnJ | M |
| SNP_mDIV_A8-154_111308.CEL | TSJ/LeJ | M |
| SNP_mDIV_D7-413_012709.CEL | YBR/EiJ | M |
| SNP_mDIV_A9-155_111308.CEL | ZRDCT Rax<ey1>/ChUmdJ | M |

836

837 CEL files of 114 classically inbred laboratory mouse strains were downloaded from the Jackson
838 Laboratory Center for Genome Dynamics for genotyping algorithm training.

839
840
841
842 **S3 Table. Genotype summary of 44 study samples genotyped at 493,290 single nucleotide**
843 **polymorphic loci located across the genome of *Mus musculus*.**
844 Genotyping summary results for all 44 Mouse Diversity Genotyping Array data files.
845 **(Too large to display. Please see separate PDF file)**
846
847
848
849
850
851
852
853
854

855 **S4 Table. Study species evaluated with publicly available nuclear genome sequence**
856 **information.**

| Sample Name | Scientific Name | Newest Assembly |
|---|---|---|
| House Mouse | *Mus musculus* | GRCm38.p6 |
| Ryukyu Mouse | *Mus caroli* | CAROLI_EIJ_v1.1 |
| Gairdner's Shrewmouse | *Mus pahari* | PAHARI_EIJ_v1.1 |
| Sprague Dawley Rat | *Rattus norvegicus* | Rnor_6.0 |
| North American Deer Mouse | *Peromyscus maniculatus* | Pman_1.0 |

857 Genomes accessed through the NCBI Genomes FTP site of samples under study
858 (ftp.ncbi.nlm.nih.gov/genomes/).
859
860