1

2

3 **Signatures of rapid evolution in urban and rural transcriptomes of white-footed mice**

4 **(*Peromyscus leucopus*) in the New York metropolitan area**

5

6 Stephen E. Harris[1], Jason Munshi-South[2*], Craig Obergfell[3], & Rachel O'Neill[3]

7

8 [1]Program in Ecology, Evolutionary Biology, & Behavior, The Graduate Center, City University

9 of New York (CUNY), New York, NY, USA

10

11 [2] Louis Calder Center-Biological Field Station & Department of Biology, Fordham University,

12 Armonk, NY, USA

13

14 [3] Molecular & Cell Biology, University of Connecticut, Storrs, CT, USA

15

16 *Corresponding author:* Jason Munshi-South

17 E-mail: jason@NYCevolution.org

18 Tel: 1-646-318-6353

19 Fax: 1-646-660-6250

20

21

**Abstract**

Urbanization is a major cause of ecological degradation around the world, and human settlement in large cities is accelerating. New York City (NYC) is one of the oldest and most urbanized cities in North America, but still maintains 20% vegetation cover and substantial populations of some native wildlife. The white-footed mouse, *Peromyscus leucopus*, is a common resident of NYC's forest fragments and an emerging model system for examining the evolutionary consequences of urbanization. In this study, we developed transcriptomic resources for urban *P. leucopus* to examine evolutionary changes in protein-coding regions for an exemplar 'urban adapter'. We used Roche 454 GS FLX+ high throughput sequencing to derive transcriptomes from multiple tissues from individuals across both urban and rural populations. From these data, we identified 31,015 SNPs and several candidate genes potentially experiencing positive selection in urban populations of *P. leucopus*. These candidate genes are involved in xenobiotic metabolism, innate immune response, demethylation activity, and other important biological phenomena in novel urban environments. This study is one of the first to report candidate genes exhibiting signatures of directional selection in divergent urban ecosystems.

37

38

## Introduction

Urbanization dramatically alters natural habitats [1], and its speed and intensity will increase as over two-thirds of the world's human population is predicted to live in urban areas by 2050 [2]. Understanding how natural populations adapt to ecologically divergent urban habitats is thus an important and immediate goal for urban ecologists and evolutionary biologists. Few ecological and evolutionary studies are conducted in urban environments [3], but recent attitude shifts and technological advancements have removed many of the obstacles to working on urban wildlife. Multiple studies have demonstrated that urban areas are biologically diverse, productive, and viable [4], and the development of next generation sequencing (NGS) has facilitated the generation of genomic resources for uncharacterized species in natural environments [5–7]. Understanding the genetic basis of adaptation in successful urban species will aid in future conservation efforts and provide insights into the effects of other anthropogenic factors, such as global climate change and evolutionary trajectories in human-dominated environments [4,8,9].

Cities typically experience a substantial decrease in biodiversity of many taxonomic groups as urban 'avoiders' disappear, accompanied by a rise in urban 'exploiters' that are primarily non-native human commensals such as pigeons or rats. Urban 'adapters' are native species that favor disturbed edge habitats such as urban forest fragments, relying on a combination of wild-growing and human-derived resources [10–12]. This last group is of primary interest for examining genetic signatures of recent evolutionary change in novel urban environments. Severe habitat fragmentation is one of the primary impacts of urbanization and often leads to genetic differentiation between populations [1,13,14]. Introductions of new predators and competitors alter ecological interactions [15], and new or more abundant parasites or pathogens influence immune system processes [16]. Air, water, and soil pollution typically increase in local urban ecosystems, and selection may favor previously-rare genetic variants that more efficiently process these toxins [17–19]. Recent studies provide some evidence of local

64    adaptation and rapid evolution in urban patches.  Using a candidate gene approach, Mueller et al.

65    [20] found consistent genetic divergence between behavioral genes for circadian behavior, harm

66    avoidance, migratory behavior and exploratory behavior in multiple urban-rural population pairs

67    of the common blackbird, *Turdus merula*.  Examining phenotypes, Brady [21] found rapid

68    adaptation to roadside breeding pond conditions in the salamander, *Ambystoma maculatum*, and

69    Cheptou et al. [22] reported a heritable increase in production of non-dispersing seeds in the

70    weed, *Crepis sancta*, over 5-12 generations in fragmented urban tree pits. The genetic

71    architecture of the phenotypes under selection has not been described for either of these urban

72    'adapters', but outlier scans of transcriptome sequence datasets are one promising approach [23].

73          *Peromyscus* spp. are an emerging model system for examining evolution in wild

74    populations [24-26], but large-scale genomic resources are not yet widely available.  The genus

75    contains the most widespread and abundant small mammals in North America, and *Peromyscus*

76    research on population ecology, adaptation, aging, and disease has a long, productive history

77    [27–31]. An increasing number of studies have demonstrated that *Peromyscus* spp. rapidly (i.e.

78    in several hundreds to thousands of generations) adapt to divergent environments. These

79    examples include adaptation to hypoxia in high altitude environments [26] and adaptive variation

80    in pelage color on light-colored soil substrates [25, 32, 33].  Presently, *P. leucopus* is the sole

81    *Peromyscus* spp. in New York City (J. Munshi-South, unpublished data) and searches of the

82    Mammal Networked Information System (MANIS) database indicate that *P. maniculatus* has not

83    occurred in NYC for several decades.  In NYC, *P. leucopus* occupies most small patches of

84    secondary forest, shrublands, and meadows within NYC parklands [33,34]. The smallest patches

85    in NYC often contain the highest population densities of white-footed mice [35], most likely due

86    to ecological release and obstacles to dispersal [36,37].  Consistently elevated population density

87    in urban patches compared to surrounding rural populations is predicted to result in density-

88     dependent selective pressures on traits related to immunology, intraspecific competition, and

89     male-male competition for mating opportunities, among others [38,39].

90     White-footed mouse populations in NYC exhibit high levels of heterozygosity and allelic

91     diversity at neutral loci within populations, but genetic differentiation and low migration rates

92     between populations [40,41]. This genetic structure contrasts with weak differentiation reported

93     for *Peromyscus* spp. at regional scales [42], or even between populations isolated on different

94     islands for thousands of generations [43,44]. High genetic diversity within and low to

95     nonexistent migration between most NYC populations suggests that selection can operate

96     efficiently within these geographically isolated populations, either on standing genetic variation

97     or *de novo* mutations. In this study we take steps to develop *P. leucopus* as a genomic model for

98     adaptive change in urban environments.

99     Pooling mRNA from multiple individuals is an effective approach to transcriptome

100     sequencing that avoids the prohibitive cost of sequencing individual genomes [45,46]. While

101     pooling results in the loss of genetic information from individuals, the ability to identify SNPs in

102     a population increases due to the inclusion of multiple individuals in the pool [47]. By analyzing

103     SNPs within thousands of transcripts, it is feasible to identify candidate genes underlying rapid

104     divergence of populations in novel environments [5, 47-49]. Certain statistical approaches, such

105     as the ratio between non-synonymous and synonymous ($p_N/p_S$) substitutions, can be applied to

106     pooled transcriptome data to identify potential signatures of selection between isolated

107     populations [23,50,51]. If positive selection is acting on a codon, then non-synonymous

108     mutations should be more common than under neutral expectations [52,53].

109     Here, we describe the results of *de novo* transcriptome sequencing, annotation, SNP

110     discovery, and outlier scans for selection among urban and rural white-footed mouse

111     populations. We used the 454 GS FLX+ system to sequence cDNA libraries generated from

112     pooled mRNA samples from multiple tissues and populations. Several *de novo* transcriptome

113    assembly programs were used and the contribution of specific tissue types to the transcriptome

114    assembly was examined. We then identified coding region SNPs between urban and rural

115    populations, and scanned this dataset for signatures of positive selection using $p_N/p_S$ between

116    populations and McDonald-Kreitman tests between multiple species.  We report several

117    candidate genes potentially experiencing directional selection in urban environments, and

118    provide annotated transcriptome datasets for future evolutionary studies of an emerging model

119    system.

120

121    **Results**

122    *Sequencing and comparison of assembly methods*

123    454 Sequencing of four full plates of *P. leucopus* cDNA libraries made from liver, brain, and

124    gonad tissue produced 3,052,640 individual reads with an average length of 309 ± 122 bp

125    (median = 341, Interquartile Range (IQR) = 188 bp).  While the initial Newbler genomic

126    assembly and Cap3 assembly produced more contigs, the mean length and N50 for both sets of

127    contigs were lower than the Newbler cDNA assembly (Table 1).  The Cap3 assembly and the

128    genomic assembly included a much higher proportion of shorter contigs than the cDNA

129    assembly (Fig. 1). Coverage was calculated for all three assemblies, and all had similar median

130    read coverage per contig (Newbler Genomic, median = 4.7 reads, IQR = 4.6; Newbler cDNA,

131    median = 4.9 reads, IQR = 4.1; Cap3, median = 5.0 reads, IQR = 7.0, Fig. S1).

132         After filtering BLASTN searches against *Mus musculus* and *Rattus norvegicus* cDNA

133    libraries, there was an average for all assemblies of 13,443 hits to known genes. The Cap3

134    assembly and Newbler genomic assembly produced the most hits, but the average alignment

135    length was longest for the Newbler cDNA assembly (Table 2).  Of the total number of contigs

136    for each assembly, the Newbler cDNA assembly had the highest proportion (47%) of 'Gene

137    Candidates' followed by the Cap3 assembly (42%) and the Newbler genomic assembly (41%).

138    Assessments important for looking at $p_N/p_S$ (longest average length of contigs, largest N50

139    value) and for reducing false positives (largest proportion of hits to one gene with known

140    function) supported the assertion that Newbler's cDNA assembly produced the best quality

141    reference transcriptome, and all further analyses used this assembly.

142

143    *cDNA transcriptome assembly*

144    The final reference *P. leucopus* Newbler cDNA assembly produced 17,371 contigs with an

145    average length of 613 ± 507 bp. These contigs were assembled into 15,004 isotigs and 12,464

146    isogroups with a combined length of 13,421,361 bp. Isotigs were constructed from an average of

147    1.6 contigs and isogroups from an average of 1.2 isotigs. The contribution of sequence reads

148    from individual tissues to the final reference transcriptome was not equal. Liver and brain cDNA

149    libraries produced higher numbers of total reads and a greater proportion of assembled reads

150    compared to ovary and testis libraries. The average read coverage of contigs for each tissue type

151    varied, but coverage from liver sequences was highest with nearly 2X more compared to brain,

152    testes, and ovaries (Table S1). Among all contigs assembled, 70% contained reads from plate 1

153    (normalized), 57% contained reads from plate 2 (non-normalized), 79% contained reads from

154    plate 3 (non-normalized), and 89% contained reads from plate 4 (non-normalized). Comparison

155    of normalized (Plate 1) and non-normalized (Plates 2-4) cDNA libraries indicated that non-

156    normalization produced nearly twice as many total sequencing reads as compared to

157    normalization, and non-normalized plates were able to sequence rare transcripts at a similar rate

158    compared to the normalized plate (Table S1).

159

160    *Mouse and rat genome comparisons*

161    Assembled mRNA transcripts from *P. leucopus* successfully mapped to both *Mus* and *Rattus*

162    reference genomes and were distributed across all chromosomes for both references (Fig. 2).

163  There were 9,418 best BLAT hits between *P. leucopus* contigs and known *Mus* genes and 8,786

164  best hits with *Rattus* genes.  The latest cDNA references include 35,900 genes for *Mus* (mm10)

165  and 29,261 genes for *Rattus* (rn5), suggesting that full or partial coding sequence from

166  approximately one-third to one-fourth of the *P. leucopus* transcriptome was sequenced.  Given

167  that many of the 15,000 contigs we assembled from our raw sequencing data may represent

168  *Peromyscus*-specific genes not found in model rodent databases, the real proportion of the

169  sequenced transcriptome may be much higher.

170

171  *Functional annotation*

172  Among isotigs from the reference *P. leucopus* transcriptome, 11,355 (75.7%) had BLASTX hits

173  to known genes, and 6,385 (42.6%) mapped to proteins and were annotated with known

174  biological functions (GO terms) from protein databases.  Top sources for these annotations were

175  the model rodents *Cricetulus griseus* (3,686 BLASTX hits, 24.5%), *Mus musculus* (2,914

176  BLASTX hits, 19.4%), and *Rattus norvegicus* (1,671 BLASTX hits, 11.1%, Fig. S2).  For cDNA

177  assemblies of individual organs, the ovary transcriptome (1,589 isotigs) had the highest

178  proportion (73.9%) of assembled contigs with GO annotations (Fig. 3).  Liver (6,240 isotigs) and

179  testes (5,728 isotigs) produced the largest number of total assembled contigs with similar

180  proportions having GO term annotations (65.6% and 64.6%, respectively).  The brain

181  transcriptome (2,613 isotigs) included a lower number of assembled contigs and percent GO

182  annotation (56.8%; Fig. 3).

183  One-tailed Fisher's Exact tests (False Discovery Rate (FDR) $\leq$ 0.05) indicated that liver

184  had the most GO terms that were significantly over-represented compared to the other tissue

185  types (Fig. 4).  1,320 annotations in liver were overrepresented in both liver to brain and liver to

186  gonad comparisons, and there were 69 overlapping annotations in brain to gonad and brain to

187  liver comparisons (Fig. 4).  Gonads had the least number of annotations (five) commonly

188    overrepresented in both brain and liver comparisons (Fig. 4). When reduced to their most-

189    specific terms, pairwise comparisons detected 64 over-represented GO annotations for liver

190    when compared to both of the other tissues, 20 for brain, and five for gonads (Table 3). Over-

191    represented GO terms in liver were related to metabolic processes including ATP binding, GTP

192    binding, NADH dehydrogenase, and electron carrier activity. Over-represented GO terms in

193    brain included regulation of behavior, actin binding, ion channel activity, motor activity, and

194    calcium ion binding. Significantly different gonad annotations were related to reproduction,

195    cilium (for sperm locomotion), the cell cycle, transcription regulation, and epigenetic regulation

196    of gene expression (See Table 3 and Table S2 for full list of overrepresented GO annotations in

197    all pairwise comparisons).

198

199    *SNP calling and calculation of $p_N/p_S$*

200    After mapping the reads used in the assembly back to the Newbler cDNA reference

201    transcriptome, 31,015 SNPs were called in 7,625 isotigs. The distribution of SNPs per isotig

202    ranged from $1 - 78$ (mean = $4 \pm 5.4$; median = 2). ORFs were identified in 11,704 isotigs

203    comprising 5.6 Mb of sequence, and 2,655 putative ORFs contained 4,893 SNPs. Of these

204    SNPs, 1,795 (36.6%) were classified as non-synonymous and 3,098 (63.3%) were classified as

205    synonymous. Aligned ORFs were used to calculate $p_N/p_S$ between each pair of populations.

206    The majority of the ORFs did not exhibit statistical signatures of positive selection (overall mean

207    $\pm$ SE $p_N/p_S = 0.28 \pm 0.56$). For the 2,307 pairs of homologous cDNA sequences between

208    populations that contained predicted ORFs, did not contain in-frame stop codons, and had greater

209    than or equal to three SNPs, $p_N/p_S$ values for 11 (0.5%) contigs exceeded 1.0 (Table 4, Fig. 5).

210    The proportion of genes with $p_N/p_S > 1.0$ is comparable to similar studies; Sun et al. [23] found

211    that 0.4% of genes in their *Pomacea canaliculata* dataset were positively selected, Renaut et al.

212    [54] reported 0.5% in *Coregonus clupeaformis*, and Wang et al [55] reported 1.8% in *Bemisia*

213     *tabaci*. Four contigs (0.2%) exhibited $p_N/p_S$ values > 1.0 in urban to urban comparisons and 7

214     contigs (0.3%) in urban to rural population comparisons. 42 (1.8%) contigs were found with $p_N/$

215     $p_S$ between 0.5 and 1 (Table S3, Fig. 5); $p_N/p_S$ greater than 0.5 is a less conservative filter for

216     detecting positive selection, especially when using truncated ORFs [56,57].

217        Different genes showed strong ($p_N/p_S > 1$) signatures of selection when urban

218     populations were compared to other urban populations than when urban and rural populations

219     were compared. Candidate genes identified from the ORF pairs (i.e. $p_N/p_S > 1$) in urban to rural

220     comparisons were related to metabolic processes (including xenobiotic metabolism),

221     reproduction, and demethylation (Table 4). Three genes were involved in metabolic processes:

222     *cytochrome P450 2A15* (xenobiotic metabolism, HP_contig01783, $p_N/p_S = 1.89$), *camello-like 1*

223     (HP_contig00870, $p_N/p_S = 1.74$), and *aldo-keto reductase family 1, member C12* (Xenobiotic

224     metabolism, HP_contig01919, $p_N/p_S = 1.18$). Our analysis also identified a reproductive gene,

225     *histone H1-like protein in spermatids 1* (HP_contig02656, $p_N/p_S = 1.07$) that is involved in

226     transcriptional regulation during spermatogenesis. The gene *phd finger protein 8*

227     (HP_contig01778, $p_N/p_S = 1.12$), codes for a demethylase that removes methyl groups from

228     histones.

229        Candidate genes in urban to urban population comparisons were primarily involved in

230     immune system processes. One of these genes is involved in regulating the innate immune

231     response, *alpha-1-acid glycoprotein 1* (CP_contig00748, $p_N/p_S = 1.97$), by modulating innate

232     immune response while circulating in the blood. The other immune system genes are involved in

233     blood coagulation and inflammation, *serine protease inhibitor a3c* (CP_contig00256, $p_N/p_S = $

234     1.76) and *fibrinogen alpha chain* (CP_contig00473, $p_N/p_S = 1.23$). We also identified *solute*

235     *carrier organic anion transporter family member 1A5* (CP_contig01204, $p_N/p_S = 1.55$), a gene

236     that facilitates intestinal absorption of bile acids and renal uptake and excretion of uremic toxins.

9

237

238        For the 22 contigs with $p_N/p_S$ between 0.5 and 1 for urban to rural comparisons, genes

239    are primarily involved in the innate immune response, metabolic processes, and methylation

240    activity, and some of these genes are involved in the same biological pathways as genes listed

241    above that exhibited $p_N/p_S > 1$ (Tables 4, S3).  For the 20 contigs with $p_N/p_S$ between 0.5 and 1

242    for urban pairwise comparisons, genes are primarily involved with the innate immune response,

243    metabolic processes (including xenobiotic), and reproductive processes.

244        Candidate genes were scanned for evidence of recombination using a phylogenetic

245    framework.  The Genetic Algorithm Recombination Detection (GARD) analysis identified no

246    evidence of recombination in any potential candidate genes.  Would-be breakpoints were

247    identified in the genes *Translocation protein SEC62, Histone H1-like protein in spermatids 1,*

248    *Aldo-keto reductase family 1 member C12, Fibrinogen alpha chain, Solute carrier organic anion*

249    *transporter family member 1A5*, and *Serine protease inhibitor a3c*, but Kishino-Hasegawa

250    testing implemented in the Data Monkey web server found the signal most likely resulted from

251    evolutionary rate variation as opposed to recombination.

252        McDonald-Kreitman tests were then performed to examine potentially adaptive evolution

253    between species in all the identified candidate genes.  *P. leucopus* was compared to *R.*

254    *norvegicus*, and *C. griseus* when *Rattus* sequences were not available.  This approach minimized

255    the number of multiple mutations at individual sites, but results were very similar when the

256    orthologous candidate genes were compared to any rodent with available orthologous gene

257    sequence.  Excess adaptive change (diversifying selection between species) was not identified in

258    any of the candidate genes.  For four genes, *39S ribosomal protein L51, PHD finger protein 8,*

259    *Cytochrome P450 2A15,* and *Solute carrier organic anion transporter 1A*, the ratio of non-

260    synonymous to synonymous polymorphisms within *P. leucopus* was significantly higher than the

261    ratio for divergent genetic changes between species (Table 5). While there were more non-

262     synonymous polymorphisms than synonymous polymorphisms in the remaining seven genes,

263     results were not significantly different from expected neutral evolution.

264

265     **Discussion**

266     *De novo transcriptome assembly and characterization*

267     Compared to other NGS technologies, 454 transcriptome sequencing provides longer read

268     lengths ideal for *de novo* assembly [58] and is especially useful for organisms without extensive

269     genomic resources like *P. leucopus* [51,54,59–61]. We compared the relative merits of two

270     established long-read assembly programs, CAP3 and Newbler, for assembling our transcriptomes

271     [60, 61]. Despite the substantially fewer megabases per run generated by 454 FLX+ compared

272     to Illumina or SOLiD sequencing [62], we still ran into computational limitations during

273     assembly when using options for cDNA sequence. Similar to Cahais et al. [63], we had the most

274     success after compressing the raw reads into a smaller number of partially assembled sequences

275     using a genome assembler followed by another assembly method better suited for transcriptome

276     data. While the CAP3 assembly produced more contigs, the Newbler v. 2.5.3 transcriptome

277     assembly performed better based on assessments useful for downstream population genomic

278     analyses (e.g. number of long contigs and average contig length). Newbler performed well at

279     assembling full-length cDNA contigs, and our results are in line with Mundry et al.'s [64]

280     findings that Newbler outperformed other assembly programs in simulated experiments. The

281     N50 value reported here is comparable to *de novo* Newbler cDNA assemblies for other

282     organisms: N50 = 1,735 bp in *Oncopeltus fasciatus*, Ewen-Campen et al. [65]; N50 = 1,333 bp

283     in *Silene vulgaris*, Sloan et al. [51]; N50 = 1,588 bp in *Spalax galili*, Malik et al. [66]; and N50 =

284     854 bp in *Arctocephalus gazella,* Hoffman & Nichols [67].

11

285      We sequenced samples using normalized and non-normalized cDNA pools and examined

286    the influence each protocol had on gene discovery. Following sequencing of the first normalized

287    plate, we used a new protocol from Roche that excluded normalization of libraries. Surprisingly,

288    we found that normalization did not necessarily improve the number of uniquely assembled

289    contigs.  Theoretically, normalization reduces the sequencing of overly abundant transcripts and

290    increases the discovery of rare sequences [68,69], but normalization does not disproportionately

291    influence gene discovery when enough sequencing coverage is achieved [70]. We found that

292    read coverage per transcript increased for our non-normalized plates compared to the normalized

293    pilot plate.  However, Ekblom et al. [71] suggest that differences in technologies and sequencing

294    effort may ultimately affect comparisons between normalized and non-normalized cDNA

295    libraries, and any differences we identify may be due to different protocols used to extract RNA

296    and prepare pooled libraries.

297

298    *Mapping to rodent genomes*

299    The mammalian laboratory models *Mus* and *Rattus* have extensively annotated genomes that

300    provide a good substitute reference for other rodent sequencing projects. The New World

301    *Peromyscus* and Old World *Mus* and *Rattus* lineages last shared a common ancestor ~25 million

302    years ago [72].  Deep divergence and high rates of chromosome evolution across these lineages

303    [73] may have affected the percentage of identified homologous gene transcripts.  Ramsdell et al.

304    [74] found the *Peromyscus* genome to be more similar to *Rattus* than *Mus* due to an enhanced

305    level of genome rearrangement in *Mus* compared to ancestral muroids.  Our results support these

306    findings given that most *Peromyscus* transcripts mapped to different chromosomes (96.1%)

307    between *Mus* and *Rattus*.  Our homologous gene matches between *Peromyscus* and *Rattus* also

308    represented a higher proportion (30.1%) of total *Rattus* genes than homologous gene matches

309    between *Peromyscus* and *Mus* (25.7%).  Non-homologous hits and mapping differences between

12

310    reference genomes may also be due to highly variable or alternatively spliced transcripts,

311    contamination by genomic DNA, or inclusion of low-quality data [75], although our assembly

312    methods included measures to limit the influence of these artifacts.

313

314    *Functional annotation and tissue comparisons*

315    Over 75% of our assembled contigs produced significant BLASTX hits to known genes in

316    NCBI's nonredundant (nr) protein database.  This rate of annotation is similar to studies on other

317    non-model species with genomic information available from closely-related model organisms,

318    e.g. 66% in the rodent *Ctenomys sociabilis* [76] and 79.7% in the plant *Silene vulgaris* [50].

319    These rates are much higher than some other organisms with few model relatives, such as

320    19.58% in a bat, *Artibeus jamaicensis*, [77], 18% in a butterfly, *Melitaea cinxia*, [59], and 29.2%

321    in the gastropod, *Pomacea canaliculata*, [23].  Phylogenetic analyses support *Peromyscus* spp.

322    and *Cricetulus* spp. as members of a monophyletic clade that diverged separately from *Mus* and

323    *Rattus* [72], and *C. griseus* represented the highest proportion of BLASTX top-hits (Figure S2,

324    Supplementary Material).  Laboratory use of *C. griseus* is not as prevalent as *Mus* or *Rattus*, but

325    Chinese hamster ovary (CHO) cell lines are commonly used *in vitro* to produce

326    biopharmaceuticals [78], and a draft genome has also been sequenced [79].  Research on protein

327    pathways and interactions within CHO cell lines provides a future resource for investigating

328    functional consequences of divergent genes between urban and rural populations of *P. leucopus*.

329        Transcriptome studies in model rodents provide useful context for understanding how

330    much of each tissue-specific transcriptome we sequenced in this study.  Yang et al. [80] used

331    microarray analysis to identify 12,845 active genes in *Mus* liver, and RNA-Seq using an Illumina

332    HiSeq 2000 on *Rattus* liver identified 7,514 known genes [81].  Our gene discovery was between

333    40-60% of these previously reported liver transcriptomes. In brain tissue, 4,508 genes were

334    identified in *Mus* by Yang et al. [80], and Chrast et al. [82] report ~4,000 genes identified in *Mus*

335   brain tissue.  The 2,610 gene annotations from our brain cDNA libraries represent between 60-

336   65% of the full *P. leucopus* brain transcriptome. Microarray analysis of testis RNA identified up

337   to 13,812 known genes [83] in *Mus*, and 454 sequencing of cDNA libraries from *C. griseus*

338   identified 13,187 annotations in ovary [76]. UniGene [84] includes 8,946 genes for *Mus* testis,

339   5,285 for *Mus* ovaries, 4,355 for *Rattus* testis, and 5,093 for *Rattus* ovaries.  The only cDNA

340   library established in UniGene for *Peromyscus* spp. includes 635 putative genes from testis [85].

341   Our assembled libraries from gonad tissue fall within these ranges, and non-annotated transcripts

342   could represent *Peromyscus*-specific genes.  To recover 100% of each tissue transcriptome,

343   samples would need to be prepared at various developmental stages and under various

344   environmental conditions.

345        Fisher's Exact Tests allowed us to identify annotated transcripts over-represented in one

346   tissue compared to the others.  The brain transcriptome of the social rodent, *C. sociabilis*,

347   exhibited highly expressed genes involved with behavior and signal transduction [76].  Over-

348   represented GO terms in *P. leucopus* brain tissue were related to similar major functions in the

349   brain, including regulation of behavior, cellular signaling, actin binding, ion transport and

350   channel activity, motor activity, and calcium ion binding.  In liver, over-represented GO terms

351   were largely dedicated to metabolic processes including ATP binding, GTP binding, NADH

352   dehydrogenase, and electron carrier activity.  There were also several GO terms related to the

353   immune response, hematopoietic processes, and nutrient binding; these annotations are supported

354   by microarray and RNA-seq analyses of liver in mouse and rat, respectively [80,81].

355

356   *SNP discovery and characterization*

357   Without a reference genome, aligning reads to assembled transcripts and assigning mismatches

358   as SNPs [86] is an acceptable substitute for generating sequence polymorphisms for non-model

359   species [51,54,87].  Difficulties may persist in distinguishing true SNPs from false positives

14

360    created by sequencing errors, misaligned reads, or alignment of reads to paralogous genes.

361    Identifying true SNPs depends on assembly quality, filtering criteria of nucleotide mismatches

362    during alignment, and statistical models used to call nucleotide variants [88]. Incorporating a

363    probabilistic framework in SNP-calling algorithms greatly reduces false positives [89,90].

364          We used conservative filtering criteria when calling SNPs to minimize false positives.

365    SAMtools [91] excels at SNP detection with low sequence coverage by comparing multiple

366    samples simultaneously [89,90].  We also filtered variants based on thresholds of quality and

367    minimum occurrence, and restricted maximum coverage to filter out false positive SNPs from

368    paralogous genes.  Excluding transcripts with the highest coverage after mapping limits

369    problems with gene duplications [92]. The thresholds we used for minimum SNP occurrence and

370    nucleotide quality reduce error rates by several orders of magnitude for pooled data, ensuring the

371    reliability of SNP libraries for downstream analyses [93].  Our SNP library represents highly

372    confident variant calls and will serve as an important resource for future population genetic

373    studies of urban and rural populations of *P. leucopus*.  We cannot completely rule out paralogous

374    genes or misalignments in our transcriptome assemblies, and thus future work will require

375    sequencing of transcripts from multiple individuals to validate SNP calls in candidate genes of

376    particular interest.

377

378    *Positive selection and the transcriptome*

379    We used the ratio of non-synonymous to synonymous substitution rates ($p_N/p_S$) to identify

380    candidate genes that may have experienced positive selection in urban populations of *P.*

381    *leucopus*.  Using SNPs to calculate ($p_N/p_S$) ratios in ORFs from assembled transcriptomes can

382    be a fruitful method for identifying the operation of natural selection on individual loci [6,52,94].

383    This approach has recently been used to identify genes under positive or purifying selection

384    between cichlid fish lineages in Nicaragua [56], between lake whitefish species pairs [54], and

385   within an invasive gastropod [23]. Studies traditionally identify positive selection in genes with

386   $p_N/p_S > 1.0$. We used this cutoff value, but also identified sequence pairs with $p_N/p_S$ between

387   0.5 and 1.0 to avoid overlooking relevant non-synonymous substitutions in candidate genes that

388   might be of interest for individual re-sequencing projects. Lack of full-length ORFs can

389   decrease $p_N/p_S$ values when some non-synonymous substitutions are unsampled [56,57]. The $p_N$

390   $/p_S$ index can also be used when samples have been pooled prior to sequencing [95], unlike

391   summary statistics that rely on allele frequencies [51].

392   We used McDonald-Kreitman tests to further elucidate patterns of evolution in candidate

393   genes. This method can identify adaptive changes between species and primarily detects

394   selection processes occurring thousands or even millions of years in the past. We calculated a

395   neutrality index (NI) as $(p_N/p_S)/(d_N/d_S)$ to look at deviations from neutral expectations [96,97].

396   While we detected an excess of non-synonymous polymorphisms within *P. leucopus* in genes

397   with functions including demethylation, xenobiotic metabolism, and innate immunity, we did not

398   find evidence of positive selection between species. While these patterns could suggest

399   purifying selection preventing the fixation of harmful mutations [98] or indicate balancing

400   selection acting to maintain favorable alleles in different populations [26], interpretation should

401   proceed cautiously due to limitations of polymorphism data generated from pooled

402   transcriptomes. The inability to assign individual allele frequencies when identifying

403   polymorphisms leads to an ascertainment bias towards high within-species $p_N/p_S$ ratios

404   compared to interspecies ratios, and this bias may explain the lack of NI values < 1 (positive

405   selection). These results could be interpreted as the result of balancing selection whereby

406   different alleles are favored in different urban populations, however, which would seem

407   consistent with the ecology of these relatively isolated populations. Individual resequencing of

408   mice from multiple populations will remove the ascertainment bias, uncover more

409    polymorphisms, and allow the use of more powerful tests to study recent selective pressures in

410    urban populations.

411        Many ecological changes arising from urbanization may drive local adaption to novel

412    conditions in fragmented urban populations, and we made several predictions about the types of

413    adaptive traits present in urban habitats from current literature. Genes involved in divergence of

414    urban and rural populations of white-footed mice are likely associated with quantitative traits

415    affected by crowded (i.e. high population density) and polluted urban environments (life history,

416    longevity, reproduction, immunity, metabolism, thermoregulatory and / or toxicological traits).

417    We identified candidate genes ($p_N/p_S > 1$) that supported these predictions between urban and

418    rural populations of mice, but also between individual urban populations.  The urban matrix is a

419    strong enough barrier to dispersal that white-footed mouse populations in individual city parks

420    may experience highly localized selective pressures in addition to selective pressures that are

421    general to urban environments [41].

422        New predators, competitors, parasites, and pathogens can drive local adaptation of traits,

423    especially those related to immunity, in novel urban environments [15,16]. We identified

424    candidate genes involved in the innate immune system and activation of the complement

425    pathway to identify pathogens.  Additionally, two candidate genes were identified in

426    comparisons of urban populations that function in blood coagulation and inflammation. The

427    innate immune system is a biochemical pathway that removes pathogens by identifying and

428    killing target cells [99], and positive selection is found to act on pathogen recognition genes

429    within the complement activation pathway [100]. The introduction of invasive species,

430    population growth of 'urban exploiters', and increased traffic, trade, and transportation within

431    cities can introduce large numbers of novel pathogens [101], and white-footed mice in NYC may

432    be evolving to efficiently recognize them and respond immunologically.  We also identified

433    several genes involved in metabolism that were divergent between populations, and a gene

17

434  expressed during spermatogenesis that was divergent between urban and rural populations.

435  Rapid evolution has been identified in reproductive proteins between *Peromyscus* spp. affecting

436  spermatogenesis, sperm competition, and sperm-egg interactions [102], and the intensity of

437  sperm competition and reproductive conflict may be increasing in dense *P. leucopus* populations

438  in NYC.

439       Increasing air, water, and soil pollution are all typical impacts of urbanization [17–19].

440  One potential marker of increased exposure to pollutants is hypermethylation of regulatory

441  regions of the genome [17,103–105]. Positive selection may also be acting on genes involved in

442  xenobiotic metabolism.  Heavy metals including mercury, lead, and arsenic occur at increased

443  concentrations within NYC park soils (S. Harris, unpublished data), and McGuire et al. [106]

444  found lower pH and higher concentrations of heavy metals in NYC parks compared to green

445  roofs.  PCB resistance was identified in multiple populations of *Fundulus heteroclitus* [18], and

446  Wirgin et al. [107] also found rapid PCB resistance in tomcod from the Hudson River through

447  positive selection.  In urban to rural comparisons we found two potential toxicological candidate

448  genes: one gene involved in metabolizing foreign chemical compounds (i.e. xenobiotics), and a

449  demethylase that removes methyl groups from histone lysines.

450       Comparing candidate genes from all pairwise analyses with $p_N/p_S$ between 0.5 and 1

451  reveals several additional patterns.  Proteins were identified that function in the alternative

452  pathway, which acts continuously in an organism without antibody activation to clear foreign

453  pathogens [108], and supports our conclusion that positive selection  ($p_N/p_S > 1$) is acting on the

454  innate immune system in these populations. Four *cytochrome p450* genes, *2d27-like*, *family 2*

455  *subfamily B*, *subfamily polypeptide 13*, and *2a15,* exhibited $p_N/p_S$ between 0.5 and 1 in urban

456  populations and between urban and rural populations.  *Cytochrome p450 2a15* was also found to

457  have a $p_N/p_S > 1$ and McDonald-Kreitman tests found significantly more polymorphisms within

458  *P. leucopus* than between species (Table 4, Table 5). The *cytochrome p450* family of genes plays

459    a major role in xenobiotic metabolism, including detoxification in variable environments

460    [109,110].  Patterns of divergence and positive selection have been robustly identified in

461    *cytochrome p450* genes in natural populations of both *Mus musculus* when ingesting toxins

462    through their diet and *Tetrahymena thermophila* exposed to toxic environments [110,111].  *P.*

463    *leucopus* in NYC populations may be experiencing different dietary demands and exposure to

464    pollutants, leading to selective pressures on detoxifying genes like the cytochrome p450 gene

465    family.

466         Alternatively, genetic differences between urban and rural populations may result from

467    genetic drift rather than selection. We will differentiate between drift and selection in future

468    work by examining genetic divergence between multiple urban and rural populations at these

469    candidate loci and additional genes. We must also be cautious when inferring the function of

470    candidate genes after identifying statistical signatures of positive selection.

471         While a $p_N/p_S$ ratio > 1.0 can represent positive selection, it may also occur due to

472    relaxation of purifying selection, and individual codons within a gene can have an excess of non-

473    synonymous substitutions due to random biological processes [112]. However, current statistical

474    tests address these issues and are generally robust in identifying positive selection [113]. In the

475    case of a single population, $p_N/p_S > 1$ may not represent positive selection.  Kryazhimskiy &

476    Plotkin [114] demonstrated that the relationship between $p_N/p_S$ and selection is radically

477    different when samples originated from the same population; $p_N/p_S$ actually decreases in

478    response to positive selection.  To infer selection between two samples using $p_N/p_S$, samples

479    must come from reproductively isolated populations with fixed substitutions [114].  All samples

480    used to calculate $p_N/p_S$ for this study came from reproductively isolated and genetically

481    structured populations [40].  We assembled transcriptome datasets individually for each

482    population to identify fixed substitutions between populations and avoid randomly segregating

483    SNPs in $p_N/p_S$ analyses.  Indices such as $p_N/p_S$ identify genes with previously unknown

19

484  signatures of selection, but candidates still need to be studied in a controlled setting to identify

485  phenotype and function [113].

486     The ability of $p_N/p_S$ and McDonald-Kreitman tests to detect genes under positive

487  selection is limited in some situations, so it is likely that we have missed many candidate genes.

488  Additionally, such analyses do not identify adaptive variation in gene regulatory regions as

489  opposed to transcribed cDNA [115].  Ratios such as $p_N/p_S$ may also vary widely when there are

490  relatively few mutations per gene [56, 108]. Given strong selection within populations, however,

491  it is plausible that multiple substitutions may rise to high frequency or become fixed within a few

492  hundred generations (i.e. in the timeframe of divergence for urban and rural populations of

493  white-footed mice).  The candidate genes identified herein can be confirmed in future work using

494  the reference genome of *P. maniculatus* (sequenced and currently being assembled) and multiple

495  tests of selection that provide more statistical power and higher resolution when identifying types

496  and age of selection in single candidate genes [116,117]. These emerging resources will allow us

497  to validate many of our predicted polymorphisms, identify paralogous genes with greater

498  certainty, and perform more powerful tests of selection by providing genetic distances and

499  genomic coordinates for our sequenced contigs. Our ongoing work in this system uses these

500  external resources with our new transcriptomic and genomic libraries from individual mice from

501  several urban, rural, and suburban populations. These ongoing studies employ multiple outlier

502  statistics based on the allele frequency spectrum and linkage disequilibrium to examine recent

503  selection in both coding and non-coding regions of urban white-footed mouse genomes.

504

505  **Materials and Methods**

506  *Ethics statement*

507  All animal procedures were approved by the Institutional Animal Care and Use Committee at

508  Brooklyn College, CUNY (Protocol No. 247), and adhered to the Guidelines of the American

509     Society of Mammalogists for the Use of Wild Mammals in Research [118].  Field work was

510     conducted with the permission of the New York State Department of Environmental

511     Conservation (License to Collect or Possess Wildlife No. 1603) and the New York City

512     Department of Parks and Recreation.

513

514     *Study Sites and population sampling*

515     *P. leucopus* were trapped and collected from each of four urban and one rural site ($N = 20\text{-}25$ /

516     population) for sequencing and analysis (total $N = 112$; Fig. 6).  The four urban sites (Central

517     Park, Flushing Meadows-Willow Lake, New York Botanical Gardens, and the Ridgewood

518     Reservoir) were chosen due to their large area, isolation by dense urban matrix, high population

519     density of mice, substantial genetic differentiation, and genetic isolation from other populations

520     [40,41].  The rural site, Harriman State Park located ~68 km north of Central Park, is one of the

521     largest contiguous protected areas nearby and the most likely representative of a non-urban

522     population of mice in proximity to NYC.  Mice were trapped over a period of 1-3 nights at each

523     site using four 7x7 transects of 3"x3"x9" Sherman live traps.  Mice were killed by cervical

524     dislocation and immediately dissected in the field.  Livers, gonads and brains were extracted,

525     rinsed with PBS to remove any debris from the surface of the tissue, and immediately placed in

526     RNALater® (Ambion Inc., Austin, TX) on ice before transport and storage at -80°C. These tissue

527     types were chosen for initial analysis due to their wide range of expressed gene transcripts [78]

528     and potential roles in adaptation to urban conditions.

529

530     *RNA extraction and cDNA library preparation*

531     Total RNA was extracted and cDNA libraries were pooled for all five populations for four

532     multiplexed plates of 454 sequencing.  The first plate of sequencing was normalized to produce

533     equalized concentrations of all transcripts present, potentially allowing enhanced gene discovery

534    and greater overall coverage of the transcriptome [119]. However, the normalization process

535    introduces additional steps and biases in library preparation [50], and resulted in a relatively low

536    number of total high-quality 454 reads.  Thus, non-normalized libraries were prepared using a

537    modified protocol for the last three 454 plates.

538         For plate 1, total RNA was isolated from ~60 mg of liver (eight males and eight females /

539    population), ~60 mg of testis (eight males / population), and ~60 mg of ovaries (eight females /

540    population) for two populations using RNaqueous® kits (Ambion, Austin, TX).  Individual RNA

541    extracts were pooled by population and organ type and selected for mature mRNA using the

542    MicroPoly(A)Purist™ kit (Ambion, Austin, TX).  Next, mRNA pools were reverse-transcribed

543    using the SMARTer™ cDNA synthesis kit (Clontech, Mountain View, CA), and normalized

544    using the Trimmer-Direct cDNA normalization kit (Evrogen, Moscow, Russia).  Then,

545    normalized cDNA pools were sequenced with multiplex identifiers using standard 454 FLX

546    Titanium protocols. This pilot plate contained cDNA pools for Harriman State Park and Flushing

547    Meadows-Willow Lake.

548         For plates 2-4, total RNA using Trizol® reagent (Invitrogen, Carlsbad, CA) was extracted

549    from ~70 mg of brain tissue (four males and four females / population), ~70 mg of testes (eight

550    males / population), and ~15 mg of liver (four males and four females / population). After

551    DNAse treatment (Promega, Madison, WI) and pooling individual samples in equimolar

552    amounts by population and tissue, the samples were treated with the RiboMinus™ Eukaryote kit

553    (Invitrogen, Carlsbad, CA) to reduce ribosomal RNA.  RNA pools were then reverse-transcribed

554    using the Roche cDNA synthesis kit (Roche Diagnostics, Indianapolis, IN) and sequenced with

555    multiplex identifiers using standard 454 FLX Titanium protocols.  Plate 2 included brain cDNA

556    pools for all five populations, plate 3 included liver and testis cDNA for Central Park,

557    Ridgewood Reservoir, New York Botanical Gardens, and Harriman State Park, and plate 4

558    included liver cDNA pools from all five populations. All raw sequencing files have been

559    deposited in the GenBank Sequence Read Archive (SRA) under accession number SRP020005.

560

561    *Transcriptome assembly*

562    Two methods were used to assemble the best transcriptome from all four 454 plates:  Cap3 [120]

563    a long-read assembler that performs well in transcriptome assemblies [60], and Roche's

564    proprietary software, Newbler (Version 2.5.3), that was designed specifically for assembling 454

565    sequencing reads with additional features for cDNA sequence.  Newbler's cDNA options

566    assemble reads into contigs, followed by assembly into larger 'isotigs' representing alternatively-

567    spliced transcripts.  Isotigs are then clustered into larger 'isogroups' representing full-length

568    genes. Transcriptome assembly was attempted with the full set of reads using Cap3 and Newbler

569    with cDNA options, but due to computational limitations the full dataset could not be assembled

570    with either software program. We addressed this issue by first assembling sequences from all

571    four plates with Newbler using the genome assembly settings and default parameters after

572    trimming 454 adaptors and barcodes from the reads. Reads that were either 'assembled' or

573    'partially assembled' in this pilot run were filtered and used as input for cDNA assemblies in

574    Newbler or Cap3. These reads were filtered from the raw sff files using a locally-installed

575    instance of Galaxy [121].  Before the cDNA assembly, nucleotides with poor quality scores,

576    primer sequences, and long poly(A) tails were removed using cutadapt (Version 1.2.1 2012,

577    [122] and the trim-fastq.pl perl script implemented in Popoolation [93].  The filtered fastq files

578    were then used as input for Cap3 or Newbler with the cDNA assembly option, using default

579    parameters for both assemblies. These assemblies (1. genome assembly with Newbler, 2. cDNA

580    assembly with Newbler, and 3. cDNA assembly with Cap3) were compared to identify the best

581    full reference transcriptome for downstream analysis.

582        For analyses of individual tissues, separate cDNA assemblies were performed.  Tissues

583    were barcoded, and sequence reads originating from liver, gonads, or brains were parsed from

584    the raw 454 sequencing reads.  These datasets were small enough to be assembled separately as

585    tissue-specific transcriptomes in Newbler using the cDNA option with default parameters.

586    Population-specific transcriptomes were also assembled, using the same methodology, to

587    examine population-specific statistical signatures of selection.

588

589    *Alignment to model rodent genomes*

590    *Peromyscus* assemblies were initially characterized and annotated by performing two separate

591    analyses using *Mus musculus* and *Rattus norvegicus* genomic resources.  The first analysis was

592    used to determine the number of likely genes in each assembly. BLASTN searches were

593    performed against *Mus musculus* (NCBI Annotation Release 103) and *Rattus norvegicus* (NCBI

594    build 5.1) cDNA reference libraries downloaded from NCBI. BLASTN matches were considered

595    significant when sequence identity was greater than 80%, alignment length was at least 50% of

596    the total length of either the query or subject sequence, and the *e*-value was less than $10^{-5}$.  While

597    significant, these hits may not be ideal for population genomic analyses due to inclusion of

598    paralogous gene matches, matches between multi-gene families, and false positive orthologous

599    gene matches. In order to identify individual isotigs representing a single gene with known

600    function useful for statistical analysis, BLASTN results were further filtered by including query

601    hits that matched only one subject ID (i.e. gene) and *vice versa*.  These contigs were annotated as

602    'Gene Candidates'.

603        The distribution of *P. leucopus* isotigs across model rodent genomes was analzyed. All *P.*

604    *leucopus* isotigs were mapped to chromosomes in the *Mus* (GRCm38) and *Rattus* (RGSC 5.0)

605    reference genomes. Default BLAT parameters were used with an exception for aligning mRNA

606    to genomes across species (-q=rnax -t=dnax, [123]), and best BLAT hits were parsed based on

607    percent identity and score (# match − # mismatch).

608

609    *Mapping and SNP discovery*

610    To generate a SNP library for downstream population genomic analysis, 454 reads were first

611    mapped to the Newbler cDNA assembly using the BWA-SW (http://bio-bwa.sourceforge.net/)

612    alignment algorithm for long read mapping [124]. We only used trimmed reads from the final

613    assembly, removed singletons before mapping to reduce false positive SNP calls from

614    sequencing errors or duplicate reads, and included reads with a mapping quality > 20 in

615    SAMtools. The SAM file from BWA-SW was used in the SAMtools package (v. 0.1.17, [89]) to

616    call SNPs using the mpileup command with a maximum coverage cutoff of 200.  The SNP

617    calling pipeline implemented in SAMtools uses base alignment quality (BAQ) calculations to

618    generate likelihoods of genotypes, can overcome low coverage by using sequence information

619    from multiple samples to call variants, and uses Bayesian inference to make SNP calls with high

620    confidence [87-89].  In addition to the default parameters in SAMtools, we included stringent

621    additional filters by removing any potential INDELs, only including SNPs with a phred quality

622    (Q-value) ≥ 20, a minimum occurrence of two, and coverage ≤ 200 to exclude alignment

623    artifacts, duplicates, and paralogous genes [93,118,124–126].

624

625    *Functional annotation of transcriptomes*

626    The reference transcriptome was annotated by performing a BLASTX search to identify

627    homologous sequences from the NCBI non-redundant protein database, and then GO terms

628    associated with BLASTX hits were retrieved using the annotation pipeline in Blast2GO

629    [125,126].  Tissue-specific assemblies were also annotated in Blast2GO, and Fisher's Exact Test

630    was used to examine whether GO terms were over-represented between pairs of tissue types.

631 Each pairwise tissue comparison (liver, brain, gonad) was analyzed for over-representation, and

632 significant results were identified with a False Discovery Rate (FDR) $\leq 0.05$.

633

634 *Prediction of Open Reading Frames (ORFs) and $p_N/p_S$ calculations*

635 Regions containing ORFs were identified using BLASTX searches of our assembled contigs

636 against the NCBI non-redundant protein database. Only best hits with an *e*-value $\leq 10^{-5}$, and

637 when query transcripts hit only one subject sequence and *vice versa*, were kept. From these

638 results, a general feature file (GFF) was manually created indicating the start and stop

639 coordinates, strand information, and reading frame from the BLASTX results. Within these

640 protein coding regions, putative ORFs were identified when a start codon was found and the

641 reading frame was greater than 150 bp long. The Perl script, Syn-nonsyn-at-position.pl,

642 implemented in Popoolation v. 1.2.2 [93] was used to define population-specific SNPs obtained

643 from the SAMtools analysis above as either non-synonymous or synonymous.

644 The ratio of non-synonymous ($p_N$) to synonymous ($p_S$) SNP substitutions ($p_N/p_S$) was

645 calculated between individual Newbler cDNA population assemblies to identify coding

646 sequences potentially experiencing directional selection in urban *P. leucopus* populations. For

647 each population pair, the fastaFromBed command in bedtools [127] was used to filter contigs

648 and generate a fasta file of putative ORFs (identified above) for each population assembly. The

649 USEARCH (http://www.drive5.com/usearch/) clustering and alignment software for genomic

650 datasets [128] was used to create pairwise alignments between all population ORFs using an *e*-

651 value $\leq 0.001$. Signatures of selection between aligned ORFs were identified using

652 KaKs_Calculator1.2 (Model GY) [129] to calculate the ratio of non-synonymous ($p_N$) to

653 synonymous ($p_S$) SNPs in each population pair. Only transcripts with at least three SNPs, an

654 ORF length greater than 150 bp, and no in-frame stop codons were included. The mean number

655 of SNPs per ORF was $1.4 \pm SE = 2.9$. A three SNP threshold was chosen to avoid bias as Ka /

656 Ks calculations lose statistical power as the number of substitutions per ORF decreases [130].

657 The maximum likelihood method was used that accounts for evolutionary characteristics (i.e.

658 ratio of transition / transversion rates, nucleotide frequencies) of our transcriptome datasets.

659 Contigs with elevated $p_N/p_S$ ratios were then annotated in Blast2GO as above.

660 Candidate genes were screened for evidence of recombination, and additional signatures

661 of natural selection were examined using McDonald-Kreitman tests. We used BLASTN

662 searches to find orthologous mRNA sequences from multiple species for each candidate gene.

663 For recombination analysis, multiple mammals were used and always included *Cricetulus*

664 *griseus*, *Rattus norevegicus*, or *Mus musculus*. Orthologous sequences were codon-aligned using

665 MACSE [131] and then scanned for evidence of recombination using a GARD analysis

666 implemented in the Data Monkey webserver [132,133]. For McDonald-Kreitman tests,

667 orthologous genes between *Peromyscus leucopus* and *Rattus norvegicus* or *Cricetulus griseus*

668 were codon aligned with MACSE [131]. Non-overlapping datasets of polymorphisms within *P.*

669 *leucopus* and fixed genetic changes between species were then generated. The McDonald-

670 Kreitman test was performed with these data using DnaSP v.5.10.1 [134]. Fasta files of

671 assembled contigs / isotigs, vcf files of SNP marker data, BLAST2GO files of functional

672 annotations, and output files from population genetics tests are available on the Dryad digital

673 repository (doi: 10.5061/dryad.r8ns3).

674

675 **Acknowledgements**

680     MacManes, three anonymous reviewers and the Associate Editor provided many comments that

681     substantially improved the manuscript.

682

683  **References**

684  1.  Shochat E, Warren PS, Faeth SH, McIntyre NE, Hope D (2006) From patterns to

685    emerging processes in mechanistic urban ecology. Trends in Ecology and Evolution 21:

686    186–191. doi:10.1016/j.tree.2005.11.019.

687  2.  United Nations (2011) World Urbanization Prospects The 2011 Revision. UN Department

688    of Economic and Social Affairs'.

689  3.  Martin LJ, Blossey B, Ellis E (2012) Mapping where ecologists work: biases in the global

690    distribution of terrestrial ecological observations. Frontiers in Ecology and the

691    Environment 10: 195–201. doi:10.1890/110154.

692  4.  Pickett ST, Cadenasso ML, Grove JM, Boone CG, Groffman PM, et al. (2011) Urban

693    ecological systems: scientific foundations and a decade of progress. Journal of

694    Environmental Management 92: 331–362. doi:10.1016/j.jenvman.2010.08.022.

695  5.  Rice AM, Rudh A, Ellegren H, Qvarnström A (2010) A guide to the genomics of

696    ecological speciation in natural animal populations. Ecology Letters: 9–18.

697    doi:10.1111/j.1461-0248.2010.01546.x.

698  6.  Hohenlohe PA, Phillips PC, Cresko WA (2011) Using Population Genomics To Detect

699    Selection in Natural Populations: Key Concepts and Methodological Considerations.

700    International Journal of Plant Sciences 171: 1059–1071. doi:10.1086/656306.USING.

701  7.  Storz J, Hoekstra H (2007) The study of adaptation and speciation in the genomic era.

702    Journal of Mammalogy 88: 1–4.

703  8.  White J, Antos M, Fitzsimons J, Palmer G (2005) Non-uniform bird assemblages in urban

704    environments: the influence of streetscape vegetation. Landscape and Urban Planning 71:

705    123–135. doi:10.1016/j.landurbplan.2004.02.006.

706    9.     Grimm NB, Faeth SH, Golubiewski NE, Redman CL, Wu J, et al. (2008) Global change

707           and the ecology of cities. Science (New York, NY) 319: 756–760.

708           doi:10.1126/science.1150195.

709    10.    Blair RB (2001) Birds and butterflies along urban gradients in two ecoregions of the U.S.

710           In: Biotic Homogenization. Lockwood JL, McKinney ML, editors Norwell, MA: Kluwer

711           Academic Publishers.

712    11.    McKinney ML (2002) Urbanization, biodiversity, and conservation. Bioscience 52: 883–

713           890.

714    12.    McKinney ML (2006) Urbanization as a major cause of biotic homogenization. Biological

715           Conservation 127: 247–260. doi:10.1016/j.biocon.2005.09.005.

716    13.    Bjorklund M, Ruiz I, Senar JC (2010) Genetic differentiation in the urban habitat: the

717           great tits ( Parus major ) of the parks of Barcelona city. Biological Journal of the Linnean

718           Society 99: 9–19. doi:10.1111/j.1095-8312.2009.01335.x.

719    14.    Wandeler P, Funk SM, Largiadèr CR, Gloor S, Breitenmoser U (2003) The city-fox

720           phenomenon: genetic consequences of a recent colonization of urban habitat. Molecular

721           Ecology 12: 647–656.

722    15.    Peluc SI, Sillett TS, Rotenberry JT, Ghalambor CK (2008) Adaptive phenotypic plasticity

723           in an island songbird exposed to a novel predation risk. Behavioral Ecology 19: 830–835.

724           doi:10.1093/beheco/arn033.

725    16.    Sih A, Ferrari MCO, Harris DJ (2011) Evolution and behavioural responses to human-

726           induced rapid environmental change. Evolutionary Applications 4: 367–387.

727           doi:10.1111/j.1752-4571.2010.00166.x.

728    17.    Yauk C, Polyzos A, Rowan-Carroll A, Somers CM, Godschalk RW, et al. (2008) Germ-

729           line mutations, DNA damage, and global hypermethylation in mice exposed to particulate

730        air pollution in an urban/industrial location. Proceedings of the National Academy of

731        Sciences of the United States of America 105: 605–610. doi:10.1073/pnas.0705896105.

732    18.   Whitehead A, Triant D, Champlin D, Nacci D (2010) Comparative transcriptomics

733        implicates mechanisms of evolved pollution tolerance in a killifish population. Molecular

734        Ecology 19: 5186–5203. doi:10.1111/j.1365-294X.2010.04829.x.

735    19.   Francis RA, Chadwick MA (2012) What makes a species synurbic? Applied Geography

736        32: 514–521. doi:10.1016/j.apgeog.2011.06.013.

737    20.   Mueller JC, Partecke J, Hatchwell BJ, Gaston KJ, Evans KL (2013) Candidate gene

738        polymorphisms for behavioural adaptations during urbanization in blackbirds. Molecular

739        Ecology 22: 3629–3637. doi:10.1111/mec.12288.

740    21.   Brady SP (2012) Road to evolution? Local adaptation to road adjacency in an amphibian

741        (Ambystoma maculatum). Scientific Reports 2. doi:10.1038/srep00235.

742    22.   Cheptou P-O, Carrue O, Rouifed S, Cantarel A (2008) Rapid evolution of seed dispersal

743        in an urban environment in the weed Crepis sancta. Proceedings of the National Academy

744        of Sciences of the United States of America 105: 3796–3799.

745        doi:10.1073/pnas.0708446105.

746    23.   Sun J, Wang M, Wang H, Zhang H, Zhang X, et al. (2012) De novo assembly of the

747        transcriptome of an invasive snail and its multiple ecological applications. Molecular

748        Ecology Resources 12: 1133–1144. doi:10.1111/1755-0998.12014.

749    24.   Weber JN, Peterson BK, Hoekstra HE (2013) Discrete genetic modules are responsible for

750        complex burrow evolution in Peromyscus mice. Nature 493: 402–405.

751        doi:10.1038/nature11816.

752    25.   Linnen CR, Kingsley EP, Jensen JD, Hoekstra HE (2009) On the origin and spread of an

753        adaptive allele in deer mice. Science (New York, NY) 325: 1095–1098.

754        doi:10.1126/science.1175826.

755  26.  Storz JF, Sabatino SJ, Hoffmann FG, Gering EJ, Moriyama H, et al. (2007) The molecular

756      basis of high-altitude adaptation in deer mice. PLoS Genetics 3: e45.

757      doi:10.1371/journal.pgen.0030045.

758  27.  Ungvari Z, Krasnikov BF, Csiszar A, Labinskyy N, Mukhopadhyay P, et al. (2008)

759      Testing hypotheses of aging in long-lived mice of the genus Peromyscus: association

760      between longevity and mitochondrial stress resistance, ROS detoxification pathways, and

761      DNA repair efficiency. Age 30: 121–133. doi:10.1007/s11357-008-9059-y.

762  28.  O'Neill R, Szalai G, Gibbs R, Weinstock G (1998) Sequencing the genome of

763      Peromyscus. White paper proposal: 14.

764  29.  Vessey SH, Vessey KB (2007) Linking behavior, life history and food supply with the

765      population dynamics of white-footed mice (Peromyscus leucopus). Integrative Zoology 2:

766      123–130.

767  30.  Metzger LH (1971) Behavioral Population Regulation in the Woodmouse, Peromyscus

768      leucopus. American Midland Naturalist 86: 434–448.

769  31.  Wang G, Wolff JO, Vessey SH, Slade NA, Witham JW, et al. (2008) Comparative

770      population dynamics of Peromyscus leucopus in North America: influences of climate,

771      food, and density dependence. Population Ecology 51: 133–142. doi:10.1007/s10144-008-

772      0094-4.

773  32.  Linnen CR, Hoekstra HE (2009) Measuring natural selection on genotypes and

774      phenotypes in the wild. Cold Spring Harbor Symposia on Quantitative Biology 74: 155–

775      168. doi:10.1101/sqb.2009.74.045.

776  33.  Mullen LM, Hoekstra HE (2008) Natural selection along an environmental gradient: a

777      classic cline in mouse pigmentation. Evolution; International Journal of Organic Evolution

778      62: 1555–1570. doi:10.1111/j.1558-5646.2008.00425.x.

779    34.    Puth LM, Burns CE (2009) New York's nature: a review of the status and trends in

780           species richness across the metropolitan region. Diversity and Distributions 15: 12–21.

781           doi:10.1111/j.1472-4642.2008.00499.x.

782    35.    Ekernas LS, Mertes KJ (2007) The influence of urbanization, patch size, and habitat type

783           on small mammal communities in the New York metropolitan region: a preliminary

784           report. Transactions of the Linnaean Society of New York 10: 239–264.

785    36.    Barko VA, Feldhamer GA, Nicholson MC, Davie DK (2003) Urban Habitat: a

786           Determinant of White-Footed Mouse (Peromyscus Leucopus) Abundance in Southern

787           Illinois. Southeastern Naturalist 2: 369–376. doi:10.1656/1528-

788           7092(2003)002[0369:UHADOW]2.0.CO;2.

789    37.    Nupp TE, Swihart RK (1996) Effect of forest patch area on population attributes of white-

790           footed mice (Peromyscus leucopus) in fragmented landscapes. Canadian Journal of

791           Zoology 74: 467–472.

792    38.    Lankau R (2010) Rapid Evolution and Mechanisms of Species Coexistence. Annual

793           Review of Ecology, Evolution, and Systematics 42: 335–354. doi:10.1146/annurev-

794           ecolsys-102710-145100.

795    39.    Lankau RA, Strauss SY (2011) Newly rare or newly common: evolutionary feedbacks

796           through changes in population density and relative species abundance, and their

797           management implications. Evolutionary Applications 4: 338–353. doi:10.1111/j.1752-

798           4571.2010.00173.x.

799    40.    Munshi-South J, Kharchenko K (2010) Rapid, pervasive genetic differentiation of urban

800           white-footed mouse (Peromyscus leucopus) populations in New York City. Molecular

801           Ecology 19: 4242–4254. doi:10.1111/j.1365-294X.2010.04816.x.

802  41.  Munshi-South J (2012) Urban landscape genetics: canopy cover predicts gene flow

803       between white-footed mouse (Peromyscus leucopus) populations in New York City.

804       Molecular Ecology 21: 1360–1378. doi:10.1111/j.1365-294X.2012.05476.x.

805  42.  Yang D-S, Kenagy GJ (2009) Nuclear and mitochondrial DNA reveal contrasting

806       evolutionary processes in populations of deer mice (Peromyscus maniculatus). Molecular

807       Ecology 18: 5115–5125. doi:10.1111/j.1365-294X.2009.04399.x.

808  43.  Degner JF, Stout IJ, Roth JD, Parkinson CL (2007) Population genetics and conservation

809       of the threatened southeastern beach mouse (Peromyscus polionotus niveiventris):

810       subspecies and evolutionary units. Conservation Genetics 8: 1441–1452.

811       doi:10.1007/s10592-007-9295-1.

812  44.  Ozer F, Gellerman H, Ashley M V (2011) Genetic impacts of Anacapa deer mice

813       reintroductions following rat eradication. Molecular Ecology: 3525–3539.

814       doi:10.1111/j.1365-294X.2011.05165.x.

815  45.  Boitard S, Schlötterer C, Nolte V, Pandey RV, Futschik A (2012) Detecting Selective

816       Sweeps from Pooled Next-Generation Sequencing Samples. Molecular Biology and

817       Evolution 29: 2177–2186. doi:10.1093/molbev/mss090.

818  46.  Futschik A, Schlötterer C (2010) Massively Parallel Sequencing of Pooled DNA Samples-

819       -The Next Generation of Molecular Markers. Genetics 218: 207–218.

820       doi:10.1534/genetics.110.114397.

821  47.  Gompert Z, Buerkle CA (2011) A hierarchical bayesian model for next-generation

822       population genomics. Genetics 187: 903–917. doi:10.1534/genetics.110.124693.

823  48.  Stinchcombe JR, Hoekstra HE (2008) Combining population genomics and quantitative

824       genetics: finding the genes underlying ecologically important traits. Heredity 100: 158–

825       170. doi:10.1038/sj.hdy.6800937.

826    49.    Bonin A (2008) Population genomics: a new generation of genome scans to bridge the gap

827            with functional genomics. Molecular Ecology 17: 3583–3584. doi:10.1111/j.1365-

828            294X.2008.03854.x.

829    50.    Ungerer MC, Johnson LC, Herman MA (2008) Ecological genomics: understanding gene

830            and genome function in the natural environment. Heredity 100: 178–183.

831            doi:10.1038/sj.hdy.6800992.

832    51.    Sloan DB, Keller SR, Berardi AE, Sanderson BJ, Karpovich JF, et al. (2012) De novo

833            transcriptome assembly and polymorphism detection in the flowering plant Silene vulgaris

834            (Caryophyllaceae). Molecular Ecology Resources 12: 333–343. doi:10.1111/j.1755-

835            0998.2011.03079.x.

836    52.    Nielsen R, Yang Z (1998) Likelihood models for detecting positively selected amino acid

837            sites and applications to the HIV-1 envelope gene. Genetics 148: 929–936.

838    53.    Ellegren H (2008) Sequencing goes 454 and takes large-scale genomics into the wild.

839            Molecular Ecology 17: 1629–1631. doi:10.1111/j.1365-294X.2008.03699.x.

840    54.    Renaut S, Nolte AW, Bernatchez L (2010) Mining transcriptome sequences towards

841            identifying adaptive single nucleotide polymorphisms in lake whitefish species pairs

842            (Coregonus spp. Salmonidae). Molecular Ecology 19 Suppl 1: 115–131.

843            doi:10.1111/j.1365-294X.2009.04477.x.

844    55.    Wang X-W, Zhao Q-Y, Luan J-B, Wang Y-J, Yan G-H, et al. (2012) Analysis of a native

845            whitefly transcriptome and its sequence divergence with two invasive whitefly species.

846            BMC Genomics 13: 529. doi:10.1186/1471-2164-13-529.

847    56.    Elmer KR, Fan S, Gunter HM, Jones JC, Boekhoff S, et al. (2010) Rapid evolution and

848            selection inferred from the transcriptomes of sympatric crater lake cichlid fishes.

849            Molecular Ecology 19: 197–211. doi:10.1111/j.1365-294X.2009.04488.x.

850 57. Swanson WJ, Wong A, Wolfner MF, Aquadro CF (2004) Evolutionary expressed

851 sequence tag analysis of Drosophila female reproductive tracts identifies genes subjected

852 to positive selection. Genetics 168: 1457–1465. doi:10.1534/genetics.104.030478.

853 58. Metzker ML (2010) Sequencing technologies - the next generation. Nature reviews

854 Genetics 11: 31–46. doi:10.1038/nrg2626.

855 59. Vera JC, Wheat CW, Fescemyer HW, Frilander MJ, Crawford DL, et al. (2008) Rapid

856 transcriptome characterization for a nonmodel organism using 454 pyrosequencing.

857 Molecular Ecology 17: 1636–1647. doi:10.1111/j.1365-294X.2008.03666.x.

858 60. Meyer E, Aglyamova G V, Wang S, Buchanan-Carter J, Abrego D, et al. (2009)

859 Sequencing and de novo analysis of a coral larval transcriptome using 454 GSFlx. BMC

860 Genomics 10: 219. doi:10.1186/1471-2164-10-219.

861 61. Santure AW, Gratten J, Mossman JA, Sheldon BC, Slate J (2011) Characterisation of the

862 transcriptome of a wild great tit Parus major population by next generation sequencing.

863 BMC Genomics 12: 283. doi:10.1186/1471-2164-12-283.

864 62. Glenn TC (2011) Field guide to next-generation DNA sequencers. Molecular Ecology

865 Resources 11: 759–769. doi:10.1111/j.1755-0998.2011.03024.x.

866 63. Cahais V, Gayral P, Tsagkogeorga G, Melo-Ferreira J, Ballenghien M, et al. (2012)

867 Reference-free transcriptome assembly in non-model animals from next-generation

868 sequencing data. Molecular Ecology Resources 12: 834–845. doi:10.1111/j.1755-

869 0998.2012.03148.x.

870 64. Mundry M, Bornberg-Bauer E, Sammeth M, Feulner PGD (2012) Evaluating

871 characteristics of de novo assembly software on 454 transcriptome data: a simulation

872 approach. PloS One 7: e31410. doi:10.1371/journal.pone.0031410.

873     65.    Ewen-Campen B, Shaner N, Panfilio KA, Suzuki Y, Roth S, et al. (2011) The maternal

874            and early embryonic transcriptome of the milkweed bug Oncopeltus fasciatus. BMC

875            Genomics 12: 61. doi:10.1186/1471-2164-12-61.

876     66.    Malik A, Korol A, Hübner S, Hernandez AG, Thimmapuram J, et al. (2011)

877            Transcriptome Sequencing of the Blind Subterranean Mole Rat, Spalax galili: Utility and

878            Potential for the Discovery of Novel Evolutionary Patterns. PLoS ONE 6: e21227.

879            doi:10.1371/journal.pone.0021227.

880     67.    Hoffman J, Nichols H (2011) A novel approach for mining polymorphic microsatellite

881            markers in silico. PloS ONE 6(8):e23283.

882     68.    Christodoulou DC, Gorham JM, Herman DS (2011) Construction of normalized RNA-seq

883            libraries for next- generation sequencing using the crab duplex-specific nuclease. Current

884            Protocols in Molecular Biology: 1–14. doi:10.1002/0471142727.mb0412s94.Construction.

885     69.    Davey JW, Hohenlohe PA, Etter PD, Boone JQ, Catchen JM, et al. (2011) Genome-wide

886            genetic marker discovery and genotyping using next-generation sequencing. Nature

887            Reviews Genetics 12: 499–510. doi:10.1038/nrg3012.

888     70.    Vijay N, Poelstra J, Künstner A, Wolf J (2012) Challenges and strategies in transcriptome

889            assembly and differential gene expression quantification. A comprehensive in silico

890            assessment of RNA-seq experiments. Molecular Ecology 46: 620–634.

891            doi:10.1111/mec.12014.

892     71.    Ekblom R, Slate J, Horsburgh GJ, Birkhead T, Burke T (2012) Comparison between

893            Normalised and Unnormalised 454-Sequencing Libraries for Small-Scale RNA-Seq

894            Studies. Comparative and Functional Genomics 2012: 281693. doi:10.1155/2012/281693.

895     72.    Steppan S, Adkins R, Anderson J (2004) Phylogeny and divergence-date estimates of

896            rapid radiations in muroid rodents based on multiple nuclear genes. Systematic Biology

897            53: 533–553. doi:10.1080/10635150490468701.

898   73.   Mlynarski EE, Obergfell CJ, O'Neill MJ, O'Neill RJ (2010) Divergent patterns of

899          breakpoint reuse in Muroid rodents. Mammalian genome : official journal of the

900          International Mammalian Genome Society 21: 77–87. doi:10.1007/s00335-009-9242-1.

901   74.   Ramsdell CM, Lewandowski A a, Glenn JLW, Vrana PB, O'Neill RJ, et al. (2008)

902          Comparative genome mapping of the deer mouse (Peromyscus maniculatus) reveals

903          greater similarity to rat (Rattus norvegicus) than to the lab mouse (Mus musculus). BMC

904          Evolutionary Biology 8: 65. doi:10.1186/1471-2148-8-65.

905   75.   Ferreira de Carvalho J, Poulain J, Da Silva C, Wincker P, Michon-Coudouel S, et al.

906          (2013) Transcriptome de novo assembly from next-generation sequencing and

907          comparative analyses in the hexaploid salt marsh species Spartina maritima and Spartina

908          alterniflora (Poaceae). Heredity 110: 181–193. doi:10.1038/hdy.2012.76.

909   76.   MacManes MD, Lacey EA (2012) The Social Brain: Transcriptome Assembly and

910          Characterization of the Hippocampus from a Social Subterranean Rodent, the Colonial

911          Tuco-Tuco (Ctenomys sociabilis). PLoS ONE 7: e45524.

912          doi:10.1371/journal.pone.0045524.

913   77.   Shaw TI, Srivastava A, Chou W-C, Liu L, Hawkinson A, et al. (2012) Transcriptome

914          Sequencing and Annotation for the Jamaican Fruit Bat (Artibeus jamaicensis). PLoS ONE

915          7: e48472. doi:10.1371/journal.pone.0048472.

916   78.   Becker J, Hackl M, Rupp O, Jakobi T, Schneider J, et al. (2011) Unraveling the Chinese

917          hamster ovary cell line transcriptome by next-generation sequencing. Journal of

918          Biotechnology 156: 227–235. doi:10.1016/j.jbiotec.2011.09.014.

919   79.   Xu X, Nagarajan H, Lewis NE, Pan S, Cai Z, et al. (2011) The genomic sequence of the

920          Chinese hamster ovary (CHO)-K1 cell line. Nature Biotechnology 29: 735–741.

921          doi:10.1038/nbt.1932.

922    80.    Yang X, Schadt EE, Wang S, Wang H, Arnold AP, et al. (2006) Tissue-specific

923            expression and regulation of sexually dimorphic genes in mice. Genome Research 16:

924            995–1004. doi:10.1101/gr.5217506.

925    81.    Chapple R (2012) The developmental liver transcriptome of Rattus norvegicus University

926            of Missouri.

927    82.    Chrast R, Scott H, Papasavvas M (2000) The Mouse Brain Transcriptome by SAGE:

928            Differences in Gene Expression between P30 Brains of the Partial Trisomy 16 Mouse

929            Model of Down Syndrome (Ts65Dn) and Normals. Genome Research 10: 2006–2021.

930            doi:10.1101/gr.158500.

931    83.    Shima JE, McLean DJ, McCarrey JR, Griswold MD (2004) The murine testicular

932            transcriptome: characterizing gene expression in the testis during the progression of

933            spermatogenesis. Biology of Reproduction 71: 319–330.

934            doi:10.1095/biolreprod.103.026880.

935    84.    Pontius JU, Wagner L, Schuler GD (2003) UniGene: a unified view of the transcriptome.

936            The NCBI Handbook: Bethesda (MD): National Center for Biotechnology Information.

937    85.    Glenn JLW, Chen C-F, Lewandowski A, Cheng C-H, Ramsdell CM, et al. (2008)

938            Expressed sequence tags from Peromyscus testis and placenta tissue: analysis, annotation,

939            and utility for mapping. BMC Genomics 9: 300. doi:10.1186/1471-2164-9-300.

940    86.    Barbazuk WB, Emrich SJ, Chen HD, Li L, Schnable PS (2007) SNP discovery via 454

941            transcriptome sequencing. The Plant Journal : For Cell and Molecular Biology 51: 910–

942            918. doi:10.1111/j.1365-313X.2007.03193.x.

943    87.    Collins L, Biggs P, Voelckel C, Joly S (2008) An Approach to Transcriptome Analysis of

944            Non-Model Organisms Using Short-Read Sequences. Genome Informatics 21: 3–14.

945    88.    De Wit P, Pespeni MH, Ladner JT, Barshis DJ, Seneca F, et al. (2012) The simple fool's

946            guide to population genomics via RNA-Seq: an introduction to high-throughput

947    sequencing data analysis. Molecular Ecology Resources 12: 1058–1067.

948    doi:10.1111/1755-0998.12003.

949    89.    Nielsen R, Paul JS, Albrechtsen A, Song YS (2011) Genotype and SNP calling from next-

950    generation sequencing data. Nature Reviews Genetics 12: 443–451. doi:10.1038/nrg2986.

951    90.    Altmann A, Weber P, Bader D, Preuß M, Binder EB, et al. (2012) A beginners guide to

952    SNP calling from high-throughput DNA-sequencing data. Human Genetics: 1541–1554.

953    doi:10.1007/s00439-012-1213-z.

954    91.    Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, et al. (2009) The Sequence

955    Alignment/Map format and SAMtools. Bioinformatics (Oxford, England) 25: 2078–2079.

956    doi:10.1093/bioinformatics/btp352.

957    92.    McCormack JE, Hird SM, Zellmer AJ, Carstens BC, Brumfield RT (2011) Applications of

958    next-generation sequencing to phylogeography and phylogenetics. Molecular

959    Phylogenetics and Evolution. doi:10.1016/j.ympev.2011.12.007.

960    93.    Kofler R, Orozco-terWengel P, De Maio N, Pandey RV, Nolte V, et al. (2011)

961    PoPoolation: a toolbox for population genetic analysis of next generation sequencing data

962    from pooled individuals. PloS One 6: e15925. doi:10.1371/journal.pone.0015925.

963    94.    Oleksyk TK, Smith MW, O'Brien SJ (2010) Genome-wide scans for footprints of natural

964    selection. Philosophical transactions of the Royal Society of London Series B, Biological

965    sciences 365: 185–205. doi:10.1098/rstb.2009.0219.

966    95.    Baldo L, Santos ME, Salzburger W (2011) Comparative transcriptomics of eastern

967    African cichlid fishes shows signs of positive selection and a large contribution of

968    untranslated regions to genetic diversity. Genome Biology and Evolution 3: 443–455.

969    doi:10.1093/gbe/evr047.

970    96.    McDonald JH, Kreitman M (1991) Adaptive protein evolution at the Adh locus in

971    Drosophila. Nature 351: 652–654. doi:10.1038/351652a0.

972     97.     Smith NGC, Eyre-Walker A (2002) Adaptive protein evolution in Drosophila. Nature 415:

973             1022–1024. doi:10.1038/4151022a.

974     98.     Andersen KG, Shylakhter I, Tabrizi S, Grossman SR, Happi CT, et al. (2012) Genome-

975             wide scans provide evidence for positive selection of genes implicated in Lassa fever.

976             Philosophical Transactions of the Royal Society of London Series B, Biological Sciences

977             367: 868–877. doi:10.1098/rstb.2011.0299.

978     99.     Kosiol C, Vinar T, da Fonseca RR, Hubisz MJ, Bustamante CD, et al. (2008) Patterns of

979             positive selection in six Mammalian genomes. PLoS Genetics 4: e1000144.

980             doi:10.1371/journal.pgen.1000144.

981     100.    Sackton TB, Lazzaro BP, Schlenke TA, Evans JD, Hultmark D, et al. (2007) Dynamic

982             evolution of the innate immune system in Drosophila. Nature Genetics 39: 1461–1468.

983             doi:10.1038/ng.2007.60.

984     101.    Bradley CA, Altizer S (2007) Urbanization and the ecology of wildlife diseases. Trends in

985             Ecology & Evolution 22: 95–102. doi:10.1016/j.tree.2006.11.001.

986     102.    Turner LM, Chuong EB, Hoekstra HE (2008) Comparative analysis of testis protein

987             evolution in rodents. Genetics 179: 2075–2089. doi:10.1534/genetics.107.085902.

988     103.    Janssens TKS, Roelofs D, van Straalen NM (2009) Molecular mechanisms of heavy metal

989             tolerance and evolution in invertebrates. Insect Science 16: 3–18. doi:10.1111/j.1744-

990             7917.2009.00249.x.

991     104.    Somers CM, Yauk CL, White P a, Parfett CLJ, Quinn JS (2002) Air pollution induces

992             heritable DNA mutations. Proceedings of the National Academy of Sciences of the United

993             States of America 99: 15904–15907. doi:10.1073/pnas.252499499.

994     105.    Somers CM, Cooper DN (2009) Air pollution and mutations in the germline: are humans

995             at risk? Human Genetics 125: 119–130. doi:10.1007/s00439-008-0613-6.

996    106.   McGuire KL, Payne SG, Palmer MI, Gillikin CM, Keefe D, et al. (2013) Digging the New

997            York City Skyline: Soil Fungal Communities in Green Roofs and City Parks. PLoS ONE

998            8: e58020. doi:10.1371/journal.pone.0058020.

999    107.   Wirgin I, Roy NK, Loftus M, Chambers RC, Franks DG, et al. (2011) Mechanistic basis

1000           of resistance to PCBs in Atlantic tomcod from the Hudson River. Science (New York,

1001           NY) 331: 1322–1325. doi:10.1126/science.1197296.

1002   108.   Carroll MC (2004) The complement system in regulation of adaptive immunity. Nature

1003           immunology 5: 981–986. doi:10.1038/ni1113.

1004   109.   Su T, Ding X (2004) Regulation of the cytochrome P450 2A genes. Toxicology and

1005           Applied Pharmacology 199: 285–294. doi:10.1016/j.taap.2003.11.029.

1006   110.   Büntge A (2010) Tracing Signatures of Positive Selection in Natural Populations of the

1007           House Mouse Christian-Albrechts-Universität, Kiel.

1008   111.   Fu C, Xiong J, Miao W (2009) Genome-wide identification and characterization of

1009           cytochrome P450 monooxygenase genes in the ciliate Tetrahymena thermophila. BMC

1010           Genomics 10: 208. doi:10.1186/1471-2164-10-208.

1011   112.   Hughes AL (2007) Looking for Darwin in all the wrong places: the misguided quest for

1012           positive selection at the nucleotide sequence level. Heredity 99: 364–373.

1013           doi:10.1038/sj.hdy.6801031.

1014   113.   Zhai W, Nielsen R, Goldman N, Yang Z (2012) Looking for Darwin in Genomic

1015           Sequences--Validity and Success of Statistical Methods. Molecular Biology and Evolution

1016           29: 2889–2893. doi:10.1093/molbev/mss104.

1017   114.   Kryazhimskiy S, Plotkin JB (2008) The population genetics of dN/dS. PLoS Genetics 4:

1018           e1000304. doi:10.1371/journal.pgen.1000304.

1019  115.  Prud'homme B, Gompel N, Carroll SB (2007) Emerging principles of regulatory

1020        evolution. Proceedings of the National Academy of Sciences of the United States of

1021        America 104: 8605–8612. doi:10.1073/pnas.0700488104.

1022  116.  Grossman SR, Shylakhter I, Karlsson EK, Byrne EH, Morales S, et al. (2010) A

1023        composite of multiple signals distinguishes causal variants in regions of positive selection.

1024        Science (New York, NY) 327: 883–886. doi:10.1126/science.1183863.

1025  117.  Li J, Li H, Jakobsson M, Li S, Sjödin P, et al. (2012) Joint analysis of demography and

1026        selection in population genetics: where do we stand and where could we go? Molecular

1027        Ecology 28: 28–44. doi:10.1111/j.1365-294X.2011.05308.x.

1028  118.  Sikes RS, Gannon WL (2011) Guidelines of the American Society of Mammalogists for

1029        the use of wild mammals in research. Journal of Mammalogy 92: 235–253.

1030        doi:10.1644/10-MAMM-F-355.1.

1031  119.  Babik W, Stuglik M, Qi W, Kuenzli M, Kuduk K, et al. (2010) Heart transcriptome of the

1032        bank vole (Myodes glareolus): towards understanding the evolutionary variation in

1033        metabolic rate. BMC Genomics 11: 390. doi:10.1186/1471-2164-11-390.

1034  120.  Huang X, Madan A (1999) CAP3: A DNA sequence assembly program. Genome

1035        Research 9: 868–877.

1036  121.  Blankenberg D, Von Kuster G, Coraor N, Ananda G, Lazarus R, et al. (2010) Galaxy: a

1037        web-based genome analysis tool for experimentalists. Current Protocols in Molecular

1038        Biology 89: 19.10.1–19.10.21. doi:10.1002/0471142727.mb1910s89.

1039  122.  Martin M (2011) Cutadapt removes adapter sequences from high-throughput sequencing

1040        reads. EMBnet 17: 10–12.

1041  123.  Kent WJ (2002) BLAT---The BLAST-Like Alignment Tool. Genome Research 12: 656–

1042        664. doi:10.1101/gr.229202.

1043    124.    Li H, Durbin R (2010) Fast and accurate long-read alignment with Burrows-Wheeler

1044            transform. Bioinformatics (Oxford, England) 26: 589–595.

1045            doi:10.1093/bioinformatics/btp698.

1046    125.    Conesa A, Götz S, García-Gómez JM, Terol J, Talón M, et al. (2005) Blast2GO: a

1047            universal tool for annotation, visualization and analysis in functional genomics research.

1048            Bioinformatics (Oxford, England) 21: 3674–3676. doi:10.1093/bioinformatics/bti610.

1049    126.    Götz S, García-Gómez JM, Terol J, Williams TD, Nagaraj SH, et al. (2008) High-

1050            throughput functional annotation and data mining with the Blast2GO suite. Nucleic Acids

1051            Research 36: 3420–3435. doi:10.1093/nar/gkn176.

1052    127.    Quinlan AR, Hall IM (2010) BEDTools: a flexible suite of utilities for comparing

1053            genomic features. Bioinformatics (Oxford, England) 26: 841–842.

1054            doi:10.1093/bioinformatics/btq033.

1055    128.    Edgar RC (2010) Search and clustering orders of magnitude faster than BLAST.

1056            Bioinformatics (Oxford, England) 26: 2460–2461. doi:10.1093/bioinformatics/btq461.

1057    129.    Zhang Z, Li J, Zhao X-Q, Wang J, Wong GK-S, et al. (2006) KaKs_Calculator:

1058            calculating Ka and Ks through model selection and model averaging. Genomics,

1059            Proteomics & Bioinformatics 4: 259–263. doi:10.1016/S1672-0229(07)60007-2.

1060    130.    Montoya-Burgos JI (2011) Patterns of positive selection and neutral evolution in the

1061            protein-coding genes of Tetraodon and Takifugu. PLoS ONE 6: e24800.

1062            doi:10.1371/journal.pone.0024800.

1063    131.    Ranwez V, Harispe S, Delsuc F, Douzery EJP (2011) MACSE: Multiple Alignment of

1064            Coding SEquences accounting for frameshifts and stop codons. PloS One 6: e22594.

1065            doi:10.1371/journal.pone.0022594.

1066    132.    Kosakovsky Pond SL, Posada D, Gravenor MB, Woelk CH, Frost SDW (2006)

1067            Automated phylogenetic detection of recombination using a genetic algorithm. Molecular

1068            Biology and Evolution 23: 1891–1901. doi:10.1093/molbev/msl051.

1069    133.    Delport W, Poon AFY, Frost SDW, Kosakovsky Pond SL (2010) Datamonkey 2010: a

1070            suite of phylogenetic analysis tools for evolutionary biology. Bioinformatics (Oxford,

1071            England) 26: 2455–2457. doi:10.1093/bioinformatics/btq429.

1072    134.    Librado P, Rozas J (2009) DnaSP v5: a software for comprehensive analysis of DNA

1073            polymorphism data. Bioinformatics (Oxford, England) 25: 1451–1452.

1074            doi:10.1093/bioinformatics/btp187.

1075

1076

**Figure Legends**

1077    **Figure Legends**

1078    **Figure 1. Frequency of contig lengths for three transcriptome assembly methods.** Inset:

1079    Zoomed-in view of frequency of longer assembled contigs from 1,500-3,000 bp. Blue line =

1080    Newbler cDNA, Red line = Newbler genome, Green line = Cap3.

1081

1082    **Figure 2. Transcriptome alignment to reference rodent genomes.** Number and distribution of

1083    contigs from *P. leucopus* transcriptome (Newbler cDNA assembly) that aligned to each

1084    chromosome of the. (a) *Rattus norvegicus*. Blue = total number of genes per chromosome for

1085    *Rattus*.  Red = number of aligned *Peromyscus* isotigs per *Rattus* chromosome. (b) *Mus musculus*.

1086    Blue = total number of genes per chromosome for *Mus*.  Red = number of aligned *Peromyscus*

1087    isotigs per *Mus* chromosome.

1088

1089    **Figure 3. Annotation of final reference transcriptome.** Number of assembled *P. leucopus*

1090    contigs from four different tissue types that had significant hits with known proteins on

1091    BLASTX, and GO term annotations from reference databases using Blast2Go; Blue = Total

1092    number of contigs, Red = BLASTX hits, Green = number of annotated contigs.

1093

1094    **Figure 4. Over-represented GO terms from pairwise tissue comparisons (FDR ≤0.05).** (a)

1095    Comparison of brain transcriptome to liver and gonad.  (b) Comparison of liver to brain and

1096    gonad.  (c) Comparison of gonad to liver and brain.

1097

1098    **Figure 5. Non-synonymous ($p_N$) SNP substitutions plotted vs. synonymous ($p_S$) substitutions**

1099    **for 354 genes.**  Each circle represents one unique assembled contig. (a) Pairwise comparisons for

1100    all urban populations. (b) Pairwise comparisons for urban to rural populations. The dashed line

46

1101 denotes $p_N/p_S = 1$, and circles above the line ($p_N/p_S > 1$) indicate candidates for positive

1102 selection. The solid line shows the slope for $p_N/p_S = 0.5$.

1103

1104 **Figure 6. Location and number of individuals collected from five populations in the NYC**

1105 **metropolitan area.** Urban populations are in shades of blue; light blue = male; dark blue =

1106 female. Rural population in orange and brown; orange = male; brown = female. Areas shaded

1107 red on the map indicate degree of urbanization (i.e. impermeable surface cover such as roads and

1108 rooftops) and green areas indicate vegetation cover from the 2006 National Landcover Database.

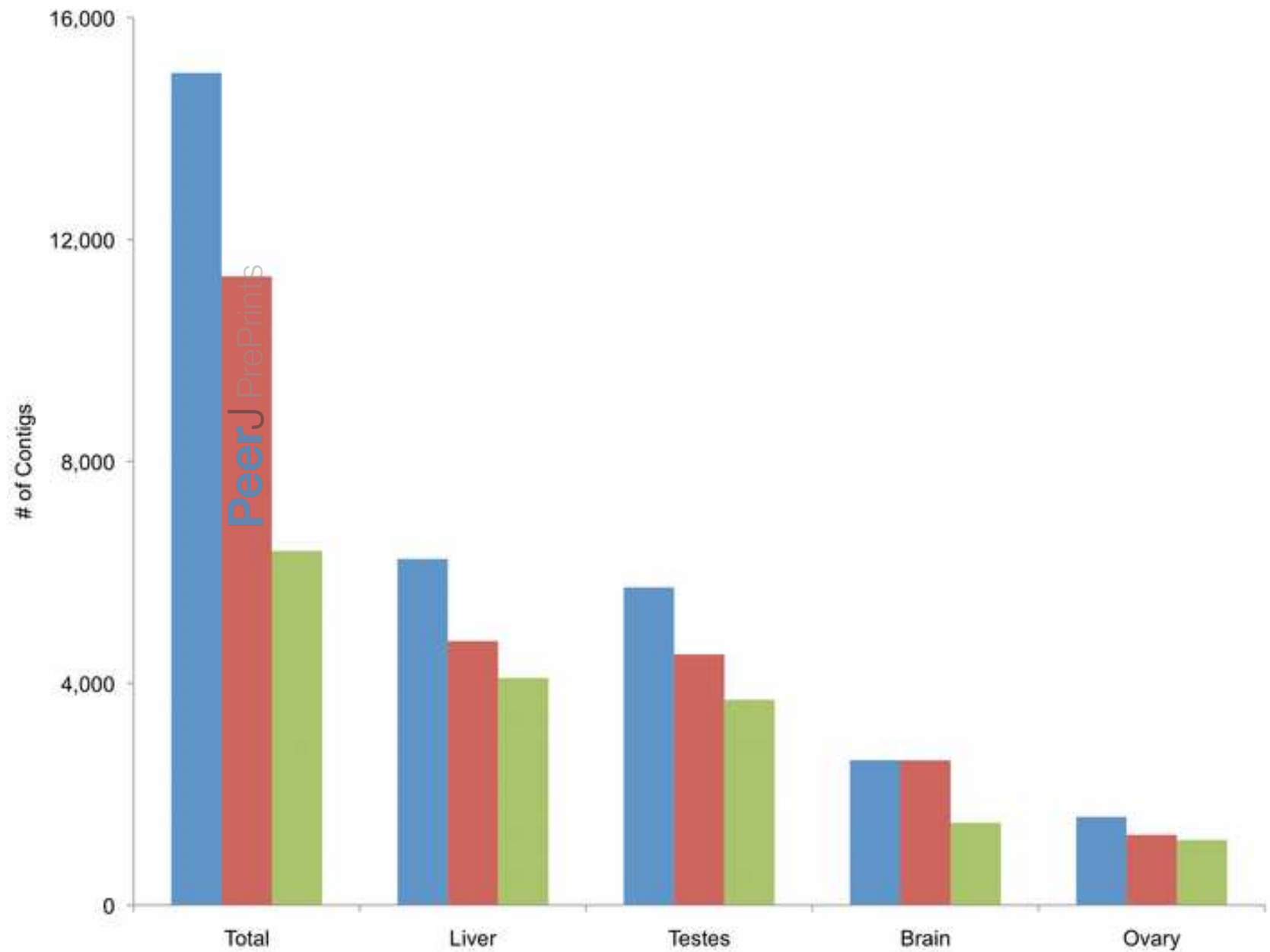1109 (CP = Central Park; NYBG = New York Botanical Gardens; RR = Ridgewood Reservoir; FM =

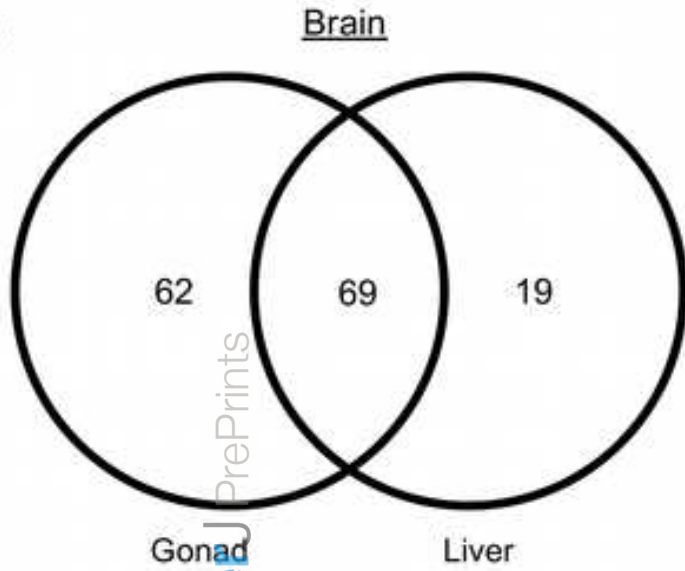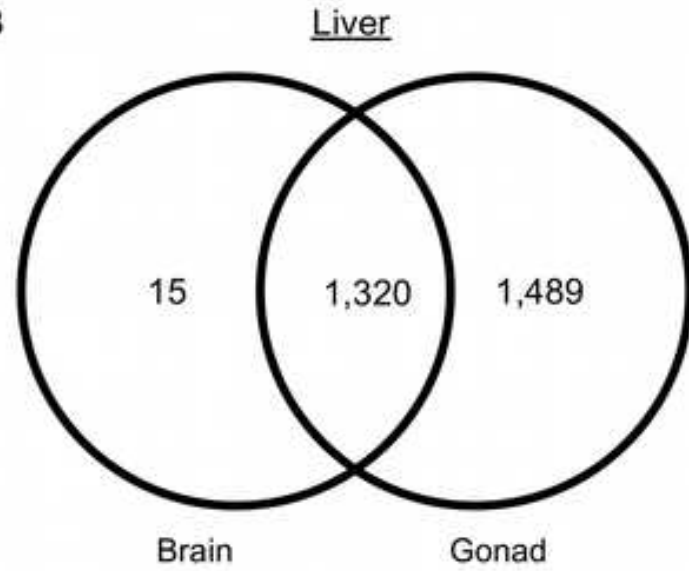1110 Flushing Meadows-Willow Lake; HP = Harriman State Park).

1111

1112
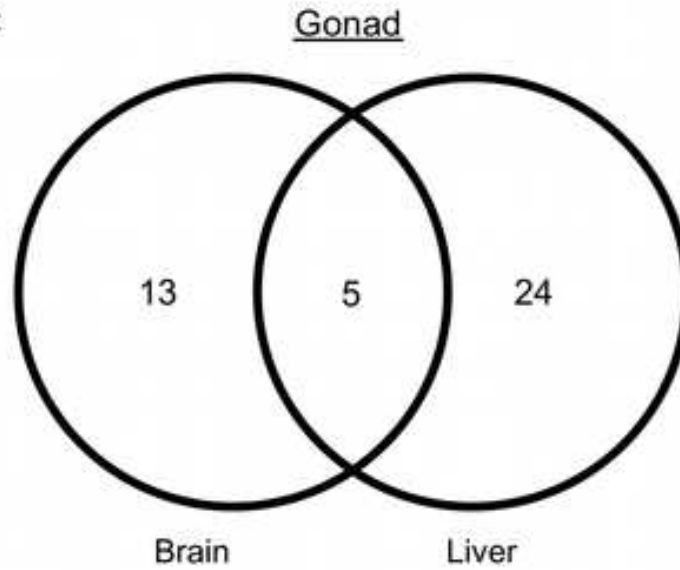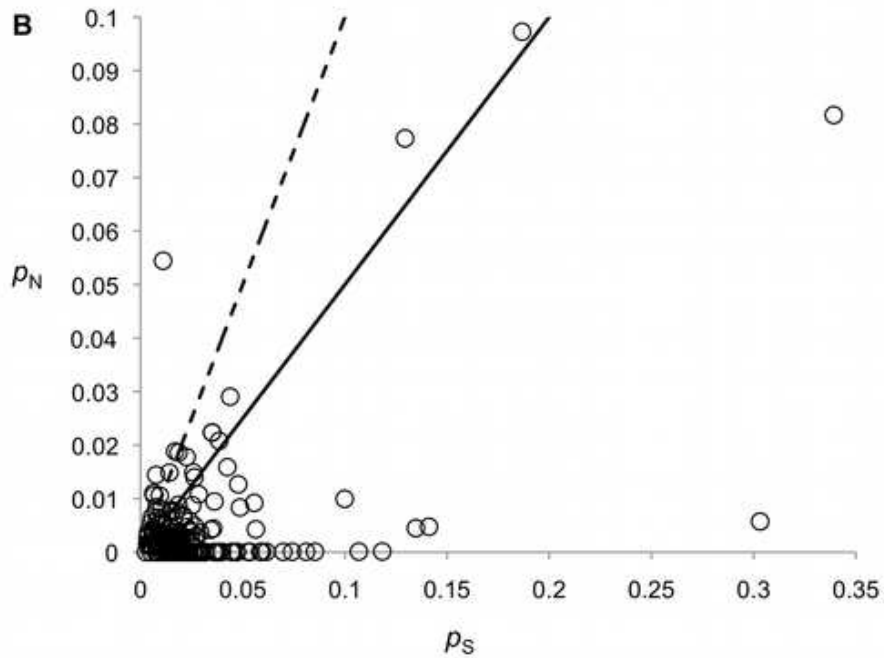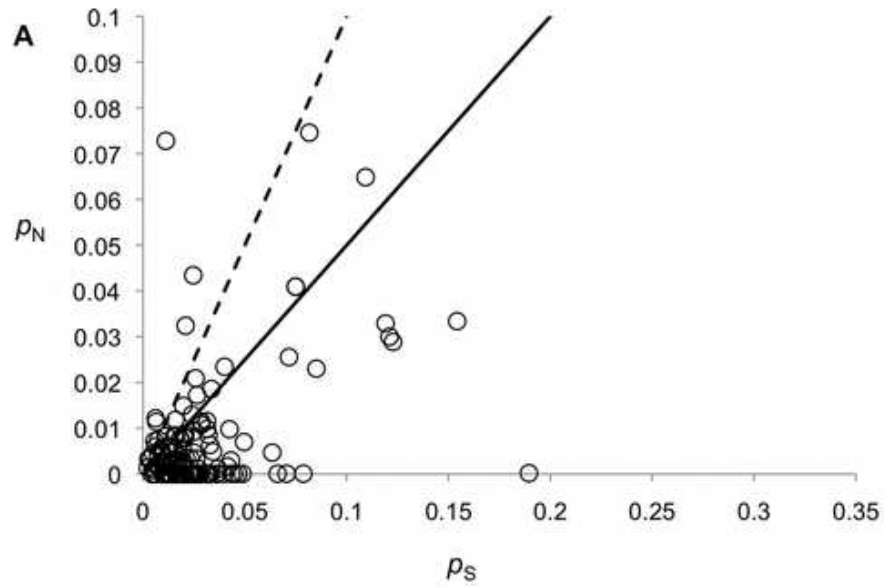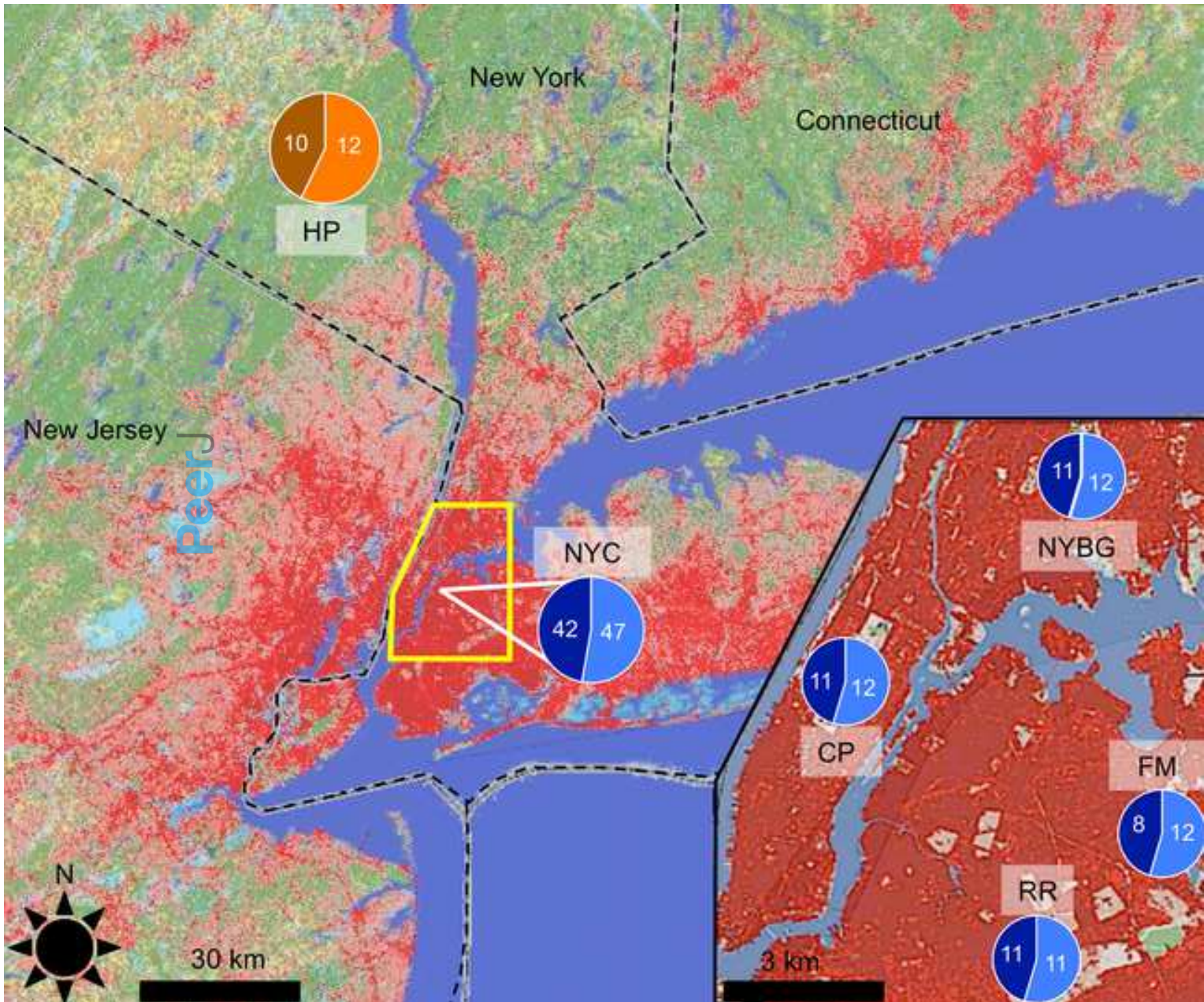
1113

## Tables

**Table 1. Results of transcriptome assembly using three different approaches.**

| Assembly Method | No. Contigs | Mean Contig Length (bp) | Median Contig Length (bp) | N50[*] | Length (Mb)[**] |
|---|---|---|---|---|---|
| Newbler genome[a] | 20,570 | 630 ± 504 | 516 | 830 | 12.95 |
| Cap3[b] | 27,497 | 653 ± 380 | 566 | 732 | 17.95 |
| Newbler cDNA[c] | 15,004 (Isotigs) | 895 ± 752 | 683 | 1,039 | 13.42 |

1116

1117   [a]Newbler v. 2.5.3 large genomic assembly of total set of raw sequencing reads

1118   [b]Cap3 assembly using 'assembled' or 'partially assembled' reads from Newbler genome

1119   assembly

1120   [c]Newbler v. 2.5.3 cDNA assembly using 'assembled' or 'partially assembled' reads from

1121   Newbler genome assembly

1122   [*]N50, The value where half the assembly is represented by contigs of this size or longer

1123   [**]Total assembly length in Megabases.

1124

1125

1126

**Table 2. BLASTN search results of three *P. leucopus* transcriptome assemblies against**

**reference cDNA libraries from *Mus* and *Rattus*.**

| Assembly Method | Total Significant Hits; *Mus* | Total Significant Hits; *Rattus* | Gene Candidates, *Mus*([*]) | Gene Candidates, *Rattus*([*]) |
|---|---|---|---|---|
| Newbler genome | 12,932 | 12,807 | 8,568 (708 bp) | 8,080 (714 bp) |
| Cap3 | 17,333 | 16,792 | 11,662 (623 bp) | 10,938 (638 bp) |
| Newbler cDNA | 10,699 | 10,094 | 7,048 (823 bp) | 6,814 (847 bp) |

1129    * = Average alignment length in base pairs

1130    Total significant hits represent sequence identity $\geq$ 80%, alignment length $\geq$ 50% of the total
1131    length of either the query or subject sequence, and *e*-value $\leq 10^{-5}$.  Gene candidates represent
1132    significant hits where one query sequence matches one subject gene and *vice versa*.

1133

1134

**Table 3. Over-represented GO terms for individual tissue types from Fisher's Exact tests**
**(FDR ≤ 0.5) in Blast2Go.**

| GO term | FDR | # Sequences |
|---|---|---|
| **Liver** | | |
| ATP binding | 5.31E-24 | 184 |
| zinc ion binding | 5.93E-20 | 154 |
| transcription factor complex | 3.91E-19 | 148 |
| electron carrier activity | 8.53E-18 | 251 |
| structural constituent of ribosome | 5.51E-15 | 117 |
| soluble fraction | 2.35E-12 | 97 |
| microsome | 1.53E-10 | 83 |
| protein homodimerization activity | 2.75E-10 | 81 |
| oxygen binding | 1.97E-09 | 93 |
| perinuclear region of cytoplasm | 9.92E-09 | 69 |
| GTP binding | 7.64E-08 | 62 |
| GTPase activity | 2.82E-05 | 42 |
| ubiquitin-protein ligase activity | 2.82E-05 | 42 |
| NADH dehydrogenase (ubiquinone) activity | 5.01E-05 | 40 |
| drug binding | 6.65E-05 | 39 |
| sequence-specific DNA binding | 6.65E-05 | 39 |
| double-stranded DNA binding | 8.90E-05 | 38 |
| mitochondrial respiratory chain complex I | 1.18E-04 | 37 |
| transcription coactivator activity | 1.18E-04 | 37 |
| catalytic step 2 spliceosome | 1.58E-04 | 36 |
| **Brain** | | |
| protein complex | 1.27E-06 | 569 |
| plasma membrane | 4.30E-92 | 567 |
| signal transduction | 2.15E-39 | 525 |
| cytosol | 1.79E-08 | 411 |
| cell differentiation | 5.07E-28 | 372 |
| anatomical structure morphogenesis | 1.89E-30 | 291 |
| cell death | 1.78E-06 | 247 |
| cell-cell signaling | 2.79E-61 | 232 |
| ion transport | 3.12E-17 | 209 |
| cytoplasmic membrane-bounded vesicle | 1.33E-22 | 197 |
| golgi apparatus | 1.51E-10 | 168 |
| cytoskeleton organization | 9.13E-13 | 145 |
| cellular homeostasis | 9.82E-16 | 134 |
| behavior | 6.72E-28 | 133 |
| calcium ion binding | 7.69E-13 | 109 |
| actin binding | 3.54E-15 | 93 |
| response to abiotic stimulus | 4.97E-08 | 88 |
| protein kinase activity | 1.61E-03 | 77 |
| ion channel activity | 5.21E-17 | 62 |
| motor activity | 8.38E-06 | 48 |
| **Gonads** | | |
| nucleic acid binding | 1.87E-08 | 1101 |
| nuclear chromosome | 9.86E-06 | 119 |
| reproduction | 1.92E-06 | 680 |
| RNA binding | 6.70E-04 | 637 |
| viral reproduction | 1.74E-02 | 339 |

1137

1138 GO terms have been reduced to their most specific terms. Only common GO terms over
1139 represented for one tissue compared to the other two tissues are shown. The top 20 terms are
1140 shown, see Table S2 for full list of GO annotations.

1141

**Table 4. Candidate loci exhibiting $p_N / p_S > 1$.**

| | Sequence name | $p_N / p_S$ | Gene name | Gene function |
|---|---|---|---|---|
| Pairwise Urban:Rural Comparisons | | | | |
| | HP_contig01773 | 1.01 | Translocation protein SEC62 | Post-translational protein translocation into the endoplasmic reticulum; plasma membrane protein |
| | HP_contig02632 | 1.05 | 39S ribosomal protein L51 | Part of mitochondrial ribosomal large subunit (39S); involved in protein translation |
| | HP_contig02656 | 1.07 | Histone H1-like protein in spermatids 1 | Transcriptional regulation and / or chromatin remodeling through DNA binding during spermatogenesis |
| | HP_contig01778 | 1.12 | PHD finger protein 8 | Removal of methyl groups from histones |
| | HP_contig01919 | 1.18 | Aldo-keto reductase family 1, member C12 | Xenobiotic metabolism; oxidation-reduction process |
| | HP_contig00870 | 1.74 | Camello-like 1 | Metabolic process; mitochondrial inner membrane protein |
| | HP_contig01783 | 1.89 | Cytochrome P450 2A15 | Metabolic process; testosterone 7a-hydroxylase activity |
| Pairwise Urban:Urban Comparisons | | | | |
| | CP_contig00473 | 1.23 | Fibrinogen alpha chain | Glycoprotein circulating in the blood; functions in blood coagulation and part of the most abundant component of blood clots |
| | CP_contig01204 | 1.55 | Solute carrier organic anion transporter family member 1A5 | Membrane protein; transports hormones; facilitates intestinal absorption of bile acids and renal uptake of indoxyl sulfate |
| | CP_contig00256 | 1.76 | Serine protease inhibitor a3c | Bind to proteases and inhibit proteolysis; often involved in blood coagulation and inflammation |
| | CP_contig00748 | 1.97 | Alpha-1-acid glycoprotein 1 | Transport protein in the blood stream; binds and distributes synthetic drugs throughout body; modulates innate immune response |

51

1143　Table 5. McDonald-Kreitman tests for candidate genes with $p_N/p_S > 1$. Comparison of the
1144　amount of polymorphisms in candidate ORFs to that of the divergence in orthologous
1145　genes between *Peromyscus* and *Rattus norvegicus*. P-values were generated from Fisher's
1146　Exact Test.

1147

| Gene Name | Polymorphisms Non-synonymous ($Pn$) | Synonymous ($Ps$) | Ratio ($Pn/Ps$) | Divergence Non-synonymous ($Dn$) | Synonymous ($Ds$) | Ratio ($Dn/Ds$) | Neutrality Index | P-value |
|---|---|---|---|---|---|---|---|---|
| Translocation protein SEC62 | 2 | 1 | 2 | 18 | 28 | 0.64 | 3.11 | 0.55 |
| 39S ribosomal protein L51 | 4 | 1 | 4 | 15 | 32 | 0.47 | 8.53 | 0.05 |
| Histone H1-like protein in spermatids 1 | 2 | 1 | 2 | 20 | 12 | 1.67 | 1.20 | 1.00 |
| PHD finger protein 8 | 9 | 3 | 3 | 36 | 51 | 0.71 | 4.25 | 0.03 |
| Aldo-keto reductase family 1, member C12 | 3 | 1 | 3 | 18 | 37 | 0.49 | 2.67 | 0.08 |
| Camello-like 1 | 4 | 1 | 4 | 41 | 23 | 1.78 | 2.24 | 0.65 |
| Cytochrome P450 2A15[*] | 6 | 1 | 6 | 13 | 28 | 0.46 | 12.92 | 0.01 |
| Fibrinogen alpha chain | 3 | 1 | 3 | 101 | 93 | 1.08 | 2.76 | 0.62 |
| Solute carrier organic anion transporter 1A5 | 9 | 3 | 3 | 21 | 37 | 0.57 | 5.29 | 0.02 |
| Serine protease inhibitor a3c[*] | 4 | 1 | 4 | 25 | 19 | 1.32 | 1.27 | 0.65 |
| Alpha-1-acid glycoprotein 1 | 4 | 1 | 4 | 68 | 44 | 1.55 | 2.59 | 0.65 |

1148　* = McDonald Kreitman test used *Cricetulus griseus*

1149

1150

1151

**Supporting Information**

1153     **Figure S1. Frequency distribution of depth of coverage (reads / contig).** (a) The Newbler

1154     cDNA assembly.  Red line indicates median coverage = 4.9 reads, Interquartile range (IQR) =

1155     4.1. (b) The Newbler genomic assembly, median = 4.7 reads, IQR = 4.6.  (c) The Cap3 assembly,

1156     median = 5.0 reads, IQR = 7.0.

1157

1158     **Figure S2. Distribution of species with the most top-hit BLASTX results in Blast2Go using**

1159     **the Newbler cDNA assembly as the query.**

1160

1161     **Table S1. Sequencing and assembly statistics for Newbler cDNA transcriptome assembly**

1162     **by tissue type and 454 sequencing plate.**

1163

1164     **Table S2. Full list of over represented GO terms for all tissue pairwise comparisons from**

1165     **Fisher's Exact Test (FDR ≤ 0.5). (a) Liver. (b) Brain. (c) Gonads.**

1166

1167     **Table S3. Candidate loci with $p_N/p_S$ between 0.5 and 1.**

1168