

Characterization of the transcriptome, nucleotide sequence polymorphism, and natural selection in the desert adapted mouse *Peromyscus eremicus*

Matthew D. MacManes¹, Michael B. Eisen²,

¹ Department of Molecular, Cellular and Biomedical Sciences, University of New Hampshire, Durham NH, USA

² HHMI and University of California, Berkeley, Berkeley, CA, USA

* E-mail: macmanes@gmail.com, @PeroMHC

1 Abstract

2 As a direct result of intense heat and aridity, deserts are thought to be amongst the most harsh of
 3 environments, particularly for their mammalian inhabitants. Given that osmoregulation can be
 4 challenging for these animals, with failure resulting in death, strong selection should be observed
 5 on genes related to the maintenance of water and solute balance. One such animal, *Peromyscus*
 6 *eremicus*, is native to the desert regions of the southwest United States, and may live its entire
 7 life without oral fluid intake. As a first step towards understanding the genetics that underlie
 8 this phenotype, we present here a characterization of the *P. eremicus* transcriptome. We as-
 9 say four tissues (kidney, liver, brain, testes) from a single individual, and supplement this with
 10 population level renal transcriptome sequencing from 15 additional animals. We identified a set
 11 of transcripts undergoing both purifying and balancing selection based on estimates of Tajima's
 12 D. In addition, we used the branch-site test to identify a transcript, Slc2a9 likely related to
 13 desert osmoregulation, undergoing enhanced selection in *P. eremicus* relative to a set of related
 14 non-desert rodents.

16 Introduction

17 Deserts are widely considered one of Earth's most harsh environments. Animals living in desert
 18 environments are forced to endure intense heat and drought, and as a result, species living in
 19 these environments are likely to possess specialized mechanisms to deal with them. While living
 20 in deserts likely involves a large number of adaptive traits, the ability to osmoregulate – to
 21 maintain the proper water and electrolyte balance – appears to be paramount [1]. Indeed, the
 22 maintenance of water balance in animals is one of the most important physiologic processes for
 23 all organisms, whether they be desert inhabitants or not. Most animals are exquisitely sensitive
 24 to changes in osmolality, with slight derangement eliciting physiologic compromise. When the
 25 loss of water exceeds dietary intake, dehydration - and in extreme cases, death - can occur.

Thus there has likely been strong selection for mechanisms supporting optimal osmoregulation in species that live where water is limiting. Understanding these mechanisms will significantly enhance our understanding of the physiologic processes underlying osmoregulation in extreme environments, having implications for studies of human health, conservation, and climate change.

The genes and structures responsible for the maintenance of water and electrolyte balance are well characterized in model organisms such as mice [2], rats [3–5], and humans [6–8]. These studies, many of which have been enabled by newer sequencing technologies, serve as a foundation for studies of renal genomics in non-model organisms. In particular, because researchers have long been interested in desert adaptation, a number of studies have looked at the morphology or expression of single genes in the renal tissues of desert adapted rodents *Phyllotis darwini* [9], *Psammomys obesus* [10], and *Perognathus penicillatus* [11]. More recently, full renal transcriptomes have been generated for *Dipodomys spectabilis* and *Chaetodipus baileyi* [12] as well as *Abrothrix olivacea* [13].

These studies provide a rich context for the current and future work, aimed at developing a synthetic understanding of the the genetic and genomic underpinnings of desert adaptation in rodents. As a first step, we have sequenced, assembled, and characterized the transcriptome (using four tissue types - liver, kidney, testes, brain), of a desert adapted cricetid rodent endemic to the Southwest United States [14], *Peromyscus eremicus*. These animals have a lifespan typical of small mammals, and therefore an individual may live it’s entire life without ever drinking water. These rodents have distinct advantage over other desert animals (e.g. *Dipodomys*) in that they breed readily in captivity, which enables future laboratory studies of the phenotype of interest. In addition the focal species is positioned in a clade of well known animals (e.g. *P. californicus*, *P. maniculatus* and *P. polionotus*) [15] with growing genetic and genomic resources [16–18] which together suggest that future comparative studies are possible.

While the elucidation of the mechanisms underlying adaptation to desert survival is beyond the scope of this manuscript, we aim here to lay the groundwork by characterizing the transcriptome from four distinct tissues (brain, liver, kidney, testes). These data will be included in current larger effort aimed at sequencing the entire genome. Further, via sequencing the renal tissue of a total of 15 additional animals, we characterize nucleotide polymorphism and genome wide patterns of natural selection. Together, these investigations will aid in our overarching goal – to understand the genetic bases of adaptation to deserts in *P. eremicus*.

Materials and Methods

Animal Collection and Study Design

To begin to understand how genes may underlie desert adaptation, we collected 16 individuals from a single population *P. eremicus* over a two year time period (2012-2013). These individuals were captured in live traps, then euthanized using isoflurane overdose and decapitation. Immediately post-mortem, the abdominal and pelvic organs were removed, cut in half (in the case of kidney), placed in RNAlater and flash frozen in liquid Nitrogen. Removal of the brain, with similar preservation techniques, followed that. Time from euthanasia to removal of all organs never exceeded five minutes. Samples were transferred to a -80C freezer at a later date. These procedures were approved by the University of California Berkeley Animal Care and Use Committee and follow guidelines established by the American Society of Mammalogy for the use of wild animals in research [19].

RNA extraction and Sequencing

Total RNA was extracted from each tissue using a TRIzol extraction (Invitrogen) following the manufacturers instructions. Because preparation of an RNA library suitable for sequencing is dependent on having high quality, intact RNA, a small aliquot of each total RNA extract was analyzed on a Bioanalyzer 2100 (Agilent, Palo Alto, CA, USA). Following confirmation of sample quality, the reference sequencing libraries were made using the TruSeq stranded RNA prep kit (Illumina), while an unstranded TruSeq kit was used to construct the other sequencing libraries. A unique index was ligated to each sample to allow for multiplexed sequencing. Reference libraries (n=4 tissue types) were then pooled to contain equimolar quantities of each individual library and submitted for Illumina sequencing using two lanes of 150nt paired end sequencing using the rapid-mode of the HiSeq 2500 sequencer at The Hubbard Center for Genome Sciences (University of New Hampshire). The remaining 15 libraries were similarly multiplexed and sequenced in a mixture of 100nt paired and single end sequencing runs across several lanes of an Illumina HiSeq 2000 at the Vincent G. Coates Genome Center (University of California, Berkeley).

Sequence Data Preprocessing and Assembly

The raw sequence reads corresponding to the four tissue types were error corrected using the software *bleed* [20], using *kmer*=25, based on the developers default recommendations. The error corrected sequence reads were adapter and quality trimmed following recommendations from MacManes [21] and Mbandi [22]. Specifically, adapter sequence contamination was removed, and low quality nucleotides (defined as *PHRED* <2) were removed using the program

93 Trimmomatic version 0.32 [23]. Reads from each tissue were assembled using the Trinity ver-
 94 sion released 17 July 2014 [24]. We used flags indicating the stranded nature of sequencing
 95 reads and set maximum allowable physical distance between read pairs to 999nt. The assembly
 96 was conducted on a linux workstation with 64 cores and 512Gb RAM. To filter the raw se-
 97 quence assembly, we downloaded *Mus musculus* cDNA and ncRNA datasets from Ensembl ([ftp:](ftp://ftp.ensembl.org/pub/release-75 fasta/mus_musculus/)
 98 [//ftp.ensembl.org/pub/release-75 fasta/mus_musculus/](ftp://ftp.ensembl.org/pub/release-75 fasta/mus_musculus/)), and the *Peromyscus manica-*
 99 *tus* reference transcriptome from NCBI ([ftp://ftp.ncbi.nlm.nih.gov/genomes/Peromyscus_](ftp://ftp.ncbi.nlm.nih.gov/genomes/Peromyscus_maniculatus_bairdii/RNA/)
 100 [maniculatus_bairdii/RNA/](ftp://ftp.ncbi.nlm.nih.gov/genomes/Peromyscus_maniculatus_bairdii/RNA/)). We used a blastN procedure (default settings, eval set to
 101 10^{-10}) to identify contigs in the *P. eremicus* dataset that are likely biological in origin. This
 102 procedure, when a reference dataset is available, retains more putative transcripts than a strat-
 103 egy employing expression-based filtering (remove if TMP < 1, e.g. [25]) of the raw assembly. We
 104 then concatenated the filtered assemblies from each tissue into a single file, reducing redundancy
 105 using the software cd-hit-est [26] using default setting except that sequences were clustered based
 106 on 95% sequence similarity. This multi-fasta file was used for all subsequent analyses including
 107 annotation and mapping.

108

109 Assembled Sequence Annotation

110 The filtered assemblies were annotated using default settings of the blastN algorithm [27]
 111 against the Ensembl cDNA and ncRNA datasets described above, downloaded on 1 August
 112 2014. Amongst other things, the Ensembl transcript identifiers were used in the analysis
 113 of gene ontology, conducted in the PANTHER package [28]. Next, because rapidly evolv-
 114 ing nucleotide sequences may evade detection by blast algorithms, we used HMMER3 [29]
 115 to search for conserved protein domains contained in the dataset using the Pfam database
 116 [30]. Lastly, we extracted putative coding sequences using Transdecoder version 4Jul2014
 117 (<http://transdecoder.sourceforge.net/>)

118

119 To identify patterns of gene expression unique to each tissue type, we mapped sequence reads
 120 from each tissue type to the reference assembly using bwa-mem [31]. We estimated expression
 121 individually for the four tissues using default settings of the software eXpress [32]. Interesting
 122 patterns of expression, including instances where expression was limited to a single tissue type
 123 were identified and visualized.

124

Population Genomics

In addition to the reference individual sequenced at four different tissue types, we sequenced 15 other conspecific individuals from the same population, located in Palm Desert, California. Sequence data were mapped to the reference transcriptome using bwa-mem. The alignments were sorted and converted to BAM format, then passed to the program ANGSD version 0.610, which was used for calculating the folded site frequency spectrum (SFS) and Tajima's D [33].

Natural Selection

To characterize natural selection on several genes related to water and ion homeostasis, we identified several of the transcripts identified as experiencing positive selection in a recent work on desert-adapted *Dipodomys* rodents. The coding sequence corresponding to these genes, Solute Carrier family 2 member 9 (Slc2a9) and the Vitamin D3 receptor (Vdr), were extracted from the dataset, aligned using the software MACSE [34] to homologous sequences in *Mus musculus*, *Rattus norvegicus*, *Peromyscus maniculatus*, and *Homo sapiens* identified by the conditional reciprocal best blast procedure (CRBB, [35]). An unrooted gene tree was constructed using the online resource Clustal-Omega, and together the tree and alignment were analyzed using the branch-site model (model=2, nsSites=2, fix_omega=0 versus model=2, nsSites=2, fix_omega=1, omega=1) implemented in PAML version 4.8 [36,37]. Significance was evaluated via use of the likelihood ratio test.

Results and Discussion

RNA extraction, Sequencing, Assembly, Mapping

RNA was extracted from the hypothalamus, renal medulla, testes, or liver from each individual using sterile technique. TRIzol extraction resulted in a large amount of high quality (RIN \geq 8) total RNA, which was used as input. Libraries were constructed as per the standard Illumina protocol, and were sequenced as described above. The number of reads per library varied from 56 million strand-specific paired-end reads in Peer360 kidney, to 9 million single-end reads in Peer321 (Table 1, available on the Short Read Archive accession XXX). Adapter sequence contamination and low-quality nucleotides were eliminated, which resulted in a loss of <2% of reads. These trimmed reads served as input for all downstream analyses.

Table 1

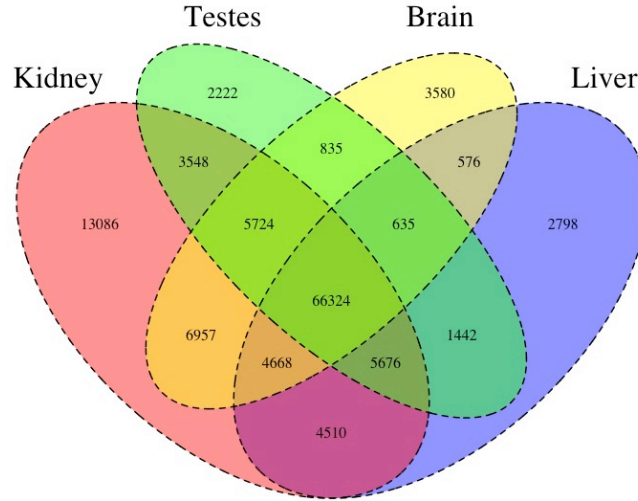
| DATASET | NUM. RAW READS |
|----------------|----------------|
| PEER360 TESTES | 32M PE/SS |
| PEER360 LIVER | 53M PE /SS |
| PEER360 KIDNEY | 56M PE/SS |
| PEER360 BRAIN | 23M PE/SS |
| PEER305 | 19M PE |
| PEER308 | 15M PE |
| PEER319 | 14M PE |
| PEER321 | 9M SE |
| PEER340 | 16M PE |
| PEER352 | 14M PE |
| PEER354 | 9M SE |
| PEER359 | 14M PE |
| PEER365 | 16M PE |
| PEER366 | 16M PE |
| PEER368 | 14M PE |
| PEER369 | 14M PE |
| PEER372 | 17M SE |
| PEER373 | 23M SE |
| PEER380 | 16M SE |
| PEER382 | 14M SE |

Table 1. The number of sequencing reads per sample, which is indicated by Peernumber.

PE=paired end, SS=strand specific, SE=single end sequencing.

Transcriptome assembly for each tissue type was accomplished using the program Trinity [24]. The raw assembly for brain, liver, testes and kidney contained 185425, 222096, 180233, and 514091 assembled sequences respectively. This assembly was filtered using a blastN procedure against the *Mus* cDNA and ncRNA, and *P. maniculatus* cDNAs which resulted in a final dataset containing 68331 brain-specific transcripts, 71041 liver-specific transcripts, 67340 testes-specific transcripts, and 113050 kidney-specific transcripts. Mapping the error-corrected adapter/quality trimmed reads to these datasets resulted in mapping 94.98% (87.01% properly paired) of brain-derived reads to the brain transcriptome, 96.07% (88.13% properly paired) of liver-derived reads to the liver transcriptome, 96.81% (85.10% properly paired) of testes-derived reads to the testes transcriptome, and 91.87% (83.77% properly paired) of kidney-derived reads to the kidney transcriptome. Together, these statistics suggest that the tissue-specific transcriptomes are of extremely high quality. All tissue-specific assemblies to be made available on Dryad.

173 **Figure 1**



174

175 Figure 1. The Venn Diagram.

176 We then estimated gene expression on each of these tissue-specific datasets, which allowed
 177 us to understand expression patterns in the multiple tissues. Specifically, we constructed a Venn
 178 diagram (Figure 1), which allowed us to visualize the proportion genes whose expression was
 179 limited to a single tissue, and those where expression was ubiquitous. 66324 transcripts are ex-
 180 pressed on all tissue types, while 13086 are uniquely expressed in kidney, 2222 in the testes, 3580
 181 in the brain, and 2798 in the liver. The kidney appears to an outlier in the number of unique
 182 sequences, though this could be the result of the recovery of more lowly expressed transcripts
 183 that may be the result of deeper sequencing.

184

185 In addition to this, we estimated mean TMP (number of transcripts per million) for all tran-
 186 scripts. Table 2 consists of the 10 genes whose mean TMP was the highest. Several genes in this
 187 list are present predominately in a single tissue type. For instance Transcript_126459, Albumin
 188 is very highly expressed in the liver, but less so in the other tissues. It should be noted, however,
 189 that making inference based on uncorrected values for TPM is not warranted. Statistical testing
 190 for differential expression was not implemented due to the fact that no replicates are available.

191

192 After expression estimation, the filtered assemblies were concatenated together, and after
 193 removal of redundancy with cd-hit-est, 123,123 putative transcripts remained (To be available
 194 on Genbank). From this filtered concatenated dataset, we extracted 71626 putative coding se-
 195 quences (72Mb, to be available on Dryad). Of these 71626 sequences, 38221 were complete

exons (containing both start and stop codons), while other were either truncated at the 5-prime end (20239 sequences), 3-prime end (6445 sequences), or were internal (6721 sequencing having neither stop nor start codon). The results of a Pfam search conducted on the predicted amino acid sequences will be found on Dryad.

Table 2

| Transcript ID | Testes | Liver | Kidney | Brain | Genbank ID | Gene ID |
|-------------------|----------|----------|----------|----------|----------------|---------|
| Transcript_83842 | 2.05E+03 | 6.40E+03 | 1.03E+04 | 5.47E+03 | DQ073446.1 | COX2 |
| Transcript_126459 | 1.43E+01 | 2.22E+04 | 2.77E+01 | 6.73E+00 | XM_006991665.1 | Alb |
| Transcript_128937 | 4.39E+00 | 1.91E+04 | 4.74E+02 | 2.23E+00 | XM_007627625.1 | Apoa2 |
| Transcript_81233 | 1.71E+03 | 5.23E+03 | 6.11E+03 | 3.08E+03 | XM_006993867.1 | Fth1 |
| Transcript_94125 | 3.67E+01 | 1.08E+04 | 2.09E+03 | 2.75E+00 | XM_006977178.1 | CytP450 |
| Transcript_119945 | 5.03E+03 | 1.15E+03 | 1.33E+03 | 3.71E+03 | XM_008686011.1 | Ubb |
| Transcript_5977 | 4.95E+00 | 1.01E+04 | 3.05E+02 | 3.58E+02 | XM_006978668.1 | Tf |
| Transcript_4057 | 2.62E+01 | 9.32E+03 | 1.34E+02 | 8.38E+01 | XM_006994871.1 | Apoc1 |
| Transcript_112523 | 4.07E+02 | 7.36E+03 | 7.78E+02 | 9.54E+02 | XM_006994872.1 | Apoe |
| Transcript_98376 | 1.98E+00 | 8.66E+03 | 1.02E+00 | 2.68E+00 | XM_006970208.1 | Ttr |

Table 2. The 10 transcripts with the highest mean TPM (transcripts per million).

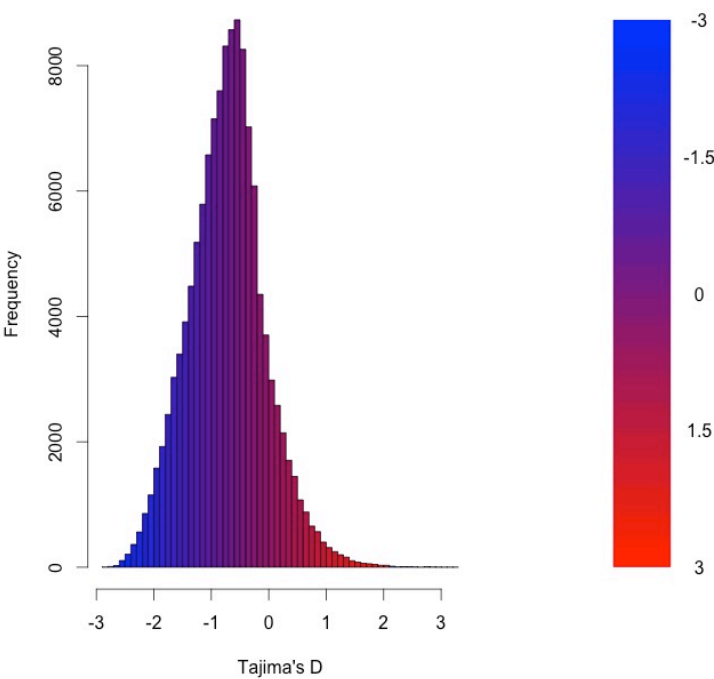
Population Genomics

As detailed above, the RNAseq data from 15 individuals were mapped to the reference transcriptome with the resulting BAM files being used as input to the software package ANGSD. The Tajima's D statistic was calculated for all transcripts covered by at least 14 of the 15 individuals. In brief, a negative Tajima's D, a result of lower than expected average heterozygosity, is often associated with purifying or directional selection, recent selective sweep or population bottleneck. In contrast, a positive value for Tajima's D represents higher than expected average heterozygosity, often associated with balancing selection.

The distribution of the estimates of Tajima's D for all assembled transcripts is shown in Figure 2. The distribution is skewed towards negative values (mean=-0.89, variance=0.58), which is likely the result of purifying selection, a model of evolution commonly invoked for coding DNA sequences [38]. Table 3 presents the 10 transcripts whose estimate of Tajima's D is greatest, while Table 4 presents the 10 transcripts whose estimate of Tajima's D is least. The former list of genes is likely to contain transcripts experiencing balancing selection in the studied population. This list includes, interestingly, genes obviously related to solute and water balance (e.g.

221 Clcnkb and a transmembrane protein gene), and those related to immune function (a interferon-
222 inducible GTPase and a Class 1 MHC gene). The latter group, containing the transcripts whose
223 estimates of Tajima's D are the smallest are likely experiencing purifying selection. Many of
224 these transcripts are involved in core regulatory functions where mutation may have strongly
225 negative fitness consequences.
226

227 **Figure 2**



228

229 Figure 2. The distribution of Tajima's D for all putative transcripts.

230 **Table 3**

231

| Transcript ID | GenBank ID | Description | Tajima's D |
|-------------------|----------------|---|------------|
| Transcript_49049 | XM.006533884.1 | heterogeneous nuclear ribonucleoprotein H1 (Hnrnp1) | 3.26 |
| Transcript_38378 | XM.006522973.1 | Son DNA binding protein (Son) | 3.19 |
| Transcript_126187 | NM.133739.2 | transmembrane protein 123 (Tmem123) | 3.02 |
| Transcript_70953 | XM.006539066.1 | chloride channel Kb (Clcnkb) | 2.96 |
| Transcript_37736 | XM.006997718.1 | h-2 class I histocompatibility antigen | 2.92 |
| Transcript_21448 | XM.006986148.1 | zinc finger protein 624-like | 2.84 |
| Transcript_47450 | NM.009560.2 | zinc finger protein 60 (Zfp60) | 2.82 |
| Transcript_122250 | XM.006539068.1 | chloride channel Kb (Clcnkb) | 2.81 |
| Transcript_78367 | XM.006496814.1 | CDC42 binding protein kinase alpha (Cdc42bpa) | 2.78 |
| Transcript_96470 | XM.006987129.1 | interferon-inducible GTPase 1-like | 2.77 |

Table 3. The 10 transcripts with the highest values for Tajima's D, which is suggestive of balancing selection.

234 **Table 4**

235

| Transcript ID | GenBank ID | Description | Tajima's D |
|-------------------|----------------|---|------------|
| Transcript_84359 | XM.006991127.1 | nuclear receptor coactivator 3 (Ncoa3) | -2.82 |
| Transcript_87121 | XM.006970128.1 | methyl-CpG binding domain protein 2 (Mbd2) | -2.82 |
| Transcript_125755 | EU053203.1 | alpha globin gene cluster | -2.78 |
| Transcript_87128 | XM.006976644.1 | membrane-associated ring finger (March5) | -2.76 |
| Transcript_55468 | XM.006978377.1 | Vpr binding protein (Vprbp) | -2.75 |
| Transcript_116042 | XM.006980811.1 | membrane associated guanylate kinase (Magi3) | -2.75 |
| Transcript_18966 | XM.006982814.1 | ubiquitin protein ligase E3 component n-recognin 5 (Ubr5) | -2.75 |
| Transcript_122204 | XM.008772511.1 | zinc finger protein 612 (Zfp612) | -2.75 |
| Transcript_100550 | XM.006971297.1 | bromodomain adjacent to zinc finger domain, 1B (Baz1b) | -2.74 |
| Transcript_33267 | XM.006975561.1 | pumilio RNA-binding family member 1 (Pum1) | -2.75 |

Table 4. The 10 transcripts with the lowest values for Tajima's D, which is suggestive of purifying or directional selection.

239 Natural Selection

240 To begin to test the hypothesis that selection on transcripts related to osmoregulation is en-
 241 hanced in the desert adapted *P. eremicus*, we implemented the branch-site test as described
 242 above, setting the sequence corresponding to *P. eremicus* for both Slc2a9 and Vdr as the fore-
 243 ground lineages in 2 distinct program executions. These two transcripts were chose specifically

because they, the former significantly, were recently linked to osmoregulation in a desert rodent [12]. The test for *Slc2a9* was highly significant ($2\Delta\text{Lnl}=51.4$, $\text{df}=1$, $p=0$), indicating enhanced selection in *P. eremicus* relative to the other lineages. The branch site test for positive selection conducted on the *Vdr* gene was non-significant ($2\Delta\text{Lnl}=0.68$, $\text{df}=1$, $p=1$). This limited analysis of selection is to be followed up by an analysis of genome wide patterns of natural selection.

250

251 Conclusions

As a direct result of intense heat and aridity, deserts are thought to be amongst the most harsh of environments, particularly for its mammalian inhabitants. Given the osmoregulation can be challenging for these animals, with failure resulting in death, strong selection should be observed on genes related to the maintenance of water and solute balance. This study aimed to characterize the transcriptome of a desert-adapted rodent species, *P. eremicus*. Specifically, we characterized the transcriptome of four tissue types (liver, kidney, brain, testes) from a single individual, and supplement this with population level renal transcriptome sequencing from 15 additional animals. We identified a set of transcripts undergoing both purifying and balancing selection based on Tajima's D. In addition, we used a branch site test to identify a transcript, likely related to desert osmoregulation, undergoing enhanced selection in *P. eremicus* relative to a set of non-desert rodents.

263

264 Acknowledgments

MDM would like to acknowledge the support of his wife and children. This work was significantly improved by comments received based on a version of the manuscript posted on biorxiv, and more generally by supporters of the Open Access and Science movements.

268 References

- 269 1. Walsberg G (2000) Small mammals in hot deserts: Some generalizations revisited. *Bio-*
270 *science* 50: 109–120.
- 271 2. Tatum R, Zhang Y, Salleng K, Lu Z, Lin JJ, et al. (2009) Renal salt wasting and chronic
272 dehydration in claudin-7-deficient mice. *AJP: Renal Physiology* 298: F24–F34.

- 273 3. Romero DG, Plonczynski MW, Welsh BL, Gomez-Sanchez CE, Zhou MY, et al. (2007)
274 Gene expression profile in rat adrenal zona glomerulosa cells stimulated with aldosterone
275 secretagogues. *Physiological Genomics* 32: 117–127.
- 276 4. Rojek A, Rojek A, Fuchtbauer E, Fuchtbauer E, Kwon T, et al. (2006) Severe urinary
277 concentrating defect in renal collecting duct-selective AQP2 conditional-knockout mice.
278 *Proceedings of The National Academy of Sciences of The United States of America* 103:
279 6037–6042.
- 280 5. Nielsen S, Chou C, Marples D, Christensen E, Kishore B, et al. (1995) Vasopressin
281 Increases Water Permeability of Kidney Collecting Duct by Inducing Translocation of
282 Aquaporin-Cd Water Channels to Plasma-Membrane. *Proceedings of The National*
283 *Academy of Sciences of The United States of America* 92: 1013–1017.
- 284 6. Mobasher A, Marples D, Young IS, Floyd RV, Moskaluk CA, et al. (2007) Distribution
285 of the AQP4 Water Channel in Normal Human Tissues: Protein and Tissue Microarrays
286 Reveal Expression in Several New Anatomical Locations, including the Prostate Gland
287 Seminal Vesicles. *Channels* 1: 30–39.
- 288 7. Bedford JJ, Leader JP, Walker RJ (2003) Aquaporin expression in normal human kidney
289 and in renal disease. *Journal of the American Society of Nephrology : JASN* 14: 2581–
290 2587.
- 291 8. Nielsen S, Kwon T (1999) Physiology and Pathophysiology of Renal Aquaporins. *Journal*
292 *of the ...*
- 293 9. Gallardo PA, Cortés A, Bozinovic F (2005) Phenotypic flexibility at the molecular and
294 organismal level allows desert-dwelling rodents to cope with seasonal water availability.
295 *Physiological and Biochemical Zoology* 78: 145–152.
- 296 10. Kaissling B, De Rouffignac C, Barrett JM, Kriz W (1975) The structural organization
297 of the kidney of the desert rodent *Psammomys obesus*. *Anatomy and embryology* 148:
298 121–143.
- 299 11. Altschuler EM, Nagle RB, Braun EJ, Lindstedt SL, Krutzsch PH (1979) Morphological
300 study of the desert heteromyid kidney with emphasis on the genus *Perognathus*. *The*
301 *Anatomical record* 194: 461–468.
- 302 12. Marra NJ, Romero a, DeWoody Ja (2014) Natural selection and the genetic basis of
303 osmoregulation in heteromyid rodents as revealed by RNA-seq. *Molecular Ecology* 23:
304 2699–2711.

- 305 13. Giorello FM, Feijoo M, a GD, Valdez L, Opazo JC, et al. (2014) Characterization of the
306 kidney transcriptome of the South American olive mouse *Abrothrix olivacea* 15: 1–10.
- 307 14. Veal R, Caire W (2001) *Peromyscus eremicus*. Mammalian Species 118: 1–6.
- 308 15. Feng BJ, Sun LD, Soltani-Arabshahi R, Bowcock AM, Nair RP, et al. (2007) Toward
309 a Molecular Phylogeny for *Peromyscus*: Evidence from Mitochondrial Cytochrome- b
310 Sequences. Journal of Mammalogy 88: 1146–1159.
- 311 16. Shorter KR, Owen A, anderson V, Hall-South AC, Hayford S, et al. (2014) Natural Ge-
312 netic Variation Underlying Differences in *Peromyscus* Repetitive and Social/Aggressive
313 Behaviors. Behavior genetics .
- 314 17. Panhuis TM, Broitman-Maduro G, Uhrig J, Maduro M, Reznick DN (2011) Analysis of
315 Expressed Sequence Tags from the Placenta of the Live-Bearing Fish *Poeciliopsis* (Poe-
316 ciliidae). Journal of Heredity 102: 352–361.
- 317 18. Shorter KR, Crossland JP, Webb D, Szalai G, Felder MR, et al. (2012) *Peromyscus* as a
318 Mammalian Epigenetic Model. Genetics Research International 2012: 1–11.
- 319 19. Sikes RS, Gannon WL, Animal Care and Use Committee of the American Society of
320 Mammalogists (2011) Guidelines of the American Society of Mammalogists for the use of
321 wild mammals in research. Journal of Mammalogy 92: 235–253.
- 322 20. Heo Y, Wu XL, Chen D, Ma J, Hwu WM (2014) BLESS: bloom filter-based error cor-
323 rection solution for high-throughput sequencing reads. Bioinformatics (Oxford, England)
324 30: 1354–1362.
- 325 21. MacManes MD (2014) On the optimal trimming of high-throughput mRNA sequence
326 data. Frontiers in Genetics 5.
- 327 22. Christoffels A (2014) A glance at quality score: implication for *de novo* transcriptome
328 reconstruction of Illumina reads : 1–5.
- 329 23. Lohse M, Bolger AM, Nagel A, Fernie AR, Lunn JE, et al. (2012) RobiNA: a user-friendly,
330 integrated software solution for RNA-Seq-based transcriptomics. Nucleic Acids Research
331 40: W622–7.
- 332 24. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, et al. (2013) *De novo*
333 transcript sequence reconstruction from RNA-seq using the Trinity platform for reference
334 generation and analysis. Nature protocols 8: 1494–1512.

- 335 25. **MacManes MD**, Lacey EA (2012) The Social Brain: Transcriptome Assembly and Char-
 336 acterization of the Hippocampus from a Social Subterranean Rodent, the Colonial Tuco-
 337 Tuco (*Ctenomys sociabilis*). PLOS ONE 7: e45524.
- 338 26. Li W, Godzik A (2006) Cd-hit: a fast program for clustering and comparing large sets of
 339 protein or nucleotide sequences. Bioinformatics (Oxford, England) 22: 1658–1659.
- 340 27. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, et al. (2009) BLAST+:
 341 architecture and applications. BMC Bioinformatics 10: 421.
- 342 28. Mi H (2004) The PANTHER database of protein families, subfamilies, functions and
 343 pathways. Nucleic Acids Research 33: D284–D288.
- 344 29. Wheeler TJ, Eddy SR (2013) nhmmer: DNA homology search with profile HMMs. Bioin-
 345 formatics (Oxford, England) 29: 2487–2489.
- 346 30. Punta M, Coggill PC, Eberhardt RY, Mistry J, Tate J, et al. (2012) The Pfam protein
 347 families database. Nucleic Acids Research 40: D290–301.
- 348 31. Li H (2013) Aligning sequence reads, clone sequences and assembly contigs with BWA-
 349 MEM .
- 350 32. Roberts A, Pachter L (2013) Streaming fragment assignment for real-time analysis of
 351 sequencing experiments. Nature Methods 10: 71–73.
- 352 33. Korneliussen TS, Moltke I, Albrechtsen A (2013) Calculation of Tajima’s D and other
 353 neutrality test statistics from low depth next-generation sequencing data. BMC
- 354 34. Ranwez V, Harispe S, Delsuc F, Douzery EJP (2011) MACSE: Multiple Alignment of
 355 Coding SEquences Accounting for Frameshifts and Stop Codons. PLOS ONE 6: e22594.
- 356 35. Aubry S, Kelly S, Kämpers BMC, Smith-Unna RD, Hibberd JM (2014) Deep Evolutionary
 357 Comparison of Gene Expression Identifies Parallel Recruitment of Trans-Factors in Two
 358 Independent Origins of C4 Photosynthesis. PLOS Genetics 10: e1004365.
- 359 36. Yang Z, dos Reis M (2011) Statistical Properties of the Branch-Site Test of Positive
 360 Selection. Molecular Biology and Evolution 28: 1217–1228.
- 361 37. Yang Z (2007) PAML 4: Phylogenetic Analysis by Maximum Likelihood. Molecular
 362 Biology and Evolution 24: 1586–1591.
- 363 38. Chamary JV, Parmley JL, Hurst LD (2006) Hearing silence: non-neutral evolution at
 364 synonymous sites in mammals. Nature Reviews Genetics 7: 98–108.