

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

Matthew D. MacManes<sup>1</sup>, Michael B. Eisen<sup>2</sup>,

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

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

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

## 1 Abstract

## 2 Introduction

3 For biologists interested in understanding the relationship between fitness, genotype, and phe-  
4 notype, modern sequencing technologies provide for an unprecedented opportunity to gain a  
5 deep understanding of genome level processes that together, underlie adaptation. Unlike more  
6 traditional approaches

7 One interesting example of adaptation lies in animals ability to survive desert conditions.  
8 Here, heat *and* drought provide for powerful selective forces, testing animals' ability to osmoreg-  
9 ulate and thus to survive.

10

11 Specifically, the maintenance of water balance in animals is one of the most important phys-  
12 iologic processes, and is critical to desert survival [?]. Indeed, mammals are exquisitely sensitive  
13 to changes in osmolality, with slight derangement eliciting physiologic compromise. When the  
14 loss of water exceeds dietary intake, dehydration - and in extreme cases, death - can occur. Un-  
15 like most mammals, animals living in desert habitats are subjected to long periods of extreme  
16 heat and intense drought. As a result, desert animals have evolved mechanisms through which  
17 physiologic homeostasis is maintained despite severe and prolonged dehydration.

18

19 One such desert-adapted rodent, a cricetid rodent endemic to the Southwest United States  
20 [?], is a novel model for the study of adaptation to desert environments. They have a lifespan  
21 typical of small mammals, and therefore an individual may live it's entire life without ever drink-  
22 ing water. These rodents have distinct advantage over other desert animals (e.g. *Dipodomys*)  
23 in that they breed readily in captivity, which enables laboratory studies of the phenotype of  
24 interest. In addition the focal species is positioned in a clade of well known animals (e.g. *P.*  
25 *californicus*, *P. maniculatus* and *P. polionotus*) [1]. There are growing genetic and genomic

resources available [2–4].

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 in *P. eremicus*.

## Materials and Methods

### Animal Collection and Study Design

To begin to understand how genes may underlie desert adaptation, I collected 15 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 [?].

### 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 14 libraries were similarly multiplexed and

59 sequenced in a mixture of 100nt paired and single end across two lanes of an Illumina HiSeq  
 60 2000 at the Vincent G. Coates Gnome Center (University of California, Berkeley).

## 61 **Sequence Data Preprocessing and Assembly**

62 Following recommendations from MacManes [5] and Mbandi [6], adapter sequence contamina-  
 63 tion was removed, and low quality nucleotides (defined as PHRED <2) were removed from the  
 64 dataset using the program Trimmomatic version 0.32 [7]. We concatenated sequence data from  
 65 each reference tissue type and assembled them jointly using the Trinity beta version released  
 66 16 March 2014 [8]. We used flags indicating the stranded nature of sequencing reads and set  
 67 maximum allowable physical distance between read pairs to 999nt. The assembly was conducted  
 68 on the XSEDE computer resource Blacklight. To filter the raw sequence assembly, I estimated  
 69 TPM for each assembled sequence using bwa-mem version 0.75 [9] and eXpress version 1.51 [10],  
 70 removing all contigs whose expression was less than TPM=1 [8].

## 72 **Assembled Sequence Annotation**

73 From the filtered assembly, I extracted putative coding sequences using Transdecoder ver-  
 74 sion 16Jan2014 (<http://transdecoder.sourceforge.net/>). These putative protein coding  
 75 sequences were annotated using default settings of the blastx algorithm [11] against the Swis-  
 76 sProt database downloaded on 1 March 2014. Because transcriptome assemblies typically con-  
 77 tain non-coding elements (e.g. ncRNA) in addition to protein coding sequence, we annotated  
 78 the entire filtered dataset using the NCBI nt dataset, downloaded on 1 March 2014. Lastly,  
 79 because rapidly evolving nucleotide sequences may evade detection by blast algorithms, we used  
 80 HMMER3 [12] to search for conserved protein domains contained in the Pfam database [13].

82 To identify sequences unique to each tissue type, I mapped sequence reads from each tissue  
 83 type to the reference assembly using bwa-mem. We estimated expression individually for the  
 84 four tissues. Interesting patterns of expression, including instances where expression was limited  
 85 to a single tissue type were identified.

## 87 **Population Genomics**

88 In addition to the reference individual sequenced at four different tissue types, we sequenced  
 89 15 other conspecific individuals from the same population, located in Palm Desert, California.  
 90 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, which was used for calculating the folded site frequency spectrum (SFS) and Tajima's D [14].

93

## 94 Results

### 95 RNA extraction, Sequencing, Assembly, Mapping

96 RNA was extracted from the hypothalamus, renal medulla, testes, or liver from each individual  
 97 using sterile technique. TRIzol extraction resulted in a large amount of high quality ( $RIN \geq 8$ )  
 98 total RNA, which was used as input. Libraries were constructed as per the standard Illumina  
 99 protocol, and were sequenced as described above. The number of reads per library varied from  
 100 56 million strand-specific paired-end reads in Peer360 kidney, to 9 million single-end reads in  
 101 Peer321. Adapter sequence contamination and low-quality nucleotides were eliminated, which  
 102 resulted in a loss of <2% of reads.

103

104 Transcriptome assembly was accomplished using the program Trinity. The raw assembly  
 105 contained 743314 assembled sequences measuring 418Mb. This assembly was filtered using  
 106  $TM \geq 1$  as a threshold. The filtered assembly contained 130764 sequences measuring 149Mb.  
 107 From this filtered dataset, I extracted 64355 putative coding sequences (60Mb). Of these 64355  
 108 sequences, 37960 were complete exons (containing both start and stop codons), while others were  
 109 either truncated at the 5-prime end (16880 sequences), 3-prime end (4203 sequences), or were  
 110 internal (5312 sequences having neither stop nor start codon).

## 111 Subsection 2

## 112 Discussion

## 113 Acknowledgments

## 114 References

- 115 1. Feng BJ, Sun LD, Soltani-Arabshahi R, Bowcock AM, Nair RP, et al. (2007) Toward  
 116 a Molecular Phylogeny for *Peromyscus*: Evidence from Mitochondrial Cytochrome- b  
 117 Sequences. Journal of Mammalogy 88: 1146–1159.

- 118 2. Shorter KR, Owen A, anderson V, Hall-South AC, Hayford S, et al. (2014) Natural Ge-  
 119 netic Variation Underlying Differences in *Peromyscus* Repetitive and Social/Aggressive  
 120 Behaviors. Behavior genetics .
- 121 3. Panhuis TM, Panhuis TM, Broitman-Maduro G, Broitman-Maduro G, Uhrig J, et al.  
 122 (2011) Analysis of Expressed Sequence Tags from the Placenta of the Live-Bearing Fish  
 123 Poeciliopsis (Poeciliidae). Journal of Heredity 102: 352–361.
- 124 4. Shorter KR, Crossland JP, Webb D, Szalai G, Felder MR, et al. (2012) *Peromyscus* as a  
 125 Mammalian Epigenetic Model. Genetics Research International 2012: 1–11.
- 126 5. MacManes MD (2014) On the optimal trimming of high-throughput mRNA sequence  
 127 data. Frontiers in Genetics 5.
- 128 6. Christoffels A (2014) A glance at quality score: implication for *de novo* transcriptome  
 129 reconstruction of Illumina reads : 1–5.
- 130 7. Lohse M, Bolger AM, Nagel A, Fernie AR, Lunn JE, et al. (2012) RobiNA: a user-friendly,  
 131 integrated software solution for RNA-Seq-based transcriptomics. Nucleic Acids Research  
 132 40: W622–7.
- 133 8. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, et al. (2013) *De novo*  
 134 transcript sequence reconstruction from RNA-seq using the Trinity platform for reference  
 135 generation and analysis. Nature protocols 8: 1494–1512.
- 136 9. Li H (2013) Aligning sequence reads, clone sequences and assembly contigs with BWA-  
 137 MEM .
- 138 10. Roberts A, Pachter L (2013) Streaming fragment assignment for real-time analysis of  
 139 sequencing experiments. Nature Methods 10: 71–73.
- 140 11. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, et al. (2009) BLAST+:  
 141 architecture and applications. BMC Bioinformatics 10: 421.
- 142 12. Wheeler TJ, Wheeler TJ, Eddy SR, Eddy SR (2013) nhmmer: DNA homology search  
 143 with profile HMMs. Bioinformatics (Oxford, England) 29: 2487–2489.
- 144 13. Punta M, Coghill PC, Eberhardt RY, Mistry J, Tate J, et al. (2012) The Pfam protein  
 145 families database. Nucleic Acids Research 40: D290–301.
- 146 14. Korneliussen TS, Moltke I, Albrechtsen A (2013) Calculation of Tajima’s D and other  
 147 neutrality test statistics from low depth next-generation sequencing data. BMC ... .

148 **Figure Legends**

149 **Tables**

150 **Table 1**

151

DATASET	NUM. RAW READS
PEER360 TESTES	32M PE
PEER360 LIVER	53M PE
PEER360 KIDNEY	56M PE
PEER360 BRAIN	23M PE
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

152

153 Table 1. The number of sequencing reads per sample