

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 Introduction

3 Deserts are widely considered one of Earth's harshest environments. Animals living in desert
4 environments are forced to endure intense heat and drought, and in turn, species having evolved
5 in these environments are likely to have evolved specialised mechanisms that may enhance
6 fitness. While living in deserts likely involves a large number of adaptive phenotypes, the abil-
7 ity to osmoregulate – to maintain the proper water and electrolyte balance – appears to be
8 paramount [?]. Indeed, the maintenance of water balance in animals is one of the most impor-
9 tant physiologic processes, and is critical to desert survival. Mammals are exquisitely sensitive
10 to changes in osmolality, with slight derangement eliciting physiologic compromise. When the
11 loss of water exceeds dietary intake, dehydration - and in extreme cases, death - can occur. Un-
12 derstanding these mechanisms will significantly enhance our understanding of the physiologic
13 processes underlying osmoregulation in extreme environments, having implications for studies
14 of human health, conservation, and climate change.

15

16 The genes and structures responsible for the maintenance of water and electrolyte balance
17 are well characterized in model organisms such as mice [?], rats [?, ?, ?], and humans [?, ?, ?].
18 These studies, many of which have been enabled by newer sequencing technologies, serve as
19 a foundation for studies of renal genomics in non-model organisms. In particular, because re-
20 searchers have long been interested in desert adaptation, a number of studies have looked at the
21 morphology or expression of single genes in the renal tissues of desert adapted rodents *Phyllotis*
22 *darwini* [?], *Psammomys obesus* [?], and *Perognathus penicillatus* [?]. More recently, full renal
23 transcriptomes have been generated for *Dipodomys spectabilis* and *Chaetodipus baileyi* [?] as
24 well as *Abrothrix olivacea* [?].

25

26 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 [?], *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*) [1] with growing genetic and genomic resources [2–4] 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 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

60 dependent on having high quality, intact RNA, a small aliquot of each total RNA extract was
 61 analyzed on a Bioanalyzer 2100 (Agilent, Palo Alto, CA, USA). Following confirmation of sample
 62 quality, the reference sequencing libraries were made using the TruSeq stranded RNA prep kit
 63 (Illumina), while an unstranded TruSeq kit was used to construct the other sequencing libraries.
 64 A unique index was ligated to each sample to allow for multiplexed sequencing. Reference
 65 libraries (n=4 tissue types) were then pooled to contain equimolar quantities of each individual
 66 library and submitted for Illumina sequencing using two lanes of 150nt paired end sequencing
 67 using the rapid-mode of the HiSeq 2500 sequencer at The Hubbard Center for Genome Sciences
 68 (University of New Hampshire). The remaining 14 libraries were similarly multiplexed and
 69 sequenced in a mixture of 100nt paired and single end across two lanes of an Illumina HiSeq
 70 2000 at the Vincent G. Coates Gnome Center (University of California, Berkeley).

71 Sequence Data Preprocessing and Assembly

72 The raw sequence reads were error corrected using the software `bleed` [?], using `kmer=25`, based
 73 on the developers recommendations. The error corrected adapter and quality trimmed following
 74 recommendations from MacManes [5] and Mbandi [6]. Specifically, adapter sequence contam-
 75 ination was removed, and low quality nucleotides (defined as `PHRED < 2`) were removed using
 76 the program `Trimmomatic` version 0.32 [7]. We concatenated sequence data from each reference
 77 tissue type and assembled them jointly using `Trinity` version released 17 July 2014 [8]. We used
 78 flags indicating the stranded nature of sequencing reads and set maximum allowable physical
 79 distance between read pairs to 999nt. The assembly was conducted on a linux workstation with
 80 64 cores and 512Gb RAM. To filter the raw sequence assembly, I estimated TPM for each as-
 81 sembled sequence using `bwa-mem` version 0.75 [9] and `eXpress` version 1.51 [10], removing all
 82 contigs whose expression was less than `TPM=1` [8].

84 Assembled Sequence Annotation

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

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

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, which was used for calculating the folded site frequency spectrum (SFS) and Tajima's D [14].

Results

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. Adapter sequence contamination and low-quality nucleotides were eliminated, which resulted in a loss of $<2\%$ of reads.

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

Subsection 2

Discussion

Acknowledgments

References

1. Feng BJ, Sun LD, Soltani-Arabshahi R, Bowcock AM, Nair RP, et al. (2007) Toward a Molecular Phylogeny for *Peromyscus*: Evidence from Mitochondrial Cytochrome- b Sequences. *Journal of Mammalogy* 88: 1146–1159.
2. Shorter KR, Owen A, anderson V, Hall-South AC, Hayford S, et al. (2014) Natural Genetic Variation Underlying Differences in *Peromyscus* Repetitive and Social/Aggressive Behaviors. *Behavior genetics* .
3. Panhuis TM, Panhuis TM, Broitman-Maduro G, Broitman-Maduro G, Uhrig J, et al. (2011) Analysis of Expressed Sequence Tags from the Placenta of the Live-Bearing Fish *Poeciliopsis* (Poeciliidae). *Journal of Heredity* 102: 352–361.
4. Shorter KR, Crossland JP, Webb D, Szalai G, Felder MR, et al. (2012) *Peromyscus* as a Mammalian Epigenetic Model. *Genetics Research International* 2012: 1–11.
5. MacManes MD (2014) On the optimal trimming of high-throughput mRNA sequence data. *Frontiers in Genetics* 5.
6. Christoffels A (2014) A glance at quality score: implication for *de novo* transcriptome reconstruction of Illumina reads : 1–5.
7. Lohse M, Bolger AM, Nagel A, Fernie AR, Lunn JE, et al. (2012) RobiNA: a user-friendly, integrated software solution for RNA-Seq-based transcriptomics. *Nucleic Acids Research* 40: W622–7.
8. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, et al. (2013) *De novo* transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis. *Nature protocols* 8: 1494–1512.
9. Li H (2013) Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM .
10. Roberts A, Pachter L (2013) Streaming fragment assignment for real-time analysis of sequencing experiments. *Nature Methods* 10: 71–73.

- 152 11. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, et al. (2009) BLAST+:
153 architecture and applications. *BMC Bioinformatics* 10: 421.
- 154 12. Wheeler TJ, Wheeler TJ, Eddy SR, Eddy SR (2013) nhmmer: DNA homology search
155 with profile HMMs. *Bioinformatics (Oxford, England)* 29: 2487–2489.
- 156 13. Punta M, Coggill PC, Eberhardt RY, Mistry J, Tate J, et al. (2012) The Pfam protein
157 families database. *Nucleic Acids Research* 40: D290–301.
- 158 14. Korneliussen TS, Moltke I, Albrechtsen A (2013) Calculation of Tajima’s D and other
159 neutrality test statistics from low depth next-generation sequencing data. *BMC ...*

160 **Figure Legends**

161 **Tables**

162 **Table 1**

163

| 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 |

164

165 Table 1. The number of sequencing reads per sample