



Published in final edited form as:

*Mol Ecol.* 2020 April ; 29(7): 1300–1314. doi:10.1111/mec.15401.

## Comparative and population genomics approaches reveal the basis of adaptation to deserts in a small rodent

Anna Tigano<sup>1,2,\*</sup>, Jocelyn P. Colella<sup>1,2</sup>, Matthew D. MacManes<sup>1,2</sup>

<sup>1</sup>Department of Molecular, Cellular and Biomedical Sciences, University of New Hampshire, Durham, 03824, NH, USA

<sup>2</sup>Hubbard Center for Genome Studies, University of New Hampshire, Durham, 03824, NH, USA

### Abstract

Organisms that live in deserts offer the opportunity to investigate how species adapt to environmental conditions that are lethal to most plants and animals. In the hot deserts of North America, high temperatures and lack of water are conspicuous challenges for organisms living there. The cactus mouse (*Peromyscus eremicus*) displays several adaptations to these conditions, including low metabolic rate, heat tolerance, and the ability to maintain homeostasis under extreme dehydration. To investigate the genomic basis of desert adaptation in cactus mice, we built a chromosome-level genome assembly and resequenced 26 additional cactus mouse genomes from two locations in southern California (USA). Using these data, we integrated comparative, population, and functional genomic approaches. We identified 16 gene families exhibiting significant contractions or expansions in the cactus mouse compared to 17 other Myodontine rodent genomes, and found 232 sites across the genome associated with selective sweeps. Functional annotations of candidate gene families and selective sweeps revealed a pervasive signature of selection at genes involved in the synthesis and degradation of proteins, consistent with the evolution of cellular mechanisms to cope with protein denaturation caused by thermal and hyperosmotic stress. Other strong candidate genes included receptors for bitter taste, suggesting a dietary shift towards chemically defended desert plants and insects, and a growth factor involved in lipid metabolism, potentially involved in prevention of dehydration. Understanding how species adapted to deserts will provide an important foundation for predicting future evolutionary responses to increasing temperatures, droughts and desertification in the cactus mouse and other species.

### Keywords

*Peromyscus eremicus*; thermal stress; hyperosmotic stress; selective sweeps; ribosomal protein; bitter taste receptor

---

\* corresponding author: anna.tigano@unh.edu.

#### Author contributions

AT and MDM conceived the study, collected the data, assembled the cactus mouse genome and performed analyses. JPC provided input on the interpretation of results. AT wrote the first version of the paper and AT, JPC and MDM reviewed and edited the paper.

## Introduction

For decades, researchers have been intrigued by adaptation, or the process by which organisms become better fitted to their environments. To this end, scientists have devoted substantial efforts to this issue and have successfully elucidated how natural selection has shaped organismal phenotypes in response to environmental pressures (Berner & Salzburger, 2015; Cooke et al., 2013; Linnen et al., 2009; Nachman, Hoekstra, & D'Agostino, 2003; Savolainen, Lascoux, & Merilä, 2013). Given their influence on metabolism, water availability and ambient temperature are environmental factors relevant to all organisms and are also of growing concern within the context of anthropogenically-induced global climate change and increasing desertification (IPCC, 2018). Studying how animals that are currently living in hot and dry environments have adapted to those conditions is one approach for helping to predict the potential impacts of increasing temperatures and aridity (Hoelzel, 2010; Somero, 2010).

Despite the challenging conditions, a wide variety of organisms have evolved adaptations to live in hot deserts. These adaptations include changes in behavior to avoid dehydration, excessive solar radiation, and heat (e.g., nocturnal life and sheltering in burrows) and a suite of anatomical modifications to dissipate heat (e.g., long body parts and pale colors). Some of the most striking adaptations are at the physiological level and help to either minimize water loss through efficient excretion and reabsorption of water (Knut Schmidt-Nielsen, 1964; K. Schmidt-Nielsen & Schmidt-Nielsen, 1952) or compensate for lack of environmental water via enhanced production of metabolic water from nutrient oxidation (Takei, Bartolo, Fujihara, Ueta, & Donald, 2012; Walsberg, 2000). While these adaptations to desert life have been described in several species (for a review on small mammals see Walsberg, 2000) and are important under current climate predictions (IPCC, 2018), the genetic underpinnings of these traits are less well known.

Genomic studies on camels (*Camelus bactrianus*) have provided substantial evidence related to the genomic basis of adaptation to deserts. For example, analysis of the camel genome showed an enrichment of fast-evolving genes involved in lipid and carbohydrate metabolism, potentially linked to energy production and storage in a food-scarce environment (Bactrian Camels Genome Sequencing and Analysis Consortium et al., 2012; Wu et al., 2014). Transcriptome analysis of renal cortex and medulla in control and water-restricted camels showed a strong response to dehydration in genes involved in water reabsorption and glucose metabolism (Wu et al., 2014). Overall, genes in the arachidonic acid pathway seem to play a role in desert adaptation in both camels and desert sheep (*Ovis aries*; Bactrian Camels Genome Sequencing and Analysis Consortium et al., 2012; Yang et al., 2016). This pathway regulates water retention and reabsorption in the kidney, primarily through changes in reno-vascular tone. Aquaporins, transmembrane water channel proteins, are also involved in water reabsorption and urine concentration, and changes in their expression levels have been associated with dry environments in kangaroo rats (*Dipodomys* spp.; Marra, Eo, Hale, Waser, & DeWoody, 2012; Marra, Romero, & DeWoody, 2014) and the Patagonian olive mouse (*Abrothrix olivacea*; Giorello et al., 2018).

The cactus mouse (*Peromyscus eremicus*) is native to the deserts of southwestern North America and displays a suite of adaptations to this extreme environment. Cactus mice have behavioral and anatomical adaptations for heat avoidance and dissipation, such as a nocturnal lifestyle, larger ears, and aestivation (Macmillen, 1965). They have also evolved lower metabolic rates, which result in a reduction in water loss, and resistance to heat stress compared to other generalist *Peromyscus* spp. (Murie, 1961). For example, while several desert rodents produce concentrated urine (Knut Schmidt-Nielsen, 1964; K. Schmidt-Nielsen & Schmidt-Nielsen, 1952), the cactus mouse is essentially anuric (Kordonowy et al., 2017), which indicates its extreme efficiency of renal water reabsorption. Kordonowy et al. (2017) showed through experimental manipulation of water availability that captive cactus mice were behaviorally and physiologically intact after three days of severe acute dehydration. Gene expression profiling of kidneys highlighted a starvation-like response at the cellular level in dehydrated mice, despite access to food, and strong differential expression of *Cyp4* genes, which are part of the arachidonic acid metabolism pathway (MacManes, 2017). Although these results indicate some degree of convergent evolution with other desert-adapted mammals (Bactrian Camels Genome Sequencing and Analysis Consortium et al., 2012; Takei et al., 2012; Yang et al., 2016), they are limited to expressed genes under particular experimental conditions and in one tissue type only.

Deserts in southwest North America formed relatively recently, only after the retreat of the Pleistocene ice sheets that covered most of the continent during the Last Glacial Maximum approximately 10,000 years ago (Pavlik, 2008). Because the ability to detect genomic signatures of selection depends on coalescent time and effective population size (Nielsen et al., 2005), and given the recent emergence of North American deserts, the footprint of these recent adaptations should continue to be evident in contemporary cactus mouse genomes. Whole genome analyses allow us to detect signatures of selection associated with life in the desert across the complete set of cactus mouse genes, regardless of expression patterns. Further, they allow for analysis of intergenic areas, and for characterization of genomic features that may promote or hinder adaptive evolution, such as the distribution of standing genetic variation and repetitive elements.

To identify genes associated with desert adaptation and to investigate the factors affecting adaptation using the cactus mouse as a model, we first generated a chromosome-level genome assembly and then integrated comparative, population, and functional genomics approaches. As dehydration is a primary challenge desert animals face, we expected to identify signatures of selection associated with metabolism and solute-water balance (i.e. adaptations that either enhance production of metabolic water or prevent fluid loss via excretion) in line with previous studies in the cactus mouse and other desert-adapted species (Bactrian Camels Genome Sequencing and Analysis Consortium et al., 2012; Giorello et al., 2018; Marra et al., 2014; Takei et al., 2012; Wu et al., 2014; Yang et al., 2016). Our analyses of gene family evolution and selective sweeps point instead to regulation of protein synthesis and degradation as the main target of selection. While we find strong support for an evolutionary response in perception of bitter taste and lipid metabolism, we do not identify an extensive signal of selection at genes linked to solute-water balance at the whole-genome level.

## Methods

### Ethics Statement

All sample collection procedures were approved by the Animal Care and Use Committee located at the University of California, Berkeley (2009 samples, protocol number R224) and University of New Hampshire (2018 samples, protocol number 130902) as well as the California Department of Fish and Wildlife (permit number SC-008135) and followed guidelines established by the American Society of Mammalogy for the use of wild animals in research (Sikes & Animal Care and Use Committee of the American Society of Mammalogists, 2016).

### Genome assembly and annotation

We extracted DNA from the liver of a female cactus mouse (ENA sample ID: SAMEA5799953) captured near Palm Desert, CA, USA using a Qiagen Genomic Tip kit (Qiagen, Hilden, Germany). We built two short-insert Illumina libraries (300 bp and 500 bp inserts) using an Illumina Genomic DNA TruSeq kit, following manufacturer recommendations. For scaffolding, we added four mate pair libraries (3 kb, 5 kb, 7 kb, 8 kb) prepared using a Nextera Mate Pair Library Prep kit (Illumina, San Diego, CA, USA). Each library was sequenced on an Illumina HiSeq 2500 sequencer by Novogene (Sacramento, CA, USA) at a depth of approximately 30x for each short-insert library, and 5x for each mate pair library. After adapter trimming with *Trimmomatic* (Bolger, Lohse, & Usadel, 2014), libraries were assembled using the program *ALLPATHS* (Butler et al., 2008). The resulting assembly was gap-filled using a PacBio library (Pacific Biosciences of California, Inc., Menlo Park, CA, USA), constructed from the same DNA extraction and sequenced at ~5x coverage, using *PBJelly* (English et al., 2012). We error-corrected the resulting assembly with short-insert Illumina data and the *Pilon* software package (Walker et al., 2014).

To improve the draft assembly, we used proximity-ligation data (Hi-C) to further order and orient draft scaffolds. We prepared a Hi-C library using the Proximo Hi-C kit from Phase Genomics (Seattle, WA, USA). We used ~200 µg of liver from a second wild-caught female animal and proceeded with library preparation following the protocol for animal tissues. The Hi-C library was sequenced at Novogene using one lane of 150 bp paired-end reads on an Illumina HiSeq 4000 platform. To arrange the draft scaffolds in chromosomes we used the program *Juicer* in an iterative fashion (Durand, Shamim, et al., 2016). Following each run, we loaded the *.map* and *.assembly* files generated by *Juicer* into *Juicebox* (Durand, Robinson, et al., 2016), the accompanying software developed to visualize crosslinks, and corrected misassemblies manually. We ran *Juicer* until no well-supported improvements in the assembly were observed. We thus obtained 24 chromosome-sized scaffolds plus 6,785 unplaced short scaffolds. We calculated assembly statistics with the *assemblathon\_stats.pl* script from the Korf Lab ([https://github.com/KorfLab/Assemblathon/blob/master/assemblathon\\_stats.pl](https://github.com/KorfLab/Assemblathon/blob/master/assemblathon_stats.pl)) and assessed assembly completeness with *BUSCO v3* (Simão, Waterhouse, Ioannidis, Kriventseva, & Zdobnov, 2015) and the Mammal gene set.

To standardize chromosome naming and enable future comparative analyses, we used the genome assembly of the deer mouse (*Peromyscus maniculatus bairdii*; NCBI Bioproject PRJNA494228) to name and orient the cactus mouse chromosomes. We used *mummer4* (Marçais et al., 2018) to align the cactus mouse genome to the deer mouse genome with the function *nucmer* and the options *--maxgap 2000* and *--minclust 1000*. We filtered alignments smaller than 10 kb with *delta-filter* and plotted the alignment using *mummerplot*. These genome alignments also allowed us to test for synteny between the two species, diverged ~9 million years ago ([timetree.com](http://timetree.com) based on Fabre, Hautier, Dimitrov, & Douzery, 2012; Fritz, Bininda-Emonds, & Purvis, 2009; León-Paniagua, Navarro-Sigüenza, Hernández-Baños, & Morales, 2007; Schenk, Rowe, & Steppan, 2013), and to assess the degree of structural divergence between the two genomes.

We identified transposable elements and other repetitive elements using *RepeatMasker v.4.0* (Smit, Hubley, & Green, 2015) and the Rodentia dataset. The masked genome was annotated using *Maker v.2.3.1* (Cantarel et al., 2008) and the *Mus musculus* reference protein dataset.

### Whole genome resequencing

Little is known about population structure in the cactus mouse (Riddle et al. 2000), therefore we included samples from two different locations and sampling times to begin characterizing geographic and temporal structuring in this species. We sequenced the genomes of an additional 26 cactus mice collected from two locations in Southern California: Motte Rimrock and Boyd Deep Canyon Reserves (both belonging to the University of California Natural Reserve System; Table 1). The Motte Rimrock Reserve (Motte hereafter) sits on a broad, rocky plateau and supports both coastal and desert habitats as it is located equidistant between the Pacific coast and the Colorado Desert. We captured cactus mice in xeric areas characterized by rocky outcrops, which constitutes their typical habitat. The Boyd Deep Canyon Reserve (Deep Canyon hereafter) is a large natural reserve extending from low to high elevation (290-2657 m). We sampled at two locations in the lower elevation part of the Deep Canyon Reserve: Deep Canyon (290 m a.s.l.), the driest and hottest location, with average monthly temperatures between 10-40 °C and mean annual rainfall of 15 cm, and Agave Hill (820 m a.s.l.), with average monthly temperatures between 9-35 °C and mean annual rainfall of 18 cm. Samples from Deep Canyon were collected in 2009 and 2018 (Table 1).

Genomic libraries were prepared at the Biotechnology Resource Center at Cornell University (Ithaca, NY, USA) using the Illumina Nextera Library Preparation kit and a modified protocol for low-coverage whole genome resequencing ('skim-seq'). Individually barcoded libraries were sequenced at Novogene using 150 bp paired-end reads from one lane on the Illumina NovaSeq S4 platform. We conducted an analysis of sequencing read quality and trimmed adapters from raw sequencing data with *fastp* (Chen, Zhou, Chen, & Gu, 2018). We mapped sequences from each of the 26 individuals the cactus mouse reference genome using *bwa mem* (H. Li & Durbin, 2009) and removed duplicates with *Samblaster* (Faust & Hall, 2014). The resulting BAM files were sorted and indexed using *Samtools* (H. Li et al., 2009).

As sequencing depth was variable among individuals (raw coverage between ~2-17X), we called variants in *ANGSD* (Korneliussen, Albrechtsen, & Nielsen, 2014) as its algorithm takes into account genotype uncertainty associated with low-coverage data. To identify a list of high-confidence variable sites, we ran a global variant calling analysis including all 26 individuals using the genotype likelihood model from *Samtools* (-GL 1; Li et al. 2009). To be included in our analyses, a site had to satisfy the following criteria: p-value below  $10^{-6}$ , minimum sequencing and mapping qualities above 20, minimum depth and number of individuals equal to half of the number of individuals included in the analysis (13 out of 26), and a minor allele frequency (MAF) above 1%.

### Genomic differentiation across space and time

To estimate the effects of temporal and spatial distance on levels of genomic differentiation among individuals, we first ran an Analysis of Principal Components (PCA) of genetic data using the *ngsCovar* program from *ngsTools* (Fumagalli, Vieira, Linderth, & Nielsen, 2014). We used all high-quality SNPs (defined above) called at the species level and ran the PCA using genotype posterior probabilities, rather than called genotypes, as input, to take into account the genotype uncertainty associated with low-coverage data (Korneliussen et al. 2014). To estimate genetic distance between individuals controlling for the effect of varying depth of sequencing across individuals, we downsampled to a single base for each site included in the high-quality SNPs list and performed multidimensional scaling (MDS) analysis in *ANGSD*. These preliminary PCA and MDS analysis (Supplementary Figure 1 and Figure 1) revealed the presence of an individual from Motte that grouped with neither population: this is particularly evident in the PCA (Supplementary Figure 1), and in the MDS this individual was the only one not laying along one of the two axes describing intrapopulation variation (Figure 1). We retained a second individual from Motte that appeared distant in the MDS statistical space from the rest of the individuals as it laid on the variation axis describing Motte, which ran parallel to the variation axis describing Deep Canyon, and grouped with the rest of the individuals from Motte in the PCA (Figure 1 and Supplementary Figure 1). This individual was subsequently removed and variants were re-called.

After outlier exclusion, we reran *ANGSD* to estimate allele frequencies in Motte and Deep Canyon separately. We provided a list of high-confidence SNPs for use in downstream analyses and applied the same filters as in the global SNP calling, excluding SNP p-value and MAF thresholds, with the major allele fixed across runs. We used the sample allele frequencies (.mafs file) from each population to calculate the 2D Site Frequency Spectrum (SFS), which we also used as prior for estimating  $F_{ST}$ , a measure of genetic differentiation. We calculated average  $F_{ST}$  across autosomes and across the X chromosome separately, and investigated broad patterns of differentiation across the genome in 50 kb sliding windows.

### Sequence and structural standing genetic variation at the species level

To estimate levels and patterns of standing genetic variation at the species level, we analyzed samples from both cactus mouse populations together. We calculated the overall proportion of polymorphic sites and nucleotide diversity ( $\pi$ ) in 50 kb sliding windows to characterize broad patterns. Estimates of  $\pi$  for each polymorphic site were based on the maximum

likelihood of the SFS calculated with *realSFS* in *ANGSD*. To obtain accurate estimates of diversity, we corrected global and window estimates of  $\pi$  by the number of variant and invariant sites covered by data. We estimated genome coverage, total and per window, by rerunning *ANGSD* (including all 25 samples) using the same filtering parameters we used for the global calling variant but without the SNP p-value and the MAF filters. We then divided the sum of per-site  $\pi$  by the number of variant and invariant sites in a given window.

To investigate the distribution of structural variation in the cactus mouse genome, we sequenced the genome of an additional individual from the Deep Canyon site within the Deep Canyon Reserve (sampled in 2009) using 10X Genomics (Pleasanton, CA, USA). This method is based on linked-reads technology and enables phasing and characterization of structural variation using synthetic long-range information. We used the program *Long Ranger v.2.2.2* from 10X Genomics and ran it in whole genome mode with *longranger wgs* using *Freebayes* (Garrison & Marth, 2012) as variant caller. Finally, we tested whether chromosome size was a predictor of either sequence or structural standing variation using linear models in *R* (R Core Team 2019).

### Gene family evolution in the cactus mouse

To investigate gene family contractions and expansions in the cactus mouse, we analyzed the genomes of 25 additional species (Supplementary Table 1) within the Myodonta clade (Order: Rodentia), which includes rats, mice, and jerboas, for which genome assemblies were publicly available from NCBI in June 2019. To avoid potential biases due to different gene annotation strategies, we re-annotated all 25 genomes with the same strategy used for the cactus mouse (see above). Genome quality was evaluated using *BUSCO*. We identified groups of orthologous sequences (orthogroups) in all species using the package *Orthofinder2* (Emms & Kelly, 2018) with *Diamond* as protein aligner (Buchfink, Xie, & Huson, 2015). In a preliminary run, we observed that fewer orthogroups and fewer genes per orthogroup were identified in species with lower genome assembly quality. Given this pattern could bias the results of analysis of gene family expansion or contraction, we took a conservative approach and filtered assemblies that had less than 70% complete benchmarking universal single-copy orthologs (BUSCOs) and thus retained 18 species (Supplementary Table 1). 70% was chosen as a compromise between including as many genomes as possible while eliminating the potential for biases.

We analyzed changes in gene family size accounting for phylogenetic history with the program *CAFE v4.2.1* (De Bie, Cristianini, Demuth, & Hahn, 2006). We filtered both invariant orthogroups and those that varied across species by more than 25 genes to obtain meaningful likelihood scores for the death-birth parameter ( $\lambda$ ) as recommended by the *CAFE* developers. We used the rooted species tree inferred by *Orthofinder2* and ran the analysis using a single value for  $\lambda$  estimated in *CAFE* for the whole tree. Finally, we summarized the results with the python script *cafetutorial\_report\_analysis.py* from the *CAFE* developers (<https://hahnlab.github.io/CAFE/manual.html>).

## Detection of signatures of selection

We used an integrative approach to detect signatures of selection from the population genomics data. First, we used *Sweepfinder2* (DeGiorgio, Huber, Hubisz, Hellmann, & Nielsen, 2016; Nielsen et al., 2005) to detect recent selective sweeps from the SFS as it is compatible with low-coverage whole genome data. We ran these analyses using all 25 individuals, based on the rationale that potential differences between the two populations (Deep Canyon and Motte) due to local adaptation or demography would be swamped by signatures common in the species across populations (Miller et al. 2020). As recommended by Huber, DeGiorgio, Hellmann, & Nielsen (2016), we included both variant and invariant sites, but we could not reconstruct the ancestral state of these sites due to lack of data from closely related species (all species had > 9 million years divergence). The X chromosome was excluded from this analysis. We converted allele frequencies estimated in *ANGSD* to allele counts, and estimated the SFS from the autosomes only in *SweepFinder2*. We then tested for sweeps using *SweepFinder2* with the *-l* setting, i.e. using the pre-computed SFS, and calculated the composite likelihood ratio (CLR) and  $\alpha$  every 10,000 sites. This window size of 10 kb was selected as a trade-off between computational time and resolution. Only the peaks with CLR values above the 99.9<sup>th</sup> percentile of the empirical distribution of CLR values were considered under selection. We functionally annotated the closest genes to these outlier peaks and ran a Gene Ontology (GO) enrichment analyses to test whether genes under putative selection were enriched for a particular function or pathway. We performed this analysis on the [geneontology.org](http://geneontology.org) webpage using the program *Panther* (Mi et al., 2017) and the *Mus musculus* gene set as reference. As GO terms are hierarchical, we summarized these results with the software package *REVIGO* (Supek, Bošnjak, Škunca, & Šmuc, 2011), which uses a clustering algorithm based on semantic similarities, setting the similarity threshold at 0.5.

Finally, we generated a list of *a priori* candidate genes potentially involved in desert adaptation. These included the genes that were most differentially expressed in response to experimental dehydration, including 11 *Cyp4* genes from the arachidonic acid metabolism pathway, and the sodium carrier gene *Slc8a1* (MacManes, 2017). We also included the 9 aquaporins, which are important in water reabsorption in the kidney but were not differentially expressed in hydrated versus dehydrated cactus mice (MacManes, 2017). We integrated this list of candidate genes with the genomic areas showing strong signatures of selective sweeps in the *Sweepfinder2* analysis. As decreases in  $\pi$  and Tajima's D can also be indicative of selective sweeps, we zoomed in around these candidate regions, plus an additional 10 kb flanking on each side, and calculated these two statistics in 1 kb windows .

## Results

### Genome assembly and annotation

Illumina and PacBio reads yielded a draft genome assembly of 2.7 Gbp and a scaffold N50 of 1.3 Mbp. Scaffolding with Hi-C data increased contiguity 100-fold and yielded 24 chromosome-sized scaffolds for a total assembly size of 2.5 Gbp, plus 173 Mbp of unplaced scaffolds. The final assembly contained 92.9% complete BUSCOs, with 1.2% of the genes duplicated and 3.7% missing. Together these statistics indicate that the cactus mouse



genome assembly has high contiguity, high completeness and low redundancy (Table 2). We annotated 18,111 protein-coding genes. *Repeatmasker* masked 35% of the genome as repetitive. LINE1 and LTR elements alone constituted 21% of the repeats. Total proportion of repeats and relative proportion of different repeats classes were similar across eight *Peromyscus* species (Supplementary Table 2).

Whole genome alignment to the *P. maniculatus* genome revealed the presence of several intrachromosomal differences between the two species, but no large inversions, translocations, or interchromosomal rearrangements were evident at the resolution granted by *mummer4* (Supplementary Figure 2). This, in combination with a conserved number of chromosomes supported by both karyotype characterization (Smalec, Heider, Flynn, & O'Neill, 2019) and genome assemblies, indicates that genome structure is highly conserved between these *Peromyscus* species.

### Genomic differentiation across space and time

Our PCA showed that the first two principal components (PCs) explained 13.32% of the variation present across 26 cactus mice. Both the PCA (Supplementary Figure 1) and the MDS (Figure 1A) clearly separated individuals from Motte and Deep Canyon Reserve. Within the Deep Canyon Reserve, no differentiation at the temporal or microspatial scale was observed (Supplementary Figure 1, Figure 1A). One individual from Motte appeared distinct from other individuals included in the analysis. As we could not ascertain the reason for such behavior, e.g., technical artifact, taxonomic misidentification, hybridization, etc., this individual was excluded from further analyses.

Differentiation between Motte and Deep Canyon Reserve populations was high, with an average  $F_{ST}$  value of 0.19 and 0.14 across the 23 autosomes and the X chromosome, respectively.  $F_{ST}$  calculated in 50 kb windows ranged from 0.06 to 0.46 with 95% of the windows ranging from 0.12 and 0.28 (Figure 1B).

### Sequence and structural standing genetic variation

A total of 1,875,915,109 variant and invariant sites, representing 75% of the genome, were included in our analyses. We identified 43,695,428 SNPs with high-confidence (one every 43 bp – 2.3% of all sites). Global genome-wide  $\pi$  was  $6 \times 10^{-3}$ .  $\pi$  was lowest on the X chromosome and seemed to increase from chromosome 1 to 23 (Figure 1C). In fact, chromosome length was a strong negative predictor of nucleotide diversity at each autosome ( $R^2 = 0.65$ ,  $F_{(1,21)} = 39.58$ ,  $p < 0.001$ ; Figure 2).

A large area of elevated nucleotide diversity (~17 Mbp long) was evident at the beginning of chromosome 1. Although this signature can indicate the collapse of paralogous sequence, conserved synteny with other *Peromyscus* species and unequivocal support from the Hi-C contact map strongly indicated that the assembly was correct for chromosome 1. To begin to understand the genome-level processes that may have generated this pattern, we calculated depth of coverage in chromosomes 1 and 2 - a reference chromosome that did not show similar large regions of elevated  $\pi$  - using a subset of the shotgun data and the number and proportion of repetitive elements in 50 kb windows (Supplementary Figure 3). Depth of sequencing in this unusual area of chromosome 1 was higher than in the rest of chromosome

1 (up to 7x higher) and compared to chromosome 2 (10% higher overall), and showed a similar peak as the one shown for nucleotide diversity (Supplementary Figure 3). The number and proportion of repetitive elements were both higher in this area relative to other parts of chromosome 1, and all of chromosome 2 (Supplementary Figure 3). Whether these repeats have been collapsed or reads from other similar repeats spread throughout the genome have mapped to this area is hard to disentangle, even though the inflated nucleotide diversity would suggest the latter. Together, these analyses indicate that this area is highly repetitive, rather than containing a misassembled large duplication.

Analysis of the 10x Genomics data using *LongRanger* resulted in an estimated mean DNA molecule length of only 13,620 bp (ideal is > 40,000 bp), and the number of linked reads per molecule was six, much lower than the ideal threshold of 13. As short molecules can negatively impact the detection of large structural variants and generate many false positives, we adopted a conservative approach by reporting only short indels (41-29,527 bp). We identified a total of 87,640 indels between the reference and an individual from the same population. Indels affected 101 Mbp of the total genome assembly, which represents 4% of the total sequence. Number of indels per chromosome was strongly positively correlated with chromosome size ( $R^2 = 0.95$ ,  $F_{(1,21)} = 438.7$ ,  $p < 0.001$ ; Figure 2).

### Gene family evolution

*Orthofinder2* grouped protein sequences from 18 Myodontine rodents into a total of 23,020 orthogroups. After removing orthogroups with either high (> 25 genes) or no variation in the number of genes across species, the dataset was reduced to 21,347 orthogroups. On average, all species included in our analysis showed gene family contraction, albeit of varying magnitude. *Mus musculus* had the highest number of significant changes ( $p < 0.01$ ) in gene family size (92 orthogroups), while *Dipodomys ordii* had no significant changes (Figure 3). The number of gene families with significant contractions or expansions varied between 16 and 24 among *Peromyscus* species, with more expansions than contractions except in the cactus mouse, which exhibited four gene family expansions and 12 contractions from the closest node in the tree (Figure 3). Four of these gene families contained genes associated with sperm motility (three contractions, one expansion), four included ribosomal proteins (two contraction, one expansion), three were associated with immune response (three contractions), and two included genes in the ubiquitin-like (ubl) conjugation pathway (two contractions). Other functions included pheromone reception (one contraction), cytoskeletal protein binding (one contraction), and prohibitin (one expansion; Figure 3).

### Identification of selective sweeps

Analysis of the signatures of selective sweeps yielded a total of 232 sites under selection. Of these, 119 clustered in 44 larger regions that included two or more adjacent CLR outliers (Supplementary Figure 4). By retrieving the genes closest to each peak (one in both up- and down-stream directions), we compiled a list of 186 genes associated with selective sweeps (Supplementary Table 3). Fourteen of these genes, including many putative olfactory and vomeronasal receptors, were not matched with a corresponding GO term. Ribosomes were overrepresented among 'cellular components', with eight GO terms pointing to this organelle ( $p < 0.001$ , after Bonferroni correction). In addition to this, 279 biological

processes and 89 molecular functions were significantly overrepresented (before correction for multiple tests; full list in Supplementary Table 4 and 5, respectively). GO terms clustering in *REVIGO* showed that terms with the lowest p-values under ‘biological processes’ included ‘membrane organization’, ‘cellular amide metabolism process’, ‘translation’, ‘ribosome assembly’, and ‘detection of chemical stimulus involved in sensory perception of bitter taste’ (Figure 4, Supplementary Table 3); while under ‘molecular functions’ they included ‘structural constituent of ribosome’, ‘binding’, ‘olfactory receptor binding’, and ‘mRNA binding’ (Figure 4, Supplementary Table 4). Note that although the selective sweeps identified in the area of inflated  $\pi$  of chromosome 1 should be considered with caution, most of the genes associated with these areas were either vomeronasal receptors that were not included in the GO enrichment analysis (see above) or genes whose function was not enriched overall.

Contrary to predictions, mean  $\pi$  and Tajima’s D were significantly higher across the candidate areas for selective sweeps when compared to genome-wide means across all samples combined (Wilcoxon test,  $p < 0.001$  for both  $\pi$  and Tajima’s D; Figure 5) and within each population (Supplementary Figure 5). Among the candidate genes from previous studies, 8 (all *Cyp4* genes, *SLC8A1*, and *aqp4*, *aqp5*, *aqp8*, and *aqp12*) and 6 genes (the *Cyp4a* gene cluster, *Cyp4v2*, *SLC8A1*, and *aqp5*, *aqp9*, and *aqp12*) showed significant deviations from genome-wide average in  $\pi$  and Tajima’s D, respectively ( $p < 0.05$  after Benjamini-Yekutieli correction for multiple testing), but not always in the predicted direction (Supplementary Figure 6). Among the *Cyp4* genes, *Cyp4f* showed a modest decrease in  $\pi$ ; among aquaporins, *aqp8* showed a decrease in  $\pi$ , and *aqp9* showed a decrease in Tajima’s D; and *Slc8a1* showed the greatest reduction in  $\pi$  and Tajima’s D overall (Supplementary Figure 5).

Number of sweeps in each chromosome did not correlate with mean  $\pi$  ( $p > 0.05$ ). A correlation with either chromosome size or number of indels ( $p < 0.01$ ) was entirely driven by the outlier behaviour of chromosome 1, and it did not hold when chromosome 1 was removed from the dataset ( $p > 0.05$ ).

## Discussion

### A chromosome-level assembly for the cactus mouse

A high-quality chromosome-level assembly of the cactus mouse genome allowed us to investigate genomic patterns of variation, differentiation, and other genomic features (i.e. genes, repetitive elements, number and size of chromosomes), and to identify regions of the genome that may be associated with desert adaptations. As the number of publicly available *Peromyscus* genome assemblies increases (Colella, Tigano, & MacManes, 2019), the cactus mouse genome will provide additional insights into adaptation, speciation, and genome evolution when analyzed in a comparative framework. For example, our comparison of the cactus mouse and the deer mouse genomes revealed higher than expected genome stability considering the divergence time between the two species (–9 million years ago), and their large effective population sizes and short generation times (Bromham, 2009; Charlesworth, 2009). Our synteny analysis confirms at the sequence level what is reported from karyotypes of several *Peromyscus* spp. (Smalec et al., 2019), i.e. between *P. eremicus* and *P.*

*maniculatus* there is no variation in chromosome number and no interchromosomal rearrangements but abundant intrachromosomal variation. Genome stability among *Peromyscus* species (Long et al., 2019) contrasts variation in the Muridae family (Order Rodentia) where the number of chromosomes varies dramatically, even within the *Mus* genus, and large chromosomal rearrangements are abundant (Thybert et al., 2018).

### High genetic diversity and differentiation

Population differentiation between cactus mouse populations inhabiting the Motte and Deep Canyon Reserves in Southern California is high despite being separated by only 90 km. Differentiation across the genome was high, but without distinguishable  $F_{ST}$  peaks (Figure 1B). The distribution of allele frequency changes is suggestive of prolonged geographical isolation, which is consistent with the Peninsular Range mountains acting as a dispersal barrier. Although this is not surprising given the limited dispersal ability of *Peromyscus* mice (e.g., Ribble 1992; Krohne et al. 1984), the close proximity of these two sampling locations suggests that population structure across the species range, which spans more than 2500 km from Nevada (USA) to San Luis Potosì (Mexico), is likely to be strong. Previous analyses based on a single mitochondrial marker split the *P. eremicus* species complex into three species – *P. eva*, *P. fraterculus*, and *P. merriami* – plus West and East *P. eremicus* clades (Riddle, Hafner, & Alexander, 2000). Our analyses suggest that population structure could be pronounced even within each *P. eremicus* clade, which warrants further investigation to elucidate the taxonomic status of these species and to reveal potential differences in adaptation to local desert conditions.

As standing genetic variation is the main source of adaptive genetic variation (Barrett & Schluter, 2008), characterizing levels and distribution of sequence and structural variation can help understand how and where in the genome adaptations evolve. With more than 43 million high-quality SNPs and ~87,000 indels, the cactus mouse exhibits high levels of standing genetic variation, which is consistent with large effective population sizes and comparable to what has been reported for the congeneric white-footed mouse (*P. leucopus*; Long et al., 2019). While SNPs are the main, and often the only, type of variation screened in genomic studies of adaptation (Wellenreuther, Mérot, Berdan, & Bernatchez, 2019), here we show that small- to mid-sized indels are common and cover ~4% of the genome. Given that our analysis of structural variation leveraged sequence data from only one individual, the level of standing structural variation in the population is likely much higher. Pezer et al. (2015) found that indels covered ~2% of a wild house mouse (*M. musculus domesticus*) genome compared to a reference, likely an underestimation considering that the analysis was based on variation of read depth only. While the inclusion of our high coverage 10X Genomics dataset allowed us to characterize structural variation in a single individual, future work at the population level will allow us to characterize the role of structural variation in the adaptive evolution of the cactus mouse.

Nucleotide variation and number of indels were significantly correlated with chromosome size in the cactus mouse. While a positive, linear relationship between chromosome size and number of indels is expected, a negative correlation between chromosome size and

nucleotide diversity may be explained by recombination rate, which is higher in shorter chromosomes (Kaback, Guacci, Barber, & Mahon, 1992).

Sex chromosomes generally harbor lower nucleotide diversity (Wilson Sayres, 2018) and higher differentiation (Presgraves, 2018) when compared to autosomes due to their reduced effective population size, different mode of inheritance, and their role in the evolution of reproductive barriers. We did observe lower  $\pi$  in comparison to autosomes (45–71% of mean  $\pi$  for each autosome), slightly lower than neutral expectations assuming equal sex ratio (75%). Demographic processes and/or selection could be implicated in this additional reduction, but their relative roles were not tested here. However, contrary to our expectations, analysis of genomic differentiation based on  $F_{ST}$  showed that the X chromosome was less differentiated than the autosomes in the interpopulation comparison. In fact, sex chromosomes are more differentiated than autosomes in 95% of the studies for which this information is available (Presgraves, 2018). In two cases regarding mammals, domestic pigs and wild cats, lower X chromosome differentiation was ascribed to hybridization and introgression (Ai et al., 2015; G. Li, Davis, Eizirik, & Murphy, 2016). Hybridization has been reported among several *Peromyscus* species (Barko & Feldhamer, 2002; Leo & Millien, 2017) and represents a viable hypothesis given that the cactus mouse is sympatric with the canyon mouse (*Peromyscus crinitus*) and numerous other *Peromyscus* species throughout its range. In the future, population genomic data from additional *Peromyscus* species will help assess how common this pattern is within and across species, to test the hybridization hypothesis, to identify potential donor and recipient species, and to test the potential role of hybridization and introgression in desert adaptation.

Neither sequence nor structural variation at the chromosome level was a strong predictor of the number of selective sweeps in a chromosome. However, mean  $\pi$  and Tajima's D were significantly higher in the areas affected by selective sweeps than across the whole genome. Theory predicts that a selective sweep should remove variation from the adaptive site and its surroundings, thus resulting in a localized reduction in  $\pi$  and lower Tajima's D relative to the ancestral level of variation (Kim & Stephan, 2002; Smith & Haigh, 1974). However, the signature of a selective sweep, and the ability to detect it, depends on the strength of selection, the recombination rate around the selected site, and whether the sweep is hard, soft, or incomplete (i.e. whether a single or multiple haplotype carries the beneficial allele, or the allele hasn't reached fixation yet; Messer & Neher, 2012). *Sweepfinder2* is best suited to detect hard selective sweeps and has limited power to identify soft or incomplete sweeps (DeGiorgio et al., 2016; Huber et al., 2016; Nielsen et al., 2005). Therefore, if we do not observe a general reduction of diversity compared to the genome-wide average around adaptive sites, it could be for one or a combination of the following reasons. a) High recombination rates due to a large effective population size may cause rapid LD decay among neighboring sites, thus reducing the size of the typical, diagnostic dip in diversity to a point of non-detectability. b) Due to computational limitations (*Sweepfinder2* is not able to parallelize, leading to long run times) we estimated the CLR of a site every 10 kb, which may not be dense enough to pinpoint the exact location of the sweep and reveal narrow reductions in  $\pi$ . c) If soft sweeps are indistinguishable from hard sweeps when selection is strong (Harris, Garud, & DeGiorgio, 2018), selective sweeps may preferentially occur in areas of high standing genetic variation. Finally, d) even if a reduction in diversity occurs

relative to ancestral levels,  $\pi$  may still not drop under the genome-wide average, especially if diversity was originally high. On chromosome 9, for example, where three consecutive *Sweepfinder2* outlier regions extend over 440 kb, the reduction of  $\pi$  and Tajima's D was drastic, suggesting that coarse resolution may prevent the detection of diversity dips around selected sites if the genomic area affected is small.

### **Lost in translation: pervasive signature of selection in genes associated with protein synthesis and degradation**

Together, the analyses of gene family evolution and selective sweeps indicate that traits associated with the synthesis and degradation of proteins have evolved under the influence of natural selection. Four ribosomal protein families are either expanded or contracted, and gene ontology analysis demonstrated an enrichment of terms associated with ribosomes (e.g., ribosome assembly and translation, structural constituents of ribosomes, mRNA binding, and unfolded protein binding). We also report a significant contraction of a gene family associated with the ubl conjugation pathway, which was similarly identified in an analysis of selective sweeps. Ubiquitin and ubl-proteins function either as a tag on damaged proteins to be degraded or as regulators of interactions among proteins (Hochstrasser, 2009). Cactus mice face many stressors including high temperatures and lack of water. Heat causes cellular stress directly, via thermal stress, and indirectly, by exacerbating the negative effects of dehydration due to lack of water and rapid water loss (e.g., respiratory water loss, evaporative cooling). Thermal and hyperosmotic stress can suppress the transcription and translation machinery, increase DNA breaks and protein oxidation, and cause cell cycle arrest, and eventually apoptosis and cell death (Burg, Ferraris, & Dmitrieva, 2007; Kampinga, 1993). However, the strongest and most immediate effect of thermal and hyperosmotic stress is protein denaturation (Burg et al., 2007; Kampinga, 1993; Lamitina, Huang, & Strange, 2006). Our results are consistent with the expected cellular response to both thermal and hyperosmotic stress, which have similar physiological effects even though the underlying mechanisms may differ. A meta-analysis of genomics and transcriptomics studies investigating the evolutionary response to different thermal environments in metazoans, including invertebrates to mammals, highlighted 'translation', 'structural constituents of ribosomes', and 'ribosome' as the gene ontology terms most commonly enriched (Porcelli, Butlin, Gaston, Joly, & Snook, 2015), in line with our results. Similarly, many genes mediating the cellular response to hyperosmotic stress are involved in the regulation of protein translation and the elimination of denatured proteins in *Caenorhabditis elegans* (Lamitina et al., 2006). These analyses suggest that selection has acted strongly on genes responsible for protection against thermal and/or hyperosmotic stress or for efficiently removing damaged proteins and resuming translation after acute stress (Kampinga, 1993). Additionally, as the volume of dehydrated cells decreases causing rearrangements in the cytoskeleton (Burg et al., 2007), the significant contraction of a cytoskeletal protein gene family could also point to additional adaptations to hyperosmotic stress in the cactus mouse. Acute dehydration experiments on captive cactus mice also found limited tissue damage and apoptosis in the kidneys of dehydrated individuals (MacManes, 2017), consistent with our hypothesis. Negative regulation of cell death was one of the most significant GO terms, suggesting that these genes may be under selection to avoid tissue necrosis during acute or chronic stress.

## Life in the desert involves dietary and metabolic adaptations

The GO analysis of genes associated with selective sweeps indicated an enrichment for bitter taste receptors. The perception of bitter taste has evolved to allow organisms to avoid toxic compounds found in many plants and insects (Garcia & Hankins, 1975; Glendinning, 1994). Although herbivorous and insectivorous animals generally have a larger repertoire of bitter taste receptors compared to their carnivorous counterparts (D. Li & Zhang, 2014; Wang & Zhao, 2015), they are also less sensitive to bitterness (Glendinning, 1994). The cactus mouse is omnivorous, with a diet predominantly based on seeds, insects, and green vegetation with proportions varying according to seasonal availability (Bradley & Mauer, 1973; Meserve, 1976). We hypothesize that repeated signal of selective sweeps at bitter taste receptor genes may have increased the frequency of alleles that decrease bitter sensitivity, thus making a greater variety of food palatable to the cactus mouse in an environment that is characterized by scarcity of resources and an abundance of bitter-tasting plants and insects.

Chromosome 9 showed the largest and strongest selective sweep in the genome (Supplementary Figure 4). This area was associated with *Gdf10* (growth/differentiation factor 10), the only annotated gene of known function in the region, which is involved in osteogenesis and adipogenesis. Overexpression of *Gdf10* in the adipose tissues of mice prevents weight gain under a high-fat diet and affects their metabolic homeostasis, including oxygen consumption and energy expenditure (Hino et al., 2017). The drastic loss of weight and the starvation-like response reported in experimentally dehydrated cactus mice suggests that lipid metabolism has a role in the adaptive response to dehydration (MacManes, 2017). Experimental water deprivation induced higher food consumption and loss of body fat in the spinifex hopping mouse (*Notomys alexis*), a desert-specialist rodent (Takei et al., 2012). In camels, accelerated evolution of genes associated with lipid metabolism was associated with food scarcity in the desert (Wu et al., 2014). The strong signature of selection around *Gdf10* in the cactus mouse therefore warrants further investigation, as adaptive changes in lipid metabolism may be pivotal for survival in the desert.

Among the candidate genes we selected from previous studies, only *Slc8a1* – the sodium carrier gene – show significant reduction in both  $\pi$  and Tajima's D, consistent with a selective sweep. However, *Cyp4v2* – one of the genes in the arachidonic acid pathway – was in proximity of a selective sweep on chromosome 17. This gene shows similar catalytic properties to other *Cyp4* genes in the arachidonic acid pathway and is commonly expressed in retinal, kidney, lung, and liver tissue of humans (Nakano, Kelly, & Rettie, 2009). Known to be strongly associated with ocular disease, its role, however, has not been investigated in the context of solute-water balance in the kidneys. The discrepancy between the strong changes in gene expression in *Cyp4* genes between hydrated and dehydrated mice (MacManes, 2017) and the results presented here suggest that these genes affect kidney physiology predominantly via gene expression and potentially through changes in regulatory regions that we have not targeted explicitly. Future comparative analyses of sequences from additional rodents, including other desert-adapted species, will help us understand the relative role of gene expression regulation versus coding changes in adaptation to desert environments.

Comparisons with studies on other desert-adapted mammals highlight a combination of convergent and idiosyncratic adaptations to life in the desert. Although the arachidonic acid pathway showed signatures of selection in camels, sheep, and the cactus mouse (Bactrian Camels Genome Sequencing and Analysis Consortium et al., 2012; MacManes, 2017; Yang et al., 2016), the specific genes involved differed among species and between the mechanisms by which they putatively affected adaptive phenotypes: gene family contractions and expansions in camel, selective sweeps in desert sheep, and changes in gene expression under acute dehydration in the cactus mouse. In addition to the stress imposed by heat and lack of water, comparative genomics analyses suggest that camels have unique adaptations to avoid the deleterious effects of dust ingestion and intense solar radiation (Wu et al., 2014). The cactus mouse avoids solar radiation altogether with a nocturnal lifestyle and shows strong signatures of selection at receptors for bitter taste, consistent with adaptation to a diet based on bitter-tasting desert plants and insects. Desert woodrats (*Neotoma lepida*) are highly specialized to bitter and toxic plants, such as juniper (*Juniperus monosperma*) and creosote bush (*Larrea tridentata*), and have evolved several adaptations to consume them, including detoxifying gut microbiomes (Kohl, Weiss, Cox, Dale, & Dearing, 2014) and hepatic enzymes (Skopec & Dearing, 2011). Although the repertoire of bitter taste receptors, or their sensitivities to bitter taste, has not been investigated in desert rodents, nor our genomic analyses of cactus mice highlighted detoxification genes, it is evident that many desert species have evolved different strategies to cope with bitter, toxic plants in the absence of more palatable options.

### Signatures of selection potentially involved in reproductive isolation

We reported significant evolutionary changes linked to genes for sperm motility, spermatogenesis, and pheromone reception that may lend support to a role of sexual selection in the evolution of the cactus mouse genome. The comparison of sperm morphology and behaviour in *Peromyscus* species has revealed a link between sperm traits and reproductive strategies. Sperm of the promiscuous *P. maniculatus*, for example, can aggregate on the basis of relatedness, thus increasing motility and providing a competitive advantage against other males, whereas sperm of the monogamous *P. polionotus* lacks these adaptations (Fisher, Giomi, Hoekstra, & Mahadevan, 2014; Fisher & Hoekstra, 2010). The four gene families associated with sperm motility showed a conspicuous contraction in the cactus mouse when compared to other *Peromyscus* species (9 versus 26, 17 and 21 in *P. leucopus*, *P. maniculatus*, and *P. polionotus*, respectively). Although these results may suggest a correlation between reduction in the number of sperm motility genes and a monogamous reproductive strategy, *P. maniculatus* also has fewer genes than *P. polionotus* in these gene families (17 versus 21 in total). Nonetheless, these candidate genes represent interesting targets for future studies on sperm competition within the cactus mouse and among *Peromyscus*.

### Conclusions

The high-quality assembly of the cactus mouse genome and the candidate genes identified in this study build on the growing body of genomic resources available to further understand the genomic and physiological basis of desert adaptation in the cactus mouse and other



species. Taken together, our results indicate that the strongest signatures of selection in the cactus mouse genome are consistent with adaptations to life in the desert, which are mostly, but not solely, associated with high temperatures and dehydration. Contrary to expectations, we did not find a pervasive signature of selection at genes involved in solute-water balance in the kidneys. However, this does not necessarily discount the relative role of these organs under thermal and hyperosmotic stress, as we have not yet tested when and where in the body the expression of these candidate genes is beneficial. Our analyses also show that signatures of selection are widespread across the cactus mouse genome, with all autosomes showing selective sweeps, and that they are not affected by chromosome-level patterns of standing genetic variation, sequence or structural. Sweeps seem to be associated with high local  $\pi$ , instead. Dynamic gene families and enrichment of several GO terms associated with selective sweeps indicate that the genetic basis of at least some desert-adapted traits may be highly polygenic. In the future, evolutionary and physiological genomics work stemming from these results will allow us to better characterize the phenotypes and genotypes associated with desert adaptations in the cactus mouse, and to understand how they evolved.

### Data availability

All read data for the genome assembly are housed on ENA under project ID PRJEB33593. Specifically, genome assembly (ERZ1195825), 300 bp PE (ERR3445708), 500 bp PE (ERR3446161), 8 kb mate pair (ERR3446162), 5 kb mate pair (ERR3446317), 3 kb mate pair (ERR3446318), 7 kb mate pair (ERR3446319), Hi-C (ERR3446437), 10X Genomics (ERR3447855). Whole genome resequencing data for 26 cactus mice are housed on ENA under project ID PRJEB35488. Scripts for the genome assembly and all other analyses can be found at [https://github.com/atigano/Peromyscus\\_eremicus\\_genome/](https://github.com/atigano/Peromyscus_eremicus_genome/)

### Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

### Acknowledgements

We would like to thank Christopher Tracy for access to the Boyd Deep Canyon Reserve, Adam Stuckert and Douglas Kelt for help with field sampling, the Biotechnology Resource Center at Cornell University for preparation of the whole genome resequencing libraries, the MacManes Lab and the Evolutionary Genomics Journal Group at the University of New Hampshire (UNH) for useful comments on earlier versions of the manuscript, and three anonymous reviewers for their constructive feedback. All analyses were performed on the UNH Premise Cluster. This work was funded by the National Institute of Health National Institute of General Medical Sciences to MDM (1R35GM128843)

### References

- Ai H, Fang X, Yang B, Huang Z, Chen H, Mao L, ... Huang L (2015). Adaptation and possible ancient interspecies introgression in pigs identified by whole-genome sequencing. *Nature Genetics*, 47(3), 217–225. [PubMed: 25621459]
- Bactrian Camels Genome Sequencing and Analysis Consortium, Jirimutu, Wang Z, Ding G, Chen G, Sun Y, ... Meng H. (2012). Genome sequences of wild and domestic bactrian camels. *Nature Communications*, 3, 1202.
- Barko VA, & Feldhamer GA (2002). Cotton mice (*Peromyscus gossypinus*) in southern Illinois: evidence for hybridization with white-footed mice (*Peromyscus leucopus*). *The American Midland Naturalist*, 147(1), 109–116.

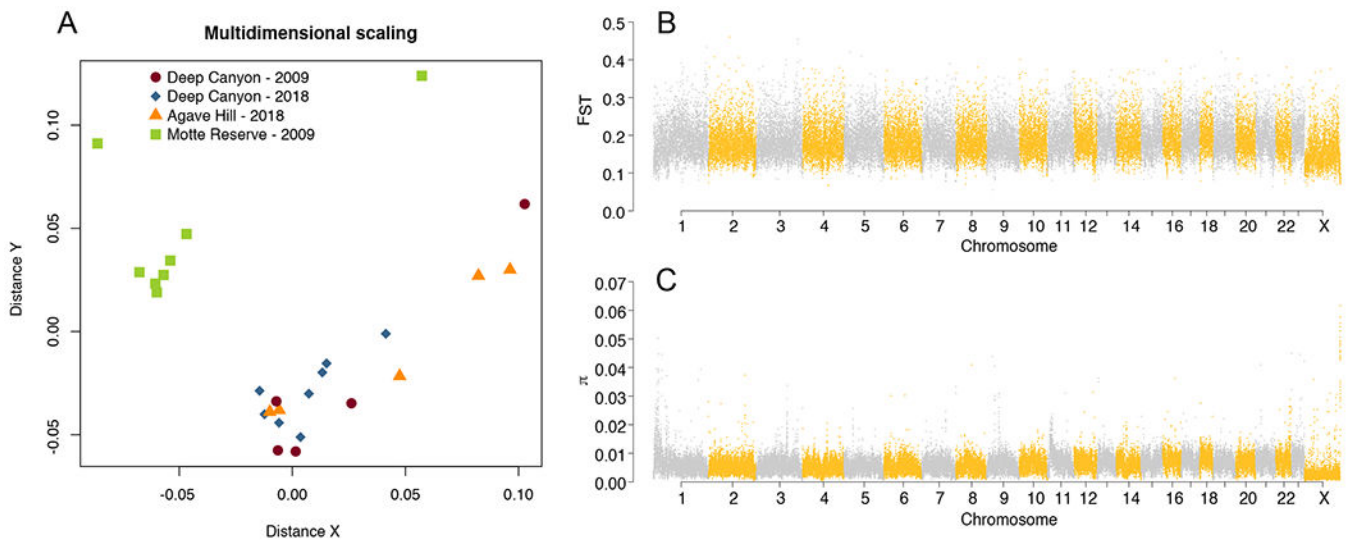
- Barrett RDH, & Schluter D (2008). Adaptation from standing genetic variation. *Trends in Ecology & Evolution*, 23(1), 38–44. [PubMed: 18006185]
- Berner D, & Salzburger W (2015). The genomics of organismal diversification illuminated by adaptive radiations. *Trends in Genetics*, 31(9), 491–499. [PubMed: 26259669]
- Bradley WG, & Mauer RA (1973). Rodents of a creosote bush community in southern Nevada. *The Southwestern Naturalist*, 17(4), 333–344.
- Bromham L (2009). Why do species vary in their rate of molecular evolution? *Biology Letters*, 5(3), 401–404. doi: 10.1098/rsbl.2009.0136 [PubMed: 19364710]
- Buchfink B, Xie C, & Huson DH (2015). Fast and sensitive protein alignment using DIAMOND. *Nature Methods*, 12(1), 59–60. [PubMed: 25402007]
- Burg MB, Ferraris JD, & Dmitrieva NI (2007). Cellular response to hyperosmotic stresses. *Physiological Reviews*, 87(4), 1441–1474. [PubMed: 17928589]
- Butler J, MacCallum I, Kleber M, Shlyakhter IA, Belmonte MK, Lander ES, ... Jaffe DB (2008). ALLPATHS: de novo assembly of whole-genome shotgun microreads. *Genome Research*, 18(5), 810–820. [PubMed: 18340039]
- Cantarel BL, Korf I, Robb SMC, Parra G, Ross E, Moore B, ... Yandell M (2008). MAKER: an easy-to-use annotation pipeline designed for emerging model organism genomes. *Genome Research*, 18(1), 188–196. [PubMed: 18025269]
- Charlesworth B (2009). Effective population size and patterns of molecular evolution and variation. *Nature Reviews. Genetics*, 10(3), 195–205.
- Chen S, Zhou Y, Chen Y, & Gu J (2018). fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics*, 34(17), i884–i890. [PubMed: 30423086]
- Colella JP, Tigano A, & MacManes MD (2019). Higher quality de novo genome assemblies from degraded museum specimens: a linked-read approach to museomics. *bioRxiv*, 716506. doi: 10.1101/716506
- Cooke SJ, Sack L, Franklin CE, Farrell AP, Beardall J, Wikelski M, & Chown SL (2013). What is conservation physiology? Perspectives on an increasingly integrated and essential science. *Conservation Physiology*, 1(1). doi: 10.1093/conphys/cot001
- De Bie T, Cristianini N, Demuth JP, & Hahn MW (2006). CAFE: a computational tool for the study of gene family evolution. *Bioinformatics*, 22(10), 1269–1271. [PubMed: 16543274]
- DeGiorgio M, Huber CD, Hubisz MJ, Hellmann I, & Nielsen R (2016). SweepFinder2: increased sensitivity, robustness and flexibility. *Bioinformatics*, 32(12), 1895–1897. [PubMed: 27153702]
- Durand NC, Robinson JT, Shamim MS, Machol I, Mesirov JP, Lander ES, & Aiden EL (2016). Juicebox provides a visualization system for Hi-C contact maps with unlimited zoom. *Cell Systems*, 3(1), 99–101. [PubMed: 27467250]
- Durand NC, Shamim MS, Machol I, Rao SSP, Huntley MH, Lander ES, & Aiden EL (2016). Juicer provides a one-click system for analyzing loop-resolution Hi-C experiments. *Cell Systems*, 3(1), 95–98. [PubMed: 27467249]
- Emms DM, & Kelly S (2018). OrthoFinder2: fast and accurate phylogenomic orthology analysis from gene sequences. *bioRxiv*, 466201. doi: 10.1101/466201
- English AC, Richards S, Han Y, Wang M, Vee V, Qu J, ... Gibbs RA (2012). Mind the gap: upgrading genomes with Pacific Biosciences RS long-read sequencing technology. *PLoS One*, 7(11), e47768. [PubMed: 23185243]
- Fabre P-H, Hautier L, Dimitrov D, & Douzery EJP (2012). A glimpse on the pattern of rodent diversification: a phylogenetic approach. *BMC Evolutionary Biology*, 12, 88. [PubMed: 22697210]
- Faust GG, & Hall IM (2014). SAMBLASTER: fast duplicate marking and structural variant read extraction. *Bioinformatics*, 30(17), 2503–2505. [PubMed: 24812344]
- Fisher HS, Giomi L, Hoekstra HE, & Mahadevan L (2014). The dynamics of sperm cooperation in a competitive environment. *Proceedings of the Royal Society B: Biological Sciences*, 281(1790), 20140296.
- Fisher HS, & Hoekstra HE (2010). Competition drives cooperation among closely related sperm of deer mice. *Nature*, 463(7282), 801–803. [PubMed: 20090679]

- Fritz SA, Bininda-Emonds ORP, & Purvis A (2009). Geographical variation in predictors of mammalian extinction risk: big is bad, but only in the tropics. *Ecology Letters*, 12(6), 538–549. [PubMed: 19392714]
- Fumagalli M, Vieira FG, Linderoth T, & Nielsen R (2014). ngsTools: methods for population genetics analyses from next-generation sequencing data. *Bioinformatics*, 30(10), 1486–1487. [PubMed: 24458950]
- Garcia J, & Hankins W (1975). The evolution of bitter and the acquisition of toxiphobia Olfaction and Taste: 5th Symposium, pp. 39–45. doi: 10.1016/b978-0-12-209750-8.50014-7
- Garrison E, & Marth G (2012). Haplotype-based variant detection from short-read sequencing. Retrieved from <http://arxiv.org/abs/1207.3907>
- Giorello FM, Feijoo M, D'Elía G, Naya DE, Valdez L, Opazo JC, & Lessa EP (2018). An association between differential expression and genetic divergence in the Patagonian olive mouse (*Abrothrix olivacea*). *Molecular Ecology*, 27(16), 3274–3286. doi: 10.1111/mec.14778
- Glendinning JI (1994). Is the bitter rejection response always adaptive? *Physiology & Behavior*, 56(6), 1217–1227. [PubMed: 7878094]
- Harris AM, Garud NR, & DeGiorgio M (2018). Detection and classification of hard and soft sweeps from unphased genotypes by multilocus genotype identity. *Genetics*, 210(4), 1429–1452. [PubMed: 30315068]
- Hino J, Nakatani M, Arai Y, Tsuchida K, Shirai M, Miyazato M, & Kangawa K (2017). Overexpression of bone morphogenetic protein-3b (BMP-3b) in adipose tissues protects against high-fat diet-induced obesity. *International Journal of Obesity*, 41(4), 483–488. [PubMed: 28104917]
- Hochstrasser M (2009). Origin and function of ubiquitin-like proteins. *Nature*, 458(7237), 422–429. [PubMed: 19325621]
- Hoelzel AR (2010). Looking backwards to look forwards: conservation genetics in a changing world. *Conservation Genetics*, 11(2), 655–660.
- Huber CD, DeGiorgio M, Hellmann I, & Nielsen R (2016). Detecting recent selective sweeps while controlling for mutation rate and background selection. *Molecular Ecology*, 25(1), 142–156. [PubMed: 26290347]
- IPCC. (2018). Global Warming of 1.5° C: An IPCC special report on the impacts of global warming of 1.5° C above pre-industrial levels and related global greenhouse gas emission pathways, in the context of strengthening the global response to the threat of climate change, sustainable development, and efforts to eradicate poverty. Intergovernmental Panel on Climate Change.
- Kaback DB, Guacci V, Barber D, & Mahon JW (1992). Chromosome size-dependent control of meiotic recombination. *Science*, 256(5054), 228–232. [PubMed: 1566070]
- Kampinga HH (1993). Thermotolerance in mammalian cells. Protein denaturation and aggregation, and stress proteins. *Journal of Cell Science*, 104, 11–17. [PubMed: 8449990]
- Kim Y, & Stephan W (2002). Detecting a local signature of genetic hitchhiking along a recombining chromosome. *Genetics*, 160(2), 765–777. [PubMed: 11861577]
- Kohl KD, Weiss RB, Cox J, Dale C, & Dearing MD (2014). Gut microbes of mammalian herbivores facilitate intake of plant toxins. *Ecology Letters*, 17(10), 1238–1246. [PubMed: 25040855]
- Kordonowy L, Lombardo KD, Green HL, Dawson MD, Bolton EA, LaCourse S, & MacManes MD (2017). Physiological and biochemical changes associated with acute experimental dehydration in the desert adapted mouse, *Peromyscus eremicus*. *Physiological Reports*, 5(6). doi: 10.14814/phy2.13218
- Korneliussen TS, Albrechtsen A, & Nielsen R (2014). ANGSD: Analysis of Next Generation Sequencing Data. *BMC Bioinformatics*, 15, 356. [PubMed: 25420514]
- Lamitina T, Huang CG, & Strange K (2006). Genome-wide RNAi screening identifies protein damage as a regulator of osmoprotective gene expression. *Proceedings of the National Academy of Sciences of the United States of America*, 103(32), 12173–12178. [PubMed: 16880390]
- León-Paniagua L, Navarro-Sigüenza AG, Hernández-Baños BE, & Morales JC (2007). Diversification of the arboreal mice of the genus *Habromys* (Rodentia: Cricetidae: Neotominae) in the Mesoamerican highlands. *Molecular Phylogenetics and Evolution*, 42(3), 653–664. [PubMed: 17070711]

- Leo SST, & Millien V (2017). Microsatellite markers reveal low frequency of natural hybridization between the white-footed mouse (*Peromyscus leucopus*) and deer mouse (*Peromyscus maniculatus*) in southern Quebec, Canada. *Genome*, 60(5), 454–463. [PubMed: 28177836]
- Li D, & Zhang J (2014). Diet shapes the evolution of the vertebrate bitter taste receptor gene repertoire. *Molecular Biology and Evolution*, 31(2), 303–309. [PubMed: 24202612]
- Li G, Davis BW, Eizirik E, & Murphy WJ (2016). Phylogenomic evidence for ancient hybridization in the genomes of living cats (Felidae). *Genome Research*, 26(1), 1–11.
- Li H, & Durbin R (2009). Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*, 25(14), 1754–1760. [PubMed: 19451168]
- Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, ... 1000 Genome Project Data Processing Subgroup. (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, 25(16), 2078–2079. [PubMed: 19505943]
- Linnen CR, Linnen CR, Kingsley EP, Kingsley EP, Jensen JD, & Hoekstra HE (2009). On the origin and spread of an adaptive allele in deer mice. *Science*, 325(5944), 1095–1098. [PubMed: 19713521]
- Long AD, Baldwin-Brown J, Tao Y, Cook VJ, Balderrama-Gutierrez G, Corbett-Detig R, ... Barbour AG (2019). The genome of *Peromyscus leucopus*, natural host for Lyme disease and other emerging infections. *Science Advances*, 5(7), eaaw6441. [PubMed: 31355335]
- MacManes MD (2017). Severe acute dehydration in a desert rodent elicits a transcriptional response that effectively prevents kidney injury. *American Journal of Physiology. Renal Physiology*, 313(2), F262–F272. [PubMed: 28381460]
- Macmillen RE (1965). Aestivation in the cactus mouse, *Peromyscus eremicus*. *Comparative Biochemistry and Physiology*, 16(2), 227–248. [PubMed: 5865202]
- Marçais G, Delcher AL, Phillippy AM, Coston R, Salzberg SL, & Zimin A (2018). MUMmer4: A fast and versatile genome alignment system. *PLoS Computational Biology*, 14(1), e1005944. [PubMed: 29373581]
- Marra NJ, Eo SH, Hale MC, Waser PM, & DeWoody JA (2012). A priori and a posteriori approaches for finding genes of evolutionary interest in non-model species: osmoregulatory genes in the kidney transcriptome of the desert rodent *Dipodomys spectabilis* (banner-tailed kangaroo rat). *Comparative Biochemistry and Physiology. Part D, Genomics & Proteomics*, 7(4), 328–339.
- Marra NJ, Romero A, & DeWoody JA (2014). Natural selection and the genetic basis of osmoregulation in heteromyid rodents as revealed by RNA-seq. *Molecular Ecology*, 23(11), 2699–2711. [PubMed: 24754676]
- Meserve PL (1976). Habitat and resource utilization by rodents of a California coastal sage scrub community. *The Journal of Animal Ecology*, 45(3), 647–666.
- Messer PW, & Neher RA (2012). Estimating the strength of selective sweeps from deep population diversity data. *Genetics*, 191(2), 593–605. [PubMed: 22491190]
- Mi H, Huang X, Muruganujan A, Tang H, Mills C, Kang D, & Thomas PD (2017). PANTHER version 11: expanded annotation data from Gene Ontology and Reactome pathways, and data analysis tool enhancements. *Nucleic Acids Research*, 45(D1), D183–D189. [PubMed: 27899595]
- Murie M (1961). Metabolic characteristics of mountain, desert and coastal populations of *Peromyscus*. *Ecology*, 42(4), 723–740.
- Nachman MW, Hoekstra HE, & D'Agostino SL (2003). The genetic basis of adaptive melanism in pocket mice. *Proceedings of the National Academy of Sciences of the United States of America*, 100(9), 5268–5273. [PubMed: 12704245]
- Nakano M, Kelly EJ, & Rettie AE (2009). Expression and characterization of CYP4V2 as a fatty acid omega-hydroxylase. *Drug Metabolism and Disposition: The Biological Fate of Chemicals*, 37(11), 2119–2122. [PubMed: 19661213]
- Nielsen R, Williamson S, Kim Y, Hubisz MJ, Clark AG, & Bustamante C (2005). Genomic scans for selective sweeps using SNP data. *Genome Research*, 15(11), 1566–1575. [PubMed: 16251466]
- Pavlik BM (2008). *The California Deserts: an ecological rediscovery*. University of California Press.
- Pezer Ž, Harr B, Teschke M, Babiker H, & Tautz D (2015). Divergence patterns of genic copy number variation in natural populations of the house mouse (*Mus musculus domesticus*) reveal three

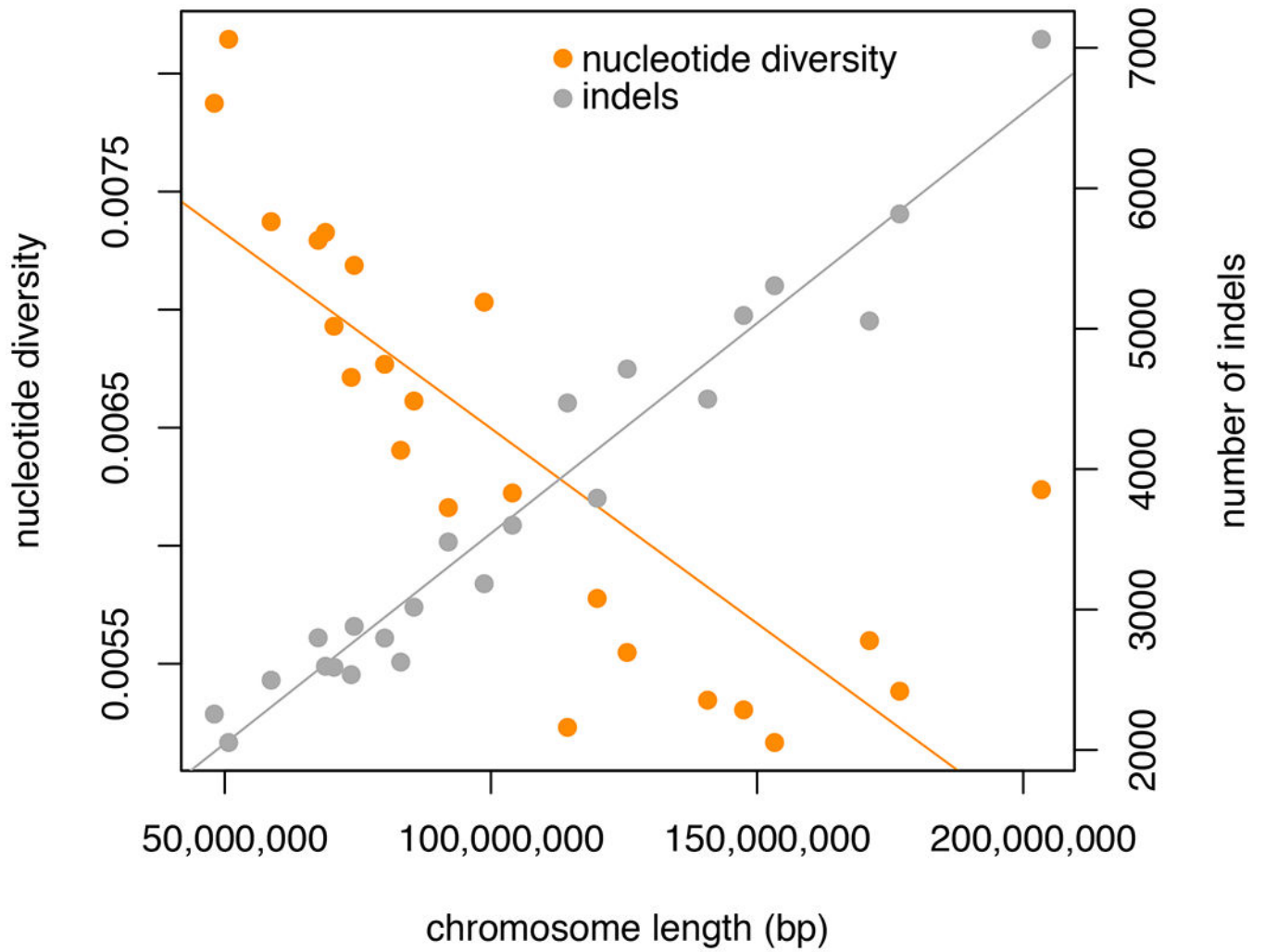
- conserved genes with major population-specific expansions. *Genome Research*, 25(8), 1114–1124. [PubMed: 26149421]
- Porcelli D, Butlin RK, Gaston KJ, Joly D, & Snook RR (2015). The environmental genomics of metazoan thermal adaptation. *Heredity*, 114(5), 502–514. [PubMed: 25735594]
- Presgraves DC (2018). Evaluating genomic signatures of “the large X-effect” during complex speciation. *Molecular Ecology*, 27(19), 3822–3830. [PubMed: 29940087]
- Riddle BR, Hafner DJ, & Alexander LF (2000). Phylogeography and systematics of the *Peromyscus eremicus* species group and the historical biogeography of North American warm regional deserts. *Molecular Phylogenetics and Evolution*, 17(2), 145–160. [PubMed: 11083930]
- Savolainen O, Lascoux M, & Merilä J (2013). Ecological genomics of local adaptation. *Nature Reviews. Genetics*, 14(11), 807–820.
- Schenk JJ, Rowe KC, & Steppan SJ (2013). Ecological opportunity and incumbency in the diversification of repeated continental colonizations by muroid rodents. *Systematic Biology*, 62(6), 837–864. [PubMed: 23925508]
- Schmidt-Nielsen K (1964). *Desert Animal: physiological problems of heat and water*. Clarendon Press.
- Schmidt-Nielsen K, & Schmidt-Nielsen B (1952). Water metabolism of desert mammals 1. *Physiological Reviews*, 32(2), 135–166. [PubMed: 14929697]
- Sikes RS, & Animal Care and Use Committee of the American Society of Mammalogists. (2016). 2016 Guidelines of the American Society of Mammalogists for the use of wild mammals in research and education. *Journal of Mammalogy*, 97(3), 663–688. [PubMed: 29692469]
- Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, & Zdobnov EM (2015). BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*, 31(19), 3210–3212. [PubMed: 26059717]
- Skopec MM, & Dearing MD (2011). Differential expression and activity of catechol-O-methyl transferase (COMT) in a generalist (*Neotoma albigula*) and juniper specialist (*Neotoma stephensi*) woodrat. *Comparative Biochemistry and Physiology. Toxicology & Pharmacology: CBP*, 154(4), 383–390. [PubMed: 21820082]
- Smalec BM, Heider TN, Flynn BL, & O’Neill RJ (2019). A centromere satellite concomitant with extensive karyotypic diversity across the *Peromyscus* genus defies predictions of molecular drive. *Chromosome Research: An International Journal on the Molecular, Supramolecular and Evolutionary Aspects of Chromosome Biology*. doi: 10.1007/s10577-019-09605-1
- Smit AFA, Hubley R, & Green P (2015). RepeatMasker Open-4.0. 2013--2015.
- Smith JM, & Haigh J (1974). The hitch-hiking effect of a favourable gene. *Genetical Research*, 23(1), 23–35. [PubMed: 4407212]
- Somero GN (2010). The physiology of climate change: how potentials for acclimatization and genetic adaptation will determine “winners” and “losers.” *The Journal of Experimental Biology*, 213(6), 912–920. [PubMed: 20190116]
- Supek F, Bošnjak M, Škunca N, & Šmuc T (2011). REVIGO summarizes and visualizes long lists of gene ontology terms. *PloS One*, 6(7), e21800. [PubMed: 21789182]
- Takei Y, Bartolo RC, Fujihara H, Ueta Y, & Donald JA (2012). Water deprivation induces appetite and alters metabolic strategy in *Notomys alexis*: unique mechanisms for water production in the desert. *Proceedings of the Royal Society B: Biological Sciences*, 279(1738), 2599–2608.
- Thybert D, Roller M, Navarro FCP, Fiddes I, Streeter I, Feig C, ... Flicek P (2018). Repeat associated mechanisms of genome evolution and function revealed by the *Mus caroli* and *Mus pahari* genomes. *Genome Research*, 28(4), 448–459. [PubMed: 29563166]
- Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, ... Earl AM (2014). Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PloS One*, 9(11), e112963. [PubMed: 25409509]
- Walsberg GE (2000). Small mammals in hot deserts: some generalizations revisited. *BioScience*, 50(2), 109–120.
- Wang K, & Zhao H (2015). Birds generally carry a small repertoire of bitter taste receptor genes. *Genome Biology and Evolution*, 7(9), 2705–2715. [PubMed: 26342138]

- Wellenreuther M, Mérot C, Berdan E, & Bernatchez L (2019). Going beyond SNPs: The role of structural genomic variants in adaptive evolution and species diversification. *Molecular Ecology*, 28(6), 1203–1209. [PubMed: 30834648]
- Wilson Sayres MA (2018). Genetic diversity on the sex chromosomes. *Genome Biology and Evolution*, 10(4), 1064–1078. [PubMed: 29635328]
- Wu H, Guang X, Al-Fageeh MB, Cao J, Pan S, Zhou H, ... Wang J (2014). Camelid genomes reveal evolution and adaptation to desert environments. *Nature Communications*, 5, 5188.
- Yang J, Li W-R, Lv F-H, He S-G, Tian S-L, Peng W-F, ... Liu M-J (2016). Whole-genome sequencing of native sheep provides insights into rapid adaptations to extreme environments. *Molecular Biology and Evolution*, 33(10), 2576–2592. [PubMed: 27401233]



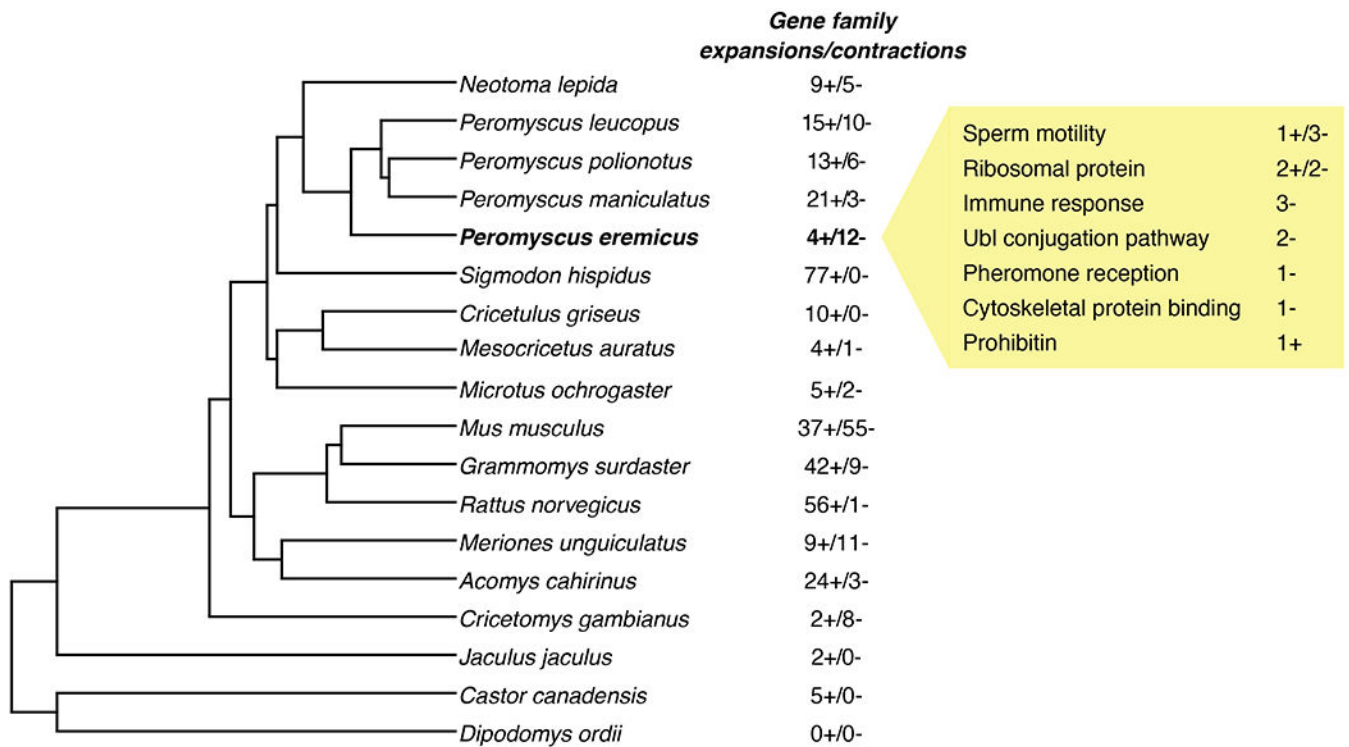
**Figure 1.**

Diversity and differentiation in the cactus mouse. a) MDS plot showing relative distance among individuals based on downsampling to a single base at 43.7 million variable sites. Note outlier from Motte on the far left side of the plot. b) Manhattan plot showing patterns of differentiation based on  $F_{ST}$  between Motte and Deep Canyon Reserves (after outlier removal). c) Manhattan plot showing patterns of nucleotide diversity  $\pi$  from all samples combined (after outlier removal).

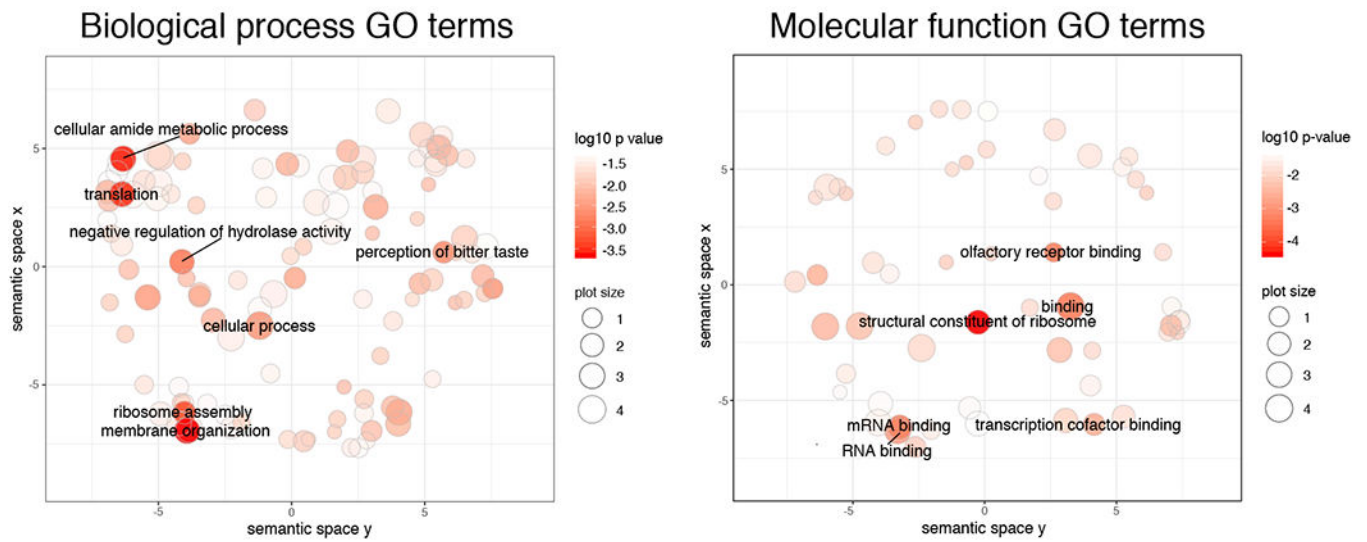


**Figure 2.** Plot showing mean nucleotide diversity and number of indels as a function of chromosome length ( $p < 0.001$  in both cases, albeit with opposite trends).

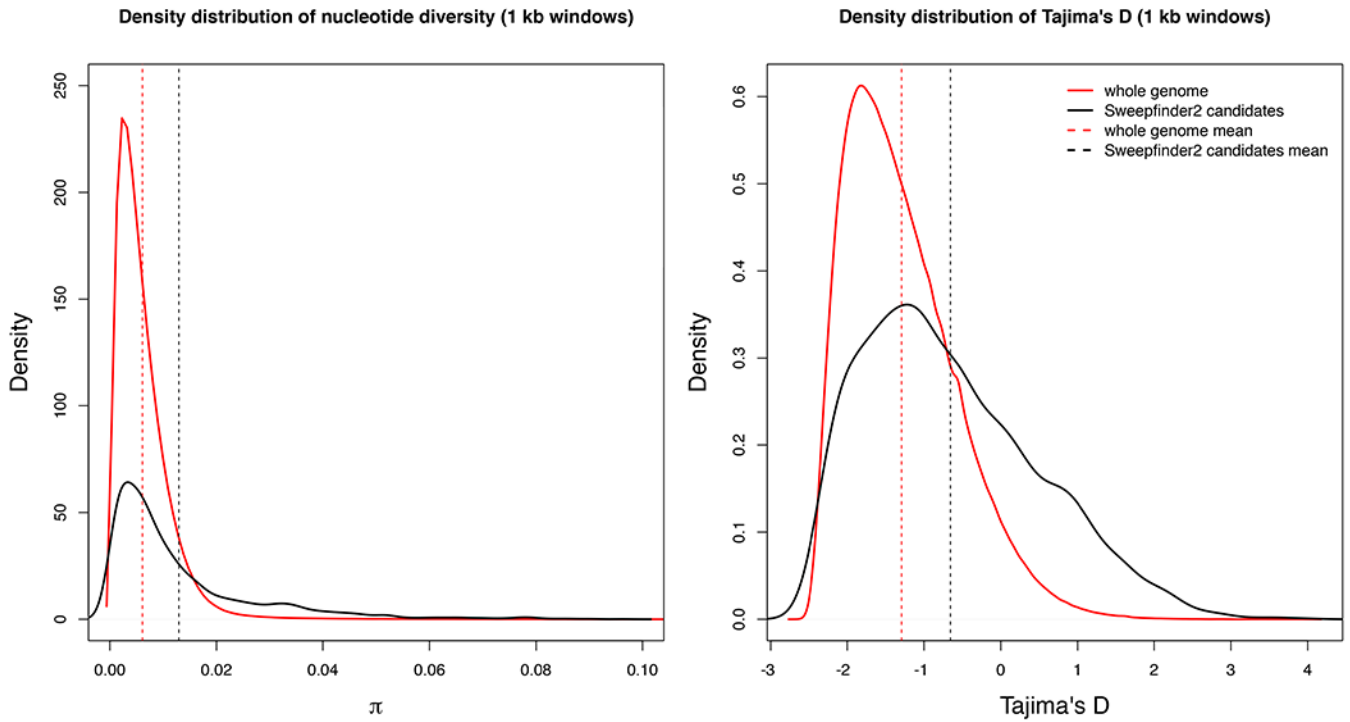




**Figure 3.** Species tree built in *Orthofinder2* from protein sequences of 18 species in the Myodonta clade (Order: Rodentia). Beside each species name are the number of gene families that underwent significant ( $p < 0.05$ ) expansions (+) or contractions (-). In the yellow box are the functions and number of expanded/contracted gene families in the cactus mouse (*Peromyscus eremicus*) relative to the closest ancestral node.



**Figure 4.** Scatterplot showing clusters representative of enriched GO terms after semantic reduction in *REVIGO* for biological process GO terms (left) and molecular function GO terms (right). Only the names of GO clusters with a p-value  $< 10^{-2.5}$  are shown for visual clarity. The full list of genes and reduced GO terms in *REVIGO* can be found in Supplementary Tables 3, 4 and 5.



**Figure 5.**

Density plots comparing distribution of  $\pi$  (left) and Tajima's D (right) across the genome (in red) and across *Sweepfinder2* candidate regions only (in black). Values are calculated in 1 kb non-overlapping windows along the genome. Dashed vertical lines show the means across the genome in red and across *Sweepfinder2* candidate regions only in black. Means across the genome and across *Sweepfinder2* candidate regions only are significantly different in both cases ( $p < 0.001$ ).

**Table 1.**

Details on sampling locations of individuals sequenced for population genomics analyses.

<b>Reserve</b>	<b>Location</b>	<b>Latitude/longitude</b>	<b>Year</b>	<b>Sample size</b>
Motte Rimrock	Motte	33°48'N/117°15'W	2009	8*
Boyd Deep Canyon	Deep Canyon	33°38'N/116°22'W	2009	8
Boyd Deep Canyon	Deep Canyon	33°38'N/116°22'W	2018	4
Boyd Deep Canyon	Agave Hill	33°38'N/116°24'W	2018	6

\* denotes that the individual outlier was sampled here, and effective sample size was reduced to 7 after outlier removal

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

**Table 2.**

Summary of assembly statistics before and after scaffolding with Hi-C data.

	<b>Pre Hi-C assembly</b>	<b>Post Hi-C assembly</b>
Number of scaffolds	7650	24 (+ 6,785 unplaced scaffolds)
Total size of scaffolds	2.7 Gbp	2.5 Gbp (+ 173 Mb unplaced sequence)
Longest scaffold	13.7 Mbp	203.4 Mbp
Scaffold N50/L50	1.3 Mbp/530	120 Mbp/9
Gaps %	5.51	3.53

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript