



Genome sequence of a diabetes-prone rodent reveals a mutation hotspot around the ParaHox gene cluster

Adam D. Hargreaves^a, Long Zhou^b, Josef Christensen^c, Ferdinand Marlétaz^{a,d}, Shiping Liu^b, Fang Li^b, Peter Gildsig Jansen^c, Enrico Spiga^e, Matilde Thye Hansen^c, Signe Vendelbo Horn Pedersen^c, Shameek Biswas^f, Kyle Serikawa^f, Brian A. Fox^f, William R. Taylor^e, John Frederick Mulley^g, Guojie Zhang^{b,h,i,1}, R. Scott Heller^{c,1}, and Peter W. H. Holland^{a,1}

^aDepartment of Zoology, University of Oxford, Oxford, OX1 3PS, United Kingdom; ^bChina National Genebank, BGI-Shenzhen, 518083, Shenzhen, Guangdong, China; ^cNovo Nordisk, Måløv, DK-2760, Denmark; ^dMolecular Genetics Unit, Okinawa Institute of Science and Technology Graduate University, Onna, Okinawa 904-0495, Japan; ^eFrancis Crick Institute, London, NW1 1AT, United Kingdom; ^fNovo Nordisk Research Centre, Seattle, WA 98109; ^gSchool of Biological Sciences, Bangor University, Bangor, LL57 2DG, United Kingdom; ^hState Key Laboratory of Genetic Resources and Evolution, Kunming Institute of Zoology, Chinese Academy of Sciences, 650223, Kunming, China; and ⁱCentre for Social Evolution, Department of Biology, University of Copenhagen, DK-2100 Copenhagen, Denmark

Edited by Richard M. Harland, University of California, Berkeley, CA, and approved June 2, 2017 (received for review February 23, 2017)

The sand rat *Psammomys obesus* is a gerbil species native to deserts of North Africa and the Middle East, and is constrained in its ecology because high carbohydrate diets induce obesity and type II diabetes that, in extreme cases, can lead to pancreatic failure and death. We report the sequencing of the sand rat genome and discovery of an unusual, extensive, and mutationally biased GC-rich genomic domain. This highly divergent genomic region encompasses several functionally essential genes, and spans the ParaHox cluster which includes the insulin-regulating homeobox gene *Pdx1*. The sequence of sand rat *Pdx1* has been grossly affected by GC-biased mutation, leading to the highest divergence observed for this gene across the Bilateria. In addition to genomic insights into restricted caloric intake in a desert species, the discovery of a localized chromosomal region subject to elevated mutation suggests that mutational heterogeneity within genomes could influence the course of evolution.

desert rodent | type 2 diabetes | homeobox | *Pdx1* | gene conversion

Arid environments impose extreme physiological demands on animals because of low food and water availability. The sand rat *Psammomys obesus* (Fig. 1A) is a member of the subfamily Gerbillinae, most species of which live in deserts and arid environments (Fig. 1B). *P. obesus* has emerged as a model for research into diet-induced type II diabetes because, if provided with high carbohydrate diets, the majority of individuals become obese and develop classic diabetes symptoms, in the most extreme cases leading to pancreatic failure and death (1–4).

In searching for the molecular basis of this unusual phenotype, attention has been paid to the *Pdx1* homeobox gene, also called *Ipf1*, *Irx1*, *Sf1*, or *Xlox* (5–9), the central and most highly conserved member of the ParaHox gene cluster (10). *Pdx1* is the only member of the Pdx gene family in tetrapods and encodes a homeodomain that has been invariant across their evolution. Mammalian *Pdx1* is expressed in pancreatic beta cells and encodes a homeodomain transcription factor that acts as a transcriptional activator of *insulin* and other pancreatic hormone genes (11, 12). A pivotal role in insulin regulation is also reflected in the association of heterozygous *Pdx1* mutations with maturity-onset diabetes of the young (*MODY4*) and type II diabetes mellitus in humans (13). Contrary to the usual conservation, several studies have reported inability to detect *Pdx1* in multiple gerbil species, including *P. obesus*, by immunocytochemistry, Western blotting, or PCR. However, *Pdx1* is readily detectable in the closely related spiny mouse, *Acomys cahirinus* (Fig. 1B), leading to the hypothesis that the gene has been lost within the Gerbillinae subfamily, contributing to the compromised ability to regulate insulin in the sand rat (14–16). Such a conclusion would raise further questions, because in addition to its adult functions, *Pdx1* is also essential for pancreatic development in the embryo. For example, targeted deletion in mice causes loss of pancreas and anterior duodenum

and is lethal (9, 17). In humans, pancreatic agenesis has been reported in a patient with a homozygous frameshift mutation before the *Pdx1* homeobox and in a compound heterozygous patient with substitution mutations in helices 1 and 2 of the homeodomain (18–20).

Results

To resolve the conundrum of a putatively absent “essential” gene, we sequenced the *P. obesus* genome by using a standard shotgun strategy (Illumina), using a combination of short and long insert libraries, initially at 85.5× coverage (*SI Appendix, SI Materials and Methods, section 1*). This assembly lacked a *Pdx1* gene, supporting the prevailing hypothesis of a loss of the *Pdx1* gene in gerbils. However, a synteny comparison between *P. obesus* and other mammals delineated a contiguous block of 88 genes (*SI Appendix, Fig. S2*) missing from the assembly including several genes essential

Significance

A core question in evolutionary biology is how mutation and selection adapt and constrain species to specialized habitats. We sequenced the genome of the sand rat, a desert rodent susceptible to nutritionally induced diabetes, and discovered an unusual chromosome region skewed toward G and C nucleotides. This region includes the *Pdx1* homeobox gene, a transcriptional activator of *insulin*, which has undergone massive sequence change, likely contributing to diabetes and adaptation to low caloric intake. Our results imply that mutation rate varies within a genome and that hotspots of high mutation rate may influence ecological adaptation and constraint. In addition, we caution that divergent regions can be omitted by conventional short-read sequencing approaches, a consideration for existing and future genome sequencing projects.

Author contributions: A.D.H., L.Z., W.R.T., G.Z., R.S.H., and P.W.H.H. designed research; A.D.H., L.Z., J.C., S.L., F.L., E.S., M.T.H., S.V.H.P., S.B., K.S., B.A.F., W.R.T., J.F.M., R.S.H., and P.W.H.H. performed research; F.M., W.R.T., and G.Z. contributed new reagents/analytic tools; A.D.H., L.Z., F.M., P.G.J., E.S., S.B., K.S., B.A.F., J.F.M., G.Z., R.S.H., and P.W.H.H. analyzed data; and A.D.H., L.Z., E.S., G.Z., R.S.H., and P.W.H.H. wrote the paper.

Conflict of interest statement: J.C., P.G.J., M.T.H., S.V.H.P., S.B., K.S., B.A.F., and R.S.H. are current or former employees of Novo Nordisk.

This article is a PNAS Direct Submission.

Freely available online through the PNAS open access option.

Data deposition: The data reported in this paper have been deposited in the National Center for Biotechnology Information Short Read Archive (accession nos. [SRA502705](#), [SRR5084169](#), [SRR5084170](#), [SRR5092818](#), [SRR5092819](#), [SRR5092820](#), and [SRR5429486](#)) and [DBJ/ENA/GenBank](#) (accession nos. [NESX000000000](#) and [NESX010000000](#)).

¹To whom correspondence may be addressed. Email: peter.holland@zoo.ox.ac.uk, zhanggj@genomics.cn, or richardscottheller@gmail.com.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1702930114/-DCSupplemental.

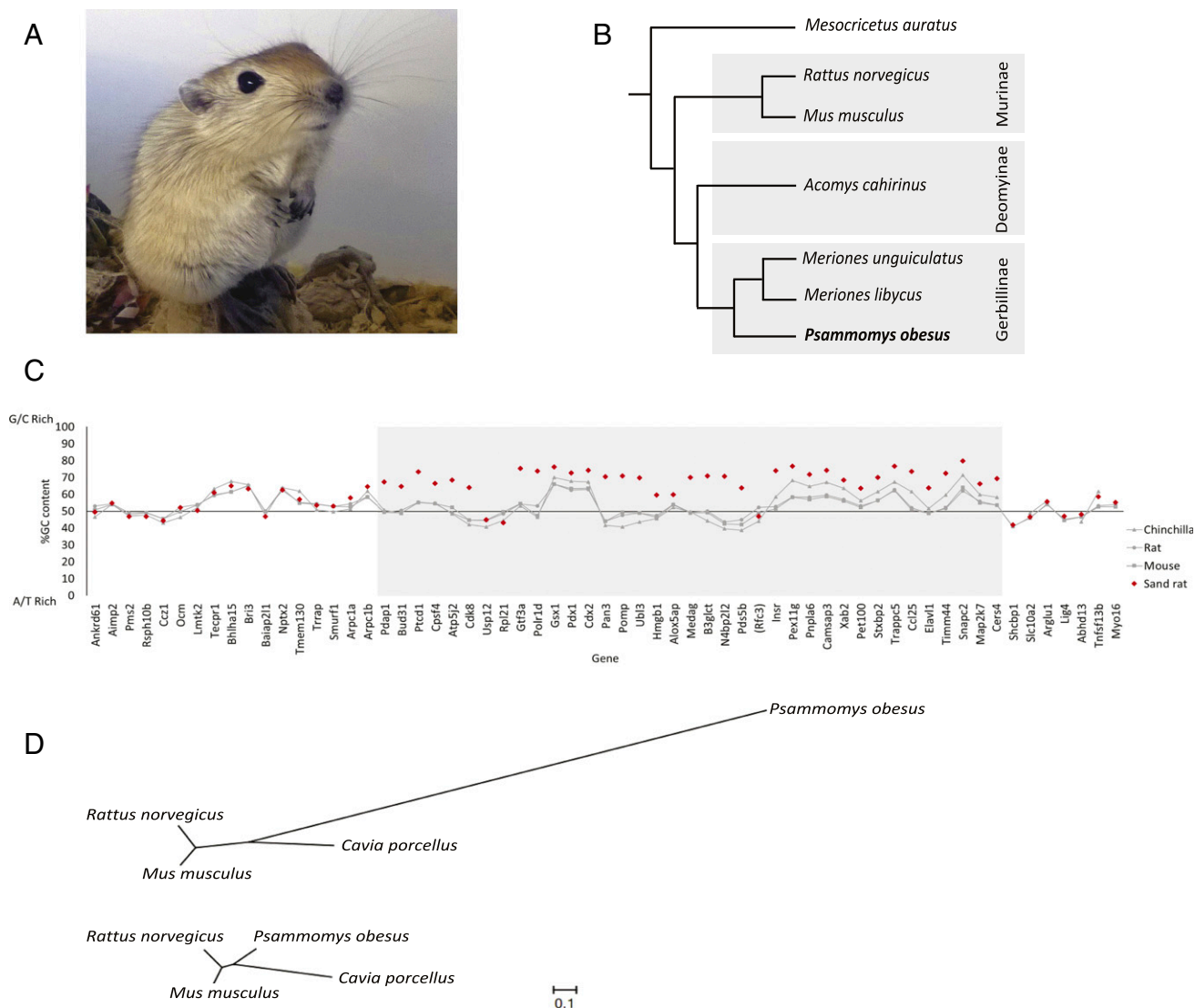


Fig. 1. The sand rat and its genomic hotspot of mutation. (A) Juvenile sand rat *P. obesus*. (B) Cladogram of representative murid rodents indicating the phylogenetic position of sand rat. (C) GC content of genes around the ParaHox cluster of sand rat and other rodents (*Mus musculus*, *Rattus norvegicus*, *Chinchilla lanigera*) revealing a chromosomal hotspot of GC skew in sand rat (shaded in gray). Genes shown in inferred ancestral gene order; parentheses around *Rfc3* indicate this gene has been transposed to a different genomic location in sand rat. Sand rat GC values based on transcriptome and genome sequences; when partial only alignable sequence is compared. (D) Unrooted phylogenetic trees inferred from synonymous changes (d5) only from concatenated alignments of 26 genes in the mutational hotspot (Upper) and 100 random genes (Lower).

to basic cellular functions, such as *Brca2* and *Cdk8*, in addition to *Pdx1*. This finding led us to suspect that standard short read sequencing may have given an incomplete genome assembly, even at high coverage. To resolve whether the gene absence reflected a large-scale deletion or an unusual genomic region, we sequenced the transcriptomes of *P. obesus* liver, pancreatic islets, and duodenum, which contained transcripts for many of the missed genes (SI Appendix, Tables S4–S6). Furthermore, these transcripts show unusually high GC content in most cases, indicating that a large contiguous stretch of elevated GC had either been underrepresented in initial sequencing data or had failed to assemble correctly, most likely due to nucleotide compositional bias. We term such cryptic or hidden sequence “dark DNA.” We therefore isolated GC-rich *P. obesus* genomic DNA by cesium chloride gradient centrifugation, sequenced this fraction after limited amplification by using Illumina MiSeq overlapping paired-end reads, and reassembled the genome incorporating this longer-read sequence data

(SI Appendix, SI Materials and Methods, section 1.5). This approach gave a refined assembly with a total size of 2.38 Gb and a scaffold N50 of 10.4 Mb (Table 1 and SI Appendix, SI Materials and Methods, sections 1, 3, 4, and 6), including much of the dark DNA region in several scaffolds, and containing genes syntenic to a region of chromosome 12 in rat and a region of chromosome 5 and the subtelomeric region of chromosome 8 in mouse. Analysis indicates that the region was initially omitted by standard genome assembly methods because of lower read coverage of GC regions coupled with short sequence read lengths. Comparison of GC content between species demonstrates that sand rat genes are elevated in GC content across this chromosomal region, syntenic to 12 Mb of the rat genome (Fig. 1C and SI Appendix, SI Materials and Methods, section 9). This large region encompasses a 250-kb repeat-rich scaffold containing the sand rat ParaHox cluster and its well-characterized genomic neighbors. We inferred a high W (weak, A/T) to S (strong, G/C) allelic mutation rate in this region of the

Table 1. Metrics of sand rat raw genomic sequencing data and final genome assembly

Genome sequencing and assembly	Value
Total no. of paired-end reads	724,377,486
Total no. of mate-pair reads	1,780,436,140
Total bases sequenced	394,396,928,120
Estimated sequencing coverage, x	87.6
No. of scaffolds >2 kb	1,737
Total length of assembly, bp	2,381,209,849
Longest scaffold, bp	54,616,910
Mean scaffold length, bp	15,794
Scaffold N50, bp	10,461,538
Scaffold L50	63
Contig N50, bp	83,904
Percentage of assembly in scaffolds, %	98.6

Coverage was calculated by using an estimated genome size of 2.51 Gb based on a k-mer analysis (*SI Appendix, SI Materials and Methods*, section 1.3) and is based on paired-end sequencing data only.

P. obesus genome compared with randomly selected genomic regions or homologous regions in other species of rodent (Fig. 1*D* and *SI Appendix, SI Materials and Methods*, section 12 and Tables S11 and S12). The existence of a localized GC-biased stretch of the *P. obesus* genome is striking and of far-reaching importance, and implies the existence of elevated and biased mutational pressure, acting in one region of a mammalian genome. Gene conversion, caused by the nonreciprocal exchange of information during meiosis, is the best characterized process known to cause GC-biased mutation (21).

The full coding sequence of the *P. obesus Pdx1* gene was deduced from the refined genome and transcriptome assemblies, and the gene was found to be expressed in sand rat pancreatic islets and duodenum (*SI Appendix, SI Materials and Methods*, section 7). The 60-aa homeodomain of Pdx1 shows 100% conservation across other mammals for which data are available; however, in *P. obesus*, there are remarkable 15-aa differences in the homeodomain, making it the most divergent *Pdx1* gene discovered in the Bilateria (Fig. 2*A*). All but one of the amino acid changes are caused by A/T to G/C mutation. The N-terminal and C-terminal regions are also divergent with numerous deletions, although the hexapeptide motif used in heterodimer formation with TALE proteins is conserved (Fig. 2*B*). Additional RNA sequencing of Mongolian jird (*Meriones unguiculatus*) duodenum reveals that extensive sequence divergence due to GC-biased mutation in *Pdx1* is not unique to sand rat (Fig. 2*A*). Analysis of synonymous and nonsynonymous mutations in *Pdx1* across vertebrates reveals a dN/dS ratio of 2.6 (dN = 39; dS = 15) in the lineage leading to *P. obesus* and *M. unguiculatus* (*SI Appendix, Fig. S10*). High dN/dS ratios are often taken as evidence for positive selection, but can be skewed by mutational processes such as GC-biased gene conversion (22). Despite its radical divergence, Pdx1 is the closest homeodomain by BLASTP, and phylogenetic analysis places it as a rodent Pdx1 on a long branch (Fig. 2*B*); extensive synteny with the ParaHox region of mouse and rat confirms it is the true and single *Pdx1* ortholog (*SI Appendix, Table S9*). Evidence that the locus is functional includes expression in pancreas and duodenum, and the fact that extensive polymorphism is found in the 3' untranslated region but is limited in the coding sequence (*SI Appendix, Fig. S11*), indicating that the coding region is under functional constraint despite extensive mutation. Extreme deviation from the expected sequence explains why antibodies and PCR failed to detect *Pdx1* in sand rat, Mongolian jird, and, potentially, other gerbil species (14–16).

These findings indicate that GC-biased mutation has driven radical changes in an otherwise highly conserved homeobox

gene; these changes could be maladaptive and constrain the physiological capability of the sand rat, or adaptive enhancing ability to live in arid regions. To test whether the extent of sequence divergence is unusual for sand rat proteins, we calculated a “protein deviation index” (PDI) (*SI Appendix, SI Materials and Methods*, section 5) for all 1:1 mammalian orthologs by dividing mouse-human protein sequence identity by mouse-sand rat sequence identity (Fig. 2*C*). This analysis is distinct from identifying the fastest evolving proteins and specifically identifies proteins that have undergone uncharacteristic divergence in sand rat. We find the majority of sand rat proteins are highly similar to mouse or human (mode PDI = 1.0); in contrast, Pdx1 is unusually divergent (mouse-sand rat 54.82%, mouse-human 91.37%; PDI = 1.67). To test whether other genes implicated in glucose metabolism or pancreatic function are also divergent, we compiled a list of 45 candidates from human studies including all genes implicated in monogenic diabetes (23) and genes for which coding sequence variants have been strongly associated with type 2 diabetes (24). Of the 33 genes with clear 1:1 orthologs between human, mouse, and sand rat, 32 lie between position 225 and 10,195 in our PDI ranking, indicating that they are not unusually divergent in sand rat. Pdx1 is ranked first and is the most unusually divergent protein identified in the sand rat predicted proteome (*SI Appendix, Materials and Methods*, section 8 and Tables S8 and S10). Strikingly, 7 of the top 10 highest PDI results correspond to genes located within the mutational hotspot (*SI Appendix, Table S8*), indicating that GC-biased mutation is contributing to coding sequence divergence across this region.

The mutations fixed in sand rat *Pdx1* gene do not cause frameshifts or truncations in known domains, and molecular modeling reveals that the sand rat Pdx1 homeodomain has the ability to form all three helices required for DNA binding (Fig. 3*A*). To examine whether these mutations have resulted in subtle effects on the stability of DNA binding, we deployed molecular dynamics simulations with atomistic representation of Pdx1 homeodomains, DNA target, and solvent. From the post-processing of the molecular dynamics simulations, we estimated the enthalpy of binding between sand rat and mouse (or other mammal) Pdx1 and monomer DNA binding sites by using the Molecular Mechanics Poisson Boltzmann Surface Area (MM-PBSA) method (*SI Appendix, Materials and Methods*, section 10). Target DNA sequences used were core Pdx1-binding sites of the mouse *insulin* A1 promoter and its sand rat ortholog. From 200-ns molecular dynamics simulations, the enthalpy of binding for protein–DNA interaction was calculated to be lower for sand rat than for mouse Pdx1 (mean –140 kcal/mol vs. mean –122 kcal/mol), indicative of sand rat Pdx1 binding DNA more “tightly” than is normal for the mammalian Pdx1 protein (Fig. 3*B*). One amino acid change was responsible for much of the difference: a Leu-to-Arg substitution in alpha helix 1 (homeodomain position 13), leading to the positive side chain of Arg making a new indirect contact with the phosphate backbone of DNA. A second substitution, Val to Arg in alpha helix 2 (homeodomain position 36), makes a smaller contribution (Fig. 3*C*). We also detect modifications to specific base interactions, with sand rat residues Met54 and Arg58 making new contacts to A and T bases within the TAAT core. Hence, stronger DNA binding is most likely driven by increased contacts with the backbone of DNA, coupled with decreased sequence specificity of DNA interaction. These results suggest that sand rat Pdx1 is divergent in DNA-binding affinity and specificity. Conserved Pdx1-binding sites in well-characterized promoters of three downstream target genes encoding pancreatic hormones (*insulin*, *somatostatin*, and *glucokinase*) show negligible divergence in sand rat compared with mouse, rat, and human (*SI Appendix, Materials and Methods*, section 11), indicating that Pdx1 divergence alone is

