



RESEARCH ARTICLE SUMMARY

ZOOMOMIA

Evolutionary constraint and innovation across hundreds of placental mammals

Matthew J. Christmas† and Irene M. Kaplow† et al.

INTRODUCTION: A major challenge in genomics is discerning which bases among billions alter organismal phenotypes and affect health and disease risk. Evidence of past selective pressure on a base, whether highly conserved or fast evolving, is a marker of functional importance. Bases that are unchanged in all mammals may shape phenotypes that are essential for organismal health. Bases that are evolving quickly in some species, or changed only in species that share an adaptive trait, may shape phenotypes that support survival in specific niches. Identifying bases associated with exceptional capacity for cellular recovery, such as in species that hibernate, could inform therapeutic discovery.

RATIONALE: The power and resolution of evolutionary analyses scale with the number and diversity of species compared. By analyzing genomes for hundreds of placental mammals, we can detect which individual bases in the genome are exceptionally conserved (constrained) and likely to be functionally important in both coding and noncoding regions. By including species that represent all orders of placental mammals and aligning genomes using a method that does not require designating humans as the reference species, we explore unusual traits in other species.

RESULTS: Zoonomia's mammalian comparative genomics resources are the most comprehensive

and statistically well-powered produced to date, with a protein-coding alignment of 427 million bases and a whole-genome alignment of 240 placental mammals representing all orders. We estimate that at least 10.7% of the human genome is evolutionarily conserved relative to neutrally evolving repeats and identify about 101 million significantly constrained single bases (false discovery rate < 0.05). We cataloged 4552 ultraconserved elements at least 20 bases long that are identical in more than 98% of the 240 placental mammals.

Many constrained bases have no known function, illustrating the potential for discovery using evolutionary measures. Eighty percent are outside protein-coding exons, and half have no functional annotations in the Encyclopedia of DNA Elements (ENCODE) resource. Constrained bases tend to vary less within human populations, which is consistent with purifying selection. Species threatened with extinction have few substitutions at constrained sites, possibly because severely deleterious alleles have been purged from their small populations.

By pairing Zoonomia's genomic resources with phenotype annotations, we find genomic elements associated with phenotypes that differ between species, including olfaction, hibernation, brain size, and vocal learning. We associate genomic traits, such as the number of olfactory receptor genes, with physical phenotypes, such as the number of olfactory turbinates. By comparing hibernators and nonhibernators, we implicate genes involved in mitochondrial disorders, protection against heat stress, and longevity in this physiologically intriguing phenotype. Using a machine learning-based approach that predicts tissue-specific cis-regulatory activity in hundreds of species using data from just a few, we associate changes in noncoding sequence with traits for which humans are exceptional: brain size and vocal learning.

CONCLUSION: Large-scale comparative genomics opens new opportunities to explore how genomes evolved as mammals adapted to a wide range of ecological niches and to discover what is shared across species and what is distinctively human. High-quality data for consistently defined phenotypes are necessary to realize this potential. Through partnerships with researchers in other fields, comparative genomics can address questions in human health and basic biology while guiding efforts to protect the biodiversity that is essential to these discoveries. ■

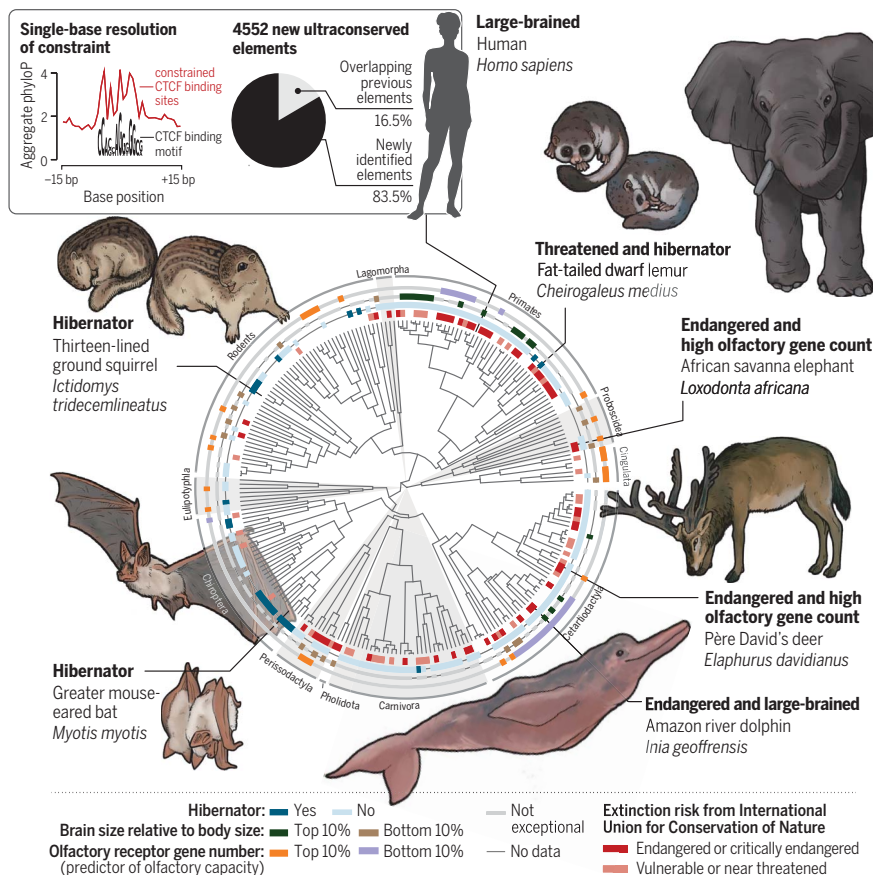
The list of author affiliations is available in the full article online.

*Corresponding author. Email: Kerstin Lindblad-Toh (kersli@broadinstitute.org); Elinor K. Karlsson (elinor.karlsson@umassmed.edu)

†These authors contributed equally to this work.

Cite this article as M. J. Christmas et al., *Science* 380, eabn3943 (2023). DOI: 10.1126/science.abn3943

READ THE FULL ARTICLE AT
<https://doi.org/10.1126/science.abn3943>



Comparing genomes from 240 species to explore the evolution of placental mammals. Our new phylogeny (black lines) has alternating gray and white shading, which distinguishes mammalian orders (labeled around the perimeter). Rings around the phylogeny annotate species phenotypes. Seven species with diverse traits are illustrated, with black lines marking their branch in the phylogeny. Sequence conservation across species is described at the top left.

RESEARCH ARTICLE

ZOOMOMIA

Evolutionary constraint and innovation across hundreds of placental mammals

Matthew J. Christmas^{1†}, Irene M. Kaplow^{2,3†}, Diane P. Genereux⁴, Michael X. Dong¹, Graham M. Hughes⁵, Xue Li^{4,6,7}, Patrick F. Sullivan^{8,9}, Allyson G. Hindle¹⁰, Gregory Andrews⁷, Joel C. Armstrong¹¹, Matteo Bianchi¹, Ana M. Breit¹², Mark Diekhans¹¹, Cornelia Fanter¹⁰, Nicole M. Foley¹³, Daniel B. Goodman¹⁴, Linda Goodman¹⁵, Kathleen C. Keough^{15,16,17}, Bogdan Kirilenko^{18,19,20}, Amanda Kowalczyk^{2,3}, Colleen Lawless⁵, Abigail L. Lind^{16,17}, Jennifer R. S. Meadows¹, Lucas R. Moreira^{4,7}, Ruby W. Redlich²¹, Louise Ryan⁵, Ross Swofford⁴, Alejandro Valenzuela²², Franziska Wagner²³, Ola Wallerman¹, Ashley R. Brown^{2,3}, Joana Damas²⁴, Kaili Fan⁷, John Gatesy²⁵, Jenna Grimshaw²⁶, Jeremy Johnson⁴, Sergey V. Kozyrev¹, Alyssa J. Lawler^{3,4,21}, Voichita D. Marinescu¹, Kathleen M. Morrill^{4,6,7}, Austin Osmanski²⁷, Nicole S. Paulat²⁶, BaDoi N. Phan^{2,3,27}, Steven K. Reilly²⁸, Daniel E. Schäffer², Cynthia Steiner²⁹, Megan A. Supple³⁰, Aryn P. Wilder²⁹, Morgan E. Wirthlin^{2,3,31}, James R. Xue^{4,32}, Zoonomia Consortium[§], Bruce W. Birren⁴, Steven Gazal³³, Robert M. Hubley³⁴, Klaus-Peter Koepfli^{35,36,37}, Tomas Marques-Bonet^{38,39,40,41}, Wynn K. Meyer⁴², Martin Nweeia^{43,44,45,46}, Pardis C. Sabeti^{4,32,47}, Beth Shapiro^{30,48}, Arian F. A. Smit³⁴, Mark S. Springer⁴⁹, Emma C. Teeling⁵, Zhiping Weng⁷, Michael Hiller^{18,19,20}, Danielle L. Levesque¹², Harris A. Lewin^{24,50,51}, William J. Murphy¹³, Arcadi Navarro^{38,40,52,53}, Benedict Paten¹¹, Katherine S. Pollard^{16,17,54}, David A. Ray²⁶, Irina Ruf⁵⁵, Oliver A. Ryder^{29,56}, Andreas R. Pfenning^{2,3}, Kerstin Lindblad-Toh^{1,4*†}, Elinor K. Karlsson^{4,7,57*†}

Zoonomia is the largest comparative genomics resource for mammals produced to date. By aligning genomes for 240 species, we identify bases that, when mutated, are likely to affect fitness and alter disease risk. At least 332 million bases (~10.7%) in the human genome are unusually conserved across species (evolutionarily constrained) relative to neutrally evolving repeats, and 4552 ultraconserved elements are nearly perfectly conserved. Of 101 million significantly constrained single bases, 80% are outside protein-coding exons and half have no functional annotations in the Encyclopedia of DNA Elements (ENCODE) resource. Changes in genes and regulatory elements are associated with exceptional mammalian traits, such as hibernation, that could inform therapeutic development. Earth's vast and imperiled biodiversity offers distinctive power for identifying genetic variants that affect genome function and organismal phenotypes.

Placental mammals, the evolutionary lineage that includes humans, are exceptionally diverse, with more than 6100 extant species (1), from the 2-g bumblebee bat to the 150,000-kg blue whale (2, 3). Over the past 100 million years, mammals have adapted to almost every habitat on Earth (Fig. 1A) (4). Zoonomia is the largest comparative genomics resource for mammals produced to date, with whole genomes aligned for 240 diverse species [2.3-fold more families and 3.9-fold more species than the mammals included in the earlier 100 Vertebrates alignment (5)] and protein-coding sequences aligned for 427 species (6). Using this resource, we can find elements that are conserved in the genomes of all placental mammals, elements that are changing unusually quickly in particular lineages, and elements that are associated with particular traits. All three approaches address a primary challenge in genomics: identifying genomic elements that affect genome function and organismal phenotypes (7).

Species evolve through selection on both small, sequence-level mutations and larger structural changes to the genome (e.g., translocation of transposable elements, inversions,

deletions, and duplications), as well as through hybridization with other species (8–10). Mutations are assumed to arise by random chance and then rise and fall in frequency within a population as a consequence of both neutral drift and selection. Mutations that disrupt characteristics that are essential for survival tend to be lost, whereas those conferring an advantage are more likely to be retained, eventually resulting in genetic differences that differentiate species.

By aligning the genomes of many different species, we can measure whether mutations at a given position in the genome are retained more or less often than expected under neutral drift (11–13). Fewer differences between species than expected suggests evolutionary constraint (dearth of variation due to purifying selection; also referred to as conservation), whereas more differences than expected in some lineages suggests acceleration (rapid evolution that may be clade-specific) (12, 13). Both metrics indicate that the given position has a role in molecular function. Measures of constraint and acceleration do not vary with cell type or developmental time point sampled, which simplifies sample collection and data generation. They are complementary to

methods for annotating the functional genome (14, 15).

Previous studies have used comparative genomics analyses to associate protein-coding changes with specific adaptations (16), such as diet type (17), echolocation (18), and subterranean habitation (19). However, these studies included few species relative to Zoonomia. As a result, they lacked the power and resolution required to investigate changes in genes and noncoding regulatory elements on a genome-wide level. Studying the evolution of regulatory elements, which make up much of the functional sequence in the genome, is particularly challenging because they tend to evolve more quickly and be less strongly conserved than coding elements (15, 20, 21). By substantially increasing the number and diversity of species in our comparative genomic analyses, we increase the sensitivity and specificity of methods used for detecting evolutionary signals and associating these signals with species-level phenotypes (22, 23).

Evolutionary constraint is a powerful tool for determining which genomic variants are causally implicated in human diseases. We explore this in detail in our companion paper (24), where we show that constrained positions are enriched for variants that explain common disease heritability more than any other functional annotation and that using the Zoonomia constraint scores improves polygenic risk scoring and fine-mapping of candidate disease loci.

Here, we use the new comparative genomics resources produced by Zoonomia to explore placental mammal evolution, including the origins of exceptional traits. We also synthesize the discoveries described by the compendium of papers in the Zoonomia package.

Evolutionary constraint and acceleration in mammals

We selected species for inclusion in Zoonomia to maximize the evolutionary branch length represented and thereby increase the power to detect constraint (4). The updated 241-way reference-free Cactus alignment with 240 species (domestic dog has two representatives) overcomes limitations of reference-based alignments (table S1) (4, 11). It includes genomic elements lost in humans, allows detection of multiple-orthology relationships, and captures complex rearrangements and copy-number variation. We observed 3.6 million perfectly conserved sites, which is 19,000-fold more than expected by chance, assuming a uniform substitution rate (4), and is consistent with purifying selection on functional positions in the genome.

We measured constraint across the human, chimpanzee, mouse, dog, and little brown bat reference genomes by projecting the Cactus alignment onto each species and then measuring sequence constraint with phyloP (Fig. 2, A

and B, and table S2) (11, 12). The chimpanzee-referenced alignment supports the investigation of bases deleted in only humans. Mouse, dog, and little brown bat have well-annotated reference genomes and represent diverse branches of the mammalian lineage, supporting comparative research in a wide range of organisms. We measured sequence constraint in the primate subset of the Cactus alignment (43 species) using PhastCons, which offers more power with fewer species by scoring multibase elements rather than single bases (24, 25).

We inferred a new phylogeny of placental mammals that we used for subsequent analyses that require a tree (26) (Fig. 1B). This phylogeny used only bases from the alignment that scored as near-neutrally evolving with phyloP ($N = 466,232$). It places interordinal diversification before the major extinction event marking the end of the Cretaceous period, addressing a long-standing debate in the field (27–30). A divergence time analysis of the phylogeny supports the “long-fuse” model of mammalian diversification, with interordinal diversification in the Cretaceous and most intraordinal diversification after the Cretaceous–Paleogene mass extinction event (31–33), and not the fossil record-derived “explosive” model, which places all inter- and intraordinal diversification after the Cretaceous–Paleogene event, or other scenarios (34–36).

At any given site in the genome, the number of species aligned can vary from just one to all 240. The variation in alignment depth distinguishes regulatory regions with differing evolutionary histories (37). In the human-referenced

alignment, 91% of the human genome aligns to at least five species, but only 11% aligns to $\geq 95\%$ (≥ 228) of species (fig. S1). Candidate cis-regulatory elements are 926,535 putative regulatory elements in the human genome defined by the Encyclopedia of DNA Elements (ENCODE) resource (14) using DNA accessibility and chromatin modification data. In the alignment at candidate cis-regulatory elements, we discern three common patterns (Fig. 2C). In highly conserved elements, most bases align in most species, including distantly related species. In actively evolving elements, most species have a partial alignment to humans. Primate-specific elements align exceptionally well in only a small number of species. Promoter-like and enhancer-like elements tend to be highly conserved. Elements that specifically bind the transcription factor CTCF or are marked by H3K4me3 (trimethylated histone H3 lysine 4) are more likely to be evolving actively, and about 20% are primate-specific (Fig. 2D).

Estimate of genome-wide constraint

We estimate that a minimum of 332 Mb (10.7%) of the human genome is under constraint through purifying selection (Fig. 2A) (12). We computed this lower-bound of the percentage under constraint by comparing the observed genome-wide phyloP score distribution to that expected in the absence of selection (modeled using ancestral repeats) (fig. S2A). Using bootstrapping, we show that the sample of ancestral repeats used had little effect on the lower-bound constraint estimate that was achieved; a 95% confidence interval spans only

1.9 mega-base pairs (Mbp). Ancestral repeats are a reasonable proxy for neutrally evolving sequence and can help account for local factors such as GC-content and mutation rate variation that might affect the phyloP score distribution (12, 38, 39). Our estimate of 10.7% falls at the upper end of previous estimates, which ranged from 3 to 12% (40). It is substantially higher than estimates of at least 5% that were calculated using similar methods but much smaller mammalian datasets (12, 13). With more species, we have more power to detect both weaker constraint across mammals and lineage-specific constraint, although these scenarios are not readily distinguished by the phyloP scores (fig. S2, B and C).

The lower-bound estimates for constraint in chimp-, mouse-, dog-, and bat-referenced projections of the alignment range from 239 Mb in the mouse (9.0%) to 359 Mb in the chimp (11.8%) (Fig. 2A and table S2). We are unable to determine whether the total amount of constraint truly varies between species. Both the species composition of the dataset and technical confounders, including differences in assembly contiguity and quality, could explain the differences observed. The amount of sequence detected as significantly constrained [false discovery rate (FDR) < 0.05] correlates with the average branch length to the nine closest species [Spearman's correlation coefficient (ρ) = -0.975 ; $p = 0.0048$], with more constraint detected in species with more closely related species in the alignment (table S3). This suggests that the amount of the genome under detectable constraint in mouse, dog, and bat will

¹Department of Medical Biochemistry and Microbiology, Science for Life Laboratory, Uppsala University, 751 32 Uppsala, Sweden. ²Department of Computational Biology, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213, USA. ³Neuroscience Institute, Carnegie Mellon University, Pittsburgh, PA 15213, USA. ⁴Broad Institute of MIT and Harvard, Cambridge, MA 02139, USA. ⁵School of Biology and Environmental Science, University College Dublin, Belfield, Dublin 4, Ireland. ⁶Morningside Graduate School of Biomedical Sciences, UMass Chan Medical School, Worcester, MA 01605, USA. ⁷Program in Bioinformatics and Integrative Biology, UMass Chan Medical School, Worcester, MA 01605, USA. ⁸Department of Genetics, University of North Carolina Medical School, Chapel Hill, NC 27599, USA. ⁹Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden. ¹⁰School of Life Sciences, University of Nevada Las Vegas, Las Vegas, NV 89154, USA. ¹¹Genomics Institute, University of California Santa Cruz, Santa Cruz, CA 95064, USA. ¹²School of Biology and Ecology, University of Maine, Orono, ME 04469, USA. ¹³Veterinary Integrative Biosciences, Texas A&M University, College Station, TX 77843, USA. ¹⁴Department of Microbiology and Immunology, University of California San Francisco, San Francisco, CA 94143, USA. ¹⁵Fauna Bio, Inc., Emeryville, CA 94608, USA. ¹⁶Department of Epidemiology and Biostatistics, University of California San Francisco, San Francisco, CA 94158, USA. ¹⁷Gladstone Institutes, San Francisco, CA 94158, USA. ¹⁸Faculty of Biosciences, Goethe-University, 60438 Frankfurt, Germany. ¹⁹LOEWE Centre for Translational Biodiversity Genomics, 60325 Frankfurt, Germany. ²⁰Senckenberg Research Institute, 60325 Frankfurt, Germany. ²¹Department of Biological Sciences, Mellon College of Science, Carnegie Mellon University, Pittsburgh, PA 15213, USA. ²²Department of Experimental and Health Sciences, Institute of Evolutionary Biology (UPF-CSIC), Universitat Pompeu Fabra, 08003 Barcelona, Spain. ²³Museum of Zoology, Senckenberg Natural History Collections Dresden, 01109 Dresden, Germany. ²⁴The Genome Center, University of California Davis, Davis, CA 95616, USA. ²⁵Division of Vertebrate Zoology, American Museum of Natural History, New York, NY 10024, USA. ²⁶Department of Biological Sciences, Texas Tech University, Lubbock, TX 79409, USA. ²⁷Medical Scientist Training Program, University of Pittsburgh School of Medicine, Pittsburgh, PA 15261, USA. ²⁸Department of Genetics, Yale School of Medicine, New Haven, CT 06510, USA. ²⁹Conservation Genetics, San Diego Zoo Wildlife Alliance, Escondido, CA 92027, USA. ³⁰Department of Ecology and Evolutionary Biology, University of California Santa Cruz, Santa Cruz, CA 95064, USA. ³¹Allen Institute for Brain Science, Seattle, WA 98109, USA. ³²Department of Organismic and Evolutionary Biology, Harvard University, Cambridge, MA 02138, USA. ³³Keck School of Medicine, University of Southern California, Los Angeles, CA 90033, USA. ³⁴Institute for Systems Biology, Seattle, WA 98109, USA. ³⁵Center for Species Survival, Smithsonian's National Zoo and Conservation Biology Institute, Washington, DC 20008, USA. ³⁶Computer Technologies Laboratory, ITMO University, St. Petersburg 197101, Russia. ³⁷Smithsonian-Mason School of Conservation, George Mason University, Fort Royal, VA 22630, USA. ³⁸Catalan Institution of Research and Advanced Studies (ICREA), 08010 Barcelona, Spain. ³⁹CNAG-CRG, Centre for Genomic Regulation, Barcelona Institute of Science and Technology (BIST), 08036 Barcelona, Spain. ⁴⁰Department of Medicine and Life Sciences, Institute of Evolutionary Biology (UPF-CSIC), Universitat Pompeu Fabra, 08003 Barcelona, Spain. ⁴¹Institut Català de Paleontologia Miquel Crusafont, Departament Autònoma de Barcelona, 08193 Cerdanyola del Vallès, Barcelona, Spain. ⁴²Department of Biological Sciences, Lehigh University, Bethlehem, PA 18015, USA. ⁴³Department of Comprehensive Care, School of Dental Medicine, Case Western Reserve University, Cleveland, OH 44106, USA. ⁴⁴Department of Vertebrate Zoology, Canadian Museum of Nature, Ottawa, Ontario K2P 2R1, Canada. ⁴⁵Department of Vertebrate Zoology, Smithsonian Institution, Washington, DC 20002, USA. ⁴⁶Narwhal Genome Initiative, Department of Restorative Dentistry and Biomaterials Sciences, Harvard School of Dental Medicine, Boston, MA 02115, USA. ⁴⁷Howard Hughes Medical Institute, Harvard University, Cambridge, MA 02138, USA. ⁴⁸Howard Hughes Medical Institute, University of California Santa Cruz, Santa Cruz, CA 95064, USA. ⁴⁹Department of Evolution, Ecology and Organismal Biology, University of California Riverside, Riverside, CA 92521, USA. ⁵⁰Department of Evolution and Ecology, University of California Davis, Davis, CA 95616, USA. ⁵¹John Muir Institute for the Environment, University of California Davis, Davis, CA 95616, USA. ⁵²BarcelonaBeta Brain Research Center, Pasqual Maragall Foundation, 08005 Barcelona, Spain. ⁵³CRG, Centre for Genomic Regulation, Barcelona Institute of Science and Technology (BIST), 08003 Barcelona, Spain. ⁵⁴Chan Zuckerberg Biohub, San Francisco, CA 94158, USA. ⁵⁵Division of Messel Research and Mammalogy, Senckenberg Research Institute and Natural History Museum Frankfurt, 60325 Frankfurt am Main, Germany. ⁵⁶Department of Evolution, Behavior and Ecology, School of Biological Sciences, University of California San Diego, La Jolla, CA 92039, USA. ⁵⁷Program in Molecular Medicine, UMass Chan Medical School, Worcester, MA 01605, USA.

*Corresponding author. Email: kersli@broadinstitute.org (K.L.-T.); elinor.karlsson@umassmed.edu (E.K.K.)

†These authors contributed equally to this work.

‡These authors contributed equally to this work.

§Zoomomia Consortium collaborators and affiliations are listed at the end of this paper.

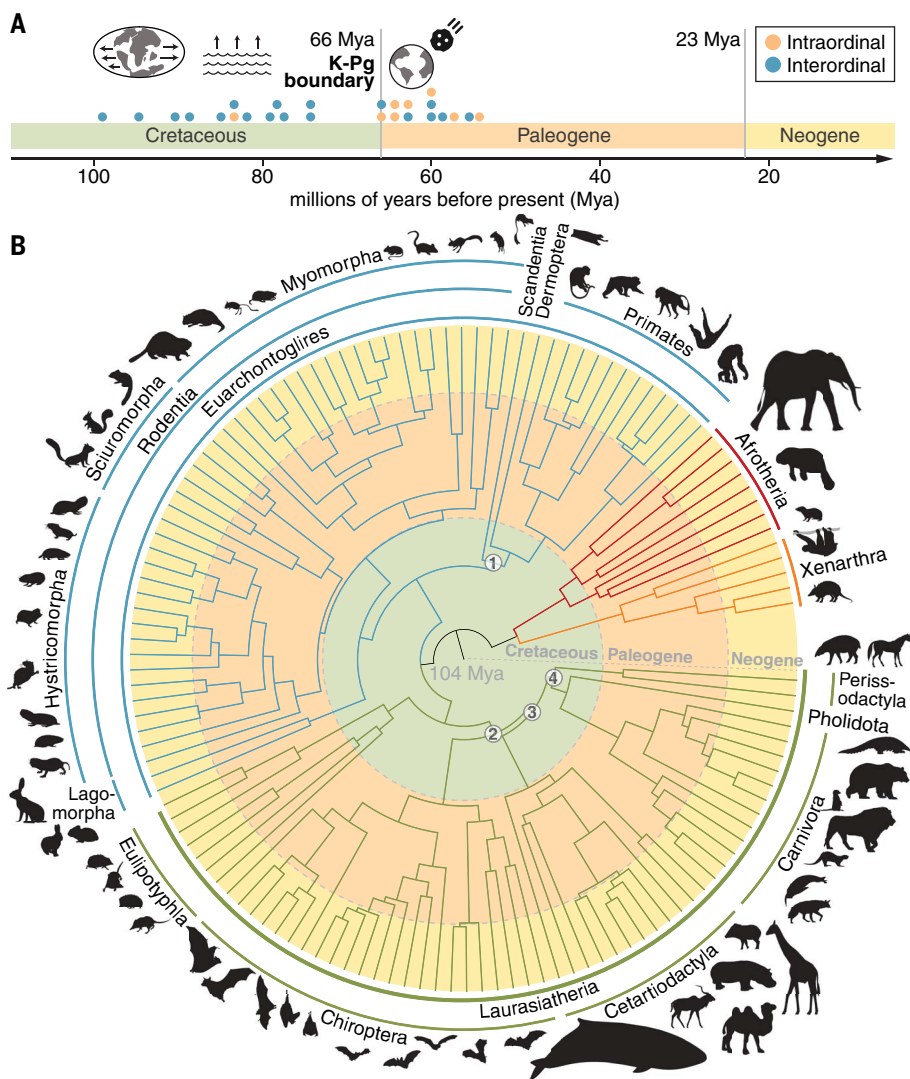


Fig. 1. New placental mammal phylogeny supports the long-fuse model of diversification. (A) Most interordinal diversification occurred in the Cretaceous, coincident with continental fragmentation and sea level changes. A pulse of intraordinal diversification occurred after the mass extinction event at the Cretaceous-Paleogene (K-Pg) boundary. Green, orange, and yellow shading bounded by gray lines demarcates different time periods. (B) A phylogeny based on divergence times estimated using ~470 kb of near-neutrally evolving sequence for 240 species resolves recalcitrant relationships in the placental mammal phylogeny (black numbers in white circles), including (1) Euarchonta (primates, colugos, and treeshrews), (2) Scrotifera [Perissodactyla (odd-toed ungulates), Cetartiodactyla (terrestrial even-toed ungulates and cetaceans), carnivorans, and bats], (3) Fereuungulata (perissodactyls, cetartiodactyls, carnivorans, pangolins), and (4) Zoomata [perissodactyls and Ferae (carnivorans and pangolins)]. [Species silhouettes are from PhyloPic]

increase as additional species are added to the alignment.

Genes enriched for constraint and acceleration

Genes with highly constrained protein-coding sequences are enriched in biological processes that function similarly across species, whereas those that are changing more quickly are enriched in processes that vary between species, consistent with previous studies (41–45). We tested the top 5% most accelerated and most conserved genes as measured by mean phyloP

score of coding sequence (data S1) against a nonredundant representative set of Gene Ontology (GO) biological processes using WebGestalt and identified overrepresented gene sets (46–48). The most constrained genes are involved in posttranscriptional regulation of gene expression (“mRNA processing”; GO:0006397; 81 of 487 genes; $p_{\text{FDR}} < 0.0002$) and embryonic development (“cell-cell signaling by wnt”; GO:0198738, 79 of 460 genes, $p_{\text{FDR}} < 0.0002$) (fig. S3A and table S4). RNA processing is essential for regulating cellular responses to environmental change (49), and

defects can cause debilitating diseases (50). “Pattern specification process” ranks third and includes all four HOX gene clusters (GO:0007389, 76 of 433; $p_{\text{FDR}} < 0.0002$). The most accelerated genes shape an animal’s interaction with its environment, including innate and adaptive immune responses, skin development, smell, and taste (fig. S3B).

We leveraged the large number of species in the Zoonomia alignments to show that a well-described gene inactivation, originally speculated to be human-specific (51), is found in 10 different lineages of mammals. The gene *CMAH* is inactivated in humans by a 92-bp frame-shifting exon deletion but is intact in other great apes (52). *CMAH* encodes an enzyme that converts the sialic acid Neu5Ac to Neu5Gc, and its loss restricts infection by pathogens dependent on Neu5Gc [e.g., malaria parasite *Plasmodium reichenowi* (53)] but increases susceptibility to viruses that bind Neu5Ac [e.g., severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) (54)]. When first observed, the loss of *CMAH* in humans was speculated to explain human-specific brain expansion (55, 56), but other mammals were subsequently shown to lack *CMAH* function (57–59). We combined the Cactus whole-genome alignment with analyses of read coverage and coding sequence alignment and found that *CMAH* has been inactivated in 40 of 239 species analyzed, representing 10 lineages (five newly discovered), including three rodent lineages and three bat lineages (fig. S4) (58). We confirm that *CMAH* loss occurred in the ancestor of all mustelids and pinnipeds using 11 species (compared with three originally) and that, among the primates, only humans and platyrrhine (New World) monkeys have lost *CMAH* (57). The role of *CMAH* in pathogen response suggests that its loss could shape the zoonotic potential of Neu5Gc-dependent pathogens, but further investigation is needed (60). Correlating *CMAH* inactivation with susceptibility to infection by SARS-CoV-2 or other viruses will require measuring infection susceptibility for a larger and more diverse set of mammals than has been studied to date.

Single-base resolution of constraint

Coding regions are the most strongly enriched for evolutionarily constrained positions, but most (80%) constrained positions are noncoding (Fig. 2E). We defined a “constrained base” as a position that has a positive phyloP score with $\text{FDR} < 5\%$. Constrained bases comprise 3.26% (101 Mb) of the human genome (Fig. 2B and table S2) and tend to cluster together, as previously described (13, 61). Most (80%) are within 5 bp of another constrained base, and 30% are in blocks ≥ 5 bp. The conservative $\text{FDR} < 5\%$ threshold limits the number of false positives but may miss weakly constrained bases or bases constrained in just a subset of

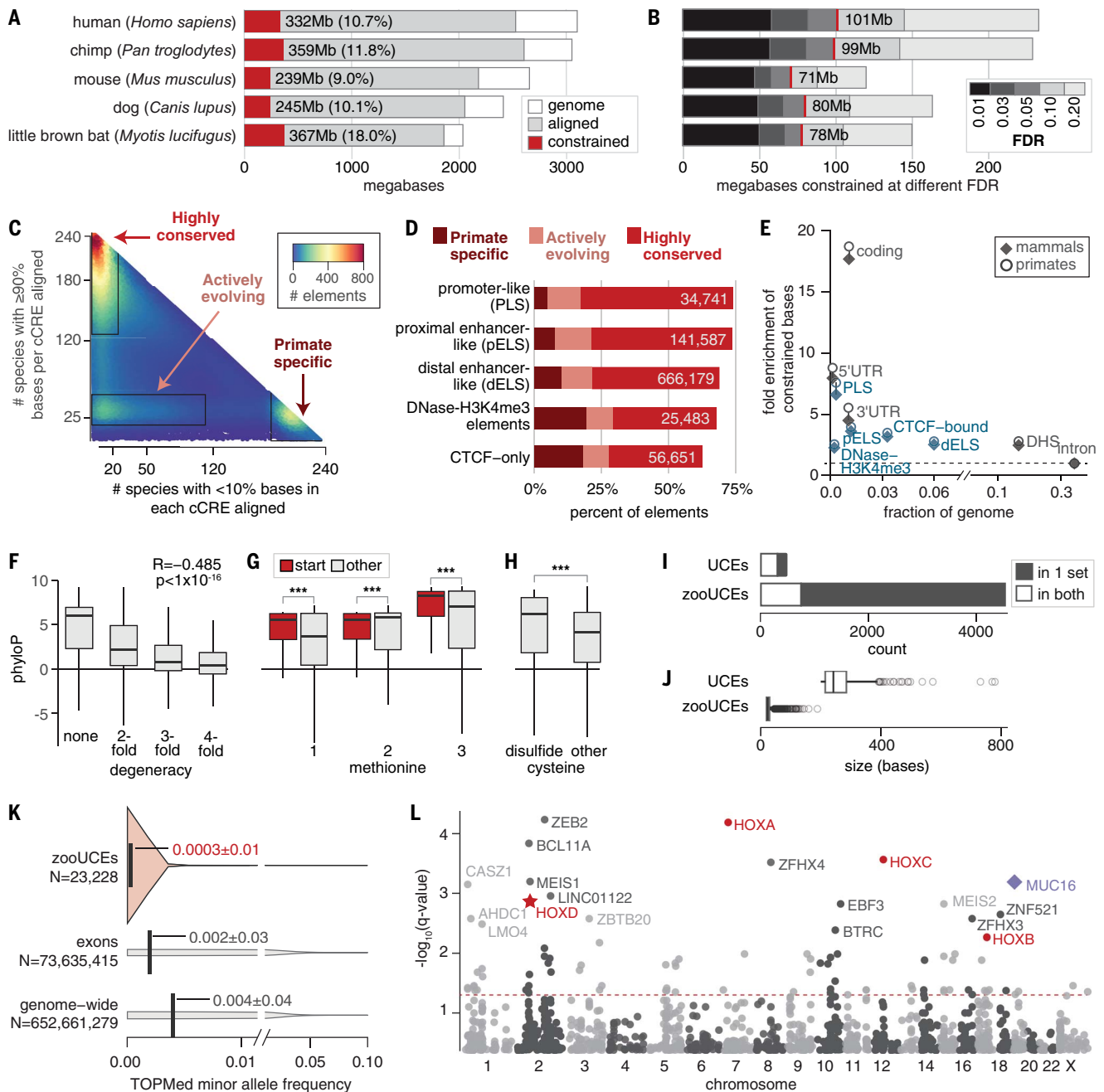


Fig. 2. Comparing 240 species resolves mammalian constraint to single bases and identifies elements under selection.

(A and B) We estimated a lower-bound on the total amount of the genome under constraint (A) and the number of single bases constrained at different FDR thresholds (B). The red lines in (B) indicate the 5% FDR threshold, with the amount of sequence below this threshold given. (C and D) Comparing the number of species with poor alignments (x axis) with those with good alignments (y axis) at 924,641 human candidate cis-regulatory elements (14) (C) reveals three clusters that are nonrandomly distributed across element types (all chi-square test $p < 2.2 \times 10^{-308}$) (D). (E) Functional elements are enriched for constraint, with candidate cis-regulatory elements in blue and other element types in black. The dashed line indicates no enrichment. DHS, DNase hypersensitivity site; 3'UTR, 3' untranslated region; 5'UTR, 5' untranslated region. (F) Constraint is negatively correlated with degeneracy across 59,504,353 protein-coding positions. (G) Methionine codons functioning as start sites in protein-coding sequence are more constrained at each of the three codon positions. (H) Cysteines in disulfide bridges are more constrained than other cysteines. In (F) to (H), the box boundaries

represent 25 and 75% quartiles, with a horizontal line at the median and the vertical line demarcating an additional 1.5 times interquartile range (IQR) above and below the box boundaries. $***p_{\text{Wilcoxon}} < 1 \times 10^{-16}$. (I) Most zooUCEs are new and do not overlap ultraconserved elements in the original set (73). (J) All zooUCEs are shorter than the original ultraconserved elements. Box and whisker parameters are the same as in (F), with outlier zooUCEs (>1.5 times IQR below or above the box boundaries) plotted as open circles. (K) Human variants in zooUCEs (light orange) have lower minor allele frequencies than they do in exons or genome-wide (gray). The vertical lines are at the means. The filled area is the distribution of allele frequencies. (L) Constraint measured in 100-kb bins genome-wide. The most constrained 100-kb bins include the HOX clusters (red). *HOXD* (red star) overlaps the longest syntenic block shared across mammals (174). Rearrangements in this locus can lead to limb malformations and other damaging outcomes. One bin containing *MUC16* (purple diamond) significantly lacks constraint. *MUC16* provides a mucosal barrier that protects epithelial cells from pathogens. The red dashed line indicates $q = 0.05$. Labeled bins have $q < 0.006$.

mammals. Using a threshold of FDR < 20% increases the estimated percentage of bases constrained from 3.26 to 7.56% (Fig. 2B and table S2).

The phyloP scores have three-base periodicity in coding sequence, consistent with the genetic code (62, 63). The Zoonomia phyloP scores are strongly correlated with the codon degeneracy at individual positions. Nondegenerate sites are far more likely to be constrained bases than fourfold degenerate sites (74.1 versus 18.5%). The median phyloP score exome-wide is 4.9 [interquartile range (IQR) = 5.8] in the first position (nondegenerate for 17 of 20 amino acids), 6.0 (IQR = 4.0) in the second (nondegenerate in 19 of 20), and 0.68 (IQR = 2.7) in the third (nondegenerate for 2 of 20) (fig. S5). The more functionally equivalent nucleotide options a coding base has in the genetic code, the weaker its phyloP score (Spearman's $\rho = -0.51$, $p < 2.2 \times 10^{-16}$) (Fig. 2F). Our ability to demonstrate expected patterns of constraint in coding sequence suggests that we have achieved sufficient power to resolve constraint to single bases in the human genome. This is unprecedented. The 29 Mammals project alignment resolved constraint to ~12 bases (13), and studies with more species examined only a subset of the genome (12). Comparing exomes for 141,456 humans achieved only gene- or exon-level resolution (64).

We discern stronger constraint at critical positions in peptides than at other protein-coding positions, supporting the utility of the Zoonomia phyloP scores for predicting functional importance. Whereas previous work had shown broadly that splice sites are often located in constrained regions (61), we discern enrichment of constraint at start codons, stop codons, and splice sites specifically (24 times, 19 times, and 25 times greater than genome-wide; chi-square test, $p < 2.2 \times 10^{-16}$). Methionine codons that function as start codons are more conserved than methionines elsewhere in the peptide (Fig. 2G). Cysteines in intrapeptide disulfide bridges, which can cause misfolding when mutated (65), are more conserved than other cysteines (Fig. 2H).

Bases constrained in mammals are less likely to be variable in humans, consistent with purifying selection (64, 66–68). Previous work showed that variants in functional positions have lower minor allele frequencies among humans in the Trans-Omics for Precision Medicine dataset (TOPMed) (69). Positions designated as evolutionarily constrained in Zoonomia similarly have lower minor allele frequencies in TOPMed, consistent with functional importance [constrained: frequency = 0.0026 ± 0.02 (\pm SD) and $N = 20,718,868$; unconstrained: 0.0040 ± 0.04 and $N = 601,458,551$; $p_{\text{Wilcoxon}} = 9.5 \times 10^{-13}$] (69). The less variable the position is in humans, the stronger its constraint across mammals (Spearman's $\rho = 0.78$, $p = 0.00014$; $N = 622,177,419$; fig. S6A).

Incorporating mammalian constraint into functional predictions will likely be particularly informative for poorly annotated positions. The correlation between the percentage of variants that are very rare in humans (minor allele frequency < 0.005 variants) and phyloP scores is strongest for positions that are scored as having unknown functional impact by SnpEff (70) (Spearman's $\rho = 0.98$, $p = 5.45 \times 10^{-7}$; $N = 608,227,093$; fig. S6B). SnpEff already considers 100-way vertebrate constraint scores in scoring variants, suggesting that constraint within mammals provides functional information that is not available through other sources.

Using versions of the reference-free Cactus alignment projected onto species other than human, we can assess constraint at positions that are deleted in the human genome and thus missing from previous resources (5, 13). We identified 10,032 human-specific deletions that overlap conserved elements and functionally assessed their regulatory effects using massively parallel reporter assays (71). Subsetting on just human-specific deletions constrained in chimp (phyloP score > 1) substantially increased concordance between measured regulatory change and predicted transcription factor binding differences [Pearson's correlation coefficient (r) increases from 0.25 ($p = 0.0037$) to 0.37 ($p = 0.00019$); Spearman's ρ increases from 0.24 ($p = 0.00614$) to 0.32 ($p = 0.00158$)].

New catalogs of conserved elements

We expanded and refined the catalog of ultraconserved elements in the human genome by 13-fold using the Cactus alignment, providing a rich new resource for exploring essential mammalian traits (72). The original set of 481 mammal ultraconserved elements consists of elements >200 bp long with identical sequence between human, mouse, and rat (73). Most are noncoding, and many function as enhancers during embryonic development (74–76). We defined Zoonomia ultraconserved elements (zooUCEs) as regions 20 bp or longer where every position is identical in at least 235 of 240 (98%) species in the alignment. Of the 4552 zooUCEs [average size 28.9 ± 13.0 bp (\pm SD)], 753 overlap 318 of the original ultraconserved elements, whereas 3799 are new (Fig. 2, I and J). Twenty-seven zooUCEs are longer than 100 bp (fig. S7A). Most of the zooUCEs are noncoding (69% are outside of protein-coding exons). Like the original ultraconserved elements, they are enriched near genes whose products are involved in transcription-related and developmental biological processes (table S5 and data S1) (73). The longest two zooUCEs (190 and 161 bp) are separated by a single base and are in an intron of *POLAI*, which encodes the catalytic subunit of DNA polymerase α .

Human TOPMed variants are rare in zooUCEs compared with the rest of the genome, suggesting purifying selection within humans

similar to the original UCEs (25, 72, 77, 78). ZooUCEs have fewer positions that are variable in humans (17.6%) than the coding sequences of genes (22.7%), which are known to be exceptionally constrained (69). When variants do occur in zooUCEs, their allele frequencies tend to be extremely low compared with those of variants that occur elsewhere in the genome. Average minor allele frequencies were 12.97 and 7.72 times lower in zooUCEs [$N = 23,228$; mean = 0.0003 ± 0.01 (\pm SD)] compared with genome-wide ($N = 652,661,279$; mean = 0.004 ± 0.04) and within exons ($N = 73,635,415$; mean = 0.002 ± 0.03), respectively (Fig. 2K).

We also cataloged constrained regions in the human genome using a phyloP score-based metric that allowed for more variability in constraint across mammals than the zooUCE criteria. Regions of contiguous constraint are regions of at least 20 bases where every individual base has a phyloP score above the FDR < 5% threshold (fig. S7B). Of the 595,536 such regions that we identified, most are short (median size = 32, IQR = 27), but 273 are longer than 500 bp and six are longer than 1 kb. The longest (1.36 kb) is in an intron of the gene *METAP1D* (chr2:172071926-172073285) and encompasses four distal enhancer-like candidate cis-regulatory elements. *METAP1D* encodes an essential mitochondrial protein that is conserved at least back to the common ancestor of human and zebrafish (79). This locus physically interacts with at least one transcription start site for each of *METAP1D* (FastHiC $q = 2.23 \times 10^{-2}$), *TLK1* (FastHiC $q = 7.62 \times 10^{-3}$), and *HAT1* (FastHiC $q = 3.92 \times 10^{-2}$) in human adult cortex Hi-C data (80–82). The synteny between these three genes is preserved in the *Xenopus* frog (83, 84). *TLK1* regulates chromatin structure (85), *HAT1* coordinates histone production and acetylation (86), and both are expressed in the cerebral cortex of 19 (*TLK1*) or 21 (*HAT1*) out of 19 or 21 mammals analyzed in a previous study, respectively (87).

We identified broad regions of unusually high constraint by scoring 100-kb nonoverlapping bins ($N = 28,218$) across the genome based on the fraction of bases that were constrained (data S2). We identified 53 bins with significantly elevated constraint ($q < 0.05$; average 17.8% constrained bases versus 3.5% for the genome; table S6). These bins are enriched for transcription-related biological processes and overlap the four *HOX* gene clusters (Fig. 2L). Five are in gene deserts, and two neighbor highly constrained developmental transcription factors (*LMO4* and *BCL11A*) (88, 89).

Constraint suggests regulatory function

Zoonomia's metrics of constraint can help detect positions likely to have regulatory function both within and outside of coding regions. In coding sequence, fourfold degenerate sites that overlap ENCODE3 transcription factor

binding sites ($N = 2,647,541$) (90) show moderately higher constraint than other fourfold degenerate sites ($N = 2,420,610$; chi-square test, $p < 2.2 \times 10^{-16}$; fig. S8). Noncoding constrained bases are enriched in regulatory elements across mammals and within primates, including at promoter-like signatures, enhancer-like signatures, sites bound by CTCF, and sites marked by H3K4me3 (Fig. 2E) (20, 91). The proportion of bases under constraint is higher in the subset of gene deserts (the longest 5% of intergenic regions) that neighbor developmental transcription factors (224 of 873 regions; $p_{\text{Wilcoxon}} = 2.15 \times 10^{-15}$) (92, 93) than in other gene deserts and is particularly high in candidate cis-regulatory elements within such gene deserts ($N = 38,065$; $p_{\text{Wilcoxon}} = 6.95 \times 10^{-280}$ compared with elements in other gene deserts; table S7).

Zoonomia constraint scores can distinguish which regulatory elements are likely to be functionally conserved across species. We identified transcription factor binding sites using convolutional neural networks and publicly available data for more than 600 ENCODE3 (14) transcription factor binding experiments spanning hundreds of cell and tissue types (37). This is a more comprehensive assessment of the regulatory landscape in mammals than was performed in previous work, which focused on two or three different transcription factors in five or six species (94, 95). We used a two-component Gaussian mixture model to classify sites as constrained or unconstrained. Of 15.6 million unique binding sites, covering 5.7% of the human genome, 1.9 million (0.8% of the genome) are constrained (table S8). Minor allele frequencies at sites variable in humans are significantly lower in constrained (mean = 0.0022, SD = 0.032) than in unconstrained (mean = 0.0036, SD = 0.041) binding sites (one-sided $p_{\text{Wilcoxon}} < 2.2 \times 10^{-16}$), consistent with strong purifying selection on these sites. The fraction of binding sites constrained varies by transcription factor and ranges from 1.5% (ZNF250) to 59.8% (YY2) (fig. S10A). The orthologs of the constrained binding sites are enriched for active histone marks [H3K4me3 and H3K27ac (acetylated histone H3 lysine 27)] in macaque, dog, mouse, and rat compared with unconstrained binding sites, suggesting that constrained sites are more likely to be functional in other species (fig. S9).

The correlation of constraint with both motif information content and functional state is evident in transcription factor binding sites for CTCF. CTCF is a highly conserved and ubiquitously expressed transcription factor that mediates genome three-dimensional (3D) structure (96–98). Overall, 14.8% of CTCF's binding sites are constrained (Fig. 3A). Motif information content for individual bases is significantly more correlated with base-level constraint in constrained sites than in uncon-

strained sites, showing that Zoonomia achieved single-base resolution constraint in noncoding regulatory elements that were missing from earlier analyses (95, 99) (Fig. 3B and fig. S10). This pattern persists across constrained binding sites for all evaluated transcription factors (Fig. 3C and fig. S10, B and C), advancing earlier work that lacked single base-level resolution (37, 95, 99). The motif logos calculated from constrained CTCF binding sites are nearly identical across species, unlike unconstrained sites (Fig. 3D), suggesting that constrained binding sites are more likely to be functional in other mammals (Fig. 3, E and F).

Unannotated constraint

Almost half of all constrained bases (48.5%) are in regions with no annotations in the

thousands of cell types, tissues, or conditions assayed by ENCODE3 (table S9) (14). We grouped constrained bases (phyloP FDR < 5%) fewer than 5 bp apart in unannotated intergenic regions (excluding repeats, centromeres, and telomeres) to define 423,586 elements, which we term unannotated intergenic constrained regions (UNICORNs) (median size = 20 bp; IQR = 23; 95th percentile = 131 bp; 0.5% of genome; Fig. 4A and fig. S7C). Most (77.0%) of these unannotated elements are within 500 kb of the transcription start site for a protein-coding gene. They tend to contain fewer variants ($p_{\text{Wilcoxon}} < 2.2 \times 10^{-16}$) with lower minor allele frequencies ($p_{\text{Wilcoxon}} < 2.2 \times 10^{-16}$) than other intergenic regions (Fig. 4B).

Many unannotated regions are likely to be functional under conditions that were not

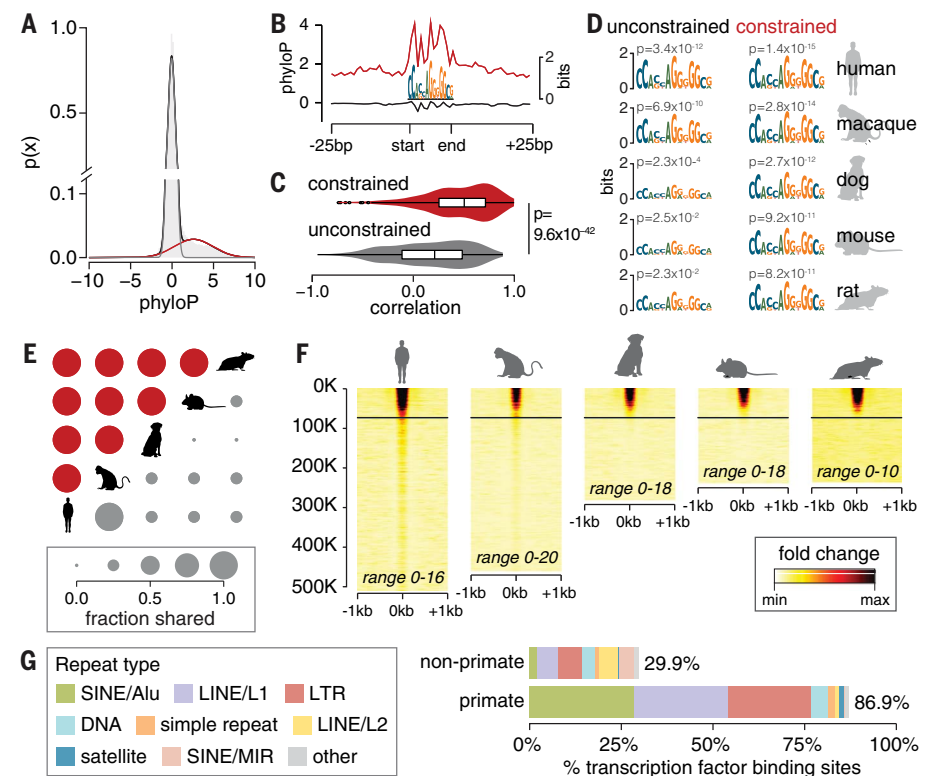


Fig. 3. Conserved function of constrained transcription factor binding sites. (A) A two-component Gaussian mixture model fit over average phyloP scores across binding sites for CTCF distinguishes the distribution for evolutionarily constrained sites (red) from others (gray). (B) At CTCF binding sites, aggregate phyloP scores are high for constrained binding sites (red, 61,832 sites) but not for unconstrained binding sites (gray, 424,177 sites). The same pattern is observed for other transcription factors (fig. S10). (C) Across all transcription factors, aggregate phyloP scores are more strongly correlated (Pearson's correlation) with binding site information content for constrained sites than for unconstrained sites. Boxes and whiskers represent 25% quartile, 75% quartile, minimum, and maximum, with a horizontal line at the median. The shading indicates the density of the data. (D) CTCF logos of constrained and unconstrained sets for four species made by lifting over human transcription factor binding sites. (E) Fraction of constrained (red) and unconstrained (gray) CTCF binding sites that are shared between pairs of species. (F) CTCF transcription factor chromatin immunoprecipitation sequencing (ChIP-seq) signal over binding sites in mammalian livers sorted by average phyloP scores. Each row is a binding site; in nonhuman species, only aligned sites are shown. The horizontal lines indicate significant constraint. Ranges give the minimum and maximum ChIP-seq fold change over input for each species. (G) Percentage of primate-specific and non-primate-specific transcription factor binding sites that are derived from individual transposable element classes. LINE, long interspersed nuclear element; LTR, long terminal repeat; MIR, mammalian-wide interspersed repeat; SINE, short interspersed nuclear element. [Species silhouettes are from PhyloPic]

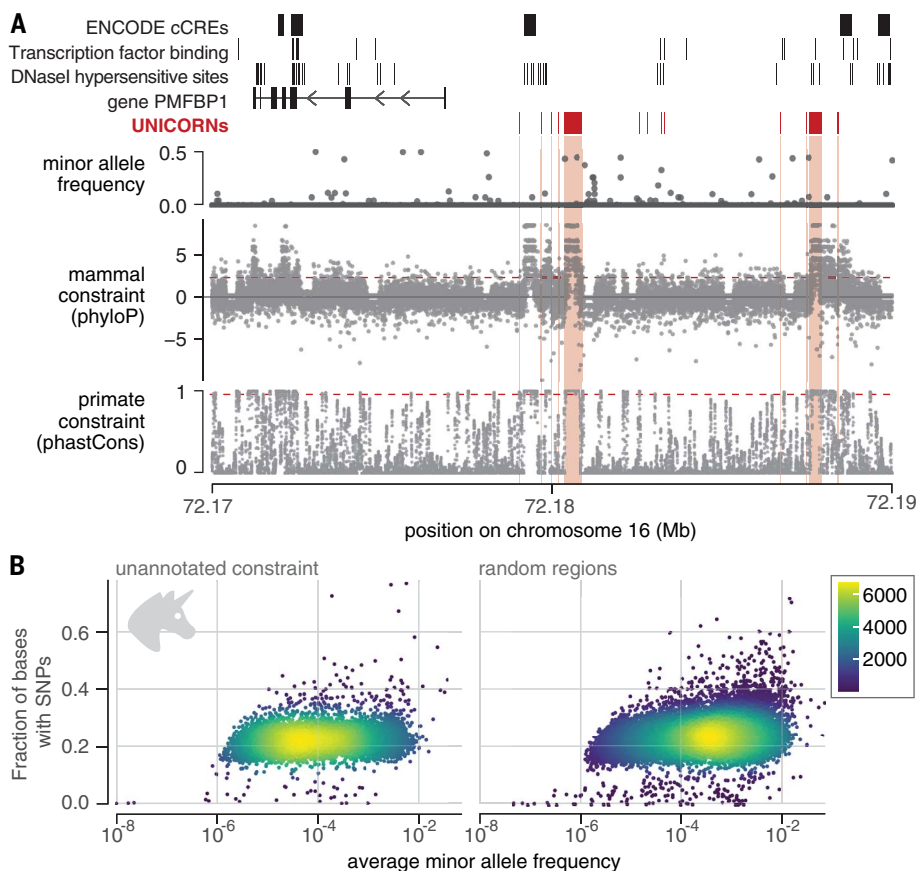


Fig. 4. Constraint highlights unannotated regions that are likely functional. (A) Example UNICORNs on human chromosome 16. The largest is 418 bp and located 3.5 kb upstream of the transcription start site of the gene *PMFBP1*; the second largest is 174 bp. Gray dots represent single bases. Red dashed lines represent the FDR < 5% threshold for phyloP and the threshold for phastCons that captures equivalent genome proportion (phastCons base score ≥ 0.961). UNICORNs lack coding or regulatory annotations in ENCODE (top track), and most have low diversity in human populations (second track). **(B)** UNICORNs contain fewer variants, and those present have lower allele frequencies than those in the random set (Wilcoxon rank sum test, $p < 2.2 \times 10^{-16}$). The fraction of bases with single-nucleotide polymorphisms (SNPs) versus mean minor allele frequency for human SNPs within UNICORNs (left) or within a random set of unannotated sequences (right) is shown. Allele frequencies were \log_{10} transformed. Human variants and allele frequencies were obtained from TOPMed data freeze 8 (69).

assayed in human ENCODE3 (table S9) (14). For example, open chromatin regions (a proxy for candidate enhancers) in developing brain tissues (100), adult motor cortical neuron cell types (101), and narrowly defined regions of young adult brain (102) overlap 8.8, 7.1, and 8.6% of UNICORNs respectively (17% collectively; 5.4, 2.7, and 4.2% are active in only developing brain, adult motor cortical neurons, and young adult brain regions, respectively). As resources like ENCODE expand to include more difficult-to-access time points, cell types, and tissues, we anticipate that the function of many UNICORNs will be elucidated.

Regions of accelerated evolution

Recent evolution in the human lineage may have occurred in part by modifying the 3D structure of the genome, which can alter gene

regulation (103). We developed an automated pipeline for identifying “accelerated” regions that are highly constrained across mammals but exceptionally variable in particular lineages (104). We found 312 regions accelerated in humans and 141 in chimpanzees, most of which are noncoding. Human (82%) and chimpanzee (86%) accelerated regions tend to have signatures of positive selection (after accounting for other factors such as GC-biased gene conversion); these accelerated regions also tend to reside near developmental and neurological genes, consistent with previous work (105–108). In domains that contain human accelerated regions, we show that the 3D genome structure is altered by human-specific structural variants, suggesting a role for enhancer hijacking in the species-specific evolution of these loci (109).

Evolution through transposable elements

We cataloged transposable elements in the genomes of 248 species (fig. S11) (110). Transposable elements are mobile DNA sequences 100 to 10,000 bp long that can accumulate to >1 million copies per genome. Despite their potential to influence genome structure and function (111, 112), they are difficult to analyze, and most studies have focused on human and mouse (113). We analyzed transposable element class, number, and distribution in 248 species (table S1). There is little variation between mammals in the fraction of the genome in transposable elements [$N = 248$; $49.0 \pm 7.5\%$ (\pm SD)], consistent with counterbalancing accumulation with DNA loss (114). Recent accumulation, especially retrotransposon accumulation, is positively correlated with genome size [hierarchical Bayesian model, coefficient of determination (R^2) = 0.54 (95% high probability density 0.42, 0.64)], suggesting insufficient time to purge insertions after a surge of activity, and negatively correlated with transposable element diversity, suggesting that genomic control mechanisms may limit the repertoire of active elements (110, 115). Younger transposable element families are more likely to include insertions that are polymorphic in the species and thus may be subsequently lost. However, any family with multiple members is likely a permanent feature of the species because there is no known mechanism to target an entire family for elimination. Bats are a hotspot for horizontal transfer of DNA transposons, with more than 200 such events, compared with just 11 transferred into other lineages (table S10) (116).

Overall, about 11% of constrained human bases are in transposable elements, with constraint enriched in simple repeats and DNA transposons and depleted in short interspersed nuclear elements, long terminal repeats, and satellite repeats (fig. S12A). This likely reflects the absence of function within more recently inserted transposable elements. DNA transposons are an ancient class of repeats known to acquire functional roles, such as the transcription factor *ZBED5* (70% constrained) (117). By contrast, the repeat classes depleted in constraint have been active more recently during primate evolution and are therefore less likely to be functional (118). In simple repeats, constraint is negatively correlated with distance to the nearest gene. Simple repeats near genes, where they are more likely to influence gene expression (119), are more constrained (Spearman’s $\rho = -0.13$, $p < 2.2 \times 10^{-16}$; fig. S12B).

Most (87%) primate-specific transcription factor binding sites overlap transposable elements, unlike most non-primate-specific sites (30%) (Fig. 3G). Sites in transposable elements, and especially those in younger elements, tend to be less conserved and change more quickly (fig. S13). Our results suggest that transposable elements may be a driver of recent regulatory

innovations in primates (120–122), with the caveat that the binding sites have not been confirmed to have regulatory function (123). Transposable element-derived CTCF binding sites found only in primates are enriched near genes involved in vision, reproduction, immunity, lower extremity development, and social behavior [enrichment analysis of cis-regulatory regions with Genomic Regions Enrichment of Annotations Tool (GREAT) (108); table S11].

Connecting genotype to phenotype

The Zoonomia resource offers an unprecedented opportunity to explore the evolution of exceptional mammalian traits by associating genomic variation with species-level phenotypes in hundreds of diverse species. For many traits, phenotype annotations are sparse, limiting the application of these methods. Here, we illustrate the potential of this approach using traits that vary within multiple clades of mammals and for which we have species-level phenotypes for a large number of Zoonomia species. We apply tests for different modes of evolution, including changes in gene number, gene sequence, and gene regulation.

Olfactory ability

Mammals have widely varying olfactory abilities, reflecting adaptation to different ecological niches (124–128). Olfactory receptor gene repertoire is a proxy for olfactory ability in mammals (128). We investigated olfactory evolution by first identifying olfactory receptor genes in genome assemblies of 249 mammalian species through genome annotation by means of a set of mammalian receptor profile hidden Markov models (table S12) (127). This increases by 10-fold the number of species with olfactory gene annotations. Our annotated gene counts do not vary with genome quality, as measured by contig N50 (Spearman's $\rho = 0.065$, $p = 0.31$, $N = 249$), scaffold N50 (Spearman's $\rho = 0.0091$, $p = 0.89$, $N = 249$), or genome completeness (129) (Spearman's $\rho = 0.10$, $p = 0.11$, $N = 249$), and capture the wide variation across species [mean count = 1218 \pm 683 (\pm SD), $N = 249$] (Fig. 5A and fig. S14).

By improving representation within lineages, most notably rodents ($N=55$), cetaceans ($N = 17$), and xenarthrans ($N = 8$), we discern variation in olfaction that was missed in earlier studies (fig. S15). Rodents have more olfactory receptor genes on average than other mammals [55 rodents versus 194 others, mean = 1434 \pm 466 (\pm SD) versus 1156 \pm 721, $t = 3.4$, $p_{t\text{-test}} = 0.0008$]. The top rodent is the Central American agouti (3233 genes), which has more genes than all but three other species (Hoffmann's two-toed sloth, the nine-banded armadillo, and the African savanna elephant). Cetaceans have the narrowest variation of any order. All cetaceans (17 species) have exceptionally small olfactory receptor gene repertoires relative to other mammals (225 \pm 75

genes compared with 1290 \pm 650 genes, $t = -22.9$, $p_{t\text{-test}} = 5.8 \times 10^{-60}$). Baleen whales retain olfactory structures that were lost in toothed whales (130, 131), and, consistent with this anatomic evidence for olfactory ability, the four baleen whale species in Zoonomia have more olfactory receptor genes than the 13 toothed whales (339 \pm 36 versus 190 \pm 40, $t = -6.96$, $p_{t\text{-test}} = 0.00064$) (fig. S14).

The association of olfactory turbinal number with olfactory receptor gene repertoire across placental mammals suggests that both evolve in response to selection on olfactory capacity. Olfactory turbinals are an anatomic feature of the nasal cavity that is known to affect olfactory capacity (132–134). In 64 species that were phenotyped for both traits, the number of olfactory turbinals correlates with the number of olfactory receptor genes (Spearman's $\rho = 0.71$, $p = 5.50 \times 10^{-11}$) (Fig. 5A). This relationship remains significant after accounting for species relationships by applying a phylogenetic generalized least squares method (phyloIm coefficient = 0.014, $p = 4.31 \times 10^{-10}$) and a permutation approach that preserves the tree topology (permutation $p = 0.0013$) (fig. S16) (135–137). We also confirm earlier observations that the number of genes is negatively associated with group living (phyloIm coefficient = -0.0013 , phylogeny-aware permutation $p = 0.022$) (127, 138), possibly because social animals are less dependent on smell. The association between the number of genes and solitary living fails to reach significance (phyloIm coefficient = 0.00086, phylogeny-aware permutation $p = 0.099$).

Hibernation

Zoonomia includes the largest mammal protein-coding alignment completed to date, with 17,795 human genes aligned in up to 488 assemblies of 427 distinct species (6). This alignment complements the Cactus whole-genome alignment (4, 11). It integrates gene annotation, ortholog detection, and classification of genes as intact or inactivated and can join orthologous fragments of genes split in fragmented assemblies.

Our protein-coding alignment includes 22 deep hibernators (species capable of core temperature depression below 18°C for >24 hours) and 154 strict homeotherms (species that maintain constant body temperature), offering an opportunity to explore the genomic origins of hibernation. Forms of torpor are found in every deep mammalian lineage, suggesting that metabolic depression through heterothermy existed in some form in the ancestor of all mammals (139, 140). Modifications, including the capacity for seasonal hibernation, may be derived. Understanding the genomics of hibernation, including cellular recovery from repeated cooling and re-warming without apparent long-term harm, could inform therapeutics, critical care, and long-distance spaceflight (141, 142).

Comparing hibernators and strict homeotherms to the reconstructed ancestral mammal protein-coding sequence using generalized least squares forward genomics (23) identified 28 100-bp regions ($p_{\text{FDR}} < 0.05$) in 20 genes where hibernators are less diverged from the placental mammalian ancestor (table S13). Two of these genes, *MFN2* and *PINK1*, overlap four GO Biological Process gene sets related to depolarization and degradation of damaged mitochondria, an organelle essential for metabolic depression (table S14) (143), although the process's enrichment is only nominally significant (top geneset $p = 7.5 \times 10^{-5}$; $p_{\text{FDR}} = 0.39$). A third, *TXNIP*, also regulates mitophagy (144) and shows torpor-responsive gene expression in rodents (145–147) and bats (148).

Testing with RERconverge identified an additional 22 genes as evolving unusually fast or slow in hibernators compared with homeotherms (Fig. 5B and data S3) (149–151). RERconverge tests for associations between relative evolutionary (substitution) rates of genes and the evolution of traits. We controlled for the high proportion of hibernators in the bat lineage, a potential confounder, through a Bayes factor analysis that quantified the amount of signal arising from hibernators and from bats and excluded genes with a hibernator signal less than fivefold larger than the bat signal (fig. S17). The top-scoring genes ($p_{\text{FDR}} < 0.05$ and phylogeny-aware permutation $p_{\text{FDR}} < 0.05$) included 11 that are evolving faster and 11 that are evolving slower in hibernating species (fig. S18). Faster-evolving genes are nominally enriched in gene sets related to temperature response and immunity (fig. S18A and table S15). Among the genes that are evolving faster in hibernators are *HSPD1* [involved in stress adaptation underlying mammalian torpor (152)], the mTor pathway inhibitor *ADAMST9* [also implicated in longevity based on sequence convergence in microbats and naked mole rats (153)], and two genes connected to neurodevelopmental disorders [the voltage-gated sodium channel gene *SCN2A* (154) and the membrane K-Cl cotransporter gene *SLC12A5* (155)].

There is no overlap between the two methods in the genes that score as significant (phylogeny-aware permutation $p_{\text{FDR}} \leq 0.05$), suggesting that their distinct methodologies are sensitive to different types of sequence change. One gene (the neurodevelopmental gene *NCDN*) is nominally significant in both sets ($p < 0.05$ and permutation $p < 0.05$ in both analyses).

Neurological traits

We developed a toolkit for associating differences in cis-regulatory elements, an important driver of phenotype divergence (156–158), with differences in phenotypes that include brain size and vocal learning (159, 160). This Tissue-Aware Conservation Inference Toolkit (TACIT) does not require tissue-specific cis-regulatory

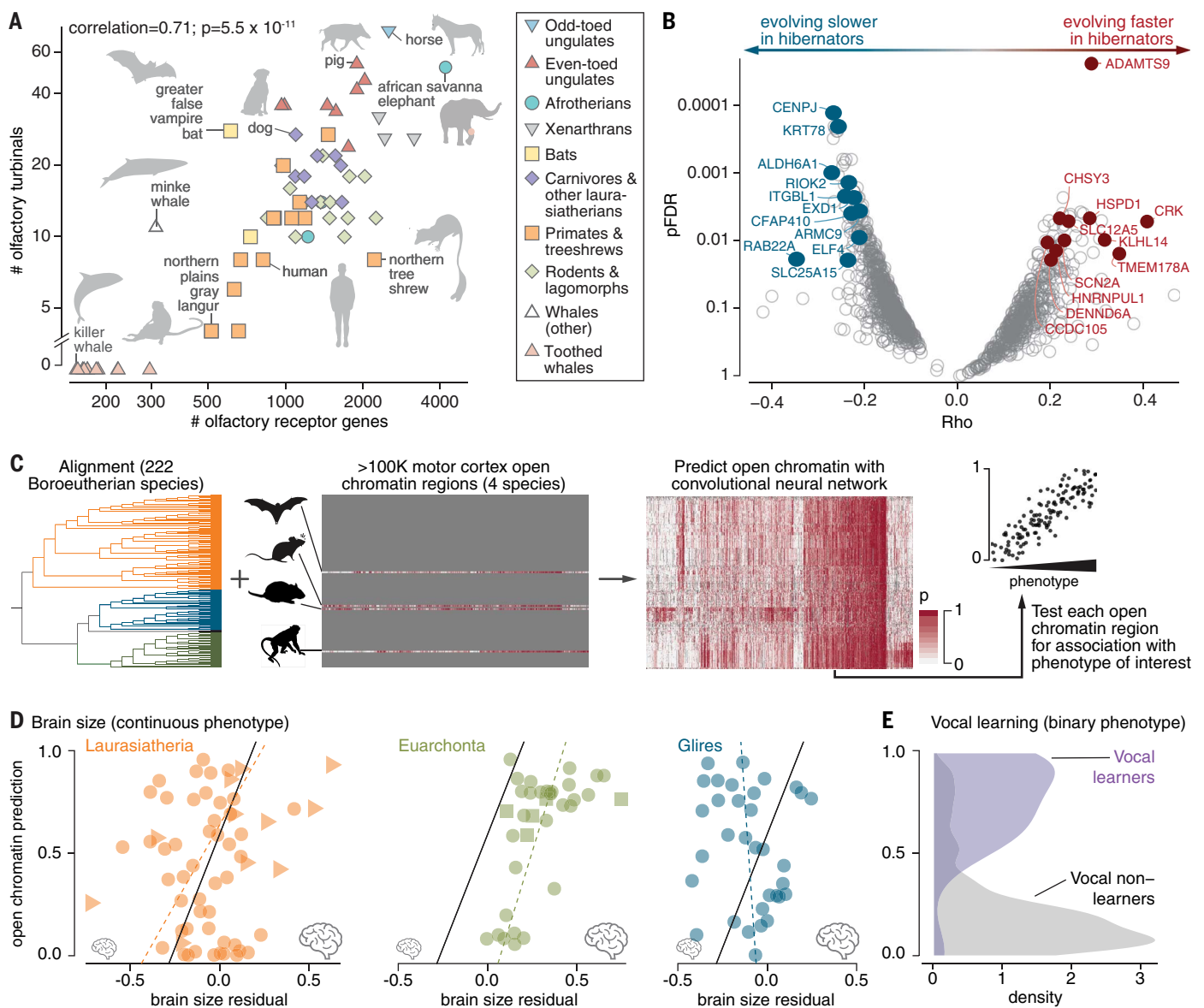


Fig. 5. Associating coding and regulatory change with species phenotypes.

(A) Olfactory receptor gene count (x axis) is associated with the number of olfactory turbinals (y axis) in 64 species. Labels and silhouettes mark outliers and species of interest. (B) Testing the coding sequence of 16,209 genes identified 341 genes that are evolving faster or slower in hibernators ($p_{\text{FDR}} < 0.05$; gray open circles), and 22 are significant after phylogeny-aware permutation testing (permutation $p_{\text{FDR}} < 0.05$; labeled), including 11 evolving faster (red filled circles) and 11 evolving slower (blue filled circles). (C) TACIT first trains a predictive classifier on sequences that underlie open chromatin regions from tissues or cell types in a few species and then predicts open chromatin in many others and tests for phenotype associations. (D) TACIT associated a motor cortex open chromatin region with brain size (a continuous-

valued trait), driven by associations within Laurasiatheria (59 species) and Euarchonta (36 species) but not within Glires (33 species). Results are for a rhesus macaque open chromatin region (chr10:48660711–48661679) near *MACROD2*. The phylom line of best fit is shown for all species [solid line; phylom coefficient (slope) = 0.45, permutation $p_{\text{FDR}} = 0.11$] and, as a visual aid, for each clade (dashed line). Triangles represent cetaceans (highest variation in brain size residual), squares represent great apes (highest variation in brain size residual within Euarchonta), and circles represent other species. (E) TACIT associated a motor cortex open chromatin region with vocal learning (a binary trait) in the *GALC* locus (phylom coefficient = 6.51, permutation $p_{\text{FDR}} = 0.045$) (137). Results are for an Egyptian fruit bat open chromatin region (PVIL01002568.1:139004–139596). [Species silhouettes are from PhyloPic]

element data from every species, which is costly and logistically challenging to obtain. Instead, it uses cis-regulatory sequence features in a tissue or cell type of interest from a few species to train machine-learning models that can be used to predict activity in that tissue or cell type at cis-regulatory element orthologs in many spe-

cies (Fig. 5C) (15). Models trained in one species can identify species- and tissue-specific cis-regulatory element activity in others, including for elements not used in training, demonstrating the feasibility of this approach (15). We then associated the predictions with phenotypes. We ran TACIT on traits that are pheno-

typed in more than 80 Zoonomia species and are proposed to involve neural cell types for which we have cis-regulatory element data from multiple species (motor cortex and parvalbumin neurons) (101, 161–163).

Brain size, measured relative to body size, is associated with predicted activity at cis-regulatory

elements that are active in the motor cortex (49 out of 98,912 elements tested, four species with training data, 158 species tested) and parvalbumin neurons (15 out of 35,034 elements tested, two species with training data, 72 species tested) (phylogeny-aware permutation $p_{FDR} < 0.15$) (159, 164–166). This includes a region near the gene *MACROD2*, a nervous system development gene implicated in microcephaly and intellectual disability in humans (Fig. 5D) (167, 168). Motor cortex cis-regulatory elements near genes previously implicated in microcephaly or macrocephaly tend to have more significant associations with brain size across mammals (one-sided $p_{Wilcoxon} = 0.013$).

In an analysis of 175 phenotyped species, both protein-coding changes and cis-regulatory changes were associated with capacity for vocal learning (160). Vocal learning is the ability to mimic noninnate sounds and likely evolved convergently in humans, bats, cetaceans, and pinnipeds (169). Our analysis of candidate cis-regulatory elements active in motor cortex ($N = 94,444$) and parvalbumin neurons ($N = 35,557$) identified motor cortex elements near *GALC* (Fig. 5E) (170), *TSHZ3* (171), and other speech disorder-related genes.

Applying genomics to biodiversity conservation

In addition to illuminating mammalian evolutionary history, Zoonomia's alignment and measures of constraint can help efforts to protect biodiversity for the future. Evolutionary constraint scores enable empirical estimation of deleterious genetic load and its demographic drivers across diverse species. We find that Zoonomia species with smaller historical effective population sizes carry higher fixed genetic load, with proportionally more missense substitutions (phyloM $p = 7.76 \times 10^{-5}$) and substitutions at constrained sites (phyloM $p = 9.63 \times 10^{-3}$). Species with a smaller historical effective population size are also more likely to be classified as threatened by the International Union for Conservation of Nature (IUCN) (phyloM $p < 3.3 \times 10^{-5}$), suggesting that historical processes are predictive of species' contemporary extinction risk status. Our analysis showed that threatened species have fewer substitutions at extremely constrained sites (phyloM $p = 0.001$), particularly in primates, whereas the opposite is true of missense substitutions, possibly because severely deleterious alleles have been purged or lost to drift (172) (Fig. 6). As the number of species with reference genomes grows, so will the power to leverage genomic data for identifying those most susceptible to the impacts of rapid environmental changes that characterize the Anthropocene.

Discussion

By aligning hundreds of mammalian genomes, Zoonomia realizes the vision of the landmark

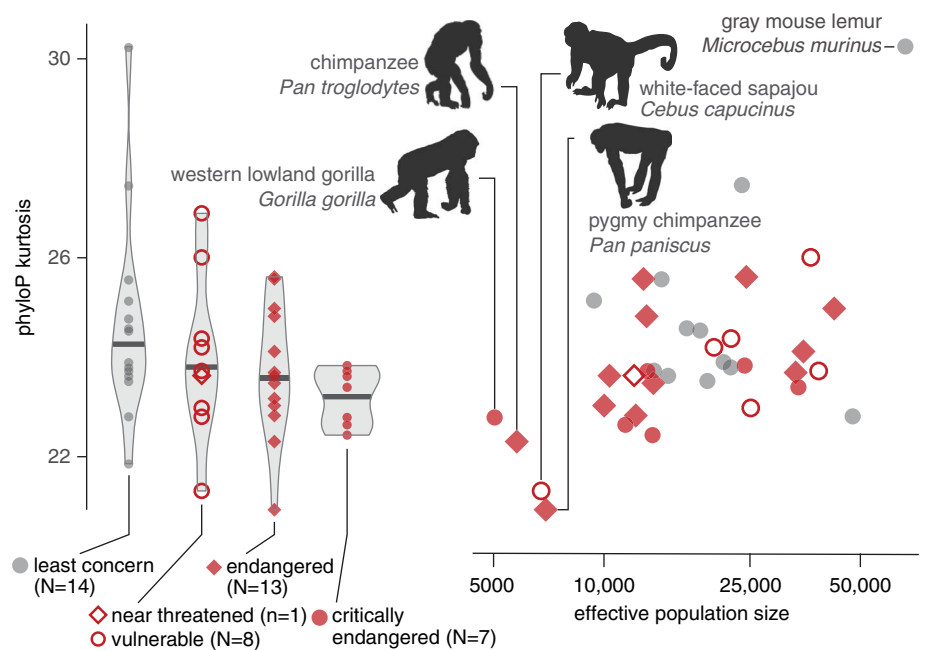


Fig. 6. Genomic metrics distinguish at-risk primate species. Primates that are categorized at increasing levels of extinction risk and with smaller effective population sizes have fewer substitutions at extremely constrained sites, measured as kurtosis (which describes the tail of the distribution) of phyloP scores (phyloM $p = 7.9 \times 10^{-4}$ and $p = 0.024$, respectively). Four at-risk species with the smallest effective population size (labeled with silhouettes) have low kurtosis (i.e., fewer phyloP outliers), and a species categorized as “least concern” with the largest effective population size has high kurtosis (gray mouse lemur; labeled). [Species silhouettes are from PhyloPic]

29 Mammals paper (13) to achieve single-base resolution of constraint across the human genome. This resource, which includes even deeper coverage of protein-coding regions (6), addresses a central goal of medical genomics: to identify genetic variants that influence disease risk and understand their biological mechanisms (7, 24, 37, 71, 173). It also opens new opportunities for exploring the evolution of mammalian genomes as species diverged and adapted to a wide range of ecological niches (15, 26, 110, 116, 160, 174) and for discovering what is distinctively human (104).

Zoonomia illustrates how new sequencing technology and analysis methods are transforming comparative genomics while underscoring the critical need for high-quality phenotype annotations. Studies into the genomic origins of exceptional mammalian traits have the potential to inform human therapeutic development (141) but are limited by sparse and inconsistent phenotype data. Here, we focus on a handful of traits for which we could define phenotypes consistently in large numbers of species, including hibernation (174 species), brain size (158 species), and vocal learning (175 species). Achieving the richer datasets that are needed to study other traits, evaluate pattern robustness, and address broader prospects requires collaborations between genomics researchers and scientists with expertise in morphology, physiology, and behavior to develop standardized phenotype definitions that apply

across species (175). It also requires proper collection, annotation, and data-handling practices that facilitate discovery, evaluation, and reuse of data (176).

Comparative genomics projects are classically motivated by the potential to advance human biomedicine, but they rely on biodiversity imperiled by human activity (177). Our analysis suggests that even a single reference genome per species may help conservation scientists identify potentially threatened populations earlier when management efforts can be more efficient and effective, but more work is needed to develop these methods (172). Through close and enduring partnerships with researchers working in biodiversity conservation, resources from Zoonomia and other comparative genomics projects can address questions in human health and basic biology while simultaneously guiding efforts to protect the biodiversity that is essential to these discoveries (178).

Methods summary

Alignment and annotation

We finalized the Zoonomia Cactus alignment by updating the initial Progressive Cactus alignment used in (11) to remove a mislabeled genome. We identified genes in Zoonomia genomes using *halLiftOver* in conjunction with the Zoonomia Cactus alignment, identifying sequences orthologous to the protein-coding sequence of human exons from ENSEMBL across each of the 241 assemblies. We also

developed an alternative reference-based approach described in our companion paper (6), which we applied to 427 species. We used a combination of two approaches using short sequencing reads and genome assemblies to determine whether the *CMAH* gene had been lost in mammalian genomes. We considered putative *CMAH* gene loss events to be cases where both these approaches indicated loss of the same part of the gene.

Constraint scoring

We used the Zoonomia alignment and a randomly selected set of ancestral repeat positions (100 kb total) to generate three different neutral models: one for autosomes and one each for the two sex chromosomes. We used PhyloFit from Phast v1.5 to estimate branch lengths. We used this same method to estimate primate-neutral models, but with the ancestral branch reconstruction based on the 43 primates from the alignment. We used phyloP (part of the PHAST v1.5 package) to calculate per-base constraint and acceleration *p* values. We calculated phyloP scores on the human-, chimpanzee-, mouse-, dog-, and bat-referenced 241-way alignments, as well as for a human-referenced, primates-only alignment (43-way). We computed a mammalian phyloP threshold by converting the *p* values corresponding to the phyloP scores into *q* values using a FDR correction. We considered any column with a resulting *q* ≤ 0.05 to be significantly evolutionarily constrained or accelerated, as determined by the sign of the score.

Analyzing constraint

Proportion of genome under constraint

We estimated lower bounds for the fraction of sites under purifying selection across the human, chimpanzee, dog, house mouse, and little brown bat genomes by comparing the empirical cumulative distribution functions of phyloP scores across each genome to the those of ancestral repeats, following the same method detailed in (12).

Constraint in functional elements

We extracted phyloP scores for all positions in protein-coding genes (GENCODE v.36) including 5' and 3' untranslated regions, and compared constraint between different positions within coding sequences. We summarized mean and standard deviation phyloP scores for positions within codons, degenerate and nondegenerate positions, methionines that act as and do not act as start codons, and cysteines that form and do not form intrapeptide disulfide bridges. We calculated constraint enrichment for several genome features (coding sequences, 5' untranslated regions, 3' untranslated regions, introns, DNase hypersensitivity sites, and the five types of cCREs [ENCODE candidate cis-regulatory regions (14)], where

we calculated constraint enrichment as the constrained fraction of the feature divided by the constrained fraction of the genome.

Highly constrained regions

We identified all positions where the number of species aligned was ≥235 and the base was the same among all species aligned at that position. We then merged neighboring positions, creating zooUCEs ranging in size from 20 to 190 bp. We assessed overlap between our zooUCEs and previously defined UCEs. We also defined regions of contiguous constraint as regions of at least 20 contiguous base pairs with phyloP scores above the FDR > 0.05 threshold and identified 100-kb bins with significantly high or low constraint.

Constraint in unannotated regions

We subsetted the human genome, removing all regions with the following annotations: GENCODE v37 exons (untranslated regions and exons for all protein-coding genes), promoters (transcription start site ±1 kb), introns, ENCODE3 cCREs, DNase hypersensitivity sites (including transcription factor binding sites), chromatin interaction analysis with paired-end tag sequencing (ChIA-PET) anchors, three promoter annotation sets, and six enhancer annotation sets (table S9). Within the remaining unannotated sequence, we identified closely located constraint positions to define a set of 423,586 UNICORNs.

Olfaction

We explored the olfactory receptor gene family across the Zoonomia species set, independently of alignment-based annotation. We mined all genomes for olfactory receptor gene sequences using the olfactory receptor assigner (179). We classified sequences as “pseudogenes” if they contained in-frame stop codons or were shorter than 650 bp and therefore not long enough to form the seven-transmembrane domain. We curated species-specific numbers of olfactory turbinals from both sides of the nasal cavity (table S12), obtaining turbinal numbers for 64 species in our sample. We tested for an association between the total number of olfactory receptor genes with the number of olfactory turbinals using phylolm (136), solitary living status, and group living status while accounting for the Zoonomia phylogenetic tree (26, 138).

Hibernation

We investigated genomic differences between mammals that we defined as hibernators and as strict homeotherms (table S1), with 22 species defined as deep hibernators and 154 species defined as strict homeotherms. We used generalized least squares forward genomics to identify genes that are more similar to the mammalian ancestor than they are to non-

hibernators as well as to identify regions conserved in hibernators relative to the placental ancestor. We also used RERconverge (149) to identify genes with significant evolutionary rate shifts in hibernating mammals versus nonhibernating mammals. Such genes are putative hibernation-related genes.

REFERENCES AND NOTES

- C. J. Burgin, J. P. Colella, P. L. Kahn, N. S. Upham, How many species of mammals are there? *J. Mammal.* **99**, 1–14 (2018). doi: [10.1093/jmammal/gyx147](https://doi.org/10.1093/jmammal/gyx147)
- K. E. Jones, K. Safi, Ecology and evolution of mammalian biodiversity. *Philos. Trans. R. Soc. London Ser. B* **366**, 2451–2461 (2011). doi: [10.1098/rstb.2011.0090](https://doi.org/10.1098/rstb.2011.0090); pmid: [21807728](https://pubmed.ncbi.nlm.nih.gov/21807728/)
- K. E. Jones et al., PanTHERIA: a species-level database of life history, ecology, and geography of extant and recently extinct mammals. *Ecological Archives* E090-184. *Ecology* **90**, 2648–2648 (2009). doi: [10.1890/08-1494.1](https://doi.org/10.1890/08-1494.1)
- Zoonomia Consortium, A comparative genomics multitool for scientific discovery and conservation. *Nature* **587**, 240–245 (2020). doi: [10.1038/s41586-020-2876-6](https://doi.org/10.1038/s41586-020-2876-6); pmid: [33177664](https://pubmed.ncbi.nlm.nih.gov/33177664/)
- University of California Santa Cruz Genomics Institute, Conservation track settings: <http://genome.ucsc.edu/cgi-bin/hgTrackUi?g=cons100way>.
- B. M. Kirilenko et al., Integrating gene annotation with orthology inference at scale. *Science* **380**, eabn3107 (2023). doi: [10.1126/science.abn3107](https://doi.org/10.1126/science.abn3107)
- T. Lappalainen, D. G. MacArthur, From variant to function in human disease genetics. *Science* **373**, 1464–1468 (2021). doi: [10.1126/science.abb8207](https://doi.org/10.1126/science.abb8207); pmid: [34554789](https://pubmed.ncbi.nlm.nih.gov/34554789/)
- S. A. Taylor, E. L. Larson, Insights from genomes into the evolutionary importance and prevalence of hybridization in nature. *Nat. Ecol. Evol.* **3**, 170–177 (2019). doi: [10.1038/s41559-018-0777-y](https://doi.org/10.1038/s41559-018-0777-y); pmid: [30697003](https://pubmed.ncbi.nlm.nih.gov/30697003/)
- H. H. Kazazian Jr., Mobile elements: Drivers of genome evolution. *Science* **303**, 1626–1632 (2004). doi: [10.1126/science.1089670](https://doi.org/10.1126/science.1089670); pmid: [15016989](https://pubmed.ncbi.nlm.nih.gov/15016989/)
- P. G. D. Feulner, R. De-Kayne, Genome evolution, structural rearrangements and speciation. *J. Evol. Biol.* **30**, 1488–1490 (2017). doi: [10.1111/jeb.13101](https://doi.org/10.1111/jeb.13101); pmid: [28786195](https://pubmed.ncbi.nlm.nih.gov/28786195/)
- J. Armstrong et al., Progressive Cactus is a multiple-genome aligner for the thousand-genome era. *Nature* **587**, 246–251 (2020). doi: [10.1038/s41586-020-2871-y](https://doi.org/10.1038/s41586-020-2871-y); pmid: [33177663](https://pubmed.ncbi.nlm.nih.gov/33177663/)
- K. S. Pollard, M. J. Hubisz, K. R. Rosenbloom, A. Siepel, Detection of nonneutral substitution rates on mammalian phylogenies. *Genome Res.* **20**, 110–121 (2010). doi: [10.1101/gr.097857.109](https://doi.org/10.1101/gr.097857.109); pmid: [19858363](https://pubmed.ncbi.nlm.nih.gov/19858363/)
- K. Lindblad-Toh et al., A high-resolution map of human evolutionary constraint using 29 mammals. *Nature* **478**, 476–482 (2011). doi: [10.1038/nature10530](https://doi.org/10.1038/nature10530); pmid: [21993624](https://pubmed.ncbi.nlm.nih.gov/21993624/)
- J. E. Moore et al., Expanded encyclopaedias of DNA elements in the human and mouse genomes. *Nature* **583**, 699–710 (2020). doi: [10.1038/s41586-020-2493-4](https://doi.org/10.1038/s41586-020-2493-4); pmid: [32728249](https://pubmed.ncbi.nlm.nih.gov/32728249/)
- I. M. Kaplow et al., Inferring mammalian tissue-specific regulatory conservation by predicting tissue-specific differences in open chromatin. *BMC Genomics* **23**, 291 (2022). doi: [10.1016/j.celrep.2012.08.032](https://doi.org/10.1016/j.celrep.2012.08.032); pmid: [23022484](https://pubmed.ncbi.nlm.nih.gov/23022484/)
- M. Hiller et al., A “forward genomics” approach links genotype to phenotype using independent phenotypic losses among related species. *Cell Rep.* **2**, 817–823 (2012). doi: [10.1016/j.celrep.2012.08.032](https://doi.org/10.1016/j.celrep.2012.08.032); pmid: [23022484](https://pubmed.ncbi.nlm.nih.gov/23022484/)
- F. Wagner et al., Reconstruction of evolutionary changes in fat and toxin consumption reveals associations with gene losses in mammals: A case study for the lipase inhibitor PNLI1P1 and the xenobiotic receptor NRII3. *J. Evol. Biol.* **35**, 225–239 (2022). doi: [10.1111/jeb.13970](https://doi.org/10.1111/jeb.13970); pmid: [34882899](https://pubmed.ncbi.nlm.nih.gov/34882899/)
- A. Marcovitz et al., A functional enrichment test for molecular convergent evolution finds a clear protein-coding signal in echolocating bats and whales. *Proc. Natl. Acad. Sci. U.S.A.* **116**, 21094–21103 (2019). doi: [10.1073/pnas.1818532116](https://doi.org/10.1073/pnas.1818532116); pmid: [31570615](https://pubmed.ncbi.nlm.nih.gov/31570615/)
- R. Partha et al., Subterranean mammals show convergent regression in ocular genes and enhancers, along with adaptation to tunneling. *eLife* **6**, e25884 (2017). doi: [10.7554/eLife.25884](https://doi.org/10.7554/eLife.25884); pmid: [29035697](https://pubmed.ncbi.nlm.nih.gov/29035697/)
- D. Villar et al., Enhancer evolution across 20 mammalian species. *Cell* **160**, 554–566 (2015). doi: [10.1016/j.cell.2015.01.006](https://doi.org/10.1016/j.cell.2015.01.006); pmid: [25635462](https://pubmed.ncbi.nlm.nih.gov/25635462/); doi: [10.1186/s12864-022-08450-7](https://doi.org/10.1186/s12864-022-08450-7); pmid: [35410163](https://pubmed.ncbi.nlm.nih.gov/35410163/)

Nevada Las Vegas, Las Vegas, NV 89154, USA. ¹²Biodiscovery Institute, University of Nottingham, Nottingham, UK. ¹³Department of Immunology, Genetics and Pathology, Science for Life Laboratory, Uppsala University, Uppsala 751 85, Sweden. ¹⁴Department of Biological Sciences, Texas Tech University, Lubbock, TX 79409, USA. ¹⁵Division of Vertebrate Zoology, American Museum of Natural History, New York, NY 10024, USA. ¹⁶Keck School of Medicine, University of Southern California, Los Angeles, CA 90033, USA. ¹⁷Fauna Bio Incorporated, Emeryville, CA 94608, USA. ¹⁸Baskin School of Engineering, University of California Santa Cruz, Santa Cruz, CA 95064, USA. ¹⁹Faculty of Biosciences, Goethe-University, 60438 Frankfurt, Germany. ²⁰LOEWE Centre for Translational Biodiversity Genomics, 60325 Frankfurt, Germany. ²¹Senckenberg Research Institute, 60325 Frankfurt, Germany. ²²Institute for Systems Biology, Seattle, WA 98109, USA. ²³School of Biology and Environmental Science, University College Dublin, Belfield, Dublin 4, Ireland. ²⁴Department of Experimental and Health Sciences, Institute of Evolutionary Biology (UPF-CSIC), Universitat Pompeu Fabra, Barcelona 08003, Spain. ²⁵Department of Computational Biology, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213, USA. ²⁶Neuroscience Institute, Carnegie Mellon University, Pittsburgh, PA 15213, USA. ²⁷Program in Molecular Medicine, UMass Chan Medical School, Worcester, MA 01605, USA. ²⁸Department of Epidemiology & Biostatistics, University of California San Francisco, San Francisco, CA 94158, USA. ²⁹Gladstone Institutes, San Francisco, CA 94158, USA. ³⁰Center for Species Survival, Smithsonian's National Zoo and Conservation Biology Institute, Washington, DC 20008, USA. ³¹Computer Technologies Laboratory, ITMO University, St. Petersburg 197101, Russia. ³²Smithsonian-Mason School of Conservation, George Mason University, Front Royal, VA 22630, USA. ³³Department of Biological Sciences, Mellon College of Science, Carnegie Mellon University, Pittsburgh, PA 15213, USA. ³⁴Senckenberg Research Institute and Natural History Museum Frankfurt, 60325

Frankfurt am Main, Germany. ³⁵Department of Evolution and Ecology, University of California Davis, Davis, CA 95616, USA. ³⁶John Muir Institute for the Environment, University of California Davis, Davis, CA 95616, USA. ³⁷Morningside Graduate School of Biomedical Sciences, UMass Chan Medical School, Worcester, MA 01605, USA. ³⁸Department of Genetics, Yale School of Medicine, New Haven, CT 06510, USA. ³⁹Catalan Institution of Research and Advanced Studies (ICREA), Barcelona 08010, Spain. ⁴⁰CNAG-CRG, Centre for Genomic Regulation, Barcelona Institute of Science and Technology (BIST), Barcelona 08036, Spain. ⁴¹Department of Medicine and Life Sciences, Institute of Evolutionary Biology (UPF-CSIC), Universitat Pompeu Fabra, Barcelona 08003, Spain. ⁴²Institut Català de Paleontologia Miquel Crusafont, Universitat Autònoma de Barcelona, 08193 Cerdanyola del Vallès, Barcelona, Spain. ⁴³Institute of Cell Biology, University of Bern, 3012 Bern, Switzerland. ⁴⁴Department of Biological Sciences, Lehigh University, Bethlehem, PA 18015, USA. ⁴⁵BarcelonaBeta Brain Research Center, Pasqual Maragall Foundation, Barcelona 08005, Spain. ⁴⁶CRG, Centre for Genomic Regulation, Barcelona Institute of Science and Technology (BIST), Barcelona 08003, Spain. ⁴⁷Department of Comprehensive Care, School of Dental Medicine, Case Western Reserve University, Cleveland, OH 44106, USA. ⁴⁸Department of Vertebrate Zoology, Canadian Museum of Nature, Ottawa, ON K2P 2R1, Canada. ⁴⁹Department of Vertebrate Zoology, Smithsonian Institution, Washington, DC 20002, USA. ⁵⁰Narwhal Genome Initiative, Department of Restorative Dentistry and Biomaterials Sciences, Harvard School of Dental Medicine, Boston, MA 02115, USA. ⁵¹Department of Evolutionary Ecology, Leibniz Institute for Zoo and Wildlife Research, 10315 Berlin, Germany. ⁵²Medical Scientist Training Program, University of Pittsburgh School of Medicine, Pittsburgh, PA 15261, USA. ⁵³Chan Zuckerberg Biohub, San Francisco, CA 94158, USA. ⁵⁴Division of Messel Research and Mammalogy, Senckenberg Research Institute and Natural History Museum Frankfurt, 60325 Frankfurt am Main,

Germany. ⁵⁵Conservation Genetics, San Diego Zoo Wildlife Alliance, Escondido, CA 92027, USA. ⁵⁶Department of Evolution, Behavior and Ecology, School of Biological Sciences, University of California San Diego, La Jolla, CA 92039, USA. ⁵⁷Department of Organismic and Evolutionary Biology, Harvard University, Cambridge, MA 02138, USA. ⁵⁸Howard Hughes Medical Institute, Harvard University, Cambridge, MA 02138, USA. ⁵⁹Department of Ecology and Evolutionary Biology, University of California Santa Cruz, Santa Cruz, CA 95064, USA. ⁶⁰Howard Hughes Medical Institute, University of California Santa Cruz, Santa Cruz, CA 95064, USA. ⁶¹Department of Evolution, Ecology and Organismal Biology, University of California Riverside, Riverside, CA 92521, USA. ⁶²Department of Genetics, University of North Carolina Medical School, Chapel Hill, NC 27599, USA. ⁶³Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden. ⁶⁴Iris Data Solutions, LLC, Orono, ME 04473, USA. ⁶⁵Museum of Zoology, Senckenberg Natural History Collections Dresden, 01109 Dresden, Germany. ⁶⁶Allen Institute for Brain Science, Seattle, WA 98109, USA.

SUPPLEMENTARY MATERIALS

[science.org/doi/10.1126/science.abn3943](https://doi.org/10.1126/science.abn3943)

Materials and Methods

Supplementary Text

Figs. S1 to S18

Tables S1 to S15

MDAR Reproducibility Checklist

References (181–334)

Data S1 to S3

[View/request a protocol for this paper from Bio-rotocol.](#)

Submitted 23 November 2021; accepted 16 December 2022
10.1126/science.abn3943

Evolutionary constraint and innovation across hundreds of placental mammals

Matthew J. Christmas, Irene M. Kaplow, Diane P. Genereux, Michael X. Dong, Graham M. Hughes, Xue Li, Patrick F. Sullivan, Allyson G. Hindle, Gregory Andrews, Joel C. Armstrong, Matteo Bianchi, Ana M. Breit, Mark Diekhans, Cornelia Fanter, Nicole M. Foley, Daniel B. Goodman, Linda Goodman, Kathleen C. Keough, Bogdan Kirilenko, Amanda Kowalczyk, Colleen Lawless, Abigail L. Lind, Jennifer R. S. Meadows, Lucas R. Moreira, Ruby W. Redlich, Louise Ryan, Ross Swofford, Alejandro Valenzuela, Franziska Wagner, Ola Wallerman, Ashley R. Brown, Joana Damas, Kaili Fan, John Gatesy, Jenna Grimshaw, Jeremy Johnson, Sergey V. Kozyrev, Alyssa J. Lawler, Voichita D. Marinescu, Kathleen M. Morrill, Austin Osmanski, Nicole S. Paulat, BaDoi N. Phan, Steven K. Reilly, Daniel E. Schffer, Cynthia Steiner, Megan A. Supple, Aryn P. Wilder, Morgan E. Wirthlin, James R. Xue, Zoonomia Consortium, Bruce W. Birren, Steven Gazal, Robert M. Hubley, Klaus-Peter Koepfli, Tomas Marques-Bonet, Wynn K. Meyer, Martin Nweeia, Pardis C. Sabeti, Beth Shapiro, Arian F. A. Smit, Mark S. Springer, Emma C. Teeling, Zhiping Weng, Michael Hiller, Danielle L. Levesque, Harris A. Lewin, William J. Murphy, Arcadi Navarro, Benedict Paten, Katherine S. Pollard, David A. Ray, Irina Ruf, Oliver A. Ryder, Andreas R. Pfenning, Kerstin Lindblad-Toh, Elinor K. Karlsson, Gregory Andrews, Joel C. Armstrong, Matteo Bianchi, Bruce W. Birren, Kevin R. Bredemeyer, Ana M. Breit, Matthew J. Christmas, Hiram Clawson, Joana Damas, Federica Di Palma, Mark Diekhans, Michael X. Dong, Eduardo Eizirik, Kaili Fan, Cornelia Fanter, Nicole M. Foley, Karin Forsberg-Nilsson, Carlos J. Garcia, John Gatesy, Steven Gazal, Diane P. Genereux, Linda Goodman, Jenna Grimshaw, Michaela K. Halsey, Andrew J. Harris, Glenn Hickey, Michael Hiller, Allyson G. Hindle, Robert M. Hubley, Graham M. Hughes, Jeremy Johnson, David Juan, Irene M. Kaplow, Elinor K. Karlsson, Kathleen C. Keough, Bogdan Kirilenko, Klaus-Peter Koepfli, Jennifer M. Korstian, Amanda Kowalczyk, Sergey V. Kozyrev, Alyssa J. Lawler, Colleen Lawless, Thomas Lehmann, Danielle L. Levesque, Harris A. Lewin, Xue Li, Abigail Lind, Kerstin Lindblad-Toh, Ava Mackay-Smith, Voichita D. Marinescu, Tomas Marques-Bonet, Victor C. Mason, Jennifer R. S. Meadows, Wynn K. Meyer, Jill E. Moore, Lucas R. Moreira, Diana D. Moreno-Santillan, Kathleen M. Morrill, Gerard Muntan, William J. Murphy, Arcadi Navarro, Martin Nweeia, Sylvia Ortmann, Austin Osmanski, Benedict Paten, Nicole S. Paulat, Andreas R. Pfenning, BaDoi N. Phan, Katherine S. Pollard, Henry E. Pratt, David A. Ray, Steven K. Reilly, Jeb R. Rosen, Irina Ruf, Louise Ryan, Oliver A. Ryder, Pardis C. Sabeti, Daniel E. Schffer, Aitor Serres, Beth Shapiro, Arian F. A. Smit, Mark Springer, Chaitanya Srinivasan, Cynthia Steiner, Jessica M. Storer, Kevin A. M. Sullivan, Patrick F. Sullivan, Elisabeth Sundstrm, Megan A. Supple, Ross Swofford, Joy-El Talbot, Emma Teeling, Jason Turner-Maier, Alejandro Valenzuela, Franziska Wagner, Ola Wallerman, Chao Wang, Juehan Wang, Zhiping Weng, Aryn P. Wilder, Morgan E. Wirthlin, James R. Xue, and Xiaomeng Zhang

Science, **380** (6643), eabn3943.
DOI: 10.1126/science.abn3943

View the article online

<https://www.science.org/doi/10.1126/science.abn3943>

Permissions

<https://www.science.org/help/reprints-and-permissions>

Use of this article is subject to the [Terms of service](#)

Science (ISSN) is published by the American Association for the Advancement of Science. 1200 New York Avenue NW, Washington, DC 20005. The title *Science* is a registered trademark of AAAS.

Copyright © 2023 The Authors, some rights reserved; exclusive licensee American Association for the Advancement of Science. No claim to original U.S. Government Works