

# Somatic mutation rates scale with lifespan across mammals

<https://doi.org/10.1038/s41586-022-04618-z>

Received: 17 August 2021

Accepted: 7 March 2022

Published online: 13 April 2022

Open access

 Check for updates

Alex Cagan<sup>1,15</sup>✉, Adrian Baez-Ortega<sup>1,15</sup>, Natalia Brzozowska<sup>1</sup>, Federico Abascal<sup>1</sup>, Tim H. H. Coorens<sup>1</sup>, Mathijs A. Sanders<sup>1,2</sup>, Andrew R. J. Lawson<sup>1</sup>, Luke M. R. Harvey<sup>1</sup>, Shriram Bhosle<sup>1</sup>, David Jones<sup>1</sup>, Raul E. Alcantara<sup>1</sup>, Timothy M. Butler<sup>1</sup>, Yvette Hooks<sup>1</sup>, Kirsty Roberts<sup>1</sup>, Elizabeth Anderson<sup>1</sup>, Sharna Lunn<sup>1</sup>, Edmund Flach<sup>3</sup>, Simon Spiro<sup>3</sup>, Inez Januszczak<sup>3,4</sup>, Ethan Wrigglesworth<sup>3</sup>, Hannah Jenkins<sup>3</sup>, Tilly Dallas<sup>3</sup>, Nic Masters<sup>3</sup>, Matthew W. Perkins<sup>5</sup>, Robert Deaville<sup>5</sup>, Megan Druce<sup>6,7</sup>, Ruzhica Bogeska<sup>6,7</sup>, Michael D. Milsom<sup>6,7</sup>, Björn Neumann<sup>8,9</sup>, Frank Gorman<sup>10</sup>, Fernando Constantino-Casas<sup>10</sup>, Laura Peachey<sup>10,11</sup>, Diana Bochynska<sup>10,12</sup>, Ewan St. John Smith<sup>13</sup>, Moritz Gerstung<sup>14</sup>, Peter J. Campbell<sup>1</sup>, Elizabeth P. Murchison<sup>10</sup>, Michael R. Stratton<sup>1</sup> & Iñigo Martincorena<sup>1,15</sup>✉

The rates and patterns of somatic mutation in normal tissues are largely unknown outside of humans<sup>1–7</sup>. Comparative analyses can shed light on the diversity of mutagenesis across species, and on long-standing hypotheses about the evolution of somatic mutation rates and their role in cancer and ageing. Here we performed whole-genome sequencing of 208 intestinal crypts from 56 individuals to study the landscape of somatic mutation across 16 mammalian species. We found that somatic mutagenesis was dominated by seemingly endogenous mutational processes in all species, including 5-methylcytosine deamination and oxidative damage. With some differences, mutational signatures in other species resembled those described in humans<sup>8</sup>, although the relative contribution of each signature varied across species. Notably, the somatic mutation rate per year varied greatly across species and exhibited a strong inverse relationship with species lifespan, with no other life-history trait studied showing a comparable association. Despite widely different life histories among the species we examined—including variation of around 30-fold in lifespan and around 40,000-fold in body mass—the somatic mutation burden at the end of lifespan varied only by a factor of around 3. These data unveil common mutational processes across mammals, and suggest that somatic mutation rates are evolutionarily constrained and may be a contributing factor in ageing.

Somatic mutations accumulate in healthy cells throughout life. They underpin the development of cancer<sup>9</sup> and, for decades, have been speculated to contribute to ageing<sup>10–12</sup>. Directly studying somatic mutations in normal tissues has been challenging owing to the difficulty of detecting mutations present in single cells or small clones in a tissue. Only recent technological developments, such as *in vitro* expansion of single cells into colonies<sup>13,14</sup>, microdissection of histological units<sup>8,15</sup>, single-cell sequencing<sup>16,17</sup> or single-molecule sequencing<sup>18</sup>, are beginning to enable the study of somatic mutation in normal tissues.

Over the last few years, studies in humans have started to provide a detailed understanding of somatic mutation rates and the contribution of endogenous and exogenous mutational processes across normal tissues<sup>8,13,14,19,20</sup>. These studies are also revealing how, as we

age, some human tissues are colonized by mutant cells that contain cancer-driving mutations, and how this clonal composition changes with age and disease. With the exception of some initial studies, far less is known about somatic mutation in other species<sup>1–7</sup>. Yet, comparative analyses of somatic mutagenesis would shed light on the diversity of mutagenic processes across species, and on long-standing questions regarding the evolution of somatic mutation rates and their role in cancer and ageing.

A decades-long hypothesis on the evolution of somatic mutation rates pertains to the relationship between body mass and cancer risk. Some models predict that the risk of cancer should increase proportionally to the number of cells at risk of transformation. However, there appears to be no correlation between body mass and cancer risk across

<sup>1</sup>Cancer, Ageing and Somatic Mutation (CASM), Wellcome Sanger Institute, Hinxton, UK. <sup>2</sup>Department of Hematology, Erasmus MC Cancer Institute, Rotterdam, the Netherlands. <sup>3</sup>Wildlife Health Services, Zoological Society of London, London, UK. <sup>4</sup>The Natural History Museum, London, UK. <sup>5</sup>Institute of Zoology, Zoological Society of London, London, UK. <sup>6</sup>Division of Experimental Hematology, German Cancer Research Center (DKFZ), Heidelberg, Germany. <sup>7</sup>Heidelberg Institute for Stem Cell Technology and Experimental Medicine GmbH (HI-STEM), Heidelberg, Germany. <sup>8</sup>Wellcome Trust–Medical Research Council Cambridge Stem Cell Institute, University of Cambridge, Cambridge, UK. <sup>9</sup>Department of Clinical Neurosciences, University of Cambridge, Cambridge, UK. <sup>10</sup>Department of Veterinary Medicine, University of Cambridge, Cambridge, UK. <sup>11</sup>Bristol Veterinary School, Faculty of Health Sciences, University of Bristol, Langford, UK. <sup>12</sup>Department of Pathology, Faculty of Veterinary Medicine, Universitatea de Stiinta Agricole si Medicina Veterinara, Cluj-Napoca, Romania. <sup>13</sup>Department of Pharmacology, University of Cambridge, Cambridge, UK. <sup>14</sup>European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Hinxton, UK. <sup>15</sup>These authors contributed equally: Alex Cagan, Adrian Baez-Ortega. ✉e-mail: ac36@sanger.ac.uk; im3@sanger.ac.uk

species<sup>21,22</sup>. This observation, known as Peto's paradox, suggests that the evolution of larger body sizes is likely to require the evolution of stronger cancer suppression mechanisms<sup>23,24</sup>. Whether evolutionary reduction of cancer risk across species is partly achieved by a reduction of somatic mutation rates remains unknown.

A second long-standing hypothesis on the evolution of somatic mutation rates relates to the proposed role of somatic mutations in ageing. Multiple forms of molecular damage, including somatic mutations, telomere attrition, epigenetic drift and loss of proteostasis, have been proposed to contribute to ageing, but their causal roles and relative contributions remain debated<sup>25,26</sup>. Evolutionary theory predicts that species will evolve protection or repair mechanisms against life-threatening damage to minimize death from intrinsic causes, but that selection is too weak to delay ageing far beyond the typical life expectancy of an organism in the wild (Supplementary Note 1). If somatic mutations contribute to ageing, theory predicts that somatic mutation rates may inversely correlate with lifespan across species<sup>27,28</sup>. This prediction has remained largely untested owing to the difficulty of measuring somatic mutation rates across species.

### Detection of somatic mutations across species

The study of somatic mutations with standard whole-genome sequencing requires isolating clonal groups of cells recently derived from a single cell<sup>8,13,14</sup>. To study somatic mutations across a diverse set of mammals, we isolated 208 individual intestinal crypts from 56 individuals across 16 species with a wide range of lifespans and body sizes: black-and-white colobus monkey, cat, cow, dog, ferret, giraffe, harbour porpoise, horse, human, lion, mouse, naked mole-rat, rabbit, rat, ring-tailed lemur and tiger (Supplementary Table 1). We chose intestinal crypts for several reasons. First, they are histologically identifiable units that line the epithelium of the colon and small intestine and are amenable to laser microdissection. Second, human studies have confirmed that individual crypts become clonally derived from a single stem cell and show a linear accumulation of mutations with age, which enables the estimation of somatic mutation rates through genome sequencing of single crypts<sup>8</sup>. Third, in most human crypts, most somatic mutations are caused by endogenous mutational processes common to other tissues, rather than by environmental mutagens<sup>8,18</sup>.

A colon sample was collected from each individual, with the exception of a ferret from which only a small intestine sample was available. This sample was included because results in humans have shown that the mutation rates of colorectal and small intestine epithelial stem cells are similar<sup>14,20</sup> (Extended Data Fig. 1). We then used laser microdissection on histological sections to isolate individual crypts for whole-genome sequencing with a low-input library preparation method<sup>29</sup> (Fig. 1a, Extended Data Fig. 2, Supplementary Table 2), with the exception of human crypts, for which sequencing data were obtained from a previous study<sup>8</sup>. A bioinformatic pipeline was developed to call somatic mutations robustly in all these species despite the variable quality of their genome assemblies (Methods). The distribution of variant allele fractions of the mutations detected in each crypt confirmed that crypts are clonal units in all species, enabling the study of somatic mutation rates and signatures (Extended Data Fig. 3).

We found substantial variation in the number of somatic single-base substitutions across species and across individuals within each species (Fig. 1b). For five species with samples from multiple individuals (dog, human, mouse, naked mole-rat and rat), linear regression confirmed a clear accumulation of somatic mutations with age (Fig. 1c, Extended Data Fig. 4, Supplementary Table 3). All linear regressions were also consistent with a non-significant intercept. This resembles observations in humans<sup>20</sup> and suggests that the time required for a single stem cell to drift to fixation within a crypt is a small fraction of the lifespan of a species. This facilitates the estimation of somatic mutation rates across species by dividing the number of mutations in a crypt by the age of the

individual (Supplementary Table 4). The number of somatic insertions and deletions (indels) was consistently lower than that of substitutions in all crypts (Fig. 1b), in agreement with previous findings in humans<sup>8</sup>.

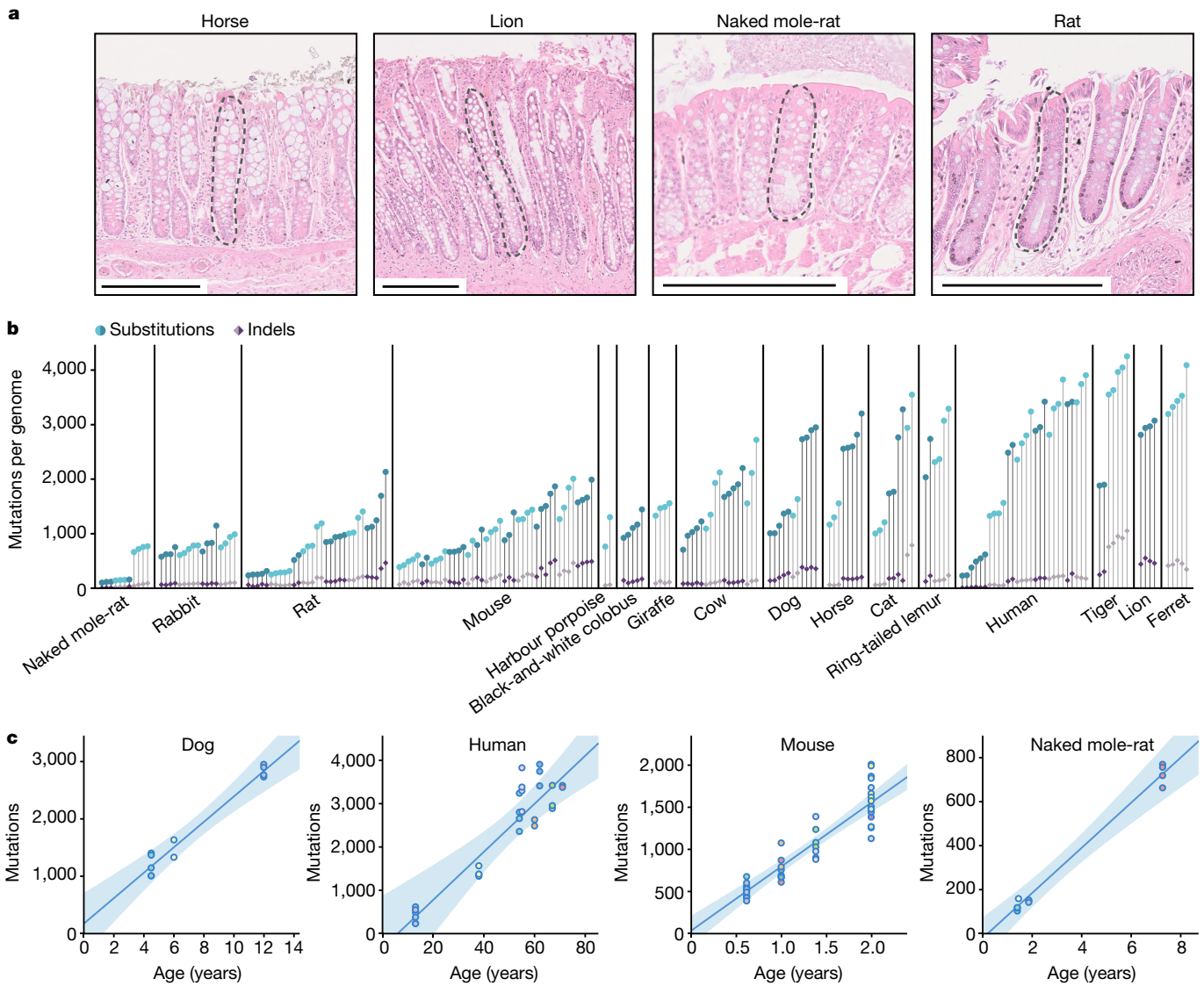
### Mutational signatures across mammals

Somatic mutations can be caused by multiple mutational processes, involving different forms of DNA damage and repair. Different processes cause characteristic frequencies of base substitution types and indels at different sequence contexts, often referred to as mutational signatures, which can be inferred from mutation data<sup>30</sup>. Across species, the mutational spectra showed clear similarities, with a dominance of cytosine-to-thymine (C>T) substitutions at CpG sites, as observed in human colon, but with considerable variation in the frequency of other substitution types (Fig. 2a). To quantify the contribution of different mutational processes to the observed spectra, we applied mutational signature decomposition<sup>8,30</sup>. We used a Bayesian model to infer mutational signatures de novo, while accounting for differences in genome sequence composition across species, and using the COSMIC human signature SBS1 (C>T substitutions at CpG sites) as a fixed prior to ensure its complete deconvolution<sup>31</sup> (Methods). This approach identified two signatures beyond SBS1, labelled SBSB and SBSC, which resemble COSMIC human signatures SBS5 and SBS18, respectively (cosine similarities 0.93 and 0.91) (Fig. 2b).

This analysis suggests that the same three signatures that dominate somatic mutagenesis in the human colon are dominant in other mammals: SBS1, which is believed to result from the spontaneous deamination of 5-methylcytosine<sup>8,32</sup>; SBSB (SBS5), a common signature across human tissues that may result from endogenous damage and repair<sup>18,33</sup>; and SBSC (SBS18), which is dominated by C>A substitutions and attributed to oxidative DNA damage<sup>30</sup>. Signature SBSC contains a minor component of T>A substitutions (resembling COSMIC SBS34), which appear to be the result of DNA polymerase slippage at the boundaries between adjacent adenine and thymine homopolymer tracts, but could also reflect assembly errors at those sites<sup>33</sup>. Although all of the species that we examined shared the three mutational signatures, their contributions varied substantially across species (Fig. 2c). SBSC was particularly prominent in mouse and ferret, and the ratio of SBS1 to SBSB/5 varied from approximately 1.2 in rat or rabbit to 6.4 in tiger. In several species with data from multiple individuals, separate linear regressions for each signature confirmed that mutations from all three signatures accumulate with age (Fig. 2d, Extended Data Fig. 5).

Although signature deconvolution identified three signatures that are active across species, we noticed some differences in the mutational profile of signature SBSB among species. To investigate this further, we inferred independent versions of SBSB from each species, while accounting for differences in genome sequence composition (Methods). This revealed inter-species variability in the mutational profile of this signature, particularly in the C>T component (Extended Data Fig. 6). Species-specific versions of SBSB showed different similarities to the related human signatures SBS5 and SBS40. For example, SBSB inferred from the human data showed a stronger similarity with the reference human signature SBS5 (cosine similarities with SBS5 and SBS40: 0.93 and 0.84), whereas SBSB from rabbit more closely resembled the reference human signature SBS40 (0.87 and 0.91). These observations are consistent with the hypothesis that SBS5 and SBS40 result from a combination of correlated mutational processes, with some variation across human tissues<sup>18,33</sup> and across species.

Analysis of the indel mutational spectra revealed a dominance of the human indel signatures ID1 and ID2, which are characterized by single-nucleotide indels at A/T homopolymers, and probably caused by strand slippage during DNA replication<sup>30</sup> (Extended Data Fig. 7a). The ratio of insertions (ID1) to deletions (ID2) appears to vary across species, possibly reflecting a differential propensity for slippage of the template and nascent DNA strands<sup>30</sup>. In addition, the indel spectra suggest a potential contribution of signature ID9 (the aetiology of which



**Fig. 1 | Somatic mutation burden in mammalian colorectal crypts. a**, Histology images of colon samples from horse, lion, naked mole-rat and rat, with one colorectal crypt marked in each. Scale bars, 250  $\mu\text{m}$ . **b**, Burden of somatic substitutions and indels per diploid genome in each colorectal crypt sample (corrected for the size of the analysable genome). Samples are grouped by individual, with samples from the same individual coloured in the same shade. Species, and individuals within each species, are sorted by mean

mutation burden. **c**, Linear regression of somatic substitution burden (corrected for analysable genome size) on individual age for dog, human, mouse and naked mole-rat samples. Samples from the same individual are shown in the same colour. Regression was performed using mean mutation burdens per individual. Shaded areas indicate 95% confidence intervals of the regression line.

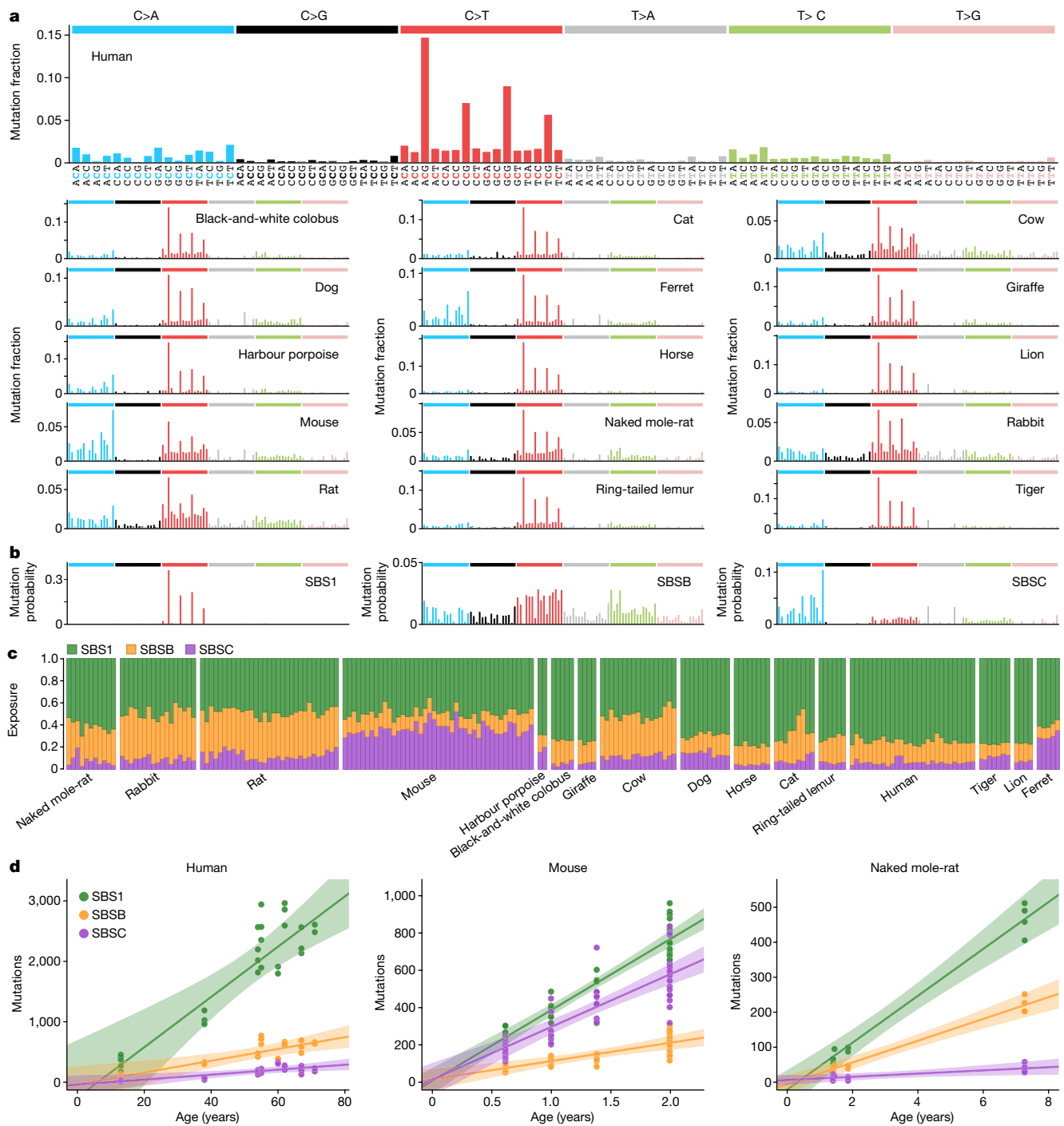
remains unknown) to human, colobus, cow, giraffe and rabbit. Analysis of indels longer than one base pair also suggested the presence of a signature of four-base-pair insertions at tetrameric repeats, which was particularly prevalent in mouse and tiger; a pattern of insertions of five or more base pairs at repeats in colobus; and a pattern of deletions of five or more base pairs at repeats, which was prominent in rabbit and resembles ID8 (a signature possibly caused by double-strand break repair through non-homologous end joining<sup>30</sup>) (Extended Data Fig. 7a).

### Other mutational processes and selection

The apparent lack of additional mutational signatures is noteworthy. A previous study of 445 colorectal crypts from 42 human donors found that many crypts were affected by a signature that was later attributed to colibactin, a genotoxin produced by *pks*<sup>+</sup> strains of *Escherichia coli*<sup>38,34,35</sup>. Analysing the original human data and our non-human data with the same methodology, we found evidence of colibactin

mutagenesis in 21% of human crypts, but only uncertain evidence of colibactin in one non-human crypt (0.6%) (Extended Data Fig. 7b, Methods). This revealed a significant depletion of colibactin mutagenesis in the non-human crypts studied (Fisher's exact test,  $P = 7 \times 10^{-14}$ ). The apparent difference in colibactin mutagenesis observed between species, or between the cohorts studied, might result from a different prevalence of *pks*<sup>+</sup> *E. coli* strains<sup>36</sup> or a different expression of colibactin by *pks*<sup>+</sup> *E. coli* across species<sup>37</sup>. Finally, we also searched for evidence of APOBEC signatures (SBS2 and SBS13), which have been reported in a small number of human crypts and are believed to be caused by APOBEC DNA-editing cytidine deaminases. We detected APOBEC signatures in 2% ( $n = 9$ ) of human crypts and found only uncertain evidence in one non-human crypt ( $P = 0.30$ ).

Beyond substitutions and indels, crypts from the eight species with chromosome-level genome assemblies were inspected for large-scale copy number changes (at least 1 Mb) (Methods). Studies in humans have found that large-scale copy number changes are relatively rare in

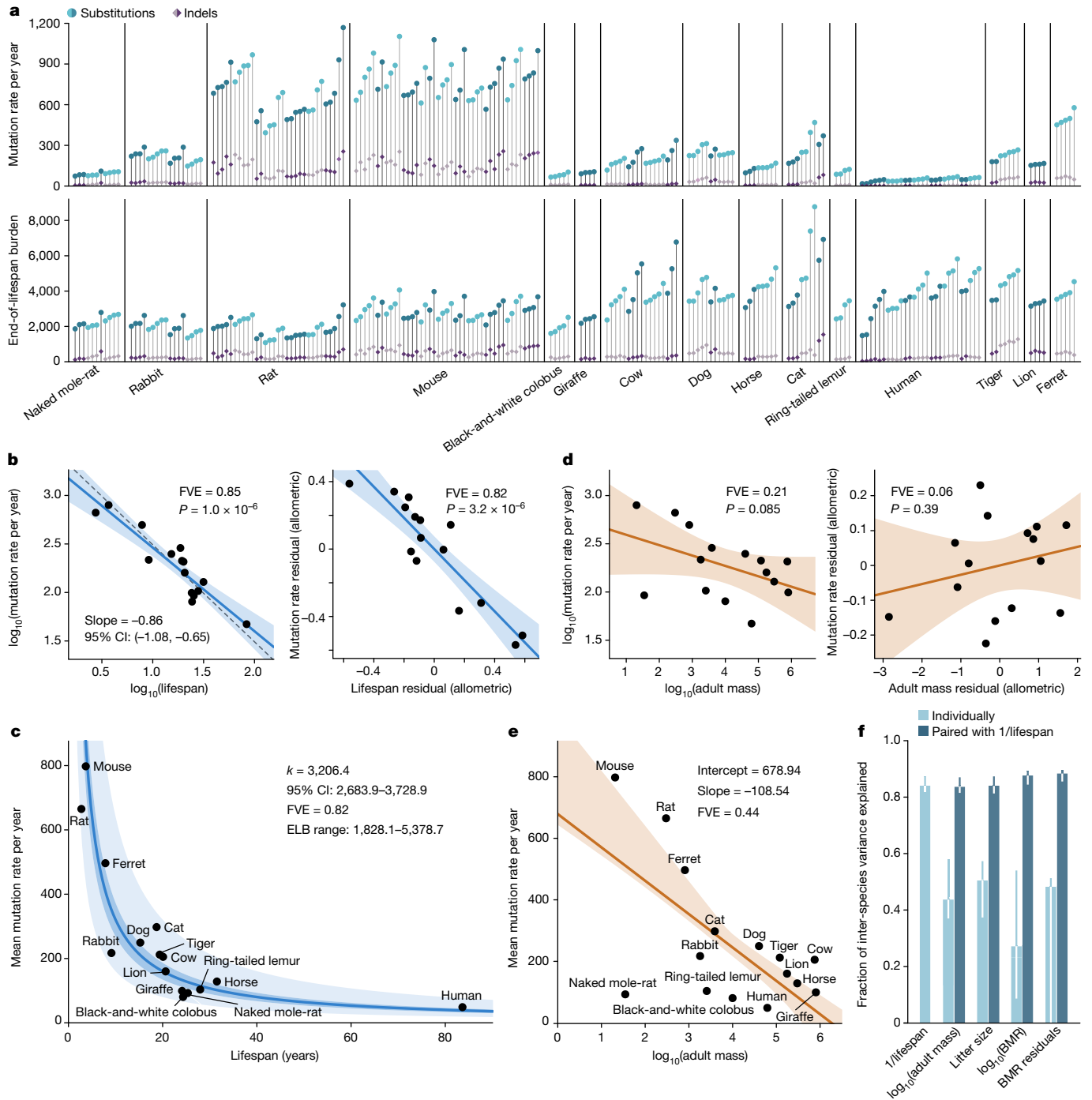


**Fig. 2 | Mutational processes in the mammalian colon. a**, Mutational spectra of somatic substitutions in each species. The x axis shows 96 mutation types on a trinucleotide context, coloured by base substitution type. **b**, Mutational signatures inferred from (SBSB, SBSC) or fitted to (SBS1) the species mutational spectra shown in **a**, and normalized to the human genome trinucleotide frequencies. The y axis shows mutation probability. **c**, Estimated

contribution of each signature to each sample. Samples are arranged horizontally as in Fig. 1b. **d**, Linear regression of signature-specific mutation burdens (corrected for analysable genome size) on individual age for human, mouse and naked mole-rat samples. Regression was performed using mean mutation burdens per individual. Shaded areas indicate 95% confidence intervals of the regression line.

normal tissues, including colorectal epithelium<sup>8</sup>. Consistent with these results, we only identified 4 large copy number changes across the 162 crypts included in this analysis: 2 megabase-scale deletions in 2 crypts from the same cow; the loss of an X chromosome in a female mouse crypt; and a 52-Mb segment with copy-neutral loss of heterozygosity in a human crypt (Extended Data Fig. 8, Methods). These results suggest that large-scale somatic copy number changes in normal tissues are also rare in other mammalian species.

Previous analyses in humans have shown that most somatic mutations in colorectal crypts accumulate neutrally, without clear evidence of negative selection against non-synonymous mutations and with a low frequency of positively selected cancer-driver mutations<sup>8</sup>. To study somatic selection in our data, we calculated the exome-wide ratio of non-synonymous to synonymous substitution rates (dN/dS) in each of the 12 species with available genome annotation. To do so and to detect genes under positive selection, while accounting for the effects of trinucleotide



**Fig. 3 | Associations between somatic mutation rates and life-history traits.**

**a**, Somatic mutation rate per year and expected end-of-lifespan mutation burden (ELB) per crypt. Samples are arranged horizontally as in Fig. 1b; harbour porpoise samples were excluded owing to the age of the sampled individual being unknown. **b**, Left, allometric regression of somatic mutation rate on lifespan. Right, regression of body-mass-adjusted residuals for somatic mutation rate and lifespan (partial correlation; Methods). Regression was performed using mean mutation rates per species. Shaded areas represent 95% confidence intervals (CI) of regression lines. FVE and  $P$  values (by  $F$ -test) are indicated (note that, for simple linear regression,  $FVE = R^2$ ). The dashed line denotes a reference slope of  $-1$ . **c**, Zero-intercept LME regression of somatic mutation rate on inverse lifespan ( $1/\text{lifespan}$ ), presented on the scale of untransformed lifespan ( $x$  axis). For simplicity, the  $y$  axis shows mean mutation rates per species, although rates per crypt were used in the regression. The

darker shaded area indicates 95% CI of the regression line, and the lighter shaded area marks a twofold deviation from the line. Point estimate and 95% CI of the regression slope ( $k$ ), FVE and range of end-of-lifespan burden are indicated. **d**, Allometric regression and linear regression of lifespan-adjusted residuals, for somatic mutation rate and body mass (elements as described in **b**). **e**, Free-intercept LME regression of somatic mutation rate on log-transformed body mass. The  $y$  axis shows mean mutation rates per species, although rates per crypt were used in the regression. Shaded area indicates 95% bootstrap interval of the regression line ( $n = 10,000$  replicates). Point estimates of the regression intercept and slope, and FVE, are indicated. **f**, FVE values for free-intercept LME models using  $1/\text{lifespan}$  or other life-history variables (alone or combined with  $1/\text{lifespan}$ ) as explanatory variables. Error bars indicate 95% bootstrap intervals ( $n = 10,000$ ).

**Table 1 | Variation in adult body mass, lifespan, somatic mutation rate and end-of-lifespan mutation burden across the 16 mammalian species surveyed**

Variable	Minimum	Maximum	Fold variation
<b>Adult mass (g)</b>	20.50	800,000.00	39,024.39
<b>Lifespan (years)</b>	2.75	83.67	30.44
<b>Mutation rate per year (substitutions per genome)</b>	47.12	796.42	16.90
<b>End-of-lifespan burden (substitutions per genome)</b>	1,828.08	5,378.73	2.94

Species-level estimates are provided in Supplementary Tables 3 and 6.

sequence context and mutation rate variation across genes, we used the dNdScv model<sup>38</sup> (Methods). Although the limited number of coding somatic mutations observed in most species precluded an in-depth analysis of selection, exome-wide dN/dS ratios for somatic substitutions were not significantly different from unity in any species, in line with previous findings in humans<sup>8</sup> (Extended Data Fig. 9). Gene-level analysis did not find genes under significant positive selection in any species, although larger studies are likely to identify rare cancer-driver mutations<sup>8</sup>.

### Correlation with life-history traits

Whereas similar mutational processes operate across the species surveyed, the mutation rate per genome per year varied widely. Across the 15 species with age information, we found that substitution rates per genome ranged from 47 substitutions per year in humans to 796 substitutions per year in mice, and indel rates from 2.5 to 158 indels per year, respectively (Fig. 3a, Supplementary Table 4, Methods).

To investigate the relationship between somatic mutation rates, lifespan and other life-history traits, we first estimated the lifespan of each species using survival curves. We used a large collection of mortality data from animals in zoos to minimize the effect of extrinsic mortality (Extended Data Fig. 10). We defined lifespan as the age at which 80% of individuals reaching adulthood have died, to reduce the effects of outliers and variable cohort sizes that affect maximum lifespan estimates<sup>39</sup> (Methods). Notably, we found a tight anticorrelation between somatic mutation rates per year and lifespan across species (Fig. 3b). A log-log allometric regression yielded a strong linear anticorrelation between mutation rate per year and lifespan (fraction of inter-species variance explained (FVE) = 0.85,  $P = 1 \times 10^{-6}$ ), with a slope close to and not significantly different from  $-1$ . This supports a simple model in which somatic mutation rates per year are inversely proportional to the lifespan of a species ( $\text{rate} \propto 1/\text{lifespan}$ ), such that the number of somatic mutations per cell at the end of the lifespan (the end-of-lifespan burden; ELB) is similar in all species.

To further study the relationship between somatic mutation rates and life-history variables, we used linear mixed-effects (LME) regression models. These models account for the hierarchical structure of the data (with multiple crypts per individual and multiple individuals per species), as well as the heteroscedasticity of somatic mutation rate estimates across species (Methods). Using these models, we estimated that the inverse of lifespan explained 82% of the inter-species variance in somatic substitution rates ( $\text{rate} = k/\text{lifespan}$ ) ( $P = 2.9 \times 10^{-9}$ ; Fig. 3c), with the slope of this regression ( $k$ ) representing the mean estimated ELB across species (3,206.4 substitutions per genome per crypt, 95% confidence interval 2,683.9–3,728.9). Of note, despite uncertainty in the estimates of both somatic mutation rates and lifespans, and despite the diverse life histories of the species surveyed—including around 30-fold variation in lifespan and around 40,000-fold variation in body mass—the estimated mutation load per cell at the end of lifespan varied by only around threefold across species (Table 1). Analogous results were obtained when repeating the analysis with estimates of

the protein-coding mutation rate, which may be a better proxy for the functional effect of somatic mutations (85% of variance explained; ELB: 31 coding substitutions per crypt) (Extended Data Fig. 11, Methods).

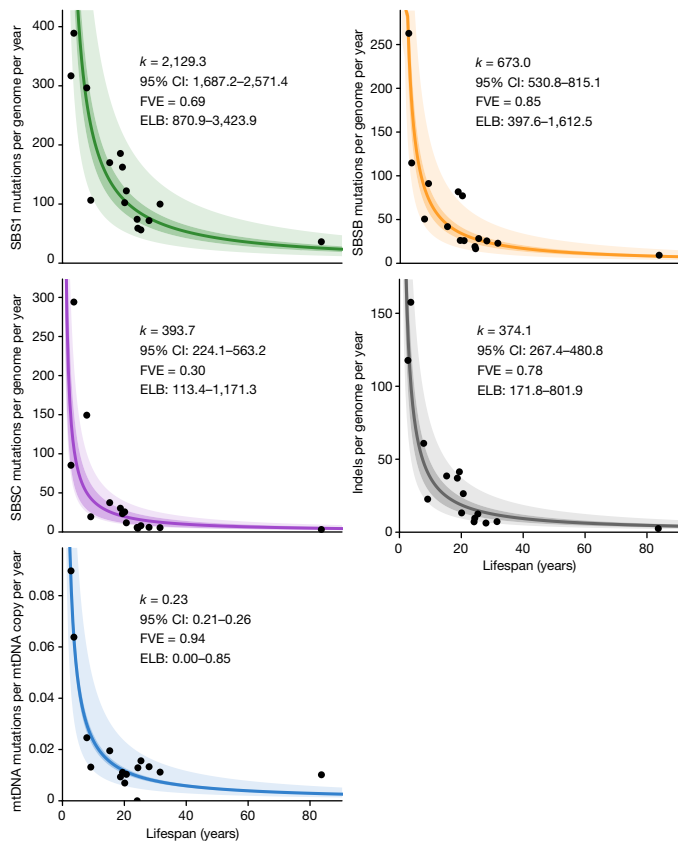
We next examined the association between somatic mutation rates and adult body mass, which is known to be a common confounder in correlations that involve lifespan<sup>40,41</sup>. An anticorrelation between somatic mutation rates and body mass may be expected if the modulation of cancer risk across species of vastly different sizes has been a major factor in the evolution of somatic mutation rates. We observed that log-transformed adult body mass was less strongly associated with somatic substitution rates than the inverse of lifespan (allometric regression FVE = 0.21, Fig. 3d; LME regression FVE = 0.44, Fig. 3e). Given that lifespan is correlated with body mass, we performed two tests to assess whether body mass explained any variation in somatic mutation rates that was not explained by lifespan. First, including both the inverse of lifespan and log-transformed adult body mass in the regression model suggested that body mass does not explain a significant amount of variance in somatic mutation rates across species after accounting for the effect of lifespan (likelihood ratio tests:  $P = 0.16$  for body mass on a model with lifespan;  $P < 10^{-4}$  for lifespan on a model with body mass; Fig. 3f, Methods). Second, partial correlation analyses using allometric regressions further confirmed that the association between somatic mutation rates and lifespan is unlikely to be mediated by the effect of body mass on both variables (lifespan residuals:  $P = 3.2 \times 10^{-6}$ , FVE = 0.82, Fig. 3b; body mass residuals:  $P = 0.39$ , FVE = 0.06, Fig. 3d; Methods).

The fact that the variation in somatic mutation rates across species appears to be dominated by lifespan rather than body size is also apparent when looking at particularly informative species. Giraffe and naked mole-rat, for instance, have similar somatic mutation rates (99 and 93 substitutions per year, respectively), in line with their similar lifespans (80th percentiles: 24 and 25 years, respectively), despite a difference of around 23,000-fold in adult body mass (Fig. 3c, e). Similarly, cows, giraffes and horses weigh much more than an average human, and yet have somatic mutation rates that are several fold higher, in line with expectation from their lifespan but not their body mass. Altogether, the weak correlation between body mass and somatic mutation rates after correction for lifespan suggests that the evolution of larger body sizes may have relied on alternative or additional strategies to limit cancer risk, as has been speculated<sup>24,42</sup> (Supplementary Note 2). Of note, the low somatic mutation rate of naked mole-rats, which is unusual for their body mass but in line with their long lifespan (Fig. 3c, e), might contribute to the exceptionally low incidence rates of cancer in this species<sup>43</sup>.

We found similar results for other life-history variables that have been proposed to correlate with lifespan, namely basal metabolic rate (BMR) and litter size<sup>44</sup> (Fig. 3f). With the caveat that estimates for these variables vary in quality, they showed weaker correlations with the somatic mutation rate as single predictors, and small non-significant increases in explanatory power when considered together with lifespan (likelihood ratio tests:  $P = 0.92$  for litter size;  $P = 0.083$  for log-BMR;  $P = 0.79$  for allometric BMR residuals; Fig. 3f, Methods). We note that the results above are robust to the use of alternative measures of the somatic mutation rate, including the rate per exome or mutations per Mb (Extended Data Fig. 11, Methods); alternative estimates of lifespan, including maximum lifespan (Extended Data Fig. 12, Methods); alternative regression models, including a Bayesian hierarchical model and a phylogenetic generalised least-squares regression, which accounts for the effect of phylogenetic relationships (Extended Data Fig. 13a, b, Methods); and bootstrapping analyses at the level of individuals or species (Extended Data Fig. 13c, Methods).

### Mutational processes and lifespan

To investigate whether a single biological process could drive the association between somatic mutation rates and lifespan, we analysed each mutational signature separately. SBS1, SBSB/5 and SBSC/18 are believed to result from different forms of DNA damage and are expected to be



**Fig. 4 | Association between mutation rate subtypes and species lifespan.** Zero-intercept LME regression of somatic rates of signature-specific substitutions, indels and mtDNA mutations on inverse lifespan ( $1/\text{lifespan}$ ), presented on the scale of untransformed lifespan (x axis). For simplicity, y axes present mean mutation rates per species, although mutation rates per crypt were used in the regressions. The darker shaded areas indicate 95% confidence intervals (CI) of the regression lines, and the lighter shaded areas mark a twofold deviation from the regression lines. Point estimates and 95% CI of the regression slope ( $k$ ), fraction of inter-species variance explained by for each model (FVE) and ranges of end-of-lifespan burden (ELB) are indicated.

subject to different DNA repair pathways<sup>18,33</sup>. They also appear to differ in their association with the rate of cell division in humans, with SBS1 being more common in fast-proliferating tissues, such as colon and embryonic or foetal tissues, and SBS5 dominating in post-mitotic cells in the absence of cell division<sup>14,18,20</sup>. Overall, we found clear anticorrelations between mutation rates per year and lifespan for the three substitution signatures and for indels, suggesting that a single biological process or DNA repair pathway is unlikely to be responsible for this association (Fig. 4). The total mutation burden also appears to show a closer fit with lifespan than individual mutational processes, as measured by the range of end-of-lifespan burden for each process across species (Fig. 4). This might be expected if the observed anticorrelation were the result of evolutionary pressure on somatic mutation rates.

DNA damage and somatic mutations in the mitochondrial genome have also attracted considerable interest in the ageing field<sup>45</sup>. Our whole-genome sequencing of individual crypts provided high coverage of the mitochondrial genome, ranging from 2,188- to 29,691-fold. Normalized against the nuclear coverage, these data suggest that colorectal crypts contain on the order of around 100–2,000 mitochondrial genomes per cell (Extended Data Fig. 14a). Using a mutation-calling algorithm that is sensitive to low-frequency variants, we found a total of 261 mitochondrial mutations across 199 crypts (Extended Data Fig. 14a, Methods). The mutational spectra across species appeared broadly consistent with that observed in humans, with a dominance of C>T and A>G substitutions

that are believed to result from mitochondrial DNA (mtDNA) replication errors rather than DNA damage<sup>46</sup> (Extended Data Fig. 14b). Although the low number of mitochondrial mutations detected per species precludes a detailed analysis, the estimated number of somatic mutations per copy of mtDNA also appears to show an anticorrelation with lifespan. Across species, we obtained an average of 0.23 detectable mutations per copy of the mitochondrial genome by the end of lifespan (Fig. 4, Methods)—a considerable burden given the coding-sequence density and the functional relevance of the mitochondrial genome.

## Discussion

Using whole-genome sequencing of 208 colorectal crypts from 56 individuals, we provide insights into the somatic mutational landscape of 16 mammalian species. Despite their different diets and life histories, we found considerable similarities in their mutational spectra. Three main mutational signatures explain the spectra across species, albeit with varying contributions and subtle variations in the profile of signature SBSB. These results suggest that, at least in the colorectal epithelium, a conserved set of mutational processes dominate somatic mutagenesis across mammals.

The most notable finding of this study is the inverse scaling of somatic mutation rates with lifespan—a long-standing prediction of the somatic mutation theory of ageing<sup>11,27</sup>. Considering evolutionary and mechanistic models of ageing together provides a framework for discussing the possible implications of these results for ageing (see Supplementary Note 1). Jointly, these models predict ageing to be multifactorial, with multiple forms of molecular and cellular damage contributing to organismal ageing owing to evolutionary limits to selection acting on the rates of these processes. The inverse scaling of somatic mutation rates and lifespan is consistent with somatic mutations contributing to ageing and with somatic mutation rates being evolutionarily constrained, although we discuss alternative explanations below. This interpretation is also supported by studies reporting more efficient DNA repair in longer-lived species<sup>47,48</sup>. Somatic mutations could contribute to ageing in different ways. Traditionally, they have been proposed to contribute to ageing through deleterious effects on cellular fitness<sup>11,49</sup>, but recent findings question this assumption (Supplementary Note 1). Instead, the discovery of widespread clonal expansions in ageing human tissues<sup>19,50–52</sup> raises the possibility that some somatic mutations contribute to ageing by driving clonal expansions of functionally altered cells at a cost to the organism<sup>49,53,54</sup>. Examples include the possible links between clonal haematopoiesis and cardiovascular disease<sup>54</sup>; between mutations in liver disease and insulin resistance<sup>55</sup>; and between driver mutations in cavernomas and brain haemorrhages<sup>49,53,56</sup>. Detailed studies on the extent and effect of somatic mutations and clonal expansions on age-related diseases and ageing phenotypes may help to clarify the precise role—if any—of somatic mutations in ageing. Even if clear causal links between somatic mutations and ageing are established, ageing is likely to be multifactorial. Other forms of molecular damage involved in ageing could be expected to show similar anticorrelations with lifespan and, indeed, such anticorrelations have been reported for telomere shortening and protein turnover<sup>57,58</sup>.

Alternative non-causal explanations for the observed anticorrelation between somatic mutation rates and lifespan need to be considered. One alternative explanation is that cell division rates could scale with lifespan and explain the observed somatic mutation rates. Available estimates of cell division rates, although imperfect and limited to a few species, do not readily support this argument (Methods). More importantly, studies in humans have shown that cell division rates are not a major determinant of somatic mutation rates across human tissues<sup>14,18</sup>. Another alternative explanation for the observed anticorrelation might be that selection acts to reduce germline mutation rates in species with longer reproductive spans, which in turn causes an anticorrelation of somatic mutation rates and lifespan. Although selective pressure on

germline mutation rates could influence somatic mutation rates, it is unlikely that germline mutation rates tightly determine somatic mutation rates: somatic mutation rates in humans are 10–20 times higher than germline mutation rates, show variability across cell types and are influenced by additional mutational processes<sup>18,20</sup>. Overall, the strong scaling of somatic mutation rates with lifespan across mammals, despite the different rates between germline and soma and the variable contributions of different mutational processes across species, suggests that somatic mutation rates themselves have been evolutionarily constrained, possibly through selection on multiple DNA repair pathways. Alternative explanations need to be able to explain the strength of the scaling despite these differences.

Altogether, this study provides a detailed description of somatic mutation across mammals, identifying common and variable features and shedding light on long-standing hypotheses. Scaled across the tree of life and across tissues, in species with markedly different physiologies, life histories, genome compositions and mutagenic exposures, similar studies promise to transform our understanding of somatic mutation and its effects on evolution, ageing and disease.

## Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41586-022-04618-z>.

1. Garcia, A. M. et al. Age- and temperature-dependent somatic mutation accumulation in *Drosophila melanogaster*. *PLoS Genet.* **6**, e1000950 (2010).
2. Behjati, S. et al. Genome sequencing of normal cells reveals developmental lineages and mutational processes. *Nature* **513**, 422–425 (2014).
3. Milholland, B. et al. Differences between germline and somatic mutation rates in humans and mice. *Nat. Commun.* **8**, 15183 (2017).
4. Schmid-Siegert, E. et al. Low number of fixed somatic mutations in a long-lived oak tree. *Nat. Plants* **3**, 926–929 (2017).
5. Jager, M. et al. Deficiency of nucleotide excision repair is associated with mutational signature observed in cancer. *Genome Res.* **29**, 1067–1077 (2019).
6. Orr, A. J. et al. A phylogenomic approach reveals a low somatic mutation rate in a long-lived plant. *Proc. R. Soc. B* **287**, 20192364 (2020).
7. López, E. H. & Palumbi, S. R. Somatic mutations and genome stability maintenance in clonal coral colonies. *Mol. Biol. Evol.* **37**, 828–838 (2020).
8. Lee-Six, H. et al. The landscape of somatic mutation in normal colorectal epithelial cells. *Nature* **574**, 532–537 (2019).
9. Stratton, M. R., Campbell, P. J. & Futreal, P. A. The cancer genome. *Nature* **458**, 719–724 (2009).
10. Szilard, L. On the nature of the aging process. *Proc. Natl Acad. Sci. USA* **45**, 30–45 (1959).
11. Morley, A. A. The somatic mutation theory of ageing. *Mutat. Res.* **338**, 19–23 (1995).
12. Vijg, J. & Dong, X. Pathogenic mechanisms of somatic mutation and genome mosaicism in aging. *Cell* **182**, 12–23 (2020).
13. Welch, J. S. et al. The origin and evolution of mutations in acute myeloid leukemia. *Cell* **150**, 264–278 (2012).
14. Blokzijl, F. et al. Tissue-specific mutation accumulation in human adult stem cells during life. *Nature* **538**, 260–264 (2016).
15. Moore, L. et al. The mutational landscape of normal human endometrial epithelium. *Nature* **580**, 640–646 (2020).
16. Lodato, M. A. et al. Aging and neurodegeneration are associated with increased mutations in single human neurons. *Science* **359**, 555–559 (2018).
17. Zhang, L. et al. Single-cell whole-genome sequencing reveals the functional landscape of somatic mutations in B lymphocytes across the human lifespan. *Proc. Natl Acad. Sci. USA* **116**, 9014–9019 (2019).
18. Abascal, F. et al. Somatic mutation landscapes at single-molecule resolution. *Nature* **593**, 405–410 (2021).
19. Martincorena, I. et al. Tumor evolution. High burden and pervasive positive selection of somatic mutations in normal human skin. *Science* **348**, 880–886 (2015).
20. Moore, L. et al. The mutational landscape of human somatic and germline cells. *Nature* **597**, 381–386 (2021).
21. Peto, R., Roe, F. J., Lee, P. N., Levy, L. & Clack, J. Cancer and ageing in mice and men. *Br. J. Cancer* **32**, 411–426 (1975).
22. Vincze, O. et al. Cancer risk across mammals. *Nature* **601**, 263–267 (2022).
23. Peto, R. Epidemiology, multistage models, and short-term mutagenicity tests. *Int. J. Epidemiol.* **45**, 621–637 (2016).
24. Tollis, M., Boddy, A. M. & Maley, C. C. Peto's paradox: how has evolution solved the problem of cancer prevention? *BMC Biol.* **15**, 60 (2017).
25. López-Otin, C., Blasco, M. A., Partridge, L., Serrano, M. & Kroemer, G. The hallmarks of aging. *Cell* **153**, 1194–1217 (2013).

26. Schumacher, B., Pothof, J., Vijg, J. & Hoeijmakers, J. H. J. The central role of DNA damage in the ageing process. *Nature* **592**, 695–703 (2021).
27. Burnet, M. *Intrinsic Mutagenesis: a Genetic Approach to Ageing* (Springer, 1974).
28. Kirkwood, T. B. & Holliday, R. The evolution of ageing and longevity. *Proc. R. Soc. B* **205**, 531–546 (1979).
29. Ellis, P. et al. Reliable detection of somatic mutations in solid tissues by laser-capture microdissection and low-input DNA sequencing. *Nat. Protoc.* **16**, 841–871 (2021).
30. Alexandrov, L. B. et al. The repertoire of mutational signatures in human cancer. *Nature* **578**, 94–101 (2020).
31. Gori, K. & Baez-Ortega, A. sigfit: flexible Bayesian inference of mutational signatures. Preprint at *bioRxiv* <https://doi.org/10.1101/372896> (2020).
32. Lindahl, T. & Nyberg, B. Heat-induced deamination of cytosine residues in deoxyribonucleic acid. *Biochemistry* **13**, 3405–3410 (1974).
33. Zou, X. et al. A systematic CRISPR screen defines mutational mechanisms underpinning signatures caused by replication errors and endogenous DNA damage. *Nat. Cancer* **2**, 643–657 (2021).
34. Wilson, M. R. et al. The human gut bacterial genotoxin colibactin alkylates DNA. *Science* **363**, eaar7785 (2019).
35. Pleguezuelos-Manzano, C. et al. Mutational signature in colorectal cancer caused by genotoxic *pkcs* *E. coli*. *Nature* **580**, 269–273 (2020).
36. Smati, M. et al. Quantitative analysis of commensal *Escherichia coli* populations reveals host-specific enterotypes at the intra-species level. *MicrobiologyOpen* **4**, 604–615 (2015).
37. Oliero, M. et al. Oligosaccharides increase the genotoxic effect of colibactin produced by *pkcs* *Escherichia coli* strains. *BMC Cancer* **21**, 172 (2021).
38. Martincorena, I. et al. Universal patterns of selection in cancer and somatic tissues. *Cell* **171**, 1029–1041 (2017).
39. Moorad, J. A., Promislow, D. E. L., Flesness, N. & Miller, R. A. A comparative assessment of univariate longevity measures using zoological animal records. *Aging Cell* **11**, 940–948 (2012).
40. Speakman, J. R. Correlations between physiology and lifespan—two widely ignored problems with comparative studies. *Aging Cell* **4**, 167–175 (2005).
41. de Magalhães, J. P., Costa, J. & Church, G. M. An analysis of the relationship between metabolism, developmental schedules, and longevity using phylogenetic independent contrasts. *J. Gerontol. A* **62**, 149–160 (2007).
42. Vazquez, J. M. & Lynch, V. J. Pervasive duplication of tumor suppressors in Afrotherians during the evolution of large bodies and reduced cancer risk. *eLife* **10**, e65041 (2021).
43. Smith, E. S. J., Schuhmacher, L.-N. & Husson, Z. The naked mole-rat as an animal model in biomedical research: current perspectives. *Open Access Anim. Physiol.* **7**, 137–148 (2015).
44. Millar, J. S. & Zammuto, R. M. Life histories of mammals: an analysis of life tables. *Ecology* **64**, 631–635 (1983).
45. Kauppi, T. E. S., Kauppi, J. H. K. & Larsson, N.-G. Mammalian mitochondria and aging: an update. *Cell Metab.* **25**, 57–71 (2017).
46. Ju, Y. S. et al. Origins and functional consequences of somatic mitochondrial DNA mutations in human cancer. *eLife* **3**, e02935 (2014).
47. Hall, K. Y., Hart, R. W., Benirschke, A. K. & Walford, R. L. Correlation between ultraviolet-induced DNA repair in primate lymphocytes and fibroblasts and species maximum achievable life span. *Mech. Ageing Dev.* **24**, 163–173 (1984).
48. Zhang, L. et al. Maintenance of genome sequence integrity in long- and short-lived rodent species. *Sci. Adv.* **7**, eabj3284 (2021).
49. Smith, J. M. Review lectures on senescence—I. The causes of ageing. *Proc. R. Soc. B* **157**, 115–127 (1962).
50. Jaiswal, S. et al. Age-related clonal hematopoiesis associated with adverse outcomes. *N. Engl. J. Med.* **371**, 2488–2498 (2014).
51. Martincorena, I. et al. Somatic mutant clones colonize the human esophagus with age. *Science* **362**, 911–917 (2018).
52. Yokoyama, A. et al. Age-related remodelling of oesophageal epithelia by mutated cancer drivers. *Nature* **565**, 312–317 (2019).
53. Ren, A. A. et al. *PIK3CA* and CCM mutations fuel cavernomas through a cancer-like mechanism. *Nature* **594**, 271–276 (2021).
54. Jaiswal, S. & Libby, P. Clonal haematopoiesis: connecting ageing and inflammation in cardiovascular disease. *Nat. Rev. Cardiol.* **17**, 137–144 (2020).
55. Ng, S. W. K. et al. Convergent somatic mutations in metabolism genes in chronic liver disease. *Nature* **598**, 473–478 (2021).
56. Jaiswal, S. & Ebert, B. L. Clonal hematopoiesis in human aging and disease. *Science* **366**, eaan4673 (2019).
57. Swovick, K. et al. Interspecies differences in proteome turnover kinetics are correlated with life spans and energetic demands. *Mol. Cell. Proteomics* **20**, 100041 (2021).
58. Whittemore, K., Vera, E., Martínez-Navado, E., Sanpera, C. & Blasco, M. A. Telomere shortening rate predicts species life span. *Proc. Natl Acad. Sci. USA* **116**, 15122–15127 (2019).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

## Methods

### Data reporting

No statistical methods were used to predetermine sample size. The experiments were not randomized and the investigators were not blinded to allocation during experiments and outcome assessment.

### Ethics statement

All animal samples were obtained with the approval of the local ethical review committee (AWERB) at the Wellcome Sanger Institute and those at the holding institutions.

### Sample collection

We obtained colorectal epithelium and skin samples from a range of sources (Supplementary Table 1). For comparability across species an approximately 1-cm biopsy of the colorectal epithelium was taken from the terminal colon during necropsy. All necropsies occurred as soon as possible post-mortem to minimize tissue and DNA degradation. Tissue samples taken later than 24 h post-mortem typically showed extensive degradation of the colorectal epithelium, making the identification of colorectal crypts challenging. These samples were also associated with poor DNA yields and so were not included in the study. Sampled tissue was fixed in PAXgene FIX (PreAnalytiX, Switzerland), a commercially available fixative, during the necropsy. After 24 h in the fixative at room temperature, samples were transferred into the PAXgene STABILIZER and stored at  $-20^{\circ}\text{C}$  until further processing.

### Sample processing

Samples were processed using a workflow designed for detection of somatic mutations in solid tissues by laser-capture microdissection (LCM) using low-input DNA sequencing. For a more detailed description see the paraffin workflow described in another study<sup>29</sup>. In brief, PAXgene-fixed tissue samples of the colorectal epithelium were paraffin-embedded using a Sakura Tissue-Tek VIP tissue processor. Sections of  $16\ \mu\text{m}$  were cut using a microtome, mounted on PEN-membrane slides and stained with Gill's haematoxylin and eosin by sequential immersion in the following: xylene (two minutes, twice), ethanol (100%, 1 min, twice), deionized water (1 min, once), Gill's haematoxylin (10 s, once), tap water (20 s, twice), eosin (5 s, once), tap water (20 s, once), ethanol (70%, 20 s, twice) and xylene or Neo-Clear, a xylene substitute (20 s, twice).

High-resolution scans were obtained from representative sections of each species. Example images are shown in Fig. 1a, Extended Data Fig. 2. Individual colorectal crypts were isolated from sections on polyethylene naphthalate (PEN) membrane slides by LCM with a Leica LMD7 microscope. Haematoxylin and eosin histology images were reviewed by a veterinary pathologist. For some samples we also cut a section of muscle tissue from below the colorectal epithelium of the section to use as a germline control for variant calling (Supplementary Table 2). Pre- and post-microdissection images of the tissue were recorded for each crypt and muscle sample taken. Each microdissection was collected in a separate well of a 96-well plate.

Crypts were lysed using the Arcturus PicoPure Kit (Applied Biosystems) as previously described<sup>8,29</sup>. Each crypt then underwent DNA library preparation, without a quantification step to avoid loss of DNA, following a protocol described previously<sup>29</sup>. For some animals, a PAXgene fixed bulk skin biopsy was used as the germline control. For these skin samples, DNA was extracted using the DNeasy Blood & Tissue Kit (Qiagen).

### Library preparation and sequencing

Libraries from microdissected samples were prepared using enzymatic fragmentation, adapter ligation and whole-genome sequencing following a method described previously<sup>29</sup>. Libraries from skin samples were prepared using standard Illumina whole-genome library preparation.

Samples were multiplexed and sequenced using Illumina XTEN and Novaseq 6000 machines to generate 150 base pair (bp) paired-end reads. Samples were sequenced to around  $30\times$  depth (Supplementary Table 2).

### Sequence read alignment

For each species, sequences were aligned to a reference assembly (Supplementary Table 2) using the BWA-MEM algorithm<sup>59</sup> as implemented in BWA v.0.7.17-r1188, with options '-T 30 -Y -p -t 8'. The aligned reads were sorted using the bamsort tool from the biobambam2 package (v.2.0.86; gitlab.com/german.tischler/biobambam2), with options 'fixmates=1 level=1 calmdnm=1 calmdnmrecompindetonly=1 calmdnmreference=<reference\_fasta> outputthreads=7 sortthreads=7'. Duplicate reads were marked using the bammarkduplicates2 tool from biobambam2, with option 'level = 0'.

### Variant calling

Identification of somatic substitutions and short indels was divided into two steps: variant calling, and variant filtering to remove spurious calls (see 'Variant filtering' below). For human colorectal crypts, we obtained previously sequenced and mapped reads from a study in which colorectal crypts were isolated by LCM<sup>8</sup>, and processed them using the sample variant calling and filtering process that was applied to the non-human samples.

Substitutions were identified using the cancer variants through expectation maximization (CaVEMan) algorithm<sup>60</sup> (v.1.13.15). CaVEMan uses a naive Bayesian classifier to perform a comparative analysis of the sequence data from a target and control sample from the same individual to derive a probabilistic estimate for putative somatic substitutions at each site. The copy number options were set to 'major copy number = 5' and 'minor copy number = 2', as in our experience this maximizes the sensitivity to detect substitutions in normal tissues. CaVEMan identifies and excludes germline variants shared in the target (colorectal crypt) and matched normal (skin or muscle tissue) samples, and produces a list of putative somatic mutations that are present only in the target sample. CaVEMan was run separately for each colorectal crypt, using either bulk skin or muscle microdissected from the sample colorectal biopsy as the matched normal control (Supplementary Table 2). For two human donors for whom an alternative tissue was not available, a colonic crypt not included as a target sample was used as the matched normal control.

Indels were identified using the Pindel algorithm<sup>61</sup> (v.3.3.0), using a second sample from the same individual as a matched control. The indel calls produced by Pindel were subsequently re-genotyped using the vafCorrect tool (<https://github.com/cancerit/vafCorrect>), which performs a local sequence assembly to address alignment errors for indels located at the end of sequence reads, and produces corrected counts of sequence reads supporting the indel and corrected estimates of variant allele fraction (VAF; the fraction of reads supporting the alternate allele at the variant site).

### Variant filtering

A number of post-processing filters were applied to the variant calls to remove false positives (Supplementary Fig. 1a, b).

**Quality flag filter.** CaVEMan and Pindel annotate variant calls using a series of quality flags, with the 'PASS' flag denoting no quality issues affecting the call<sup>60,61</sup>. Variant calls presenting any flag other than 'PASS' were discarded.

**Alignment quality filter.** Variants were excluded if more than half of the supporting reads were clipped. The library preparation methods create short insert size libraries that can result in reads overlapping. To avoid the risk of double counting mutant reads we used fragment-based statistics. Variants without at least four high-quality fragments (alignment

# Article

score  $\geq 40$  and base Phred quality score  $\geq 30$ ) were excluded. Variants were excluded if reads supporting the variant had a secondary alignment score that was greater than the primary alignment score. This filter was not applied to indel calls.

**Hairpin filter.** To remove variants introduced by erroneous processing of cruciform DNA during the enzymatic digestion, we applied a custom filter to remove variants in inverted repeats<sup>29</sup>. This filter was not applied to indel calls.

**Chromosome and contig filter.** For species with chromosome-level assemblies, we discarded variants located in non-chromosomal contigs, including the mitochondrial genome (calling of mitochondrial variants is described in the section ‘Mitochondrial variant calling and filtering’). For males, variants on the Y chromosome were excluded for species in which the Y chromosome was annotated in the assembly.

**N-tract and contig-end filter.** To reduce artefactual calls due to read misalignment, we discarded variants located within 1 kb of a tract of 50 or more consecutive N bases in the reference assembly, as well as variants within 1 kb of the start or end of a contig (this implies discarding all variants in contigs shorter than 2 kb).

**Sequencing coverage filter.** A sample-specific read depth filter was designed to exclude sites with coverage above the 99th coverage percentile in the sample or its matched normal control, as well as sites with coverage of less than  $10\times$  in the sample or its matched normal control.

**Allelic strand bias filter.** We discarded variants without any supporting reads on either the forward or the reverse strand.

**Indel proximity filter.** We discarded variants for which the total number of reads supporting the presence of an indel within 10 bp of the variant was more than three times larger than the number of reads supporting the presence of the variant. This filter was not applied to indel calls.

**Spatial clustering filter.** Visual assessment of variant calls and mutational spectra showed spatially clustered variants to be highly enriched for artefacts. Therefore, we discarded groups of two or more variants located within 1 kb of each other.

**Beta-binomial filter.** For each crypt, an artefact filter based on the beta-binomial distribution was applied, which exploits read count information in other crypts from the same individual. More specifically, for each sample, we fitted a beta-binomial distribution to the variant allele counts and sequencing depths of somatic variants across samples from the same individual. The beta-binomial distribution was used to determine whether read support for a mutation varies across samples from an individual, as expected for genuine somatic mutations but not for artefacts. Artefacts tend to be randomly distributed across samples and can be modelled as drawn from a binomial or a lowly overdispersed beta-binomial distribution. True somatic variants will be present at a high VAF in some samples, but absent in others, and are hence best captured by a highly overdispersed beta-binomial. For each variant site, the maximum likelihood estimate of the overdispersion factor ( $\rho$ ) was calculated using a grid-based method, with values ranging between  $10^{-6}$  and  $10^{-0.05}$ . Variants with  $\rho > 0.3$  were considered to be artefactual and discarded. The code for this filter is based on the Shearwater variant caller<sup>62</sup>. We found this to be one of the most effective filters against spurious calls (Supplementary Fig. 1b).

**Minimum VAF filter.** For each sample, we discarded variants for which the VAF was less than half the median VAF of variants passing the beta-binomial filter (see above) in that sample.

**Maximum indel VAF filter.** For each sample, we discarded indels that presented a VAF of greater than 0.9, as such indels were found to be highly enriched in spurious calls in some species. This filter was not applied to substitution calls.

To validate our variant calling strategy, we used LCM to microdissect two sections from the same mouse colorectal crypt. We expected to detect a high fraction of shared somatic variants in these two sections, as their cells should be derived from the same ancestral epithelial stem cell. Both sections were submitted for independent library preparation, genome sequencing, variant calling and filtering using our pipeline. The majority of substitution variant calls (2,742 of 2,933, 93.5%) were shared between both sections (Supplementary Fig. 1c). By contrast, when comparing five separate crypts from a mouse, a maximum of two variants were shared between two crypts, and no variants were shared by three or more crypts (Supplementary Fig. 1d).

## Sample filtering

Our method for estimation of mutation rates assumes monoclonality of colorectal crypt samples. This assumption can be violated owing to several causes, including contamination from other colorectal crypts during microdissection or library preparation, contamination with non-epithelial cells located in or near the crypt, insufficient time for a stem cell to drift to clonality within the crypt, or the possibility that in some species, unlike in humans<sup>8</sup>, polyclonal crypts are the norm. Therefore, a truncated binomial mixture model was applied so as to remove crypts that showed evidence of polyclonality, or for which the possibility of polyclonality could not be excluded. An expectation–maximization (EM) algorithm was used to determine the optimal number of VAF clusters within each crypt sample, as well as each cluster’s location and relative contribution to the overall VAF distribution. The algorithm considered a range of numbers of clusters (1–5), with the optimal number being that which minimized the Bayesian information criterion (BIC). As the minimum number of supporting reads to call a variant was four, the binomial probability distribution was truncated to incorporate this minimum requirement for the number of successes, and subsequently re-normalized. The EM algorithm returned the inferred optimal number of clusters, the mean VAF (location) and mixing proportion (contribution) of each clone, and an assignment of each input variant to the most likely cluster. After applying this model to the somatic substitutions identified in each sample, sample filtering was performed on the basis of the following three criteria.

**Low mutation burden.** We discarded samples that presented fewer than 50 somatic variants, which was indicative of low DNA quality or sequencing issues.

**High mutation burden.** We discarded samples with a number of somatic variants greater than 3 times the median burden of samples from the same individual (excluding samples with fewer than 50 variants). This served to exclude a small minority of samples that presented evident sequencing quality problems (such as low sequencing coverage), but which did not fulfil the low-VAF criterion for exclusion (see below).

**Low VAF.** We discarded samples in which less than 70% of the somatic variants were assigned to clusters with VAF  $\geq 0.3$ . However, this rule was not applied to those cases in which all the samples from the same individual had primary clusters with mean VAF  $< 0.3$ ; this was done to prevent the removal of samples from individuals presenting high fractions of non-epithelial cells, but whose crypts were nonetheless dominated by a single clone.

These criteria led to the exclusion of 41 out of 249 samples. On the basis of visual assessment of sequencing coverage and VAF distributions, we decided to preserve three samples (ND0003c\_lo0004,

ND0003c\_lo0011, TIGRD0001b\_lo0010) that we considered to be clonal, but which would have been discarded on the basis of the criteria above.

### Mitochondrial variant calling and filtering

For six species whose reference genome assemblies did not include the mitochondrial sequence, mitochondrial reference sequences were obtained from the GenBank database (Supplementary Table 5). For each species, alignment to the reference genome was performed using BWA (v.0.7.17-r1188), as described above (see ‘Sequence read alignment’). Pileup files were generated for mtDNA genomes using the ‘bam2R’ function in the deepSNV (v.1.32.0) R package<sup>62,63</sup>. The mapping quality cut-off was set to 0, taking advantage of the fact that the mitochondrial genome coverage for most samples was more than 100-fold higher than the nuclear genome coverage, and hence most reads with poor mapping scores should be of mitochondrial origin. Mitochondrial variants were called using the Shearwater algorithm<sup>62</sup> (deepSNV package v.1.32.0). Multiple rounds of filtering were applied to identify and remove false positives. The first set of filters removed germline polymorphisms, applied a maximum false discovery rate (FDR) threshold of  $q > 0.01$ , required that mismatches should be supported by at least one read on both the forward and reverse strands, and merged consecutive indel calls. Further filtering steps were as follows.

**Minimum VAF filter.** Only variants with VAF  $> 0.01$  were considered for analysis, based on the quality of the mutational spectra.

**Sequencing coverage filter.** Owing to species-specific mtDNA regions of poor mappability, we discarded sites with a read coverage of less than  $500\times$ .

**D-loop filter.** Analysis of the distribution of mutations along the mitochondrial genome revealed clusters of mutations within the hypervariable region of mtDNA known as the D-loop. To obtain estimates of the mutation burden in mtDNA unaffected by hypermutation of the D-loop, mutations in the D-loop region (coordinates MT:1–576 and MT:16,024–16,569 in human) were excluded from this analysis.

**High mutation burden.** We discarded samples that had a number of somatic mtDNA variants greater than four times the mean mtDNA burden across all samples. This served to exclude a small minority of samples that were suspected of enrichment in false positive calls. Visual inspection of these samples in a genome browser confirmed the presence of high numbers of variants found on sequence reads with identical start positions and/or multiple base mismatches, suggestive of library preparation or sequencing artefacts.

We examined the mutational spectra of somatic mtDNA substitutions on a trinucleotide sequence context (Extended Data Fig. 14b). The specificity of the filtered variant calls was supported by the observation that the mutational spectra across species were broadly consistent with those previously observed in studies of human tissues<sup>46</sup>, with a dominance of C>T and T>C transversions and a strong replication strand bias.

### Mitochondrial copy number analysis

Sequence reads from each sample were separately mapped to the species-specific mtDNA reference sequence to estimate average mtDNA sequencing coverage. Excluding nuclear reference sequences from the alignment enabled even coverage to be obtained across the mitochondrial genome by preventing the mismapping of sequence reads to inherited nuclear insertions of mitochondrial DNA (known as NuMTs). Next, coverage information from individual mtDNA and whole-genome alignment (BAM) files was obtained using the genomcov tool in the bedtools suite (v.2.17.0)<sup>64</sup>. Mitochondrial copy number was calculated according to the formula

$$\text{depth}_{\text{mtDNA}} \times \text{ploidy} / \text{depth}_{\text{gDNA}}$$

where  $\text{depth}_{\text{mtDNA}}$  and  $\text{depth}_{\text{gDNA}}$  are the mean coverage values for mtDNA and the nuclear genome, respectively, and  $\text{ploidy} = 2$  (assuming normal somatic cells to be diploid). For simplicity, the sex chromosomes were excluded from the calculation of the mean nuclear genome coverage.

### Calculation of analysable genome size

To estimate the somatic mutation rate, it was first necessary to establish the size of the analysable nuclear genome (that is, the portion of the genome in which variant calling could be performed reliably) for each sample (Supplementary Table 4). For both substitutions and indels, the analysable genome of a sample was defined as the complement of the union of the following genomic regions: regions reported as ‘not analysed’ by the CaVEMan variant caller; regions failing the ‘chromosome and contig’ filter; regions failing the ‘N-tract and contig-end’ filter; and regions failing the ‘sequencing coverage’ filter (see ‘Variant filtering’). For the analysis of mitochondrial variants, the analysable genome of a sample was defined as the portion of mtDNA that satisfied the ‘sequencing coverage’ filter (see ‘Mitochondrial variant calling and filtering’), after subtracting the hypervariable region (D-loop).

### Life-history data

Obtaining accurate lifespan estimates is challenging; although point estimates of maximum lifespan are available for many species, their veracity is often difficult to assess and estimates can vary widely for the same species (Supplementary Table 6). There can be many causes for this variation, including errors in recording and real variation in longevity between populations (that is, captive versus wild). As we were interested in whether the somatic mutation burden has an association with lifespan in the absence of extrinsic mortality, we sought to obtain estimates of longevity from individuals under human care, to minimize the effect of external factors such as predation or infection.

Mortality records for 14 species were obtained from the Species360 database, authorized by Species360 research data use agreement no. 60633 (Species360 Zoological Information Management System (ZIMS) (2020), <https://zims.species360.org>). This database contains lifespan data of zoo animals from international zoo records. Using records from 1980 to the present, we excluded animals for which the date of birth or death was unknown or uncertain. To avoid infant mortality influencing the longevity estimates for each species, we removed animals that died before the age of female sexual maturity, as defined by the AnAge database<sup>65</sup>. This resulted in a mean of 2,681 animal lifespan records per species for the species in the study (minimum 309, maximum 8,403; Supplementary Table 6). For the domestic dog, we combined records for domestic dogs (*Canis lupus familiaris*) and wolves (*Canis lupus*), because of the paucity of records for domestic dogs in Species360. Although the data are curated, they are still vulnerable to the presence of inaccurate records, which can bias the lifespan estimates. To reduce the effect of these outliers, for each species lifespan was estimated as the age at which 80% of the adults from that species had died<sup>66</sup> (Supplementary Table 6).

Human longevity estimates were obtained using census birth and death record data from Denmark, (1900–2020), Finland (1900–2019) and France (1900–2018), retrieved from the Human Mortality Database (University of California, Berkeley (USA), and Max Planck Institute for Demographic Research (Germany); <https://www.mortality.org>, <https://www.humanmortality.de>). We selected these countries because they had census records going back at least 100 years. To remove the effect of infant mortality, we excluded individuals who died before the age of 13. For each country, we selected the cohort born in 1900 and calculated the age at which 80% of the individuals had died (Denmark, 87 years; Finland, 83 years; France, 81 years). We then used the mean

of the three countries as our estimate of the human 80% lifespan (83.7 years) (Supplementary Table 6).

To test the effects of different estimates of lifespan on our results, we also obtained maximum longevity estimates for each species from a range of databases<sup>67</sup> and a survey of the literature (Supplementary Table 6). Other life-history metrics were obtained from the AnAge database<sup>65</sup> (Supplementary Table 6).

### Mutational signature analysis

Mutational signatures of substitutions on a trinucleotide sequence context were inferred from sets of somatic mutation counts using the sigfit (v.2.1.0) R package<sup>31</sup>. Initially, signature extraction was performed de novo for a range of numbers of signatures ( $N = 2, \dots, 10$ ), using counts of mutations grouped per sample, per individual and per species. To account for differences in sequence composition across samples, and especially across species, mutational opportunities per sample, per individual and per species were calculated from the reference trinucleotide frequencies across the analysable genome of each sample (see 'Calculation of analysable genome size'), and supplied to the 'extract\_signatures' function in sigfit. The 'convert\_signatures' function in sigfit was subsequently used to transform the extracted signatures to a human-relative representation (Fig. 2b), by scaling the mutation probability values using the corresponding human genome trinucleotide frequencies. The best-supported number of signatures, on the basis of overall goodness-of-fit<sup>31</sup> and consistency with known COSMIC signatures (<https://cancer.sanger.ac.uk/signatures/>), was found to be  $N = 3$ . The cleanest deconvolution of the three signatures was achieved when using the mutation counts grouped by species, rather than by sample or individual. The three extracted signatures (labelled SBSA, SBSB and SBSC) were found to be highly similar to COSMIC signatures SBS1 (cosine similarity 0.96), SBS5 (0.89) and SBS18 (0.91), respectively. These signatures were independently validated using the Mutational Patterns (v.1.12.0) R package<sup>68</sup>, which produced comparable signatures (respective cosine similarities 0.999, 0.98 and 0.89).

This de novo signature extraction approach, however, failed to deconvolute signatures SBSA and SBSB entirely from each other, resulting in a general overestimation of the exposure to SBSA (Extended Data Fig. 15). To obtain more accurate estimates of signature exposure, the deconvolution was repeated using an alternative approach that combines signature fitting and extraction in a single inference process<sup>31</sup>. More specifically, the 'fit\_extract\_signatures' function in sigfit was used to fit COSMIC signature SBS1 (retrieved from the COSMIC v.3.0 signature catalogue; <https://cancer.sanger.ac.uk/signatures/>) to the mutation counts grouped by species (with species-specific mutational opportunities), while simultaneously extracting two additional signatures de novo (SBSB and SBSC). Before this operation, COSMIC SBS1 was transformed from its human-relative representation to a genome-independent representation using the 'convert\_signatures' function in sigfit. By completely deconvoluting SBS1 and SBSB, this approach yielded a version of SBSB that was more similar to COSMIC SBS5 (cosine similarity 0.93); the similarity of SBSC to COSMIC SBS18 was the same under both approaches (0.91).

Finally, the inferred signatures were re-fitted to the mutational spectra of mutations in each sample (using the 'fit\_signatures' function in sigfit with sample-specific mutational opportunities) to estimate the exposure of each sample to each signature. The fitting of the three signatures yielded spectrum reconstruction similarity values (measured as the cosine similarity between the observed mutational spectrum and a spectrum reconstructed from the inferred signatures and exposures) with median 0.98 and interquartile range 0.96–0.99. Although the purely de novo extraction approach and the 'fitting and extraction' approach yielded comparable versions of signatures SBSB and SBSC, the fixing of COSMIC SBS1 in the latter approach resulted in lower SBS1 exposures and higher SBSB exposures in most samples, owing to the cleaner deconvolution of these two signatures (Fig. 2, Extended Data Fig. 15).

To examine potential variation in the spectrum of signature SBS5 across species, the following procedure was conducted for each species: individual-specific mutation counts and mutational opportunities were calculated for each individual in the species, and the 'fit\_extract\_signatures' function was used to fit COSMIC signatures SBS1, SBS18 and SBS34 (transformed to a genome-independent representation using the 'convert\_signatures' function) to the mutational spectra of each individual, while simultaneously inferring one additional signature (corresponding to signature SBS5 as manifested in that species; Extended Data Fig. 6).

To assess the presence in non-human colorectal crypts of mutational signatures caused by APOBEC or colibactin, which have been previously observed in human crypts<sup>8</sup>, we used an expectation–maximization algorithm for signature fitting, in combination with likelihood ratio tests (LRTs). More specifically, for each non-human sample, we tested for exposure to colibactin (signature SBS88, COSMIC v.3.2) by comparing the log-likelihoods of (i) a model fitting COSMIC signatures SBS1, SBS5, SBS18, SBS34 and SBS88, and (ii) a reduced model fitting only the first four signatures. Benjamini–Hochberg multiple-testing correction was applied to the  $P$  values that resulted from the LRTs, and colibactin exposure was considered significant in a sample if the corresponding corrected  $q$ -value was less than 0.05. We followed the same approach to assess exposure to APOBEC (SBS2 and SBS13), using two separate sets of LRTs for models including either SBS2 or SBS13, in addition to SBS1, SBS5, SBS18 and SBS34. APOBEC exposure was considered significant in a sample if its  $q$ -values for the models including SBS2 and SBS13 were both less than 0.05. This analysis identified 1/180 crypts with significant exposure to each of colibactin and APOBEC (although the evidence for the presence of the relevant signatures in these two crypts was not conclusive). To test for depletion of colibactin or APOBEC exposure in non-human crypts relative to human crypts, we first applied the LRT-based method described above to a published set of 445 human colorectal crypts<sup>8</sup>, identifying 92 colibactin-positive and 9 APOBEC-positive crypts. We then compared the numbers of colibactin- and APOBEC-positive crypts in the human and non-human sets using two separate Fisher's exact tests ('fisher.test' function in R). This revealed the difference in colibactin exposure to be highly significant ( $P = 7 \times 10^{-14}$ ), unlike the difference in APOBEC exposure ( $P = 0.30$ ).

Mutational spectra of somatic indels identified in each species were generated using the 'indel.spectrum' function in the Indelwald tool for R (24/09/2021 version; <https://github.com/MaximilianStammnitz/Indelwald>).

### Selection analysis

Evidence of selection was assessed using the ratio of nonsynonymous to synonymous substitution rates (dN/dS) in the somatic mutations called in each species. The dNdScv (v.0.0.1.0) R package<sup>38</sup> was used to estimate dN/dS ratios for missense and truncating substitutions in each species separately. Reference CDS databases for the dNdScv package were built for those species with available genome annotation in Ensembl (<https://www.ensembl.org>; Supplementary Table 2), using the 'buildref' function in dNdScv. For each species, the 'dndscv' function was applied to the list of somatic substitutions called in samples of that species, after de-duplicating any substitutions that were shared between samples from the same individual to avoid counting shared somatic mutations multiple times. In addition, the analysis was restricted to genes that were fully contained in the analysable genomes of all samples from the species (a condition satisfied by the vast majority of protein-coding genes). Genome-wide and gene-specific dN/dS ratios were obtained for missense and truncating substitutions in each species; no genes with statistically significant  $dN/dS \neq 1$  were observed.

### Copy number analysis

For species with chromosome-level assemblies (cat, cow, dog, horse, human, mouse, rabbit and rat), the total and the allele-specific copy

number (CN) was assessed in each sample, adapting a likelihood model that was previously applied to the detection of subclonal CN changes in healthy human skin<sup>19</sup>. This method exploits two sources of evidence: relative sequencing coverage and B-allele fraction (BAF; the fraction of reads covering a heterozygous single-nucleotide polymorphism (SNP) that support one of the alleles). Human samples PD36813x15 and PD36813x16 were excluded from this analysis owing to the poor quality of their SNP data.

For each sample, sequencing coverage was measured in non-overlapping 100-kb bins along the reference genome of the species, using the coverageBed tool in the bedtools suite (v.2.17.0)<sup>64</sup>. For each bin, the coverage per base pair was calculated by dividing the number of reads mapping to the bin by the bin length, and multiplying the result by the read length (150 bp). A normalized sample-normal coverage ratio was then calculated for each bin by dividing the bin coverage in the sample by the corresponding coverage in its matched normal control (see 'Sample processing'). Heterozygous SNPs were isolated for each sample by selecting germline SNPs with a BAF between 0.4 and 0.6 in the matched normal sample, and a coverage of at least 15 reads in both the target sample and its matched normal sample. After assigning each SNP to its corresponding 100-kb genome bin, the bins in each sample were divided into two sets: (i) bins with coverage  $\geq 10$  in both the target sample and its matched normal, and at least one heterozygous SNP; and (ii) bins with coverage  $\geq 10$  in both the target sample and its matched normal, and no heterozygous SNPs. For the first set, estimates of total and allele-specific CN were inferred by maximizing the joint likelihood of a beta-binomial model for BAF and a negative binomial model for relative coverage, as previously described<sup>19</sup>. The most likely combination of allele CN values was obtained for each bin by conducting an exhaustive search of CN values between 0 and 4, and selecting the combination maximizing the joint likelihood (calculated on the basis of expected BAF and relative coverage values). A penalty matrix was used to penalize more complex solutions over simpler ones, as previously described<sup>19</sup>. For the second set of bins (bins without SNPs), only estimates of total CN were inferred, by maximizing the likelihood of a negative binomial model for relative coverage. The most substantial differences between these methods and the one previously published are: (i) SNPs were obtained from the variant calling output, instead of from a public database; (ii) relative coverage was calculated per 100-kb bin, rather than per SNP; (iii) SNPs were not phased within each gene, but within each bin; (iv) no reference bias was assumed (that is, the underlying BAF of heterozygous SNPs was assumed to be 0.5); (v) the minimum sample purity was raised to 0.85; (vi) putative CN changes were not subjected to significance testing, but selected according to their likelihood, and subsequently filtered by means of a segmentation algorithm (see below).

Estimates of total and allele-specific CN per bin were merged into CN segments, which were defined as contiguous segments composed of five or more bins with identical CN states. Segmentation was performed separately for total and allele-specific CN estimates in each sample. After this process, any pair of adjacent segments with the same CN assignment, and separated by a distance shorter than five bins, was merged into a single segment. Finally, within each species, segments presenting CN values other than 2 (or 1/1 for allele-specific CN), and being either shorter than 10 bins (1 Mb), or shared among two or more samples, were discarded, resulting in the removal of nearly all spurious CN changes.

### Estimation of mutation rate

For each sample, the somatic mutation density (mutations per bp) was calculated by dividing the somatic mutation burden (total number of mutations called) by the analysable genome size for the sample (see 'Calculation of analysable genome size'). The adjusted somatic mutation burden (number of mutations per whole genome) was then calculated by multiplying the mutation density by the total genome size of the species (see below). The somatic mutation rate per year

(mutations per genome per year) was obtained by dividing this adjusted mutation burden by the age of the individual, expressed in years (Supplementary Table 2). The expected ELB for each sample was calculated by multiplying the somatic mutation rate by the estimated lifespan of the species (see 'Life-history data').

The total genome size of a species was estimated as the total size of its reference genome assembly. Across species, the mean genome size was 2.67 Gb, ranging between 2.41 Gb and 3.15 Gb and with a standard deviation of 221 Mb (Supplementary Table 4). This suggests that inter-species variation in genome size should not have a substantial influence on the somatic mutation rate estimates. For an assessment of alternative measures of mutation rate, see 'Association of mutation rate and end-of-lifespan burden with lifespan'.

### Association of mutation rate with life-history traits

The association of the somatic mutation rate with different life-history traits was assessed using LME models. In particular, associations with the following traits were examined: lifespan (in years), adult mass (or adult weight, in grams), BMR (in watts), and litter size (see 'Life-history data'). Associations for lifespan, adult mass and BMR were assessed using the following transformed variables:  $1/\text{lifespan}$ ,  $\log_{10}(\text{adult mass})$  and  $\log_{10}(\text{BMR})$ . To account for the potentially confounding effect of the correlation between metabolic rate and body mass, the residuals of a fitted allometric regression model of BMR on adult mass (equivalent to a simple linear regression of  $\log_{10}(\text{BMR})$  on  $\log_{10}(\text{adult mass})$ ) were used as a mass-adjusted measure of metabolic rate, referred to as 'BMR residuals'.

For each variable, an LME model was implemented for the regression of somatic mutation rates per sample on the variable of interest, using the 'lme' function in the nlme R package (v.3.1-137; <https://cran.r-project.org/web/packages/nlme>). To account for non-independence of the samples, both at the individual level and at the species level, the model included fixed effects (intercept and slope parameters) for the variable of interest, and random effects (slope parameters) at the individual and species levels. In addition, to account for the heteroscedasticity of mutation rate estimates across species, the usual assumption of constant response variance was replaced with explicit species-specific variances, to be estimated within the model.

To determine the fraction of inter-species variance in mutation rate explained by each life-history variable individually, the LME model described above was used to produce predictions of the mean mutation rate per species; only fixed effects were used when obtaining these predictions, random effects being ignored. The variance of these predictions was then compared to the variance in observed mean mutation rates; the latter were calculated for each species as the mean of the observed mean rates per individual, to avoid individuals with larger numbers of samples exerting a stronger influence on the species mean. The fraction of inter-species variance explained by the model was calculated using the standard formula for the coefficient of determination,

$$R^2 = \text{ESS}/(\text{ESS} + \text{RSS}),$$

where ESS is the explained sum of squares, and RSS is the residual sum of squares:

$$\text{ESS} = \sum_i (\hat{y}_i - \bar{y})^2, \text{RSS} = \sum_i (y_i - \hat{y}_i)^2.$$

In this formulation,  $y_i$  and  $\hat{y}_i$  denote the observed and predicted mutation rates for species  $i$ , respectively, and  $\bar{y}$  is the overall mean rate. This definition of  $R^2$  coincides with the fraction of variance explained (FVE), defined as 1 minus the fraction of variance unexplained (FVU):

$$\text{FVE} = 1 - \text{FVU} = 1 - [\text{RSS}/(\text{ESS} + \text{RSS})] = \text{ESS}/(\text{ESS} + \text{RSS}) = R^2.$$

As the predicted and observed values correspond to mean mutation rates per species, rather than mutation rates per sample, FVE provides a

measure of the fraction of inter-species variance explained by the fixed effects of the LME model. Among the variables considered, 1/lifespan was found to have the greatest explanatory power (FVE = 0.84, using a free-intercept model).

To compare the explanatory power of variables other than 1/lifespan when considered either individually or in combination with 1/lifespan, the method described above was also applied to two-variable combinations of 1/lifespan and each of the remaining variables, using an LME model with fixed effects for both variables and random effects for 1/lifespan only. The  $R^2$  formula above was used to measure the fraction of inter-species variance explained by each model. In addition, to test whether the inclusion of a second explanatory variable was justified by the increase in model fit, LRTs between each two-variable LME model and a reduced LME model including only 1/lifespan were performed using the 'anova' function in the nlme R package.

To further assess the potential effects of body mass and lifespan on each other's association with the somatic mutation rate, allometric regression models (equivalent to simple linear models under logarithmic transformation of both variables) were fitted to the mean somatic mutation rate per species, using either adult mass or lifespan as the explanatory variable. In addition, the 'allometric residuals' of mutation rate, adult mass and lifespan (that is, the residuals of pairwise allometric regressions among these three variables) were used to examine the associations between somatic mutation rate and either body mass or lifespan, after accounting for the effect of the third variable (partial correlation analysis). For example, to account for the potential influence of body mass on the relationship between somatic mutation rate and lifespan, the residuals of an allometric regression between mutation rate and adult mass, and the residuals of an allometric regression between lifespan and adult mass, were analysed using simple linear regression. This analysis supported a strong association between somatic mutation rate and lifespan (independently of the effect of mass; FVE = 0.82,  $P = 3.2 \times 10^{-6}$ ; Fig. 3c), and a non-significant association between somatic mutation rate and body mass (independently of the effect of lifespan). Therefore, the relationship between somatic mutation rate and lifespan does not appear to be mediated by the effect of body mass on both variables. Of note, this result remains after excluding naked mole-rat: after removing this species, partial correlation analysis still reveals a strong association between somatic mutation rate and lifespan (FVE = 0.77,  $P = 4.1 \times 10^{-5}$ ), and a non-significant association between somatic mutation rate and body mass ( $P = 0.84$ ). This demonstrates that the observed relationships are not dependent on the presence of naked mole-rat in the study.

To assess the robustness of the LME regression analyses described above, we performed bootstrap analysis on each LME model, at the level of both individuals and species. More specifically, for each level we used each of the LME models to perform regression on 10,000 bootstrap replicates, produced by resampling either species or individuals with replacement. We then assessed the distributions of FVE across bootstrap replicates (Extended Data Fig. 13c). In addition, we performed a similar bootstrap analysis using a collection of maximum longevity estimates obtained from the literature (see 'Life-history data'). We applied the zero-intercept LME model described above (for regressing mutation rate on inverse lifespan) on a set of 5,000 bootstrap replicates, each of which used a set of species lifespan estimates randomly sampled from the collection of literature-derived estimates (Extended Data Fig. 12).

The results obtained with the LME models were additionally validated using an independent hierarchical Bayesian model, in which the mean somatic mutation burden of each individual was modelled as following a normal distribution with mean defined as a linear predictor containing a species-specific slope parameter and a multiplicative offset (corresponding to the individual's age; inclusion of this offset minimizes the heteroscedasticity of the mutation rate across species, which results from dividing mutation burdens by age). Species-specific slope parameters were in turn modelled as normally distributed around a global slope

parameter, equivalent to the fixed-effect slope estimated by the LME model. This hierarchical model produced very similar results to those of the LME model for all life-history variables (Extended Data Fig. 13a).

We note that samples CATD0002b\_lo0003 and MD6267ab\_lo0003 were excluded from all regression analyses, owing to the fact that each shared the most of its somatic variants with another sample from the same individual (indicating, in each case, that both samples were closely related), hence violating the assumption of independence among samples. The inclusion of these two samples, however, had no effect on the outcome of the analyses.

## Association of mutation rate and end-of-lifespan burden with lifespan

The relationship between somatic mutation rate and species lifespan was further explored by adapting the LME model described in the previous section to perform constrained (zero-intercept) regression of the adjusted mutation rate per year on the inverse of lifespan, 1/lifespan (see 'Life-history data', 'Estimation of mutation rate' and 'Association of mutation rate with life-history traits'). The use of zero-intercept regression was motivated by the prediction that, if somatic mutation is a determinant of maximum lifespan, then it would be expected for all species to end their lifespans with a similar somatic mutation burden. Indeed, this was confirmed by simple linear regression of the species mean end-of-lifespan mutation burden against species lifespan (slope  $P = 0.39$ ). Thus, if  $m$  is the mutation rate per year, and  $L$  is the species' lifespan, the expected relationship is of the form.

$$m \approx kL,$$

where  $k$  is a constant representing the typical end-of-lifespan mutation burden across species. According to this relationship, the mutation rate per year is linearly related to the inverse of lifespan,

$$m \approx k(1/L).$$

Therefore, the cross-species average end-of-lifespan burden ( $k$ ), can be estimated as the slope parameter of a zero-intercept linear regression model with the mutation rate per year ( $m$ ) as the dependent variable, and the inverse of lifespan ( $1/L$ ) as the explanatory variable. To this purpose, the LME model described in the previous section was altered by removing the fixed-effect intercept parameter, thus considering only fixed- and random-effect slope parameters for 1/Lifespan.

The zero-intercept LME model estimated a value of  $k = 3,210.52$  (95% confidence interval 2,686.89–3,734.15). The fraction of inter-species variance explained by the zero-intercept model (FVE) was 0.82, whereas the LME model described in the previous section (which estimated  $k = 2,869.98$ , and an intercept of 14.76) achieved FVE = 0.84 (see 'Association of mutation rate with life-history traits'). To test whether the increase in model fit justifies the inclusion of an intercept, both models were compared using a LRT (as implemented by the 'anova' function in the nlme R package (v.3.1-137)). This yielded  $P = 0.23$ , indicating that the free-intercept model does not achieve a significantly better fit than the zero-intercept model. Similarly, the zero-intercept model yielded lower values for both the Bayesian information criterion (BIC) and the Akaike information criterion (AIC). Notably, equivalent analyses using somatic mutation rates per megabase and per protein-coding exome (instead of per whole genome) yielded comparable results (Extended Data Fig. 11).

To investigate the possibility of phylogenetic relationships between the species sampled confounding the analysis, a phylogenetic generalized linear model was used to regress the mean mutation rate of each species on the inverse of its lifespan ( $1/L$ ), while accounting for the phylogenetic relationships among species. A phylogenetic tree of the 15 species examined was obtained from the TimeTree resource<sup>69</sup>, and the phylogenetic linear model was fitted using the 'pgls' function

in the caper R package (v.1.0.1; <https://cran.r-project.org/web/packages/caper>). The estimates produced by zero-intercept regression of mean mutation rates per species on 1/lifespan were compared between this phylogenetic generalized linear model and a simple linear model ('lm' function in R). The use of this simple model, as well as the use of mean mutation rates per species (rather than mutation rates per sample), was necessary owing to the impossibility of replicating the heteroscedastic mixed-effects structure of the LME model used for the main association analyses (see 'Association of mutation rate with life-history traits') within the phylogenetic linear model. Both the phylogenetic linear model and the simple linear model produced similar estimates (Extended Data Fig. 13b), suggesting that phylogenetic non-independence of the samples does not have a substantial effect on the association analyses.

### Cell division analysis

To investigate the extent to which differences in cell division rate could explain differences in mutation rate and burden across species, we obtained estimates of intestinal crypt cell division rates from mouse<sup>70</sup>, rat<sup>71</sup> and human<sup>72,73</sup> (Supplementary Table 7). Using these cell division rates, our lifespan estimates and the observed substitution rates, we calculated the number of cell divisions at the end of lifespan and the corresponding number of mutations per cell division expected under a simple model assuming that all mutations occur during cell division (Supplementary Table 7).

We investigated whether differences in the number of cell divisions among species could explain the observed differences in mutation burden. Although colorectal cell division rate estimates are lacking for most species, existing estimates from mouse, rat and human indicate that the total number of stem cell divisions per crypt in a lifetime varies greatly across species—for example, there are around 6- to 31-fold more divisions per intestinal stem cell in a human than in a rat over their respective lifetimes, depending on the estimate of cell division rate used (Supplementary Table 7). Mouse intestinal stem cells are estimated to divide once every 24 h (ref.<sup>70</sup>), whereas estimates of the human intestinal stem cell division rate vary from once every 48 h (ref.<sup>72</sup>) to once every 264 h (ref.<sup>73</sup>). Thus, mouse intestinal stem cells divide 2–11 times faster than human intestinal stem cells. By the end of lifespan, an intestinal stem cell is predicted to have divided around 1,351 times in a mouse, around 486 times in a rat and 2,774–15,257 times in a human (depending on the estimate of cell division rate used). Applying our somatic mutation burden and lifespan data, this implies that the somatic mutation rate per cell division in a mouse is around 1.5- to 8.4-fold higher than in a human. However, the observed fold difference in somatic mutation rate between these two species is 16.9 (Table 1). Therefore, differences in cell division rate appear unable to fully account for the observed differences in mutation rate across species. Nevertheless, we note that accurate cell division rate estimates for basal intestinal stem cells are lacking for most species.

### Reporting summary

Further information on research design is available in the Nature Research Reporting Summary linked to this paper.

### Data availability

DNA sequence data have been deposited in the European Genome-Phenome Archive (EGA; <https://ega-archive.org>) under overarching accession EGAD00001008032. Human DNA sequence data from a previous study<sup>8</sup> are deposited in the EGA (accession EGAD00001004192). Processed mutation calls and other data used in the analyses have been deposited in Zenodo (<https://doi.org/10.5281/zenodo.5554777>). Raw mortality data used to estimate

lifespan (Species360 Data Use Approval Number 60633) cannot be publicly shared, as Species360 is the custodian (not the owner) of their members' data. Raw data are accessible through Research Request applications to Species360. Once Species360 grants access to data, they are intended only for and restricted to use in the project they were approved for and for a single publication. Any email communications should be directed to [support@species360.org](mailto:support@species360.org).

### Code availability

The computer code used in the analyses has been deposited in Zenodo (<https://doi.org/10.5281/zenodo.5554801>) and GitHub (<https://github.com/baезortega/CrossSpecies2021>).

- Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. Preprint at <https://arxiv.org/abs/1303.3997> (2013).
- Jones, D. et al. cgpCaVEManWrapper: simple execution of CaVEMan in order to detect somatic single nucleotide variants in NGS data. *Curr. Protoc. Bioinformatics* **56**, 15.10.1-15.10.18 (2016).
- Raine, K. M. et al. cgpPindel: identifying somatically acquired insertion and deletion events from paired end sequencing. *Curr. Protoc. Bioinformatics* **52**, 15.7.1-15.7.12 (2015).
- Gerstung, M., Papaemmanuil, E. & Campbell, P. J. Subclonal variant calling with multiple samples and prior knowledge. *Bioinformatics* **30**, 1198-1204 (2014).
- Gerstung, M. et al. Reliable detection of subclonal single-nucleotide variants in tumour cell populations. *Nat. Commun.* **3**, 811 (2012).
- Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841-842 (2010).
- Tacutu, R. et al. Human Ageing Genomic Resources: new and updated databases. *Nucleic Acids Res.* **46**, D1083-D1090 (2018).
- Tidière, M. et al. Comparative analyses of longevity and senescence reveal variable survival benefits of living in zoos across mammals. *Sci. Rep.* **6**, 36361 (2016).
- Conde, D. A. et al. Data gaps and opportunities for comparative and conservation biology. *Proc. Natl Acad. Sci. USA* **116**, 9658-9664 (2019).
- Blokzijl, F., Janssen, R., van Bostel, R. & Cuppen, E. MutationalPatterns: comprehensive genome-wide analysis of mutational processes. *Genome Med.* **10**, 33 (2018).
- Kumar, S., Stecher, G., Suleski, M. & Hedges, S. B. TimeTree: a resource for timelines, timetrees, and divergence times. *Mol. Biol. Evol.* **34**, 1812-1819 (2017).
- Snippert, H. J. et al. Intestinal crypt homeostasis results from neutral competition between symmetrically dividing Lgr5 stem cells. *Cell* **143**, 134-144 (2010).
- Rijke, R. P., Plaisier, H. M. & Langendoen, N. J. Epithelial cell kinetics in the descending colon of the rat. *Virchows Arch. B Cell Pathol. Incl. Mol. Pathol.* **30**, 85-94 (1979).
- Potten, C. S., Kellett, M., Rew, D. A. & Roberts, S. A. Proliferation in human gastrointestinal epithelium using bromodeoxyuridine in vivo: data for different sites, proximity to a tumour, and polyposis coli. *Gut* **33**, 524-529 (1992).
- Bach, S. P., Renehan, A. G. & Potten, C. S. Stem cells: the intestinal stem cell as a paradigm. *Carcinogenesis* **21**, 469-476 (2000).

**Acknowledgements** We thank the staff of the Wellcome Sanger Institute for help in this project, including the Scientific Operations, Informatics and CASM support teams. This research was made possible by the worldwide information network of zoos and aquariums that are members of Species360, and is authorized by Species360 research data use and grant agreement 60633. We thank D. Conde and J. Staerk for help accessing Species360 data; and J. P. de Magalhaes, J. DeGregori, J. Vijg and M. Lynch for their comments on the manuscript. This research was funded by Wellcome (grant number 206194), the Dunhill Medical Trust (RPGF2002188) and the Deutsche José Carreras Leukämie-Stiftung. I.M. is funded by Cancer Research UK (C57387/A21777) and the Wellcome Trust. P.J.C. is a Wellcome Trust Senior Clinical Fellow.

**Author contributions** A.C., E.P.M., M.R.S. and I.M. conceived the project. I.M., E.P.M. and M.R.S. supervised the project. E.F., S.S., I.J., E.W., N.M., R.D., M.W.P., M.D., R.B., M.D.M., F.G., F.C.-C., L.P., D.B., E.St.J.S., B.N. and E.P.M. performed and facilitated sample collection. A.C. performed the laser capture microdissection. A.C., L.M.R.H., A.R.J.L., Y.H., K.R., E.A. and S.L. processed the samples. A.C., A.B.-O., F.A., N.B., T.H.H.C., M.A.S., D.J., R.E.A. and S.B. processed the data. A.C., A.B.-O. and N.B. led the analysis with help from F.A., T.H.H.C., M.A.S., A.R.J.L., T.M.B., T.D., H.J., E.P.M. and I.M. The manuscript was written by A.C., A.B.-O., N.B. and I.M., with input from all of the authors.

**Competing interests** The authors declare no competing interests.

### Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41586-022-04618-z>.

**Correspondence and requests for materials** should be addressed to Alex Cagan or Iñigo Martincorena.

**Peer review information** Nature thanks Kamila Naxerova and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. Peer reviewer reports are available.

**Reprints and permissions information** is available at <http://www.nature.com/reprints>.