



Reporting quality, effect sizes, and biases for aging interventions: a methodological appraisal of the DrugAge database



Austin Parish¹ ✉, John P. A. Ioannidis², Kevin Zhang³, Diogo Barardo⁴, William R. Swindell⁵ & João Pedro de Magalhães⁶

Though interest has grown significantly over the past decades in interventions that may slow the aging process, most evidence for these interventions still comes from experiments in non-human animals. These studies may suffer from design, quality and reporting issues. The quality and reporting of preclinical studies have not yet been studied systematically in anti-aging research. Here we analyzed the DrugAge database, assessing reporting study quality, bias and effect sizes across 667 anti-aging preclinical studies. We found significant shortcomings in reporting of crucial design features such as randomization and blinding, as well as large variation in reporting quality and effects across species. Only one third of non-mammal findings translated to mammals. Although anti-aging interventions may have different effects depending on when they are started, most studies began giving the intervention under investigation very early in the organism's lifespan. Our findings suggest there is substantial room for improvement in preclinical anti-aging research.

There is increasing interest in interventions targeting the aging process^{1,2}. The “geroscience hypothesis” posits that a shared pathophysiology of aging shapes most chronic diseases and interventions targeting aging will confer larger health benefits than those targeting any individual disease^{3,4}. Research into such anti-aging interventions has grown substantially, including trials repurposing commonly used drugs such as metformin⁵.

Because of the large sample sizes and long durations of trials required to demonstrate anti-aging effects, most evidence to date has come from pre-clinical experiments in non-human animals⁶. Aging is a universal pathological process in eukaryotes^{7,8} with conservation of aging pathways across organisms^{9,10}; interventions targeting aging may be more successfully translated than interventions for specific diseases which often rely on artificial disease models^{11,12}.

Given the possible substantial health benefits of slowing aging, the quality of preclinical studies in this area may be especially important. However, alongside the challenges translating results from one species to another, model organism studies have a long history of shortcomings and design flaws^{13–15}.

Here, we systematically analyzed studies from DrugAge, a curated database of preclinical experiments investigating the effects of

interventions on aging and lifespan in non-human animals¹⁶. We aimed to evaluate the quality of reporting and methodological rigor of this literature, assess the distribution of observed effect sizes, and probe for the presence of diverse biases. We also investigated how these features changed over time.

Results

Of 667 included studies, 617 included only experiments with one species and one start time; 29 summarized experiments with two species and three summarized experiments with three species. Eighteen studies had experiments with two start times, for a total of 720 experiments. Of these, 364 involved an organism that reproduces sexually; of these, 130 used only males (35.7%); 47 used only females (12.9%), 172 used both (47.3%), and 15 did not report the sex(es) used (4.1%). The median sample size across experiments was 200 animals (IQR: 105–341).

All studies were published in peer-reviewed journals (667, 100.0%) and most stated control of temperature (607, 91.0%). Randomization was mentioned in 133 studies (19.9%). Blinding to intervention was mentioned in 27 studies (4.0%), blinded assessment of outcomes in 20 (3.0%), and sample size calculations in 40 (6.0%). Following animal welfare regulations

¹Department of Emergency Medicine, Brookdale University Hospital Medical Center, Brooklyn, NY, USA. ²Meta-Research Innovation Center at Stanford (METRICS), Stanford University, and Departments of Medicine, of Epidemiology and Population Health, and of Biomedical Data Science, Stanford University School of Medicine, Stanford, CA, USA. ³Hackensack Meridian School of Medicine, Nutley, NJ, USA. ⁴NOVOS Labs, New York, NY, USA. ⁵Department of Internal Medicine, Division of Hospital Medicine, University of Texas Southwestern Medical Center, Dallas, TX, USA. ⁶Genomics of Ageing and Rejuvenation Lab, Department of Inflammation and Ageing, College of Medicine and Health, University of Birmingham, Birmingham, UK. ✉e-mail: auparish@bhmcny.org

Table 1 | Mention of CAMARADES components across studies of different species in the database, along with median CAMARADES count (sum of the 8 included CAMARADES elements, minimum 0 and maximum 8)

Species, Studies (N)	Peer Reviewed Journal (%)	Control of Temperature (%)	Randomization (%)	Blinded Intervention (%)	Blinded Assessment of outcomes (%)	Sample Size Calculation (%)	Animal Welfare Regulations (%)	Conflict of Interest Statement (%)	CAMARADES Count Median (IQR)
<i>Caenorhabditis elegans</i> (316)	316 (100%)	313 (99.1%)	16 (5.1%)	14 (4.4%)	7 (2.2%)	6 (1.9%)	4 (1.3%)	194 (61.4%)	3 (2-3)
<i>Drosophila melanogaster</i> (150)	150 (100%)	144 (96%)	24 (16%)	2 (1.3%)	2 (1.3%)	6 (4%)	10 (6.7%)	70 (46.7%)	3 (2-3)
<i>Mus musculus</i> (102)	102 (100%)	72 (70.6%)	60 (58.8%)	9 (8.8%)	10 (9.8%)	27 (26.5%)	64 (62.7%)	52 (51%)	4 (2-6)
<i>Rattus norvegicus</i> (29)	29 (100%)	19 (65.5%)	17 (58.6%)	0 (0%)	0 (0%)	1 (3.4%)	9 (31%)	6 (20.7%)	3 (2-4)
<i>Saccharomyces cerevisiae</i> (12)	12 (100%)	11 (91.7%)	1 (8.3%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	6 (50%)	2 (2-3)
<i>Musca domestica</i> (7)	7 (100%)	5 (71.4%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	2 (1.5-2)
<i>Zapionus parvittiger</i> (7)	7 (100%)	7 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (14.3%)	0 (0%)	2 (2-2)
<i>Asplanchna brightwellii</i> (5)	5 (100%)	5 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	2 (2-2)
<i>Nothobranchius guentheri</i> (5)	5 (100%)	3 (60%)	3 (60%)	0 (0%)	0 (0%)	0 (0%)	2 (40%)	4 (80%)	3 (3-4)
<i>Brachionus manjavacas</i> (4)	4 (100%)	4 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	2 (50%)	2.5 (2-3)
<i>Philodina acuticornis</i> (4)	4 (100%)	3 (75%)	2 (50%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	2 (50%)	2.5 (2-3.2)
<i>Aedes aegypti</i> (2)	2 (100%)	2 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (50%)	1 (50%)	3 (2.5-3.5)
<i>Anastrepha ludens</i> (2)	2 (100%)	2 (100%)	2 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (50%)	3.5 (3.2-3.8)
<i>Mesocricetus auratus</i> (2)	2 (100%)	0 (0%)	1 (50%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1.5 (1.2-1.8)
<i>Nothobranchius furzeri</i> (2)	2 (100%)	2 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (50%)	0 (0%)	2.5 (2.2-2.8)
<i>Paramecium tetraurelia</i> (2)	2 (100%)	2 (100%)	2 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	3 (3-3)
<i>Podospora anserina</i> (2)	2 (100%)	2 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (50%)	2.5 (2.2-2.8)
<i>Acheta domesticus</i> (1)	1 (100%)	1 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (100%)	3
<i>Aedes albopictus</i> (1)	1 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (100%)	2
<i>Aeolosoma viride</i> (1)	1 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1
<i>Anopheles stephensi</i> (1)	1 (100%)	1 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (100%)	3
<i>Apis mellifera</i> (1)	1 (100%)	1 (100%)	1 (100%)	1 (100%)	1 (100%)	0 (0%)	0 (0%)	1 (100%)	6
<i>Bombyx mori</i> (1)	1 (100%)	1 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (100%)	3
<i>Caenorhabditis briggsae</i> (1)	1 (100%)	1 (100%)	1 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (100%)	4
<i>Caenorhabditis tropicalis</i> (1)	1 (100%)	1 (100%)	1 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (100%)	4
<i>Canis lupus familiaris</i> (1)	1 (100%)	0 (0%)	1 (100%)	1 (100%)	0 (0%)	0 (0%)	1 (100%)	0 (0%)	4
<i>Daphnia pulex</i> clone TCO (1)	1 (100%)	1 (100%)	1 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	3
<i>Drosophila bipectinata</i> (1)	1 (100%)	1 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	2
<i>Drosophila mojavensis</i> (1)	1 (100%)	1 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	2
<i>Mytilina brevispina</i> (1)	1 (100%)	1 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	2
<i>Tribolium castaneum</i> (1)	1 (100%)	1 (100%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	0 (0%)	1 (100%)	3
All (667)	667 (100%)	607 (91%)	133 (19.9%)	27 (4%)	20 (3%)	40 (6%)	93 (13.9%)	347 (52%)	3 (2-3)
p values	1.0	<0.0001	<0.0001	0.14	0.12	<0.0001	<0.0001	<0.0001	<0.0001

P values for count outcomes reflect the results of 2 x 31 exact tests; p value for the median CAMARADES count is the result of the Kruskal-Wallis test.

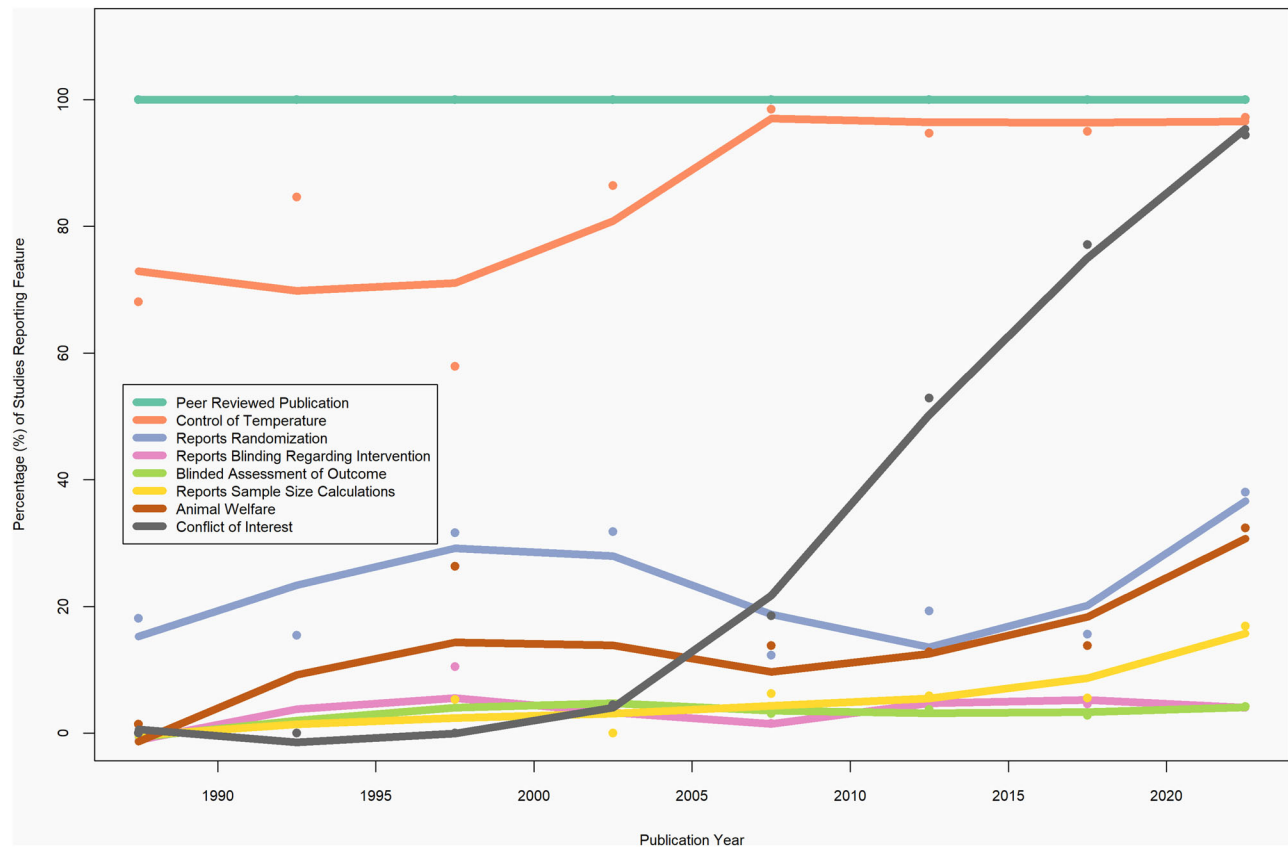


Fig. 1 | Percentage of studies with specific CAMARADES features over time (scatterplot, with local polynomial regression fit curves). Reporting of potential conflicts of interest, animal welfare regulations, control of temperature and sample size calculations increased significantly over time ($p < 0.0001$ for each).

was mentioned in 93 studies (13.9%). Conflict of interest statements were included in 347 studies (52.0%).

The median CAMARADES score across studies was 3 (IQR: 2–3) and varied significantly across species ($p < 0.0001$). 61 studies (9.1%) had CAMARADES counts > 4 . Except for peer-review publication that was ubiquitous and blinding that was rare, all CAMARADES components varied significantly across species ($p < 0.0001$, Table 1). *Caenorhabditis* and *Drosophila* studies almost always stated control of temperature but rarely reported randomization or sample size calculations. Studies of mice and rats stated control of temperature less commonly but did better on all other fronts.

Of the 667 studies, 153 reported whether the organisms included were from an inbred (genetically homogenous) line or an outbred/hybrid line; of these, 73 used inbred lines. None of the CAMARADES components differed significantly between inbred and non-inbred studies.

The earliest included study was published in 1948, the latest in 2024. Over time, there was a significant increase in reporting conflicts of interest, compliance with animal welfare regulations, control of temperature and sample size calculations (linear regression $p < 0.0001$ for each). There was no significant increase over time in reporting of randomization ($p = 0.60$) or blinding with regard to intervention ($p = 0.07$) or outcomes ($p = 0.011$) (Fig. 1). Studies had higher CAMARADES counts over time ($p < 0.0001$).

Across 720 experiments, the median percentage of average lifespan that interventions were started at was 5.7% (IQR: 5.3–13.8%) (Fig. 2), with significant variation across species (Table 2). Mammal experiments started at a relatively later point in lifespan than non-mammal experiments (26.2% vs 5.6%, $p < 0.0001$). Most experiments started “early” (before 20% of average lifespan) ($n = 592$, 82.2%). Few experiments started at 50% of average lifespan or later (58, 8.1%).

Of the 720 included experiments, most SMDs were positive (638, 88.6%), indicating a favorable effect of the intervention on lifespan. The

median SMD was 0.43 (IQR: 0.24–0.70); the random effects meta-analysis estimate was 0.57 (95% CI: 0.48–0.66, $p < 0.0001$), with significant heterogeneity ($I^2 = 95%$, 94–96%, $p < 0.0001$). As a fraction of average species lifespan, the median percentage increase in lifespan was 13.3% (IQR: 6.8–23.1%); the meta-analysis estimate was 14.3% (12.9–15.7, $p < 0.0001$). Table 3 summarizes these results.

Comparing experiments in studies with specific CAMARADES components, reporting of randomization was associated with a smaller SMD (0.38 in those reporting vs 0.45, Kruskal-Wallis $p = 0.0074$). Other CAMARADES components were not associated with significant differences: peer-reviewed publication ($p = 1.0$), control of temperature ($p = 0.094$), blinded intervention ($p = 0.35$), blinded assessment of outcome ($p = 0.84$), sample size calculations reported ($p = 0.17$), compliance with animal welfare regulations ($p = 0.48$), conflict of interest statement ($p = 0.041$). However, in multivariate linear regression adjusting also for species, publication year, sample size, there was no statistically significant association between any CAMARADES component and SMD. There was no significant difference between the 592 early start experiments and the 128 late start experiments (median SMD 0.43 vs 0.40, Kruskal-Wallis $p = 0.60$). Median SMD did not vary significantly with publication year ($p = 0.11$) (see Supplementary Fig. 1 for bubble plot). Studies with mammalian species had lower median SMDs than non-mammal studies (0.39 vs 0.44, $p = 0.040$).

There were 35 compounds that were tested in at least one mammal and at least one non-mammal experiment, allowing for comparisons within the same drug (see Supplementary Table 1). Of these, 21 showed a significant increase in lifespan ($p < 0.005$) for non-mammals. Of these, only 7 also showed a significant increase in mammal lifespan (curcumin, spermidine, epthalamin, D-glucosamine, estradiol, SKQ and taurine); additionally, two in contrast showed a significant decrease in mammal lifespan (quercetin and butylated hydroxytoluene).

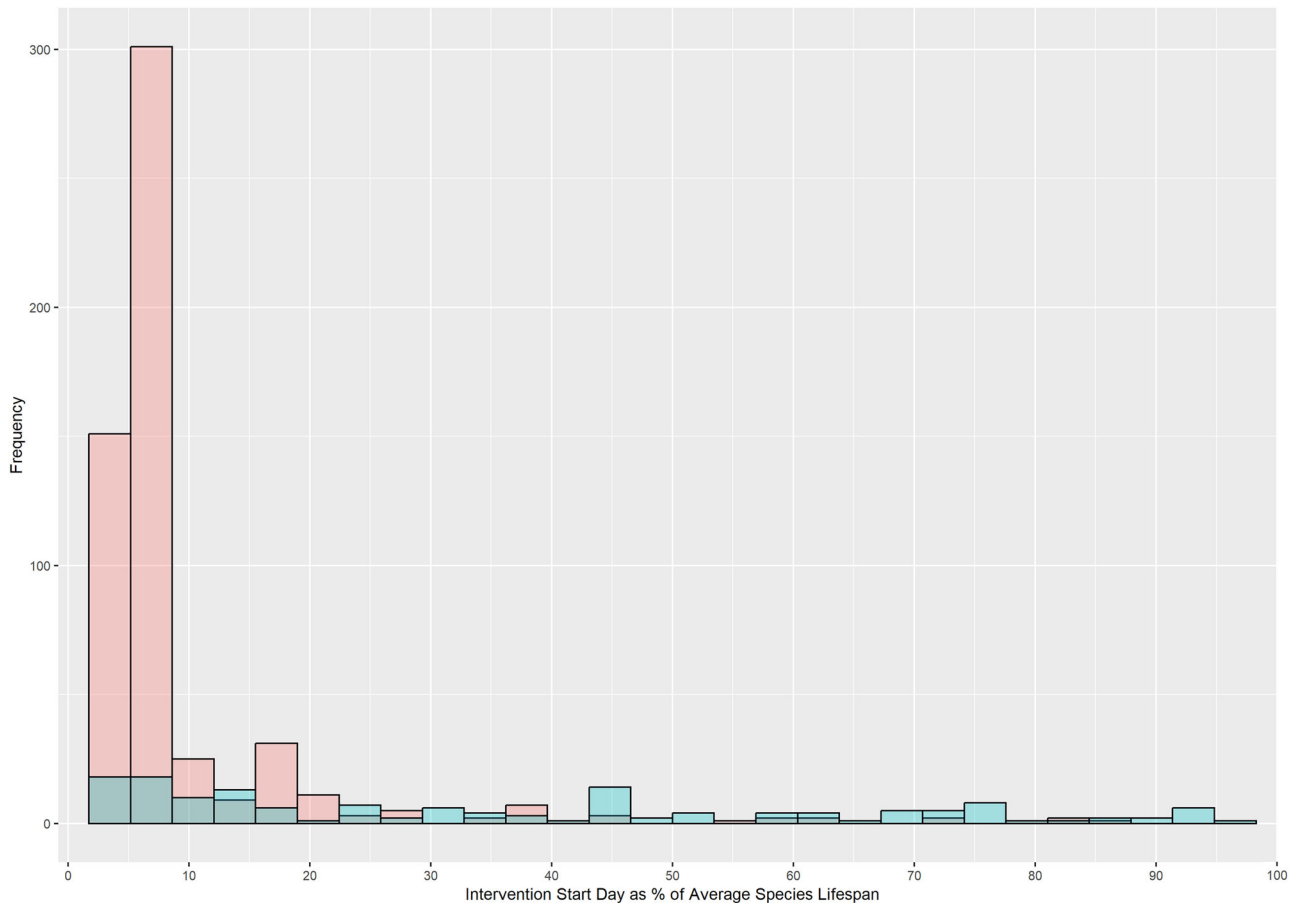


Fig. 2 | Distribution of when an intervention was started in the lifespan of an organism, expressed as a percentage of the median lifespan for that species, across 720 experiments. Blue represents 153 mammal experiments (median = 26.2%), and pink represents 567 non-mammal experiments (median = 5.6%).

Ten compounds showed a significant increase in mammal lifespan ($p < 0.005$) and most showed some increase in non-mammals as well. However, the amount of mammalian evidence for these compounds other than rapamycin was limited (total sample sizes: 2761 for rapamycin, 293 for curcumin, 360 for spermidine, 160 for melatonin, 171 for epitalamin, 44 for berberine, 146 for D-glucosamine, 370 for estradiol, 50 for SKQ, and 122 for taurine).

The absolute percent error between non-mammal and mammal effects was 83% (IQR: 49–164%). Across compounds, the median percentage increase in mammal lifespan was significantly smaller than the percentage increase in non-mammal lifespan (7.4 vs 17.5%, paired Wilcoxon $p = 0.0006$). There was no significant correlation between non-mammal and mammal SMDs or percentage increases ($r = 0.26$, $p = 0.10$ and $r = 0.25$, $p = 0.10$ respectively).

Of the 720 experiments, 638 (88.6%) were associated with an increase in lifespan and 82 (11.4%) were associated with a decrease in lifespan. Of the 638 experiments associated with increasing lifespan, 495 showed $p < 0.05$ (77.6%); of the 82 experiments associated with decreasing lifespan, 51 showed $p < 0.05$ (62.2%); Supplementary Fig. 2 shows the p value distributions.

Across the 720 experiments there was evidence of significant funnel plot asymmetry (Egger's $Z = 11.3$, $p < 0.0001$); see Fig. 3 for the contour enhanced funnel plot. The expected number of significant findings was 499 out of 720, while the observed number was 546, indicating significantly more significant results than expected (test of excess significance $\chi^2 = 14.0$, $p < 0.0001$); similarly, the proportion of statistical significance test resulted in a test statistic of $Z = 3.74$ ($p < 0.0001$)^{17,18}.

Discussion

Despite growing excitement about the possibility of anti-aging interventions impacting human healthspan and lifespan, most studies of these interventions have been conducted in non-human animals. In this review of 720 experiments from 667 such studies, we found widely varying effect sizes and reporting of factors potentially associated with study quality across species and compounds. Important design features such as randomization, blinding of intervention, blinded assessment of outcome, compliance with animal welfare regulations, and sample size calculations were infrequently reported, despite evidence that the absence of such features can bias experimental results^{19–23}. Only slightly more than half included conflict of interest statements, although all studies were published in peer-reviewed journals, and over 90% reported control of temperature.

Although reporting quality improved somewhat over time, this was mainly due to increases in reporting of compliance with animal welfare regulations or conflict of interest statements; crucial design features such as randomization and blinding did not increase substantially over time. Generally, most studies did not meet standard reporting guidelines for pre-clinical experiments²².

Preclinical studies on various diseases have also shown infrequent reporting of randomization and blinding. A review of 271 preclinical studies across different diseases found 13% of studies reported randomization and 14% blinding²⁴. In another similar review of 290 studies, 32% reported randomization and 11% blinding²⁵. Overall, our results are comparable to these, although the reporting of both randomization and blinding seems to be even less frequent in the DrugAge database.

In addition to the overall low rate of reporting factors associated with study quality, we found significant differences across species: the four most

Table 2 | Timepoint in lifespan of organisms that each experiment was started at expressed as a percentage of the median lifespan for that organism. Kruskal-Wallis test for comparison between species $p < 0.0001$

Species	Number of experiments	Percentage of average lifespan start time (median, IQR)
<i>Caenorhabditis elegans</i>	333	5.6% (5.6–5.6)
<i>Drosophila melanogaster</i>	159	2.1% (2.1–4.2)
<i>Mus musculus</i>	117	33.0% (11.2–62)
<i>Rattus norvegicus</i>	32	13.8% (6.2–35.4)
<i>Saccharomyces cerevisiae</i>	13	11.2% (11.2–11.2)
<i>Musca domestica</i>	8	4.8% (4.8–4.8)
<i>Zapionus parvittiger</i>	7	2.7% (2.7–2.7)
<i>Asplanchna brightwelli</i>	5	18.6% (18.6–18.6)
<i>Nothobranchius guentheri</i>	5	36.4% (36.4–36.4)
<i>Brachionus manjavacas</i>	4	9.2% (7.3–27.6)
<i>Philodina acuticornis</i>	4	5.5% (5.5–5.5)
<i>Caenorhabditis briggsae</i>	3	3.2% (3.2–3.2)
<i>Mesocricetus auratus</i>	3	10.8% (7.7–40.6)
<i>Aedes aegypti</i>	2	33.2% (18.7–47.7)
<i>Anastrepha ludens</i>	2	9.0% (5.3–12.7)
<i>Caenorhabditis tropicalis</i>	2	4.6% (4.6–4.6)
<i>Nothobranchius furzeri</i>	2	71.0% (67.1–75)
<i>Paramecium tetraurelia</i>	2	1.7% (1.7–1.7)
<i>Podospora anserina</i>	2	16.0% (11.2–20.8)
<i>Acheta domesticus</i>	1	15.7%
<i>Adineta vaga</i>	1	12.5%
<i>Aedes albopictus</i>	1	1.8%
<i>Aeolosoma viride</i>	1	1.6%
<i>Anopheles stephensi</i>	1	14.2%
<i>Apis mellifera</i>	1	11.1%
<i>Bombyx mori</i>	1	2.3%
<i>Canis lupus familiaris</i>	1	69.5%
<i>Daphnia pulex</i> clone TCO	1	29.0%
<i>Drosophila bipectinata</i>	1	4.3%
<i>Drosophila kikkawai</i>	1	4.0%
<i>Drosophila mojavensis</i>	1	40.7%
<i>Drosophila virilis</i>	1	1.4%
<i>Mytilina brevispina</i>	1	34.5%
<i>Tribolium castaneum</i>	1	1.2%
All species	720	5.7% (5.3–13.8)

represented species in the database (nematodes, fruit flies, mice and rats) varied widely in reporting of randomization and blinding, as well as in the average effect size found. Over half of mammal studies in the database reported randomization, while less than 10% of non-mammal studies did. The better reporting quality of mammal studies does not alleviate concerns, since even for mammal studies reporting was often suboptimal and most

studies in the database were from non-mammals. Additionally, the average effect found in non-mammal studies was significantly larger than that found in mammal studies.

For 35 compounds with both mammal and non-mammal experiments, only eight showed a significant lifespan increase in both non-mammals and mammals; the number of experiments and sample sizes for these results were limited. These results are exploratory, and the numbers are small, but they raise hesitation about the direct translation of these results to more complex organisms such as humans.

Furthermore, previous work has suggested that some interventions may have different effects if started late in an organism's lifespan rather than early^{26,27}, and there is significant interest in discovering interventions that slow aging in older adults²⁸. In our assessment, we found that most pre-clinical experiments started the anti-aging intervention early in the organism's lifespan, often prior to sexual maturity, when key senescence mechanisms may lack relevance²⁹. Although we did not find a significant difference in the effect of interventions between early and late start experiments, the sparsity of late start results makes this comparison uncertain. Our study clearly highlights the paucity of late start experiments in the literature, a deficit of evidence that needs to be remedied. We emphasize the need for whole-lifespan aging experiments with a greater diversity of start times, including more starting in middle or late life, as these better reflect the intended translational application of anti-aging interventions and likely design of future clinical trials investigating proposed interventions.

Overall, the analyzed studies have numerous significant reported results and many studies suggest sizeable effect sizes. However, the lack of methodological rigor (at least based on reported information) and the strong suggestion of bias (larger effects in smaller studies and an excess of significant results) prompt skepticism about this overall favorable picture and prospects for translation to humans. Our work has several limitations. The DrugAge database may not include some compounds that have never shown any promising results. Moreover, we did not extract quantitative data from all of the 3423 experiments in the database, but rather extracted a random experiment from each species represented in each study, as well as the earliest and latest start time experiments, obtaining quantitative data from only 21% of the experiments in the database. Although experiments were selected randomly, this method may still have led to a biased estimate of quantitative effects in some cases. Nevertheless, our selection process resulted in a dataset with largely independent observations, while the full database may have a lot of highly correlated data and may have over-represented specific experiments that were rather similar.

It is also possible that some studies that did not mention randomization still carried out randomization, and the same may apply to other design features. Nevertheless, the large variation in reporting of randomization and blinding is concerning. Also, we found the reporting of randomization did not differ significantly between studies that used genetically homogenous populations of organisms and those using more heterogenous populations. Furthermore, while we observed improvements over time in compliance with animal welfare, control of temperature and sample size calculations, it is unclear whether this represents genuine improvements in experimental practices, or simply better realization that these are features that should be reported in their publications.

Preclinical experiments investigating anti-aging interventions do not regularly follow reporting guidelines and infrequently report important design features such as randomization and blinding. There are significant differences in the average lifespan effect, as well as study quality, across different species commonly used in preclinical experiments. Non-mammal results do not seem to reliably predict mammal results, raising further concern for translation. Despite the interest in interventions able to slow aging when initiated late in human lifespan, most preclinical experiments started interventions early in organism lifespans. Our work highlights multiple concrete areas for improvement of preclinical anti-aging research, areas that may be critical for successful translation into human trial results.

Table 3 | Median and IQR of SMD for lifespan increase, as well as the median percentage increase in lifespan, for each species in the database, across 720 experiments from 667 studies

Species	Number of experiments	Median SMD (IQR)	Median Percentage Increase in Lifespan (%)
<i>Caenorhabditis elegans</i>	333	0.50 (0.29–0.79)	17.8% (16–19.6)
<i>Drosophila melanogaster</i>	159	0.35 (0.20–0.56)	12.8% (10–15.6)
<i>Mus musculus</i>	117	0.38 (0.14–0.57)	6.7% (5.2–8.2)
<i>Rattus norvegicus</i>	32	0.47 (0.10–0.79)	5.9% (2.5–9.3)
<i>Saccharomyces cerevisiae</i>	13	1.00 (0.53–2.43)	42.7% (17.4–68)
<i>Musca domestica</i>	8	–0.59 (–1.39 to –0.08)	–14.5% (–33.4 to 4.5)
<i>Zaprionus parvittiger</i>	7	0.59 (0.23–6.39)	11.2% (4.5–18)
<i>Asplanchna brightwelli</i>	5	1.24 (0.70–1.46)	13.6% (–3.5 to 30.7)
<i>Nothobranchius guentheri</i>	5	0.42 (0.42–0.53)	18.5% (7.5–29.4)
<i>Brachionus manjavacas</i>	4	0.43 (0.40–0.56)	25.5% (8.8–42.2)
<i>Philodina acuticornis</i>	4	0.62 (0.34–2.08)	72.4% (–30.5 to 175.2)
<i>Caenorhabditis briggsae</i>	3	–0.25 (–0.35 to 0.05)	–9.5% (–28.2 to 9.2)
<i>Mesocricetus auratus</i>	3	0.16 (–0.01 to 0.40)	0.5% (–5.4 to 6.5)
<i>Aedes aegypti</i>	2	0.53 (0.44–0.62)	48.1% (25.7–70.5)
<i>Anastrepha ludens</i>	2	0.65 (0.50–0.81)	37.2% (–5.8 to 80.2)
<i>Caenorhabditis tropicalis</i>	2	0.21 (–0.02 to 0.43)	2.5% (–14.6 to 19.7)
<i>Nothobranchius furzeri</i>	2	0.70 (0.58–0.82)	37.4% (–25.2 to 100.1)
<i>Paramecium tetraurelia</i>	2	0.43 (0.32–0.54)	25.6% (8.5–42.7)
<i>Podospora anserina</i>	2	0.27 (0.06–0.48)	4.6% (–7.2 to 16.4)
<i>Acheta domesticus</i>	1	0.28	67.4% (–0.1 to 134.8)
<i>Adineta vaga</i>	1	0.42	20.6% (8.3–33)
<i>Aedes albopictus</i>	1	0.43	16.2% (5.8–26.7)
<i>Aeolosoma viride</i>	1	1.75	53.1% (27.8–78.4)
<i>Anopheles stephensi</i>	1	–0.27	–11.4% (–20.1 to –2.6)
<i>Apis mellifera</i>	1	0.37	37.9% (12.8–63)
<i>Bombyx mori</i>	1	0.52	2.7% (1.2–4.3)
<i>Canis lupus familiaris</i>	1	0.52	2.5% (–0.7 to 5.7)
<i>Daphnia pulex clone TCO</i>	1	–0.13	–6.5% (–24.4 to 11.4)
<i>Drosophila bipectinata</i>	1	0.31	13.0% (5.7–20.4)
<i>Drosophila kikkawai</i>	1	0.28	16.0% (7.9–24.1)
<i>Drosophila mojavensis</i>	1	0.88	47.5% (–0.1 to 95)
<i>Drosophila virilis</i>	1	0.28	31.0% (15.3–46.7)
<i>Mytilina brevispina</i>	1	3.47	43.7% (37–50.3)
<i>Tribolium castaneum</i>	1	0.47	32.1% (7.6–56.7)
All species	720	0.43 (0.24–0.70)	13.3% (6.8–23.1)
Kruskal-Wallis p:		$p < 0.0001$	$p < 0.0001$

Methods

We downloaded the fifth build of the DrugAge database on May 1st, 2025, which contained 3423 different lifespan experiments from a total of 680 unique studies. From these, we excluded 12 studies focusing on replicative aging in the yeast *Saccharomyces cerevisiae* and one duplicate study, yielding 667 unique studies.

For the 32 studies containing experiments with more than one species, a single experiment was randomly selected for each organism. If a study contained experiments where the same compound was started at different points in the same organism’s lifespan, we included one experiment for the earliest and one experiment for the latest start time. After selecting experiments in this way, our final dataset contained 720 experiments. See Fig. 4 for a flowchart of data extraction and Data Supplement 1 for all included studies and experiments. Overall, the 720 experiments represented 568 different species-drug pairs.

For each study, we extracted eight relevant quality checklist items from CAMARADES (Collaborative Approach to Meta-Analysis and Review of Animal Data in Experimental Studies) (Macleod 2004): (1) whether the study was peer reviewed, and reporting of (2) control of temperature, (3) random allocation to treatment/control; (4) blinded intervention; (5) blinded assessment of outcomes; (6) sample size calculations; (7) adherence to animal welfare regulations; and (8) potential conflicts of interest.

We also extracted the median or mean lifespans of experimental and control groups and whenever possible their confidence intervals (CIs) or standard errors (SEs). When these were not reported, we estimated them from included Kaplan-Meier figures and corresponding log-rank test *p* values.

The difference between experimental and control groups in lifespan was calculated, as well as the standardized mean difference (SMD) and its SE. Whenever possible, the SEs reported for the experimental and control

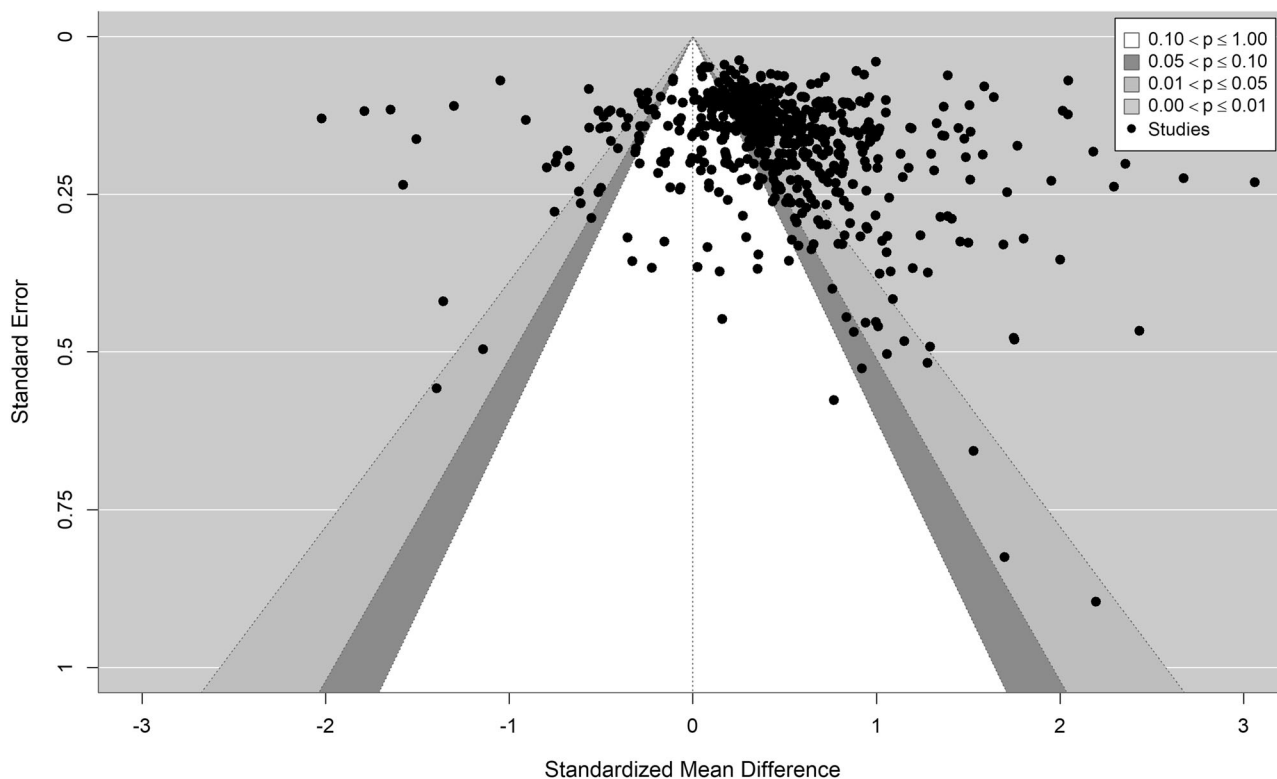


Fig. 3 | Contour-enhanced funnel plot of SMD from the 720 experiments.

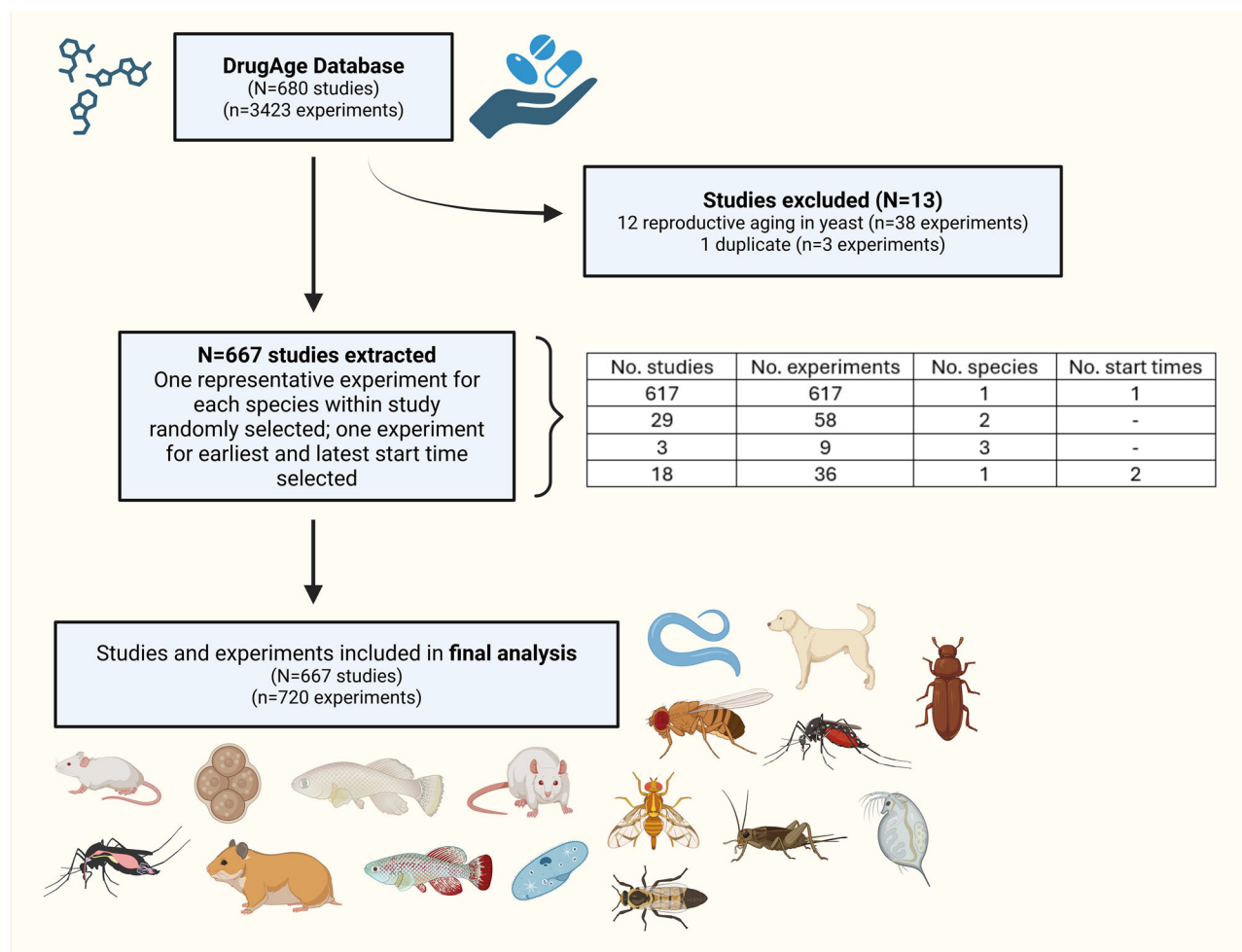


Fig. 4 | Flowchart of data extraction.

lifespans were used to calculate the SMD and its SE; when these were not reported, the log-rank p value was used to estimate the SE³⁰. We also calculated the relative increase in lifespan by dividing the difference in experimental and control lifespan by the control lifespan for each experiment, expressing the result as a percentage. For calculations of start times, we estimated the average lifespan of each included species as the median control lifespan in days across all experiments for that species in the database.

Random-effects meta-analysis was performed using the Sidik-Jonkman estimator³¹. Heterogeneity was estimated with the I^2 statistic and Q test^{32,33}. Meta-analysis calculations used the meta package in R³⁴. Contour-enhanced funnel plots and Egger's test were used to detect small study effects³⁵, considering all 720 results together in the same funnel plot. Funnel plot asymmetry indicates smaller studies showing larger effects than larger studies, possibly indicating selection biases. We also applied the test of excess significance and proportion of statistical significance test, tests that may suggest selective reporting biases^{17,18}. These various analytical tools allow estimation of publication selection bias, where selection processes favoring statistically significant results lead to biased selection of which outcomes to report, which statistical tests to use, and even which studies to publish or not¹⁷.

Results were considered statistically significant for $p < 0.005$ and possibly suggesting significance for p values between 0.05 and 0.005³⁶. The <0.005 threshold has been proposed as a way to diminish the risk of false-positive results that is high in many fields³⁶. The 4.1.0 version of the R programming language was used for all calculations³⁷.

Data availability

All data used for this study is available in the Data Supplement 1.

Code availability

The R code used for statistical analysis in this study will be made available upon request.

Received: 14 July 2025; Accepted: 22 October 2025;

Published online: 27 November 2025

References

- Burch, J. B., et al. Advances in geroscience: impact on healthspan and chronic disease. *J. Gerontol. A Biol. Sci. Med. Sci.* **69**, S1–S3 (2014).
- Guarente, L., Sinclair, D. A. & Kroemer, G. Human trials exploring anti-aging medicines. *Cell Metab.* **36**, 354–376 (2024).
- Kritchevsky, S. B. & Justice, J. N. Testing the geroscience hypothesis: early days. *J. Gerontol. A Biol. Sci. Med. Sci.* **75**, 99–101 (2020).
- Goldman, D. P. et al. Substantial health and economic returns from delayed aging may warrant a new focus for medical research. *Health Aff.* **32**, 1698–1705 (2013).
- Kulkarni, A. S. et al. Geroscience-guided repurposing of FDA-approved drugs to target aging: a proposed process and prioritization. *Aging Cell* **21**, e13596 (2022).
- Huffman, D. M. et al. Evaluating health span in preclinical models of aging and disease: guidelines, challenges, and opportunities for geroscience. *J. Gerontol. A Biol. Sci. Med. Sci.* **71**, 1395–1406 (2016).
- Cohen, A. A. Aging across the tree of life: the importance of a comparative perspective for the use of animal models in aging. *Biochim. Biophys. Acta Mol. Basis Dis.* **1864**, 2680–2689 (2018).
- Folgueras, A. R., Freitas-Rodríguez, S., Velasco, G. & López-Otín, C. Mouse models to disentangle the hallmarks of human aging. *Circ. Res.* **123**, 905–924 (2018).
- Flatt, T. & Partridge, L. Horizons in the evolution of aging. *BMC Biol.* **16**, 93 (2018).
- Kenyon, C. A conserved regulatory system for aging. *Cell* **105**, 165–168 (2001).
- Van Der Worp, H. B. et al. Can animal models of disease reliably inform human studies?. *PLoS Med.* **7**, e1000245 (2010).
- Jickling, G. C. & Sharp, F. R. Improving the translation of animal ischemic stroke studies to humans. *Metab. Brain Dis.* **30**, 461–467 (2015).
- van Luijk, J. et al. Systematic reviews of animal studies; missing link in translational research?. *PLoS ONE* **9**, e89981 (2014).
- Hooijmans, C. R. et al. SYRCLE's risk of bias tool for animal studies. *BMC Med. Res. Methodol.* **14**, 43 (2014).
- Macleod, M. R., O'Collins, T., Howells, D. W. & Donnan, G. A. Pooling of animal experimental data reveals influence of study design and publication bias. *Stroke* **35**, 1203–1208 (2004).
- Belikov, A. V., Talay, A. & de Magalhães, J. P. Sex-specific insights into drug-induced lifespan extension and weight loss in mice. *NPJ Aging* **11**, 37 (2025).
- Stanley, T. D., Doucouliagos, H., Ioannidis, J. P. A. & Carter, E. C. Detecting publication selection bias through excess statistical significance. *Res Synth. Methods* **12**, 776–795 (2021).
- Ioannidis, J. P. & Trikalinos, T. A. An exploratory test for an excess of significant findings. *Clin. Trials* **4**, 245–253 (2007).
- Schulz, K. F., Chalmers, I., Hayes, R. J. & Altman, D. G. Empirical evidence of bias. Dimensions of methodological quality associated with estimates of treatment effects in controlled trials. *JAMA* **273**, 408–412 (1995).
- Schulz, K. F. & Grimes, D. A. Allocation concealment in randomised trials: defending against deciphering. *Lancet* **359**, 614–618 (2002).
- Schulz, K. F. & Grimes, D. A. Blinding in randomised trials: hiding who got what. *Lancet* **359**, 696–700 (2002).
- Crossley, N. A. et al. Empirical evidence of bias in the design of experimental stroke studies: a meta epidemiologic approach. *Stroke* **39**, 929–934 (2008).
- Kringe, L. et al. Quality and validity of large animal experiments in stroke: a systematic review. *J. Cereb. Blood Flow. Metab.* **40**, 2152–2164 (2020).
- Kilkenny, C. et al. Survey of the quality of experimental design, statistical analysis and reporting of research using animals. *PLoS ONE* **4**, e7824 (2009).
- Bebarta, V., Luyten, D. & Heard, K. Emergency medicine animal research: does use of randomization and blinding affect the results? [published correction appears in *Acad Emerg Med.* 200310(12):1410]. *Acad. Emerg. Med.* **10**, 684–687 (2003).
- Espada, L. et al. Loss of metabolic plasticity underlies metformin toxicity in aged *Caenorhabditis elegans*. *Nat. Metab.* **2**, 1316–1331 (2020).
- Harrison, D. E. et al. Rapamycin fed late in life extends lifespan in genetically heterogeneous mice. *Nature* **460**, 392–395 (2009).
- Scott, A. J., Ellison, M. & Sinclair, D. A. The economic value of targeting aging. *Nat. Aging* **1**, 616–623 (2021).
- Williams, G. C. Pleiotropy, natural selection, and the evolution of senescence. *Evolution* **11**, 398–411 (1957).
- Altman, D. G. & Bland, J. M. How to obtain the confidence interval from a P value. *BMJ* **343**, d2090 (2011).
- Int'Hout, J., Ioannidis, J. P. & Borm, G. F. The Hartung-Knapp-Sidik-Jonkman method for random effects meta-analysis is straightforward and considerably outperforms the standard DerSimonian-Laird method. *BMC Med. Res. Methodol.* **14**, 25 (2014).
- Sedgwick, P. Meta-analyses: what is heterogeneity?. *BMJ* **350**, h1435 (2015).
- von Hippel, P. T. The heterogeneity statistic $I(2)$ can be biased in small meta-analyses. *BMC Med Res Methodol.* **15**, 35 (2015).
- Balduzzi, S., Rucker, G. & Schwarzer, G. How to perform a meta-analysis with R: a practical tutorial. *Evid. Based Ment. Health* **22**, 153–160 (2019).
- Egger, M., Davey Smith, G., Schneider, M. & Minder, C. Bias in meta-analysis detected by a simple, graphical test. *BMJ* **315**, 629–634 (1997).

36. Ioannidis, J. P. A. The proposal to lower P value thresholds to .005. *JAMA* **319**, 1429–1430 (2018).
37. R Core Team. R: The R Project for Statistical Computing. R-project.org. <https://www.r-project.org> (2021).

Author contributions

Austin Parish conceived and designed the study, collected data, performed analysis, and wrote and edited manuscript. John Ioannidis supported the analysis, wrote and edited manuscript. Kevin Zhang collected data and wrote and edited manuscript. Diogo Barardo supported the analysis, wrote and edited manuscript. William Swindell supported the analysis, wrote and edited manuscript. João Pedro de Magalhães supported the analysis, wrote and edited manuscript.

Competing interests

J.P.d.M. is CSO of YouthBio Therapeutics, an advisor/consultant for the BOLD Longevity Growth Fund and NOVOS, and the founder of Magellan Science Ltd, a company providing consulting services in longevity science. D.B. is director of R&D at NOVOS Labs. The other authors have no conflicts of interest to disclose.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41514-025-00287-0>.

Correspondence and requests for materials should be addressed to Austin Parish.

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

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-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, 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 licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence 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 licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

© The Author(s) 2025