

Supplementary Materials for

How bad is human extinction? The psychology of existential risk

Stefan Schubert, Lucius Caviola, Nadira S. Faber

Correspondence to: stefan.schubert@psy.ox.ac.uk

This PDF file includes:

Links to the Full Studies, Material, and Pre-Registrations

Supplementary results for study 1

Table S1

Supplementary results for study 2a

Table S2

Supplementary results for study 2b

Table S3

Supplementary results for study 2c

Table S4

Supplementary results for study 3

Table S5

Supplementary results for study S1

Supplementary results for study S2

Supplementary results for study S3

Supplementary results for study S4

Supplementary results for study S5

Links to the Full Studies, Material, and Pre-Registrations:

<https://osf.io/pd9ca/>

Supplementary results for study 1 - beliefs about extinction

Most participants found it unlikely that humanity will go extinct within the next 100 years (1 = *very unlikely*, 7 = *very likely*, $M = 2.51$, $SD = 1.76$), and the more likely participants found extinction to be, the less bad they found extinction to be ($r = -0.21$, $P = 0.004$). Finding extinction bad correlated with other variables as shown in Table S1. In addition, a *t*-test with

1 gender revealed that there was no significant difference between how much men and women felt
 2 extinction needed to be prevented.

3
 4 Table S1. Correlations.

	1	2	3	4	5
1. Finding extinction bad	-				
2. Importance of preventing extinction	0.75***				
3. Personal moral obligation to prevent extinction	0.66***	0.78***			
4. Prioritize preventing extinction over other causes	0.49***	0.61***	0.61***		
5. Quality of the future	0.51***	0.44***	0.42***	0.30***	
6. Probability of extinction	-0.21**	-0.17*	-0.13	-0.12	-0.11
7. Happy person into existence good	0.22**	0.25***	0.22**	0.28***	0.18*
8. Optimism - pessimism	0.32***	0.30***	0.31***	0.23**	0.30***
9. OUS - impartial beneficence subscale	0.09	0.04	0.15*	0.10	0.05
10. OUS - instrumental harm subscale	0.05	-0.04	-0.05	-0.02	0.10
11. Cognitive Reflection Test	-0.03	0.06	0.13	0.13	-0.05
12. Gender	-0.09	-0.06	-0.09	-0.13	-0.11
13. Education	0.05	0.06	0.08	0.02	0.04
14. Income	0.05	0.02	0.11	-0.02	0.04
15. Age	0.14	0.12	0.15	0.03	0.05

16. Religiosity	0.20**	0.12	0.15*	0.03	0.07
17. Economic conservatism	0.14	0.06	0.08	0.02	0.19*
18. Social conservatism	0.17*	0.07	0.07	0.05	0.18*

p-values are marked as follows: * = significant at the .05 level, ** = significant at the .01 level, *** = significant at the .001 level.

Supplementary results for study 2a - UK sample

Finding the difference between outcomes B and C to be greater than the difference between outcomes A and B (finding extinction uniquely bad) correlated with other variables as follows (Table S2).

Table S2. Logistic regression with finding extinction uniquely bad as the predicted variable.

Measure	Odds ratio	<i>p</i> -value
OUS - instrumental harm subscale	1.04	.46
OUS - impartial beneficence subscale	0.94	.32
Cognitive Reflection Test (CRT)	1.16	.01*
Gender (male = 0, female = 1)	0.91	.63
Age	0.99	.13
Education	0.93	.30
Income	1.07	.17

p-values are marked as follows: * = significant at the .05 level, ** = significant at the .01 level, *** = significant at the .001 level. CRT and OUS were analyzed in one logistic regression, demographics in another logistic regression.

The utopia condition was called “the good future condition” in the pre-registration.

1 **Supplementary results for study 2b - US sample (Mturk)**

2
3 Finding the difference between outcomes B and C to be greater than the difference between
4 outcomes A and B (finding extinction uniquely bad) correlated with other variables as follows
5 (Table S3).

6
7 Table S3. Logistic regression with finding extinction uniquely bad as the predicted variable.
8

Measure	Odds ratio	p-value
OUS - instrumental harm subscale	0.84	.008**
OUS - impartial beneficence subscale	1.00	1.000
Cognitive Reflection Test (CRT)	1.23	.005**
Gender (male = 0, female = 1)	0.66	.005**
Age	1.01	.202
Education	1.04	.569
Income	0.98	.593

9
10 *p*-values are marked as follows: * = significant at the .05 level, ** = significant at the .01 level,
11 *** = significant at the .001 level.

12
13 Again, the utopia condition was called “the good future condition” in the pre-registration.
14

15 **Supplementary results for study 2c - Oxford students**

16
17 Finding the difference between outcomes B and C to be greater than the difference between
18 outcomes A and B (finding extinction uniquely bad) correlated with other variables as follows
19 (Table S4).

20
21 Table S4. Logistic regressions with finding extinction uniquely bad as the predicted variable.

1

Measure	Odds ratio	<i>p</i> -value
OUS - instrumental harm subscale	1.29	.052
OUS - impartial beneficence	0.85	.250
Gender (male = 0, female = 1)	0.47	.036*
Age	0.98	.413

2

3 *p*-values are marked as follows: * = significant at the .05 level, ** = significant at the .01 level,
4 *** = significant at the .001 level. Consistent with our other studies, OUS and demographics
5 were analyzed in two separate logistic regressions.

6

7

8 **Supplementary materials for study 3 - Effective Altruists**

9

10 As reported in the paper, almost all participants found the difference between extinction and a
11 catastrophe to be more relevant than the difference between a catastrophe and no catastrophe.
12 We therefore did not conduct logistic regressions with this variable and instead conducted
13 correlational analyses with the mean of the three extinction prevention questions across the
14 control condition and the utopia condition (Table S5). In addition, a *t*-test with gender revealed
15 that there was a significant difference between how much men ($M = 5.84, SD = 1.46$) and
16 women ($M = 4.64, SD = 1.79$) felt extinction needed to be prevented, $t(14.15) = 2.17, P = .04, d$
17 $= 0.79$. Overall, participants strongly favored preventing extinction both in the control condition
18 ($M = 5.54, SD = 1.55$) and the utopia condition ($M = 5.74, SD = 1.60$).

19

20 Table S5. Correlations with the mean of the three extinction prevention questions.

21

Measure	Coefficient	<i>p</i> -value
OUS - instrumental harm subscale	0.32	.007**

OUS - impartial beneficence	0.03	.793
CRT	-0.03	.822
Age	-0.26	.028*

1 *p*-values are marked as follows: * = significant at the .05 level, ** = significant at the .01 level,
2 *** = significant at the .001 level.

3
4
5
6

7 **Supplementary studies**

8
9 We conducted five supplementary studies (studies S1 to S5; total N=1625). Their results are
10 broadly in line with the findings reported in the main text.

11

12 **Study S1**

13 Study S1 (pre-registered) was identical to studies 2a and 2b reported in the main text, with the
14 exception that we only had two conditions in study S1: the control condition and the utopia
15 condition. The results were in line with the effects of the utopia manipulation in studies 2a, 2b
16 and 2c. Out of our final sample of 182 participants, large majorities ranked no catastrophe as the
17 best outcome and 100% dying as the worst outcome both in the control condition (92.40%,
18 85/92) and the utopia condition (90.00%, 81/90). Subsequently, we found, in line with our pre-
19 registered hypothesis, that the proportion of participants who found the difference between 80%
20 dying and 100% dying the largest was significantly larger ($\chi^2(1) = 29.439, P < .001$) in the
21 utopia condition (69.14%, 56/81 participants) than in the control condition (25.88%, 22/85).
22 Hence, the results are convergent with those of studies 2a, 2b, and 2c that are reported in the
23 main text.

24

25 **Study S2**

1 Studies S2-S4 (all pre-registered) employed another paradigm. Here, participants were asked to
2 compare painless extinction with a catastrophe killing 50% of all humans (after which point
3 “humanity recovers to its original size in a few centuries, and then goes on to live for a very long
4 time”).

5 In Study 2 (final N=230), a majority (142 out of 230; 61.7%) of participants rated
6 painless extinction as worse than a catastrophe killing 50%. We also asked three questions
7 measuring on a seven point-scale which of the two scenarios would be more important to
8 prevent. Using a merged variable over these three questions we found that participants valued
9 preventing extinction significantly more than preventing the non-extinction catastrophe ($M =$
10 4.52 , $SD = 1.97$), $t(229) = 3.99$, $P < .001$. These results contrasted with our pre-registered
11 hypotheses: at the outset of our studies, we expected participants to find human extinction less
12 bad than they actually do. Study 1 confirmed the finding of study S2, that people find human
13 extinction very bad—however, as subsequent studies (studies 2a-2c) showed, they do not always
14 view it as *uniquely bad* in comparison to non-extinction catastrophes.

15

16 **Study S3**

17 In study S3 (final N=427), we used the same dependent variables as in study S2, but here we had
18 three conditions: one where the long-term future was said to be good, one where it was said to be
19 bad, and one with no such information (control; same as in study S2). An ANOVA indicated
20 that the merged variable from Study S2 varied significantly across conditions, $F(2, 424) = 29.00$,
21 $P < .001$. In line with our pre-registered hypothesis, a *post hoc* Tukey test showed that each
22 condition differed significantly from the others, with participants in the utopia condition ($M =$
23 5.01 , $SD = 1.74$) valuing preventing extinction more (relative to preventing the non-extinction

1 catastrophe) than control condition participants ($M = 4.46$, $SD = 2.11$, $P = .049$), who in turn
2 valued preventing extinction more than participants in the bad future condition ($M = 3.31$, $SD =$
3 1.95 , $P < .001$). This finding further supports the results from studies 2a-2c reported in the main
4 text, which show that people find extinction worse relative to non-extinction catastrophes, the
5 better the future is predicted to be (conditional on non-extinction).

6

7 **Study S4**

8 In study S4 (final $N=288$), we again used the same dependant variables as in study S2, but added
9 a condition where the catastrophe/extinction affected animals (zebras) rather than humans. In
10 line with our pre-registered hypothesis, an independent-samples t -test on the merged variable (cf.
11 study S2) showed that participants valued preventing the extinction of animals more highly
12 (relative to preventing a catastrophe killing 50% of the animal population; $M = 5.24$, $SD = 1.76$)
13 than they valued preventing human extinction (relative to preventing a catastrophe killing 50%
14 of all humans; $M = 4.74$, $SD = 1.89$), $t(286) = 2.31$, $P = .022$. Hence, the results of study S4 are
15 in line with the finding from studies 2a-2c reported in the main text, that participants find
16 extinction worse (relative to non-extinction catastrophes) when the extinction/catastrophes affect
17 animals rather than humans.

18

19 **Study S5**

20 Study S5 (pre-registered) was an initial attempt to study whether participants deem human
21 extinction uniquely bad relative to non-extinction catastrophes. We used a paradigm that we in
22 hindsight recognized as overly complex and potentially confusing to participants, and therefore
23 not well-suited to test our research questions. Rather than directly and explicitly asking

1 participants about their judgments of the *relative* difference between extinction and non-
2 extinction catastrophes (as we did in our later studies studies), in study S5 we asked participants
3 (final N=398) to rate the badness of different catastrophes (resulting in the death of 25%, 50%,
4 75% or 100% of the population) and a non-catastrophe (killing no one) in *absolute* terms, with
5 the aim to indirectly infer the relative differences from the absolute values later. We later
6 realized that this method is not reliable, because participants may not realize what claims about
7 the relative differences their judgments of the badness of the different outcomes commit them to.
8 For that reason, we rather have to ask participants directly about the relative differences, as we
9 did in studies 2a-2c and 3. Our confidence in the data yielded with this paradigm further
10 decreased due to an overly high participant drop-out rate of nearly 35% (206 out of 604).

11 In one condition of study S5, we asked participants to focus on the short-term future, and
12 in another to focus on the long-term future (cf. the salience condition from studies 2a-2c in the
13 main text). We also had a control condition where we did not give any such instruction (cf. the
14 control condition from the studies in the main text). We found that in the control condition
15 participants did not deem extinction uniquely bad. That is, the participants indicated that the
16 difference in badness between the catastrophes killing 75% and 100% was not greater than the
17 difference in badness between the catastrophes killing 50% and 75% (in line with our pre-
18 registered hypothesis and the findings reported in the main text). We also found that none of the
19 two interventions introducing a certain focus on the future (short-term or long-term) changed
20 these results (in contrast to our pre-registered hypothesis and the findings on from the “salience
21 condition” reported in the main text). We refrain from interpreting these mixed results because of
22 our doubts about the paradigm that we used.