This database was created by Michael Aird, in relation to my work with Convergence Analysis. I intend for it to be a "living document", built up collaboratively by anyone interested in existential risk reduction.

Please add "comments" to suggest additions or mention places where you think I've misinterpreted or misrepresented an estimate, or where there would just be some other context worth noting. I'll check these comments regularly, and, where relevant, copy and paste them into the spreadsheet as additions.

See this post for context on this database

**Some things worth bearing in mind:**

Some limitations of existential-risk-related estimates, or explicit probability estimates in general, noted in this post (forum.effectivealtruism.org/posts/JQQAQrunyGGhzE23a/database-of-existential-risk-estimates), Beard et al. (https://www.sciencedirect.com/science/article/pii/S0016328719303313), and this post (https://forum.effectivealtruism.org/posts/KfqFLDkoccf8NQsQe/potential-downsides-of-using-explicit-probabilities).

These estimates are unlikely to be very independent; many estimators had probably seen each others' estimates, interacted extensively with other estimators, etc.

I don't provide much context or caveats for most of the specific estimates. It may often be worth checking the appendix of Beard et al. (https://www.sciencedirect.com/science/article/pii/S0016328719303313) and/or the original source.

It's often hard to be sure precisely what is being estimated, or what other conditions are perhaps being assumed. And I may sometimes misinterpret or misrepresent this; please make a comment if you think that that's the case.

I've organised these estimates by the broad categories of what's being estimated. Most estimates within each category are **not** of exactly the same thing; for example, they may differ in the timelines over which they're estimated, or in whether they're about existential catastrophe broadly or extinction specifically.

I've converted all estimates into percentages. Where the estimator expressed their estimate in another way, I've shown how they expressed it in brackets after the percentage. All bold is added by me.

**This database was last updated on:**

20 May 2020

This Sheet includes what I'm calling "existential-risk-level estimates", meaning estimates of things like the existential risk from something, the extinction risk from something, or how much something reduces the expected value of the future.

I'm not including here estimates of "less extreme" things, like major global catastrophes that don't cause existential catastrophe, some such estimates can be found in the Sheet "Estimates of somewhat less extreme outcomes". (Of course, such events may still be very extreme by regular standards, and I don't mean to imply otherwise.)

I'm also mostly trying to include "unconditional" estimates here, and "conditional" estimates in the next Sheet. For example, I'd include here an estimate of the chance of existential catastrophe as a result of nuclear war, but not **separate** estimates of the chance that a nuclear war, **if** it occurred, would cause existential catastrophe.

| Who is the estimator? | When was the estimate made/published? | What is the estimator estimating? | What is their estimate? | Num. Years | % Probability | Annualised % | Century % | Subjective Quality Weighting | Source | Have I properly read the source myself? | Is this estimate included in at x's appendix? | Other notes |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | "Total risk" (or similar) | | | | | | | |
| Toby Ord | 2020 | "Total existential risk" by 2120 | ~17% (~1 in 6) | 100 | 17% | 0.19% | 17.00% | 1 | The Precipice | Yes | No | |
| GCR Conference | 2008 | "Overall risk of **extinction** prior to 2100" | 19% | 92 | 19% | 0.23% | 20.47% | 10 | https://www.fhi.ox | Yes | Yes | |
| Will MacAskill | 2019/2020 | Existential risk in the 21st century | 1% | 80 | 1% | 0.01% | 1.25% | 1 | https://80000hours | Yes | No | |
| Ben Todd or 80,000 Hours | 2017 | Extinction risk "in the next century" | Probably **at above** 3% | ? | | INVALID! | INVALID! | 1 | https://80000hours | Yes | No | |
| John Leslie | 1996 | Risk of **extinction** over the next two centuries | At or above 30% | 500 | 30% | 0.07% | 6.96% | 1 | Leslie, J. (2002) | No | No | |
| Martin Rees | 2003 | "Odds that our present civilization on earth will survive to the end of the present century" | 50% ("no better than fifty-fifty") | 97 | 50% | 0.71% | 51.06% | 1 | Rees, M. J. (2003) | No | No | |
| Meteulous respondents | | "there be zero living humans on planet earth on January 1, 2100" | Median: 1%. Mean: 9%. | 80 | 1% | 0.01% | 1.25% | 3 | https://www.metac | Yes | Yes | That median and mean is as of 3rd July 2019. |
| Nick Bostrom | 2002 | "existential disaster will do us in" | Probably **at or above** 25% | ? | | INVALID! | INVALID! | | https://www.nickb | No | Yes | |
| Gott III | 1993 | | 5% | ? | | INVALID! | INVALID! | | Gott III, J. R. (1993a) | | Yes | |
| Wells | 2009 | Annual probability **as of 2009** of extinction | 0.3-0.4% | 1 | 0.35% | 0.35% | 29.57% | 1 | Wells, W. (2009) | No | Yes | |
| Simpson | 2016 | "a global catastrophic risk of 0.2% per year" | 0.2% | 1 | 0.20% | 0.20% | 18.14% | 1 | Simpson, F. (2016) | No | Yes | |
| Toby Ord | 2020 | its potential: achieving something close to the best future open to us | 50% (~1 in 2) | ? | | | | | The Precipice | Yes | No | |
| | | | | | Mean | 0.22% | 19.20% | | | | | |
| | | | | | Median | 0.19% | 17.57% | | | | | |
| | | | | AI | | | | | | | | |
| Toby Ord | 2020 | Existential catastrophe by 2120 as a result of "unaligned AI" | ~10% | | | | | | The Precipice | Yes | No | |
| Global Catastrophic Risks Conference | 2008 | Human **extinction** by 2100 as a result of "superintelligent AI" | 5% | | | | | | https://www.fhi.ox | Yes | Yes | |
| Survey of "AI experts" | 2017 | "Extremely bad (e.g. extinction)" long-run impact on humanity from "high-level machine intelligence" | 5% | | | | | | https://www.nickb | Yes | No | |
| Pamlin & Armstro | 2015 | "infinite impact" from AI within the next 100 years, which "refers to the | 0-10% | | | | | | Pamlin, D. & Armstro | | Yes | |
| Rohin Shah | 2020 | Chance that AI, through "adversarial optimization against humans only", will cause existential catastrophe | ~5% | | | | | | https://www.less | Yes | No | |
| Paul Christiano | 2019 | Amount by which risk of failure to align AI (using only a narrow conception of alignment) reduces the expected value of the future | ~10% | | | | | | https://sideways-l | Yes | No | |
| Buck Schlegris | 2020 | "the probability of AI-induced existential risk" (but from context, I believe this actually meant the probability of an AI-induced existential catastrophe) | 50% | | | | | | https://futurefift | Yes | No | |
| James Fodor | 2020 | Existential risk from unaligned AI over the coming 100 years | 0.05% | | | | | | Critical Review of | Yes | No | |
| Stuart Armstrong | 2014 | Chance of humanity not surviving AI | 50, 40, or 33% | | | | | | https://www.youtu | Yes | No | |
| | | | | | Biotech | | | | | | | |
| Toby Ord | 2020 | Existential catastrophe from "engineered pandemics" by 2120 | ~3% (~1 in 30) | | | | | | The Precipice | Yes | No | |
| GCR Conference | 2008 | Human **extinction** by 2100 as a result of "the single biggest natural pandemic" | 0.05% | | | | | | https://www.fhi.ox | Yes | Yes | |
| Toby Ord | 2020 | Existential catastrophe from "**naturally** arising pandemics" by 2120 | ~0.01% (~1 in 10,000) | | | | | | The Precipice | Yes | No | |
| GCR Conference | 2008 | Human **extinction** by 2100 as a result of "single biggest engineered pandemic" | 2% | | | | | | https://www.fhi.ox | Yes | Yes | |
| Millet & Snyder-Beattie | 2017 | The annual probability of an existential catastrophe arising from a global pandemic | 0.00005% (0.0000005)% | | | | | | Millett, P., & Snyd-No | | Yes | |
| Millet & Snyder-Beattie | 2017 | The annual probability of an existential catastrophe arising from biowarfare or bioterrorism | 0.00019% (0.0000019) | | | | | | Millett, P., & Snyd-No | | Yes | |
| Pamlin & Armstro | 2015 | "infinite impact" from a global pandemic within the next 100 years, whi | 0.0001% | | | | | | Pamlin, D. & Armstro | | Yes | |
| Pamlin & Armstro | 2015 | "infinite impact" from "synthetic biology" within the next 100 years, whi | 0.01% | | | | | | Pamlin, D. & Armstro | | Yes | |
| James Fodor | 2020 | Extinction risk from engineered pandemics over the coming 100 years | 0.0002% | | | | | | Critical Review of | Yes | No | |
| | | | | | Nanotechnology | | | | | | | |
| GCR Conference | 2008 | Human **extinction** by 2100 as a result of "molecular nanotech weapons" | 5% | | | | | | https://www.fhi.ox | No | Yes | |
| GCR Conference | 2008 | Human **extinction** by 2100 as a result of "the single biggest nanotech accident" | 0.5% | | | | | | https://www.fhi.ox | No | Yes | |
| Pamlin & Armstro | 2015 | "infinite impact" from nanotechnology within the next 100 years, which | 0.0100% | | | | | | Pamlin, D. & Armstro | | Yes | |
| Toby Ord | 2020 | Existential catastrophe from "other anthropogenic risks" (which includes but is not limited to nanotechnology) by 2120 | ~2% (~1 in 50) | | | | | | The Precipice | Yes | No | |
| | | | | | Nuclear | | | | | | | |
| Toby Ord | 2020 | Existential catastrophe from nuclear war by 2120 | ~0.1% (~1 in 1000) | | | | | | The Precipice | Yes | No | |
| GCR Conference | 2008 | Human **extinction** by 2100 as a result of "nuclear wars" | 1% | | | | | | https://www.fhi.ox | Yes | Yes | |
| GCR Conference | 2008 | Human **extinction** by 2100 as a result of "acts of nuclear terrorism" | 0.03% | | | | | | https://www.fhi.ox | Yes | Yes | |
| Ben Todd or 80,0 | 2017 | The chance of "a civilization-ending nuclear war" on one that causes | Probably **at above** 0.3% | | | | | | https://80000hours | Yes | No | |
| Dave Denkenberj | 2018 | "Reduction in far future potential per year (from the risk of) full vs <0.29% | | | | | | | | https://www.getin | Sort-of | No | |
| Anders Sandberg | 2018 | "Reduction in far future potential per year (from the risk of) full vs <0.001% | | | | | | | | https://www.getin | Sort-of | No | |
| Peter McIntyre | 2016 | Amount by which "the future potential of humanity" is reduced due t | 0.4% | | | | | | https://80000hours | No | No | |
| Alexey Turchin | 2008 | The risk of extinction due to the consequences of nuclear war, or as | in the order of 1% | | | | | | Turchin, A. V. (2008a | | Yes | |
| Pamlin & Armstro | 2015 | "infinite impact" from nuclear war within the next 100 years, which "ref | 0.005% | | | | | | Pamlin, D. & Armstro | | Yes | |
| | | | | | Climate change | | | | | | | |
| Toby Ord | 2020 | Existential catastrophe from climate change by 2120 | ~0.1% (~1 in 1000) | | | | | | The Precipice | Yes | No | |
| Roman Duda | 2016 | Amount by which "the future potential of humanity" is reduced due to the risk of "extreme climate change (+4°C) in the next 100 years" | ~0.1-2% | | | | | | https://80000hours | No | No | |
| Pamlin & Armstro | 2015 | "infinite impact" from climate change within the next 100 years, which | 0.01% | | | | | | Pamlin, D. & Armstro | | Yes | |
| Toby Ord | 2020 | Amount by which existential risk till 2120 would decrease if "we could just somehow have the next century but make it so that climate change wasn't an issue" | 3.1-1 percentage points | | | | | | https://80000hours | Yes | No | |
| | | | | | Natural risks (excluding natural pandemics) | | | | | | | |
| Toby Ord | 2020 | Existential catastrophe from supervolcanic eruption by 2120 | ~0.01% (~1 in 10,000) | | | | | | The Precipice | Yes | No | |
| Pamlin & Armstro | 2015 | "infinite impact" from a super-volcano within the next 100 years, which | 0.00003% | | | | | | Pamlin, D. & Armstro | | Yes | |
| Pamlin & Armstro | 2015 | "infinite impact" from asteroid or comet impact by 2120 | ~0.0001% (~1 in 1,000,000) | | | | | | Pamlin, D. & Armstro | | Yes | |
| Pamlin & Armstro | 2015 | "infinite impact" from "a major asteroid impact" within the next 100 yea | 0.00013% | | | | | | Pamlin, D. & Armstro | | Yes | |
| Toby Ord | 2020 | Existential catastrophe from stellar explosion by 2120 | <0.00000(1)% (~1 in 1,000,000,000) | | | | | | The Precipice | Yes | No | |
| Toby Ord | 2020 | "Total natural [existential] risk" by 2120 | ~0.01% (~1 in 10,000) | | | | | | The Precipice | Yes | No | |
| Snyder-Beattie, Ord, & Bonsall | 2019 | The probability that humanity goes extinct from natural causes in any | | | | | | | | https://www.natur | No | Yes | |
| | | | | | Miscellaneous | | | | | | | |
| Toby Ord | 2020 | Existential catastrophe from "other environmental damage" (non-climat | ~0.1% (~1 in 1,000) | | | | | | The Precipice | Yes | No | |
| Pamlin & Armstro | 2015 | "infinite impact" from an "ecological catastrophe" within the next 100 ye | 0.5% | | | | | | Pamlin, D. & Armstro | | Yes | |
| Toby Ord | 2020 | Existential catastrophe from "unforeseen anthropogenic risks" by 2120 | ~5% (~1 in 30) | | | | | | The Precipice | Yes | No | |
| Toby Ord | 2020 | Existential catastrophe from "other anthropogenic risks" by 2120 (see | ~2% (~1 in 50) | | | | | | The Precipice | Yes | No | |
| Toby Ord | 2020 | "Total anthropogenic [existential] risk" by 2120 | ~17% (~1 in 6) | | | | | | The Precipice | Yes | No | |
| Toby Ord | 2020 | Amount by which existential risk till 2120 would decrease if there definitely wouldn't be a great power war during that time | ~1.7 percentage points ("something like a tenth of the total risk over that time") | | | | | | The Precipice | Yes | No | |
| GCR Conference | 2008 | Human **extinction** by 2100 as a result of "wars (including civil wars)" | 4% | | | | | | https://www.fhi.ox | Yes | Yes | |
| Pamlin & Armstro | 2015 | "infinite impact" from "an uncertain risk" within the next 100 years, whi | 0.5% | | | | | | Pamlin, D. & Armstro | | Yes | |
| Dave Denkenberj | 2018 | "**Reduction in far future potential** due (to the risk of?) 10% agricult | 0.18% | | | | | | https://www.getin | Sort-of | No | |
| Anders Sandberg | 2018 | "**Reduction in far future potential** due to (the risk of?) 10% agricult | <0.0023% | | | | | | https://www.getin | Sort-of | No | |
| | | | | | How various actions may reduce certain risks | | | | | | | |
| Ben Todd or 80,0 | 2017 | **One plausible** amount by which "$100 billion spent on reducing extin | 1 percentage point | | | | | | https://80000hours | Yes | No | |
| Paul Christiano | 2019 | Amount by which "really realing" some portion of AI safety work could | 0.03 | | | | | | https://sideways-l | Yes | No | |
| Paul Christiano | 2019 | Amount by which "a marginal person" doing some portion of AI safety | 0.003% ("one in 20,000 | | | | | | https://sideways-l | Yes | No | |

This is for conditional estimates of existential risk (or similar), such as estimates of the chance that, if a nuclear war occurs, that would cause existential catastrophe.

| Who is the estimator? | When was the estimate made/published? | What is the estimator estimating? | What is their estimate? | Source | Have I properly read the source myself? | Is this estimate included in Beard et al.'s appendix? | Other notes |
|---|---|---|---|---|---|---|---|
| | | "Total risk" (or similar) | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | AI | | | | | |
| Survey of "experts" in the AI field | 2016 | The probability that the long-run overall impact on humanity of human level machine intelligence will be "Extremely bad (existential catastrophe)", assuming HLMI will at some point exist. | 10% | Müller, V. C., & B | No | Yes | This is the mean. According to Beard et al, the question was "4. Assume for the purpose of this question that such Human Level Machine Intelligence (HLMI) will at some point exist. How positive or negative would be overall impact on humanity, in the long run?" |
| Rohin Shah | 2019 | Chance that AI, through "adversarial optimization against humans only", will cause existential catastrophe, conditional on there not being "**additional** intervention by longtermists" (or perhaps "**no** intervention from longtermists") | ~10% | https://www.less | Yes | No | This is my interpretation of some comments that may not have been meant to be taken very literally. I think he updated this in 2020 to ~15%, due to pessimism about discontinuous scenarios: https://www.lesswrong.com/posts/7dwpN46HeTbPSv2km/rohin-shah-on-reasons-for-ai-optimism?commentId=n577gwG8DvRpakBmj Rohin also discusses his estimates here: https://futureoflife.org/2020/04/15/an-overview-of-technical-ai-alignment-in-2018-and-2019-with-buck-shlegeris-and-rohin-shah/ |
| Rohin Shah | 2019 | Chance that AI, through "adversarial optimization against humans only", will cause existential catastrophe, **conditional on "discontinuous takeoff"** | ~70% (but with 'way m | https://www.less | Yes | No | |
| Toby Ord | 2020 | Chance that we don't "manage to survive that transition [to there bei | ~20% | https://80000hou | Yes | No | This may have been specifically if the transition happens in the next 100 years; it's possible Ord would estimate we'd have a different chance if this transition happened at a later time. "Basically, you can look at my [estimate that the existential risk from AI in the next 100 years is] 10% as, there's about a 50% chance that we create something that's more intelligent than humanity this century. And then there's only an 80% chance that we manage to survive that transition, being in charge of our future. If you put that together, you get a 10% chance that's the time where we lost cont [For people who would disagree, a question] is why would they think that we have much higher than an 80% chance of surviving this 'passing this baton to these other entities', but still retaining control of our future or making sure that they build a future that is excellent, surpassingly good by our own perspective? I think that the very people who are working on trying to actually make sure that artific |
| | | Biorisk | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | Nanotechnology | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | Nuclear | | | | | |
| Toby Ord | 2020 | Chance that "a full-scale nuclear war in the next century" would "be | ~2% | https://80000hou | Yes | No | "I give existential risk over the next century from nuclear war at about one in a thousand. I initially thought it would be higher than that. That's actually something that while researching the book, thought was a lower risk than I had initially thought. And how I'd break it down is to something like a 5% chance of a full-scale nuclear war in the next century and a 2% chance that that would be the end of |
| Luke Oman | 2012 | "The probability I would estimate for the global human population of zero resulting from the 150 Tg of black carbon scenario in our 2007 paper" | 0.001-0.01% ("in the ra | http://www.overc | Yes | No | **I think** that this is Oman's estimate of the chance that extinction would occur if that black carbon scenario occurred, rather than an estimate that also takes into account the low probability that that black carbon scenario occurs. I.e., I think that this estimate was conditional on a particular type of nuclear war occurring. But I'm not sure about that, and the full context doesn't make it much clearer. |
| | | Climate change | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | Natural risks (excluding natural pandemics) | | | | | |
| | | | | | | | |
| | | | | | | | |
| | | Miscellaneous | | | | | |

This Sheet is for estimates that I don't count as "existential-risk-level estimates". These may be estimates of, for example, the likelihood of Spanish-Flu-level pandemics, nuclear war, catastrophes causing over a billion deaths, etc.

Most of these events are of course still "extreme" by most standards.

I haven't yet taken the time to include many estimates of this type, because such estimates are further from the key thing I'm most concerned about (existential risk), and because it seems that it's much less hard to find those than to find what I'm considering "existential risk estim...

Sources from Beard et al.'s appendix that would be relevant here:

Other relevant sources

1 https://www.metaculus.com/questions/2568/ragnar%25C3%25B6k-question-series-results-so-far/

12 Adam Gleave gave an answer to the question "what is the chance that advanced artificial intelligence poses a significant risk of harm?". But it's very unclear to me whether this means the chance that harm occurs vs the chance that a risk is posed, somewhat unclear what the latter would mean, and very unclear what level of harm is being discussed. So I haven't included the numbers here. https://aiimpacts. org/conversation-with-adam-gleave/

13 Kevin Esvelt answered a relevant question at the 21:53 minute mark of this video: https://youtu.be/BbOHQLrVSX4?t=1313

| Who is the estimator? | When was the estimate made/published? | What is the estimator estimating? | What is their estimate? | Source | Have I properly read the source myself? | Is this estimate included in Beard et al.'s appendix? | Other notes |
|---|---|---|---|---|---|---|---|
| | | "Total risk" (or similar) | | | | | |
| | | AI | | | | | |
| Adam Gleave | 2019 | Chance that AI safety is as hard as a (caricature of) MIRI suggests | ~10% | https://aiimpacts. | Yes | No | "So, decent chance-- I think I put a reasonable probability, like 10% probability, on the hard-mode M |
| Adam Gleave | 2019 | Chance that "AI safety basically [doesn't need] to be solved, we'll just solve it by default unless we're completely completely careless" | ~20-30% | https://aiimpacts. | Yes | No | |
| Rohin Shah | 2020 | "chance that the first thing we try just works and we don't need | ~30% | https://futureoflife | Yes | No | "There's some chance that the first thing we try just works and we don't even need to solve any so |
| Buck Schlegris | 2020 | Chance we have "good competitive alignment techniques by the tim | ~30% | https://futureoflife | Yes | No | "I haven't actually written down these numbers since I last changed my mind about a lot of the inp |
| Toby Ord | 2020 | "chance that we create something that's more intelligent than human | ~50% | https://80000hou | Yes | No | "Basically, you can look at my [estimate that the existential risk from AI in the next 100 years is] 10%...  Toby Ord: With that number, I've spent a lot of time thinking about this. Actually, my first degree wa |
| | | Biorisk | | | | | |
| | | Nanotechnology | | | | | |
| | | Nuclear | | | | | |
| Toby Ord | 2020 | "chance of a full-scale nuclear war in the next century" | ~5% | https://80000hou | Yes | No | "I give existential risk over the next century from nuclear war at about one in a thousand. I initially t |
| | | Climate change | | | | | |
| | | Natural risks (excluding natural pandemics) | | | | | |
| | | Miscellaneous | | | | | |

Right-hand numbered column:

14
15
16
18 (for "infinite threshold")
20
21
22
23
24
25
26
29
31 (for "infinite threshold")
32
33
34
35
36
37 (for "infinite threshold")
40
42 (for "infinite threshold")
43
44
45
46
47 (perhaps this in fact belongs in the main sheet; I haven't read the source)
48
49 (for "infinite threshold")
50
51
52
53 (for "infinite threshold")
54 (for "infinite threshold")
55 (for "infinite threshold")
57 (for "infinite threshold")
58
59
60
61
62
63 (perhaps this in fact belongs in the main sheet; I haven't read the source)
64 (perhaps this in fact belongs in the main sheet; I haven't read the source)
65 (perhaps this in fact belongs in the main sheet; I haven't read the source)
66 (perhaps this in fact belongs in the main sheet; I haven't read the source)
67 (for "infinite threshold")

| Source | Description |
| --- | --- |
| https://wiki.lesswron | Apparently during (or from?) 2011, various people were asked questions including "What probability do you assign to the possibility of human extinction within 100 years as a result of AI capable of self-modification (that is not provably non-dangerous, if that is even possible)? P(human extinction by AI \| AI capable of self-modification and not provably non-dangerous is created)". I haven't looked into these interviews yet. |
| https://forum.effectiv | Gregory Lewis "[threw] together a guesstimate as a first-pass estimate" of "the expected value of x-risk reduction by the lights of person affecting views". But I think he was just using "reasonable to conservative" assumptions in order to build an illustrative model of what the value of x-risk reduction would be under those conditions (see footnote 23 here: https://80000hours.org/problem-profiles/global-catastrophic-biological-risks/) |
| https://foundersplege.com/research/fp-existential-risk | John Halstead gives numbers related to total existential risk this century, existential risk from nuclear war, and existential risk from advanced machine intelligence. However, from a skim, I think these are a mixture of illustrative examples and just reporting on other estimates which I mention in this post. |
| https://www.getguesstimate.com/models/1176 2 | In Denkenberger's model, he presents a modified version of an earlier model from others, which is focused on AI. I haven't included that here because I'm not sure whether Denkenberger actually endorses the values he provides in it, or if he was intentionally being "conservative" (i.e., using high estimates of AI risk) to make it harder for him to "make the case" for his own, non-AI intervention. |
| https://papers.ssrn.com/sol3/papers.cfm?abstract_id=31580 83 | This paper may have relevant estimates; I haven't read it. |
| https://forum.effectivealtruism.org/posts/dFYDzrJspMnS0zTwP/cau sal-network-model-ii-findings | This post may have relevant estimates; I haven't read it. |
| https://www.youtube.com/watch?v=fOSp19eows and footnote 2 of this: https://80000hours.org/articles/extincti on-risk/ | Spencer Greenberg surveyed users of Mechanical Turk, and some EAs, and asked about their estimates of the chance of human extinction within 50 from the time the survey was taken (around 2017, I believe). I haven't included this as I'm less interested in the general public's estimates of existential risk. But perhaps I should include it anyway, or at least the EAs' estimates. |
| https://forum.effectivealtruism.org/posts/mKwTzX5XLGs94tvBUpub lic-opinion-about-existential-risk | "An MTurk study of people in the United States (N=395) found median estimates of 1%, 5%, and 20% for the chance of human extinction in 50, 100, and 500 years, respectively. People were fairly confident in their answers and tended to think the government should prioritize preventing human extinction more than it currently does." But I haven't read the post, and, as stated above, I'm not necessarily especially interested in the general public's estimates. |
| https://forum.effectivealtruism.org/posts/XXLf6FmWukknq3E6/are -we-living-at-the-most-influential-time-in-history-1 | MacAskill writes: "Quantitatively: These considerations push me to put my posterior on [the idea that we're in the hinge of history] into something like the [1%, 0.1%] interval. But this credence interval feels very made-up and very unstable." |
| https://forum.effectivealtruism.org/posts/bYYAAa5K4nHpRCPGQ/our -current-estimates-for-likelihood-of-x-risk?commentId=wpb9gihGRoQguRS | "Anders Sandberg's Flickr account has a 2014 photo of a whiteboard from FHI containing estimates for [various relevant statements/questions". But I'm not sure how seriously these estimates were meant to be taken. |
| https://www.getguesstimate.com/models/1176 2 | Denkenberger's model also provides estimates of how both ALLFED's work so far and alternative foods planning and R&D may protect against these reductions in far future potential. Perhaps that should be included in the section for estimates of "How various actions may reduce certain risks". |
| https://80000hours.org/problem-profiles/nuclear-security/ | My interpretation is that McIntye suggests that a "major effort at the problem" could recover 30% of the 0.6-6% of the "future potential of humanity" that is currently reduced by the risk of nuclear war. In other words, I interpret him as effectively suggesting that such an effort would "increase" the "future potential of humanity" from 94-99.4% of what it would be if there was no risk of such a war to 95.8-99.58% of what it would be. But I'm unsure if he really meant that, and I'm very unsure how to interpret that anyway (among other things, what classifies as a "major effort"?). Thus, I haven't included this in the section for estimates of "How various actions may reduce certain risks". But perhaps I should. |
| https://80000hours.org/problem-profiles/climate-change/ | My interpretation is that Duda suggests that a "major effort at the problem" could recover 50% of the 0.1-2% of the "future potential of humanity" that is currently reduced by the risk of extreme climate change. In other words, I interpret him as effectively suggesting that such an effort would "increase" the "future potential of humanity" from 98-99.9% of what it would be if there was no risk of such climate change to 99-99.95% of what it would be. But I'm unsure if he really meant that, and I'm very unsure how to interpret that anyway (among other things, what classifies as a "major effort"?). Thus, I haven't included this in the section for estimates of "How various actions may reduce certain risks". But perhaps I should. |
| The Precipice (Toby | He writes that his estimates already: "incorporate the possibility that we get our act together and start taking these risks very seriously. Future risks are often estimated with an assumption of 'business as usual': that our levels of concern and resources devoted to addressing the risks stay where they are today. If I had assumed business as usual, my risk estimates would have been substantially higher. But I think they would have been misleading, overstating the chance that we actually suffer an existential catastrophe. So instead, I've made allowances for the fact that we will likely respond to the escalating risks, with substantial efforts to reduce them. The numbers therefore represent my actual best guesses of the chance the threats materialise, taking our responses into account. If we outperform my expectations, we could bring the remaining risk down below these estimates. Perhaps one could say that we were heading towards Russian roulette with two bullets in the gun, but that I think we will remove one of these before it's time to pull the trigger. And there might just be time to remove the last one too, if we really try." One could perhaps interpret this as him saying that his estimate of existential risk by 2120 conditional on continued business as usual is 2/6, and his estimate of how much existential risk by 2120 we can reduce if "we really try" is ~100% of it (meaning ~17 percentage points). But it's not clear to me that these were actually meant as estimates, rather than somewhat poetic illustrations of certain points. |
| https://scripps.ucsd | "Researchers identify a one-in-20 chance of temperature increase causing catastrophic damage or worse by 2050". In places, that page seems to imply this is a 1-in-20 chance of extinction, but I think that's just using ambiguous language in a dramatic way, and I was just skimming. |
| https://futureoflife.or | Andrew Critch: "In the same way you can predict an eclipse, you can predict that an asteroid is almost certainly not going to cause human extinction this century. I would bet, 99.99% chance, that we will not be extinct from an asteroid impact." But I think he just threw that number out there to make a point, and I wouldn't be surprised if his real estimate would be notably different. It's also unclear precisely what he's predicting (the chance of extinction from an asteroid this century? the chance of extinction from an asteroid this century conditional on us avoiding other problems, like AI catastrophe? the chance of extinction from an asteroid ever?). |