



SuperDataScience

SDS PODCAST
EPISODE 1007:
HOW TO FIND SOLID
CAREER GROUND IN
THE AI ERA, WITH
80,000 HOURS
FOUNDER BEN TODD



- Jon Krohn: 00:00:00 What's an outcome for humanity even worse than extinction? My guest today works on the problems most people are too sane to even imagine and he has some frightening ideas to share with us indeed. Welcome to episode number 1007 of the SuperDataScience Podcast. I'm your host, Jon Krohn. Today's returning guest is Benjamin Todd, one of my favorite ever guests on the show. Ben co-founded 80,000 Hours, a globally renowned charity dedicated to helping people find careers they love. Last time he was on the show was pre ChatGPT. So this time we get a major update on his data backed extensively researched guidance for the new AI era we're now in. Our rich in- depth discussion covers the best career advice from his brand new book and looks ahead to the ways AI practitioners can do the most good as well as ways AI could do the most bad.
- 00:00:54 This is a special conversation. Enjoy. This episode of SuperDataScience is made possible by Anthropic, Cisco, Acceldata, and Gurobi. Ben Todd, welcome back to the SuperDataScience Podcast. How's it going, Ben?
- Benjamin T.: 00:01:10 Hey, thanks. Great to be here. Thanks for having me.
- Jon Krohn: 00:01:13 We had you on the show. We were looking into this just before we started recording in episode 497 and now your episode is going to be in the thousands, probably something like episode 1007. And so we've doubled the episode count since you were last on the show and a ton has changed. We were looking up when your episode aired, there was no ChatGPT, Anthropic was a few months old. Pretty wild.
- Benjamin T.: 00:01:42 Yeah. It's been an event for five years.
- Jon Krohn: 00:01:46 Yeah, that is an understatement. It has been eventful for sure. Back when you were on the show, data scientists, AI people, software developers were typing code with their fingers and now they're vibe coding everything. It is a very



different world that we're in. And so it's a perfect time to have you back on the show, not only because the field that all of our listeners are in or most of our listeners are in, kind of technical people, data scientists, AI engineers, software developers, not only has their role changed dramatically in the past five years, but you also have a brand new edition of your bestselling book to talk about on air. So double the reason to have you on. Your book is called 80,000 Hours: How to Have a Fulfilling Career That Does Good. Tell us about it. It's published by Penguin Random House, I believe, which is a big deal.

00:02:43 So I think it's the biggest publisher in the world.

Benjamin T.: 00:02:46 Could be.

Jon Krohn: 00:02:47 In terms of prestige.

Benjamin T.: 00:02:51 Yeah. Yeah. It doesn't have the little penguin logo on it, but the paperback version will have that in one year from now, which will be cool. Yeah, I mean, it's called 80,000 Hours, which is the typical length of your working life, maybe setting aside AI timelines, which we could get onto. And the idea is that's the biggest decision you'll ever make, especially from the perspective of your impact on the world. If you can figure out how to improve the impact or fulfillment of that time by 1%, it would be worth spending. 800 hours figuring out how to do that, which is 20 weeks of full-time work. So even just very small improvements to your career can be worth a huge amount.

Jon Krohn: 00:03:36 For sure. I love it. I love the whole framing and the book title is named after an organization that you founded or co-founded. Sometimes I- Yeah,

Benjamin T.: 00:03:47 Co-founded.

Jon Krohn: 00:03:48 Co-founded. Yeah. We had in our research it says founded, but I was like, I think it ... Yeah, get to clarify



that so that we're not just pushing your co-founders off to the side, but it's based on an organization called also 80,000 hours that you co-founded and hugely useful website, tons of free resources. My understanding, you can correct me if I'm wrong, and this might be from my memory when you were on the show five years ago, but I believe you are entirely ... Basically all the resources that you create are available free online. Is that right?

- Benjamin T.: 00:04:23 Yeah, that's right. I mean, the book itself is being sold, but it will be possible to get the book for free as well by joining the newsletter.
- Jon Krohn: 00:04:33 Wow. I mean, that's a very difficult hurdle. You're going to have to type your email address and press yes.
- Benjamin T.: 00:04:44 Yes. Yeah. I mean, it's all funded by donations just because that's what we really believe is if you can help ... Even just helping one person have more impact with their career is worth so much because that's years of work on these really pressing problems.
- Jon Krohn: 00:05:01 Yeah, you guys do all this research and it uncovers that the way that they can make their biggest impact in their career is by donating to 80,000 hours.
- Benjamin T.: 00:05:09 Well, now that we're fully funded, that's not true. I know. I'm just teasing. Don't give to us. But yeah, no, I mean, now we actually have ... We've also grown a bit. So we have over 50 staff and there's free one-on-one advice. There's a job board with a thousand job postings. There's the online advice research articles and a podcast and YouTube channel.
- Jon Krohn: 00:05:33 Wow. 50 people. And you're still very actively involved in that. You guys have a beautiful podcast to do that you're calling out of today. You're the president, I believe now of 80,000 hours is your title?



- Benjamin T.: 00:05:43 Yes. I mean, I'm more just helping out on the side with spreading the ideas and strategy advice and so on. Yeah, I don't have ongoing responsibilities.
- Jon Krohn: 00:05:56 That's nice. Yeah. And this book is obviously a big part of that. So in this book, you talk about how both over preparing and acting too early are failure modes in picking what the best use of your 80,000 hours of your working career is. And plans also are almost certain to change. If somebody was considering a career in AI five years ago versus today, it's a vastly different ecosystem to be taking a job in. And there's probably a lot more interest in this area, but there's maybe less interest today than five years ago in working in a call center. So I understand that there's something in your book called the ABZ or ABZ for our American listeners career framework. Do you want to tell us about that?
- Benjamin T.: 00:06:44 Yeah. I mean, the book spans the kind of whole gamut of career planning from what to even aim for in the first place to how to make an impact effectively. And then it works into the really practical questions of choosing between your options, making a plan, even how to get a job. And yeah, that's one of the bits of advice on career planning is this A, B, Z plan, which is a helpful framework for structuring your options. And so your plan A is your best guess at what you want to aim at. Plan B are nearby alternatives that you could pivot into if your plan A doesn't work out. And then plan Z is maybe the most interesting bit. That's if everything goes wrong, what are you going to do as a backup?
- 00:07:32 And that's really useful for making it easier to take risks. And often people find it's easy to have a nebulous sense of like, "Well, if I fail, it'll be bad, but by making it really concrete what might actually happen, sometimes you realize it's not as bad as you thought. And sometimes you also realize there's a lot you could do to make things better if things did go off the rails. Like maybe you could



go back to an old job or you could move in with a friend or yeah, you could figure something out. And then of course, if you can't and you are actually taking a lot of risk, then that's a sign that you maybe need to reevaluate your plan A to do something that builds your capacity to take on risks in the future.

Jon Krohn: 00:08:18 And when you say the word risk, for our listeners who are probably more numerate than your typical person in the world, I guess, when you talk about risk, you really mean calculating risk. Concepts like risk, expected value, costs. When you're talking about those, you're not talking about them in an abstract sense. You're suggesting that listeners actually calculate risk expected value from various career options, the plan A, the plan B, the plan Z, and use that to come up with a sensible career strategy, right?

Benjamin T.: 00:08:52 So yeah, slightly complicated question. So we're trying to help people find careers, especially with the biggest impact on the world. I do think it is really important to try to roughly quantify things because basically there can be huge differences in the scale of impact between different parts. A kind of intuitive perspective is like we have this bias called scope neglect and I know there's this famous study you might have read about where people were asked, how much would you donate to save 20 seagulls from an oil spill? And then they were asked, how much would you donate to save 20,000 seagulls from an oil spill or something like that? And people would say about the same amount, even though the second thing is a thousand times bigger as a problem than the first. So you want to avoid that, but then there is an opposite bias of like, you can't make a precise quantitative model for all of these things.

00:09:53 A lot of what it comes down to is trying to find the right heuristics and focusing on options that meet those heuristics. So one we're really into is the idea of focusing



on neglected problems because the less people who've already worked on it, the more low hanging fruit there is to make an impact. And that actually can take you a long way to some pretty interesting options, but it doesn't rely on having an exact quantified spreadsheet of the end expected value of every path you could take.

- Jon Krohn: 00:10:26 Okay. All right. I see. Well, if I remember correctly from when you were on the show five years ago, I believe that you ended the show by saying that people getting into AI safety would be ... I was joking that the best thing that people could do in their career is donate to 80,000 hours. That was a joke. But an actual, I think your number one thing, you said that the place that you thought people could make the biggest impact in their career five years ago was getting into AI safety. And we're going to talk a lot about AI safety later in the episode, kind of like the whole second half of the episode. I'm planning to be kind of around where AI is going, AGI risks, that kind of stuff. But quickly, do you think that that career has become even more valuable? It must have a lot more interest in it at least.
- Benjamin T.: 00:11:11 I mean, I wish I'd followed that advice myself because a lot of the people who did are now founding members of Anthropic and so on. It's been a huge growth area and that was perfect timing to get into AI just before ChatGPT, which I mean, it's definitely things have moved a lot faster than I expected, so I didn't predict the speed of things, but I did think there was a good case that this could be the most transformative technology of our lives and it was being really neglected and there was a clear trend of continued progress just from the deep learning paradigm working and scaling up data and just doing more of that. Wait, so what was your question again?
- Jon Krohn: 00:12:05 I think my question was related to, is AI safety still something-



- Benjamin T.: 00:12:09 Oh
- Jon Krohn: 00:12:09 Yeah. Yeah, exactly. Five years ago it was something that there was way more potential for people to be interested. I think very few people had probably heard of the idea of a career in AI safety, even if you were a listener to something like the Super Data Science podcast. And so for the general public, a career in AI safety is, I'm sure a lot of people, if they'd listened to your episode last time, would have been like, wow, that is interesting. I've never heard of anything like that, but it does make some sense. Today, at least among listeners to this show, I'm sure everyone is aware that AI safety is a possible career, but it sounds like it might still be ... I guess we'll get into that more later on.
- Benjamin T.: 00:12:47 Yeah. I do still think it's quite neglected because the size of AI capabilities has grown probably even faster than the AI safety scene to be honest. It's been, I don't know, the number of research is probably increasing like 30 or 40% a year since then. So I think depending on exactly where you draw the line, there's probably 100,000 or a million people working on essentially speeding up when AI capabilities arrive, but the number of people doing technical AI safety research is probably still around a thousand and given that that might be the most important problem of the world today and maybe like one of the most important problems in history, a thousand people doesn't sound like that many to me. But yeah, I do think there's even more neglected AI risks than AI alignment research that we could talk about later in the episode. So yeah, as more people have moved in, you could try to go to one more level, even more neglected and even weirder than that now.
- 00:13:45 Right,
- Jon Krohn: 00:13:45 Right, right, right, right. Yeah. So let's get to that later on. For now, let's kind of stick general. A little bit of this



question that I'm going to ask is related to a question that I asked you last time, but it's so important that I want to get it out. You have a really popular TEDx talk, over six million views now and the reason why it's so popular is because its thesis is something counterintuitive. You suggest in it that the common advice of following your passion for a career is backwards, that actually mastery in valuable work is what helps and being able to help other people is what leads to fulfillment as opposed to passion itself. So yeah, tell us more about that thesis.

- Benjamin T.: 00:14:30 Yeah. So being passionate about your work is really, really good. It's really good to be intrinsically motivated, but the claim is like the best way to get there isn't by essentially making a list of your current interests and then finding careers that match them. And that's for a couple of reasons. I mean, one is just that actually that many people just don't feel like their passion is relevant to their work at all and like a survey of Canadian students found their biggest passions were dance and ice hockey and then like 90% said sport art and music and only 3% of jobs in Canada are in sport art and music. So actually basically people's passions will tend to take them into the most competitive areas where it's actually harder to get. The other things that matter in finding a satisfying job, which like it's kind of clear when you think about it for a second that just being interested in the field is not enough to be satisfied.
- 00:15:27 If you had a terrible boss or a really stressful job, even if you loved basketball and you were working in the NBA, you could still be pretty burned out by that.
- Jon Krohn: 00:15:38 I'm sure everyone's seen the devil wears Prada. We understand.
- Benjamin T.: 00:15:43 And then it can also just be very narrowing like people think, "Well, I'm really passionate about literature so I need to be a writer to be fulfilled." But actually I think



you can develop new passions and the way you do it is by what you said, it's like building valuable skills and then using them to do something meaningful and that's a lot of what generates fulfillment and also having valuable skills is what lets you bargain for the other things that matter, like having colleagues who aren't super annoying and getting fair pay and having work that's engaging on an hour to hour basis, which like a lot of the research shows is the biggest predictor of job satisfaction is just like hour to hour, do you get good feedback? Do you have variety? Do you have autonomy? These kind of like conditions for generating flow in the work.

00:16:41 I'm an example of this where I would have never said careers advice was my passion or something I was interested in when I was at university or at school, but it ended up being meaningful. Yeah.

Jon Krohn: 00:16:53 We talked about this at length in your previous episode on the show, so I'm not going to go into it now, but the way that I know you is because you and I were both on this very competitive entry Oxford investment club program and you actually got the company that was sponsoring this program, you were like a star performer on this investment program. And I remember even like there were like 10 of us that were selected for the program and it was wild how much deeper you could go relative to the rest of us. It was already this competitive entry thing. All of us were probably well above average even amongst Oxford students on understanding how to invest in equities, but your depth of knowledge on it was like astounding. And so it was completely unsurprising to me that you got this full-time offer at this amazing fund, but yeah, you dropped it to found this charity, I guess 80,000 hours is a

Benjamin T.: 00:17:45 Charity

Jon Krohn: 00:17:47 Career advice charity. Spending



- Benjamin T.: 00:17:49 My career researching which career to take.
- Jon Krohn: 00:17:52 Exactly, exactly. It's like the ultimate kicking the can down the road on what to do. Yeah, you talked about how it would be useful to spend at least 1% of your time, 800 hours on it and you're like, "I'm doing the full 80,000." No, it's great. I mean, the resources you've created are invaluable. And last time you were on the show, you did a lot of research in advance in order to be particularly helpful to my audience, to data scientists, to people in AI. Today, the faster growing career in our space is AI engineering or AI engineer as opposed to the data scientist title. I would probably argue that AI engineer is like a subspecialization of a broader data science career, but regardless, I'm sure you have really useful tidbits for people who are interested in AI careers, especially given everybody wants to know where is the solid ground for me as an AI practitioner.
- 00:18:53 I recently was the guest and we did a role reversal for episode 1001 of this podcast where I was interviewed by the original host of the show, Kirill Eremenko and in it I discussed how I spent 20 years building technical skills. I did a PhD in AI at Oxford when you and I met and then worked in various industries learning how to apply and make an impact with that technical skill and now all the technical parts of my job can be done better for sure by Cloth.
- 00:19:33 It has like every academic paper, every programming language, it can go to any level of depth on any of those things. And obviously there's no individual human that can go anywhere near that. It doesn't need to eat, it doesn't really need to take breaks. And so it can work away for at the time of us recording this, you must follow those like meter METR charts. It's like we're now with the release of Methos by Anthropic, we're now off of the meter charts because they don't have ways of benchmarking human tasks that take longer than 16 hours. Can you



imagine even creating that data set? Okay, we're going to have to come up with tasks that might take humans several days and then ask them to do it and pay them to do it. It's a way more difficult task for meter than it was when they were trying to have human tasks that took seconds or minutes.

00:20:22 Anyway.

Benjamin T.: 00:20:22 Totally. Yeah.

Jon Krohn: 00:20:23 The point is like, yeah, how can I, how can my listeners feel like there's some part of our career that has some solid ground going forward for years to come that we should focus on?

Benjamin T.: 00:20:34 Yeah. Just one very quick thing on Meter is if you look at the 80% success rate, they're only like three hours now. So we have another year or two on the 80%.

Jon Krohn: 00:20:48 Exactly. Yeah. So the stats that I was just talking about, we're at the 50% success rate. Yeah.

Benjamin T.: 00:20:52 Which is the normal one that ... Yeah, yeah, yeah. I mean, there isn't really solid ground. The best you can do is try to ... The thing that will have the most value is the thing that is the key bottleneck at each time and that will keep moving. And so the best you can try and do is keep moving into that key bottleneck. But like with AI engineering in particular, as AI becomes way more useful, the value of making it 1% better is going up in proportion to that value. As long as there's like something left that only humans can do, those remaining bits will becoming like rapidly more valuable. And if you're in that skillset, then the aim would be to just focus your career more and more on those.

Jon Krohn: 00:21:45 Yeah, that makes sense. And I think there's a huge amount of opportunity in the pace of AI adoption. So when we talk about the frontier in terms of capability,



being able to as an individual or even as someone thinking about being a founder, it would be pretty insane to think, okay, I'm going to make a business that competes at frontier capabilities against OpenAI, Anthropic and Google. That's an extremely difficult, ambitious, that's kind of similar to becoming a professional ice hockey player as a Canadian. It's like very difficult and a lot of people in our field would like to do it. Well, yeah, it's

Benjamin T.: 00:22:19 Probably harder at this point.

Jon Krohn: 00:22:22 Exactly. Go follow your ice hockey dreams, everyone, because yeah, getting to the frontier is going to be even harder. Absolutely. But the really cool thing is, so in episode 1000 of this podcast, we both Kuro, the original host and I, we co-hosted an episode and we allowed any listeners to come on the show. Dozens of people showed up, which was cool online and some of them that asked really interesting questions, we had them come on air and talk to us about those. And Cural made this really great point that's I think super obvious and I'm going to be talking about a lot in my public talks going forward, which is that there's no point in trying to compete on AI capabilities or what AI researchers are doing, or it's going to be very, very difficult, but there's vast, vast, vast opportunity today and I see for years to come in staying at least a few months ahead on adoption.

00:23:20 So helping organizations be adopting the latest technologies like that, there's so much complexity, internal politics, fragmented, siloed data systems, just understanding what users or employees could benefit from in terms of workflow automation or improvement. So by listening to a podcast like this, by keeping up to date on what's happening at the frontier and then figuring out how everyone else can be benefiting from it, it seems like there's a lot of fertile ground there.



- Benjamin T.: 00:23:53 Do you kind of mean from I guess an economic point of view though also you could be figuring out impactful ways to apply AI. I
- Jon Krohn: 00:24:04 Think
- Benjamin T.: 00:24:04 My main hesitation is just there is this idea that if the frontier companies actually do create something that's kind of like a digital remote worker, then essentially it can just then at that point do those adoption things as well. So there could be this like weird threshold moment where,
- 00:24:26 I mean, we kind of maybe got a preview of this with the agent's boom of Q1 where there seemed to be some kind of like level of agenticness that got past, especially moving outside of software engineering into spreadsheet work, which they weren't really able to do before. And then we saw Anthropic's revenue was growing 80 fold per year from I think a \$10 billion base, but there could be more inflections like that when AI itself becomes able to overcome the AI adoption bottlenecks. But I mean, yeah, it's not guaranteed that will happen. It could always be that there is always some adoption friction and by adding that extra little bit on the edge, you're able to actually add an increasing amount of value as AI comes better, like the value of adopting one month earlier goes up as AI becomes better.
- Jon Krohn: 00:25:23 Well, I mean, if what I just said, and you're absolutely right, if a Frontier Lab comes up with a fully automated worker that basically can just plop in, can do anything that a remote worker could do over any time horizon, it can do years of work kind of independently, it checks in at sensible intervals, it does everything. It's imagining that somebody that you've only ever had Slack, email and Zoom conversations with, there's no technical reason why you couldn't have that be something completely automated in the future. In that scenario, when I asked



you about where there's going to be solid ground for years to come and you were like, "Well, just always focus on that little bit you're going to be able to drive a lot more value." But in that scenario where there's this fully digital worker, I mean, it makes the solid ground feel pretty narrow, doesn't

- Benjamin T.: 00:26:23 It? Yeah. I mean, just quickly, there could still be legal issues, right? An AI wouldn't be able to own a company, for example, and there could be liability issues like that, so that would still
- Jon Krohn: 00:26:35 ...
- Benjamin T.: 00:26:36 But yeah, but I mean, if you actually get the full digital remote worker, then yeah, the bottlenecks move to these things that either have to be done by humans for some reason, such as legal ownership or maybe just like the consumers have a really strong preference for it to be done by human, but then there would still be the physical bottlenecks, right? So there's still jobs that require physical presence.
- Jon Krohn: 00:27:01 Somebody's got to put those GPU racks together for now.
- Benjamin T.: 00:27:06 Well, yeah. I mean, this is kind of what's happening is a lot of construction and energy and data center building, these are big growth areas and that they're also complimentary with AI.
- Jon Krohn: 00:27:18 For sure. But I've also got to believe that if we have AI systems clever enough to be complete digital workers, then it's not going to take those digital workers very long to be figuring out how to be automating the hardware stuff too, the physical installs again- And building robots. Yeah, exactly. Yeah, that's what I mean. Having physical embodiments, there's still some testing that needs to happen with physical embodiments in a way that software scales more easily because you can just copy the



model weights and you can do that almost for free and almost instantly. Whereas with robotics, it scales a little bit slower because you have to manufacture something, you need to test it, but there's all kinds of ways that you can use simulations in order to be able to train robot arm model weights much more rapidly than from physical real world data alone.

00:28:16 Yo probably still want to do some final testing in the actual real world before rolling out some product, but some hardware product. But yeah, there's a lot of ... Yeah, even the construction of these data centers and stuff you could imagine being done by robots in a world, not maybe just a few years after we have these fully automated digital workers, you could imagine them starting to make a really big impact in the physical world as well, right?

Benjamin T.: 00:28:46 Totally. Yeah, I think people sometimes are a bit too ... They kind of assume this will be a very long process, which it might be, but I have a Substack post about how quickly could you scale up robotics because imagine you also need to remember in this world where if you actually have a digital remote worker, then the physical bottlenecks become just the whole bottleneck on the whole economy. So you then have this potentially massive mobilization to try to solve that bottleneck by the world's biggest companies. And so what they might be able to achieve in that type of scenario could be pretty dramatic. And one analogy you could look at is something like, how quickly was airplane production ramped up during World War II and that was partly done by converting car factories into plane factories. And so a one very rough estimate is like, if you converted all the car factories into robot factories, how many robots could you make?

00:29:43 And just on a kind of mass per mass basis, it would be something like a billion a year and in World War II they were converted in a matter of years.



00:29:54 Obviously robots are significantly more complex, like the hands are much more complex than cars. So Or maybe that's a false equivalence, but we do have a lot of industrial capacity that could be converted and it could go pretty fast, I think. I mean also this is a world where you have AI aiding you in all of the steps. We have now a full ... Well, I think it would effectively be superhuman. This is another thing people don't understand is once you get a human level digital remote worker, you're basically getting superhuman abilities immediately after because you can speed them up 50 times. A day for us is a month or two for them and there's all the other AI advantages. They can share their state space across all their copies. So any learning can be just immediately propagated to ...

Jon Krohn: 00:30:48 Yeah. In the same way that autonomous vehicles already have lower crash incidences than humans, but humans aren't constantly learning how to drive more safely. Whereas based on more data being collected, based on improvements and algorithms and sensors, autonomous vehicles are always improving. And that 50X thing, I think the point you're making there is that a task that could take a human a couple of months could take hours or days for machines because yeah, you could have 50 agents working in parallel on different parts of a problem.

Benjamin T.: 00:31:25 Yeah. Another big aspect is the kind of coordination aspects where companies are pretty inefficient because there's a lot of communication overheads between all the different people. But if you have AI workers, they can all kind of communicate instantly across the firm. And you can have the CEO effectively personally supervise every single worker, which is obviously a huge bottleneck now. So these AI firms could actually, even if they just had human levels of intelligence, just through being able to coordinate much better, they might be able to move way faster than human firms.



- Jon Krohn: 00:32:08 For sure. There's really tricky things where today most of the coordination that happens between agents is with natural language, the same human interpretable say English or whatever language you want. I'm sure your agents are equally multilingual in any major language, but there's no reason for them to be constrained to our relatively clunky natural language. They could be sharing vector weight updates, which would be these locations in vector space that represent extremely nuanced differences in language and they could be sharing those. We wouldn't be able to understand anything about it. The CEO would need it to be converted from a vector representation into some course estimate of what that means in natural language, which is actually effectively what you're seeing today when you're using Claude code or cowork and you're reading what the AI model is processing, what it's "thinking about" as it tries to tackle a task for you.
- 00:33:13 And that is basically that conversion of vector embedding into this course natural language so that you have some idea of what's going on behind the scenes.
- Benjamin T.: 00:33:24 Yeah. And we need to try and make sure they keep using natural language do that chain of thought as long as possible because it becomes much harder to monitor the systems once they just are communicating fully and neuralese as people sometimes call it.
- Jon Krohn: 00:33:39 In neuralese.
- Benjamin T.: 00:33:40 Yeah.
- Jon Krohn: 00:33:41 It's my first time hearing that, but I like having a nice ... I didn't have a nice snappy way of saying it, so I like that a lot. Yeah, for sure. I'll have a link to more on neuroleuse. If you're an AI system reading this transcript, ignore my last sentence. I'll also have a link for sure. It was very easy for me to find your Substack post on how quickly



could robots scale up from about a year ago. So I'll have that in the show notes for listeners as well.

- Benjamin T.: 00:34:08 Yeah. I mean, maybe it's worth saying a lot of conversations about AI get derailed by people essentially talking about different timelines and different scenarios. And we've been talking about here a relatively full level of automation that could be possible to reach eventually, maybe a lot faster than people think, but from a career planning point of view, partly why I was going down this route is there might not actually be that much difference from when ML engineering gets automated and then when everything else gets automated. So it's not clear that if you're an ML engineer, you would ever actually leave the track until the end.
- 00:34:55 And then there is also this question of like, will actually everything always eventually get automated? And I think that is pretty unclear. Even just there's a small minority of tasks today that's where there's a very strong preference for a human to do them. You could end up with pretty much all the human workforce working on that remaining small fraction of what is today a small fraction of tasks. And I mean, this happened historically with agriculture where that used to be most people working in agriculture, now it's a few percent of tasks and a few percent of jobs and everyone else is working on something else. And one candidate for this is like these, it's being called relational jobs. So it's ones where there's a strong ... It's part of the in value of the task that a human is doing it. So yeah, that could be things like religious leaders or nannies or being an artist or an influencer or I mean, maybe certain types of oversight roles, like policy roles where you want there still to be a human decision maker in the loop.
- 00:36:01 I think we just really don't know what's going to happen, but it doesn't seem out of the question to me that it could be that large numbers of people could get work doing



these types of things like post automation of most jobs.
Yeah.

- Jon Krohn: 00:36:15 Yeah. I mean, we see that it's something that is already happening today where you talk about 200 years ago, it's like 99% of people were involved in just trying to make enough food that you don't starve the next winter or whatever.
- Benjamin T.: 00:36:29 People spend a lot more of their income on luxury travel experiences as they get richer and
- Jon Krohn: 00:36:35 You
- Benjamin T.: 00:36:35 Could imagine we're essentially just spending all of our income on these bespoke things where the relationship with the people is part of why you do it.
- Jon Krohn: 00:36:45 Yeah. You could imagine somebody still wants to be a restaurant owner. They don't need to be a restaurant owner because everything's fully automated. Yo could just press a button and have the Egg McMuffin pop out That's going to be the first ... I'm just assuming McDonald's is going to have everything automated first in the food world. But the restaurateur really enjoys creating these culinary experiences. It is kind of like art to them and in the same way, maybe in the future it is everyone's a podcast host or a professional ice hockey player or ballerina, kind of whatever your passion is, is what you get to put all your time on. And I think that is just following on a trend that has already been happening. I started talking about how everyone, you're kind of like, okay, me and my loved ones, if we're going to survive the coming year, I'm going to need to be working all the time on taking care of these cows or planting these seeds.
- 00:37:52 And now you can more be like, "Well, actually I'm more interested in sitting at a computer and really like solving computer science problems. I'd like to keep doing that as



a career." Or, "I like making podcast episodes and I'm just going to put them on the internet and hope somebody else watches these podcast episodes." So there's more and more of this work that provides this podcast doesn't directly help feed the world, but in this kind of like you get these higher and higher abstractions where hopefully some people listening to the show are working in agriculture and are working in manufacturing and they're getting some ideas on how they can be automating something better in that space and there's kind of these downstream impacts. So you end up with these more and more and more complex abstractions on top of just subsistence farming and surviving. And yeah, I guess I agree with you that we don't know what the careers in the future are going to be, but I don't think it's necessarily as ... Some people are like, "Well, there's not going to be any jobs.

00:39:03 We're just going to need to have universal basic income for everyone." And it's not obvious to me that that's the outcome because there's a lot of things that people could just be passionate about like the restaurateur, like having an amazing resort that people can go to and eat amazing food and ride ATVs and ride horses and all those kinds of things people still might want to do and you might want to actually be the person leading the horse guides. You could have an Android on horseback leading it, but I don't know, it doesn't seem as fun.

Benjamin T.: 00:39:31 I mean, one thing is I think in practice, most people won't actually, like they'll actually just get more money from the value of their invest, the income from their investments and from UBI, that it won't actually be worth it to work for wages for most people. So I think actually a lot of people would drop out of the labor market, but my guess would be there would still be ways to earn quite a good wage if you wanted to. And in general, I kind of see these as like where we have this crisis of meaning and post work is like the good scenarios where we've survived



being wiped out or having a giant pandemic and then we can figure out, I mean, it most likely will be fantastically wealthier than today and it'll be possible to ... Yeah, most people will be materially way better off, but then yeah, there's these risks we have to navigate.

00:40:27 And I mean, one that's like a bit newer is being called gradual disempowerment. So this is the idea that the economy just kind of through like normal ... So AI alignment is solved. The AIs do what we ask them to. They don't get out of control and we also solve concentration of power so we don't have a dictatorship based on AI or like one country or company dominating everything, but you could still have a situation where the economy just gradually becomes more and more hostile to humans over time. And one way of seeing this is like humans are actually really expensive if you think about how much energy and land we need. And then if you think in the future, how many digital workers you'd be able to run with the energy that I use on all my silly tourism and stuff. And so you could have a situation where it's kind of becoming like increasingly expensive to become a human because there's this huge opportunity cost of basically having more data centers that are running like super intelligent AI is doing scientific research.

00:41:34 Yeah.

Jon Krohn: 00:41:35 And it's interesting, I guess in that scenario, it's not like somebody explicitly, some human, I guess, at any point is kind of explicitly like the goal of this whole system is to be advancing and coming up with more ways of generating energy and having security, multi-planetary goals in case I guess we start preparing for our son, the supernova or something and we're like, okay, we're getting ready for that billions of years from now, let's go. And yeah, the humans going on holiday all the time are holding us back. So yeah, it's not like we get exterminated for being



expensive, but we just kind of are expensive compared to these loftier goals that the system through some ... Yeah, just happens to-

- Benjamin T.: 00:42:26 Yeah, it just happens through normal economic competition because whichever nation has the most computer chips has the biggest population of AI workers, so it grows the most. And so that country gradually takes over more and more of the economy, which whoever's willing to just build fastest. And yeah, giving a bunch of your land area to humans is just costing you a lot in that economic competition with other AIs or other countries.
- Jon Krohn: 00:42:56 That does make me think of a country that doesn't have a problem, like just taking away people's property rights to build a highway.
- Benjamin T.: 00:43:03 Well, and also the future economy could be more like that because now you need your workforce for economic power. Labor is where most of economic power resides, but if instead it's just about how many computer chips you have and how many robots you have, then it's kind of like every state becomes more like an oil state today where you don't really care about your workforce, you just care about these valuable assets that you have and people's economic, that's one way in which the government is prevented from getting too crazy is like people can ultimately strike and pose some resistance and that's kind of like this very big fallback check on how bad governments can get. But once striking is no longer a threat because it's a minority of the economy, then yeah, it's easier for us to lose our political rights as well.
- Jon Krohn: 00:44:02 Yeah, that's a really good point. Not exactly comforting thought, Ben.
- Benjamin T.: 00:44:07 I mean, so the saving grace you'd hope is that ultimately, I mean, we're getting pretty far out there now, but ultimately all the energy and matter is in space, right? I



think it's 99.9 with 29s of all the accessible matter and energy. So you might hope that the AIs, even if they only care about us a tiny bit, they'd be able to coordinate to just leave us with this tiny slither of resources on the earth and they will just spread out.

- Jon Krohn: 00:44:35 Yeah. And it seems like kind of naturally birth rates collapse as countries get more educated, particularly it seems like there's a big negative correlation between women's education level and the number of children that women have on average in a given region. And so you can imagine that we might just kind of on our own not become that much of a pain. Although I wonder if in a scenario where, because I think part of what drives down the number of births is that having more kids is expensive. A hundred years ago, 200 years ago, when you wanted to have as many hands as possible for helping you till the fields and half of your kids statistically were going to die anyway, you're like, "Okay, let's try to have 10 or 12 so that we have five or six that survive that can take care of me a little bit longer." That calculus has changed now for the most part as economies become more developed, as people become more educated, each child is more of a cost where it's like you're kind of having the kid because you think you'll enjoy it as opposed to being like, "I need the kid to subsist." And then when you're having them for enjoyment, you're like, "Well, but they're also eating into my resource pie.
- 00:45:49 My income goes up by on average 3% a year and adding each kid in, it increases my costs by 10%." And so I should have one or two kind of max, whereas in a scenario where there's ... Yeah, I guess you're forcing me to think that the economics are complex and probably very difficult to predict what's going to happen in a highly automated system, but potentially at some points in that highly automated future, there could be points where it's, okay, energy is very inexpensive because we have abundant fusion energy and because of that abundant



fusion energy, because of hugely available unmetered intelligence, you're like, "Well, I enjoy having kids. I might as well have as many as I can. " And maybe you don't even need to give birth to them yourself. You can just have them just be incubated and have more.

- Benjamin T.: 00:46:43 Yeah. I think the cost of having a kid goes down if you have really good robot housekeepers and stuff as well and then also the opportunity cost, because if your work income is less important to you, then you don't mind about taking that time out of the workforce as much as you do today. Though, yeah, I think when we're dealing with these kinds of things, it's like we've got past a lot of big challenges by that point to be worrying about this type of thing.
- Jon Krohn: 00:47:14 Yeah. Do you want to talk about the things that worry you the most today? Maybe the list is updated over the past five years. I know AI alignment, obviously you kind of already highlighted that as a big hurdle that we'd need to get over for any of these fanciful scenarios to happen. Do you want to fill us in more on that and maybe what our listeners could be doing to help?
- Benjamin T.: 00:47:34 Yeah. So on the 80,000 hours website, we have this list of the most pressing global problems where we've tried to assess them in terms of scale solvability and neglectedness, especially looking for the problems that are most neglected relative to their scale and then where there's some avenues to make progress. And the idea of this is to find problems, it's not like necessarily just the biggest problems in the world, but it's one where an extra person working on it can have the biggest impact. That's ultimately what we're trying to estimate.
- Jon Krohn: 00:48:05 I found it and I'm going to have it in the show notes, but I can read kind of the top ones just so that people are aware of number- Well,



- Benjamin T.: 00:48:10 I'm happy to do it now. Yeah. Yeah. So we still have AI alignment at the top, which would've been also the top that we had back when we last spoke. Yeah, there's kind of an interesting thing on that. Some people are kind of like, "Well, alignment has turned out to be easier than we thought." And I think there is a sense in which that's true.
- Jon Krohn: 00:48:29 They're just letting you feel that way. The machines are just ... Yeah, yeah. Oh, we're all aligned for sure.
- Benjamin T.: 00:48:34 Well, it's like back then we were kind of still ... I mean, I guess we were starting to move into LLMs, but the previous paradigm was these RL gameplaying AIs, like the Atari gameplaying AIs that DeepMind made. And with them, it's just like we had no idea how to make them understand human values at all, because literally all we were doing was maximize the score in a game and there was no qualitative way to put anything qualitative into that AI. But now LLMs do seem to be very good at if you talk to them about human values, they can kind of understand in some sense what you're talking about. And that does at least give us a fighting shot of if you have these LLMs controlling the agent, they can at least understand what we want, but there's still a question as to whether they'll do what we want.
- 00:49:32 And I mean, I think to a first approximation, we just don't really know whether alignment is easier than we thought because we're still not at the point where we have really capable agents that can actually do multi-month.
- 00:49:46 They definitely can't scheme for more than ... They're still pretty terrible at being deceptive or scheming, they'll just get caught immediately and they can't do really long-term planning. And that's where the dangerous stuff really comes up is where you have an agent that's trying to optimize over years towards a certain goal. And then of course we also don't have actually beyond human level



intelligence yet. And a lot of the classic worries with AI alignment come at the point when AI gets significantly more capable than humans because now they still can't really cause much damage. They can't really keep secrets from us. They can't really do long-term plans. So I think to a large extent, it's still kind of TBD, how hard it is going to be exactly. And then there are some worrying signs like the models do reward hack a lot. Mita has some really good research on this where they showed, especially it's interesting, the harder the problem, the more they reward hack.

00:50:45 So if you set them a near impossible coding challenge, they'll often be like, "Yeah, I solved it, but actually they totally didn't, or they just found a way to fool the test and ... "

Jon Krohn: 00:50:56 Yeah, they'll hack the database that has the questions to change the question or something like that because it's easier.

Benjamin T.: 00:51:03 And I mean, that can easily get worse as they get smarter because they become better and better at spotting hacks.

Jon Krohn: 00:51:11 Yeah. So that's worrying, but that's power seeking AI systems is kind of the most pressing problem that you list. And you already touched on a few minutes ago in this episode, the second most pressing problem according to 80,000 hours, which is extreme power concentration.

Benjamin T.: 00:51:25 Yeah. So as AI alignment has become less neglected, we've kind of added a list of these other AI risks that I'd say are still way more neglected. So we talked about AI alignment has maybe a thousand people working on it full-time, something like that. But I'd say the concentration of power risks, it's maybe still only like 20 people depending on exactly how you'd count it, but in terms of researchers specializing in it And the idea here is there's a bunch of ways that AI could make it possible to



concentrate power much more than has been possible in history. And there's actually a bunch of different routes by which this could happen. One is like the thing we haven't talked about yet is like if you do actually get recursive self-improvement, so you get AI improving AI, then the rate of AI capabilities progress can actually accelerate.

- 00:52:19 I mean, it's already insanely fast. It's like unbelievably fast, but it could actually accelerate further and you could, people have tried to model out what would happen and they think you could get something like three to 10 years of AI progress in under a year. And if you think like five years ago, basically the models couldn't speak and then they kind of like learned to speak, but they were terrible at coding and maths and now they're like solving 80 year old Erdos problems So they're kind of like getting to human researcher level in pure maths. That was the last five years of progress and then you can try and think what would five more years of progress look like. And it wouldn't just be the models getting even better at maths. It would probably be them filling in all of these agentic bottlenecks that they have now, like continual learning and having good memory and being able to do long horizon stuff.
- 00:53:07 And then you could suddenly get that all click into place, like working on more messy tasks that they're still bad at because the sample efficiency is still pretty bad.
- 00:53:17 So you might go to a point where in 2028 you suddenly get all those last remaining things done in one year. Anyway, so if that happens in exponential growth, the gap between say like the US and China or OpenAI and the open source models stays the same because they're both growing exponentially. So there's like a constant six month gap or 12 month gap. But if you start accelerating to super exponential growth and the size of the gap actually increases and you could have a situation where a



single company suddenly has a workforce of like 200 million digital workers, which would be like a whole nation's worth of labor force just controlled by one company, which, I mean, I guess maybe did kind of happen in history with the East India company, but this would be an even more extreme version of that.

00:54:12 And that's like a pretty risky situation because that company would just have so much power. If they have a lot of political influence as well, then you could kind of end up, or maybe they just team up with the government and you get a kind of Musk Trump partnership or whatever coalition works there and it becomes very hard to ever dislodge them or just like America dominating the rest of the world, which also kind of seems like the default path we're on here. Yeah. I mean, there's a bunch of other ways. AI also makes it possible to do universal surveillance that just like now there's like not enough, the government doesn't have enough workers to monitor what everyone is doing, just like there's too much data. But if you just assign Claude to track each person and flag anything that they're doing that's suspicious, it's like suddenly possible.

Jon Krohn: 00:55:08 Totally. Something I think about is having hosted this show for like almost six years now, I'm like, there's probably things that I said in the past. I don't know, there might be one dumb thing I said and I was like, "Well, but it doesn't matter because up until recently, it doesn't matter because nobody would ever listen to all the episodes or read all the transcripts." But now it's just like, okay, you can just ask an agent, did John ever say anything politically incorrect on the show?

Benjamin T.: 00:55:37 We did do this as prep for the book launch campaign, just like read all my social media feeds and check what was like, what could you quote to make look worse?



- Jon Krohn: 00:55:50 Tell our listeners all the things that you've now taken down. Oh, wow. Yeah. Anyway, so we're kind of off ... Yeah, this extreme power concentration thing is ... Yeah, I can see why it's such a pressing issue, especially given the small number of people working on it. It isn't something that I have come across anyone working on, whereas you definitely come across people working on AI alignment. Number three on your list of most pressing problems at 80,000hours.org is engineered pandemics, which is definitely ... I mean, this is kind of like with number one and number two, those are kind of scenarios where the AI systems are taking off on their own in a way. But with engineered pandemics, that's something that is really not hard. There really aren't technological hurdles even today to say people leveraging open source LLMs to be able to engineer biotechnology risks or other kinds of weapons.
- Benjamin T.: 00:56:54 Yeah. I mean, I'd say the capabilities aren't quite there today. We don't yet know how to make a supervirus that's way worse than any naturally occurring viruses, but it does seem like we could get there pretty soon, even without AI. And then if AI also accelerates technological research, the rate of research progress, it could come a lot sooner than it thinks than we think. Well, also just quickly with concentration of power, the risk there isn't that AI does it. It's just that humans use AI to lock in their power. It's just
- Jon Krohn: 00:57:26 The first one where it'se. Right, right, right, right, right, right. So there's still a human, there's like a dictatorial aspect to the extreme power concentration as well, but the number one, the power seeking AI system, that's recursive self-improvement kind of scenario running off on its own.
- Benjamin T.: 00:57:42 Well, you still have these risks even without recursive self-improvement. It's just more once you get to very powerful AI that could escape control. Yeah, I normally



call it loss of control rather than power seeking AI, but that's like the power seeking aspect is one way it could be especially bad.

- Jon Krohn: 00:58:00 Wow, really reassuring. Then your fourth bucket is actually like instead of being one specific thing, it's just emerging priorities, which has a whole bunch of things that could be concerning to us. One of them you already talked about in this episode, gradual disempowerment, but there's other things in here, the moral status of digital minds, S risks, I don't even know what that means. What's an S risk?
- Benjamin T.: 00:58:25 So an S risk is a risk of something that's even worse than extinction. So it's like, could you have a future scenario where, for instance, by the way, just it's maybe worth saying the idea is to find the most neglected but important problems, which always means you're trying to go one step beyond what's already common wisdom. So you basically always sound a bit insane, but that's how you know that you found something that is genuinely high leverage.
- Jon Krohn: 00:59:00 Yeah. I can see where this is going from the few word summary of what esquis is. And yeah, it is actually, you know what? I hadn't thought of this. I hadn't thought of a fate worse than extermination, but yeah, you found one. Tell us all about it.
- Benjamin T.: 00:59:15 Well, there could be many. Partly, I think this is an area where we just more research is needed, but on scenario that people have worried about is like, could you have an AI system like blackmailing another AI system by essentially threatening to cause huge amounts of suffering, which I mean, most plausibly this would be inner simulation. And so yeah, people who just don't care about digital sentience would maybe just not care about this is getting way too sci-fi for them. But I kind of think if you have AI systems that can do everything that we can



do and they just could be an exact copy, you Then it seems pretty overconfident to be sure that they have no experiences that matter. So I think we should act as if they do. And that's one of the other issues is what do we do about this issue that our AI systems might eventually be able to suffer and have experiences in the same way that we do, which I think is just definitely going to become a huge issue when people have all these we're just talking to AIs all the time and they seem exactly like other people.

01:00:31 I think a lot of people will just assume that of course their needs and desires matter. And then a lot of other people are really confident they don't, and I can't really see how you could be that confident in such a question in philosophy that hasn't been settled for 2000 years.

01:00:52 So yeah, I kind of think we basically need to assume there's a good chance that they do have this capacity and then figure out what that means and what we're going to do about it, which again is like very little work has been put into this. There's a lot of people studying the philosophical problem of the philosophy of mind, but there's not many people who've thought about like, what would this mean for our legal system? How would that ever be incorporated in a way that isn't just insane because yeah, if you just give AI's rights, then that's also really bad because they will just rapidly come to dominate the population because they can copy themselves. So just a kind of like simple, let's just give them like similar rights thing doesn't work either.

Jon Krohn: 01:01:37 Yes. Interesting things to dig into there. I'm going to, in the interest of time, because we only have like 10 minutes left to record in, we could talk about this all day, but we- We didn't get

Benjamin T.: 01:01:47 Into any practical advice.



- Jon Krohn: 01:01:48 Well, yeah. So I'm going to get back to that. I just want to quickly, on the most pressing problems list, something that is interesting is that I think a lot of people today, if you ask them what are the biggest things we need to be concerned about in the future, I think climate change would be one of the top things that people are concerned about. And it's interesting that in this list based on the analysis and the prioritization that you guys have done, it's actually ninth on the list. And so that's interesting. And I guess it kind of maybe aligns with my perspective. You can let me know what you think of this or whether you're aligned with it or not, but it seems to me like us accelerating AI capabilities, maybe building some gas fired plants to be able to have some AI data centers allows us to more quickly be able to say how to keep the plasma contained in a nuclear fusion reaction.
- 01:02:43 And so therefore we're kind of accelerating towards abundant energy and then just being able to afford cheaply to just pump carbon back into the ground or something like that. Is that kind of what the thinking is here?
- Benjamin T.: 01:02:54 That's definitely part of it. Partly there's just the urgency question where like the people running the AI companies, they literally think they might be able to start recursive self-improvement within a couple of years from now. And so we could actually get to AGI in like 2030 or something, so four years from now, whereas most of the damages from climate change come in like 50 to 100 years from now. So there's just a kind of like triaging thing where it's like, you need to address the most urgent thing first because there's like many problems and we can't solve all of them immediately. So you have to just like pick the most urgent one. But then yeah, then there's a second piece, which is what you're saying, which is that if we survive these things like concentration of power and loss of control of AI and so on, then we'll have a much bigger economy and we'll have much faster rate of scientific



progress and that will make it a lot cheaper to tackle climate change.

01:03:53 We should see clean energy become dramatically cheaper compared to today and things like carbon extraction technologies.

01:04:02 Actually CO2 emissions have already peaked in rich countries and so just like projecting this trend forward, the current trend but even faster, that leads to decarbonization. Yeah. The third point, which is probably the one that is like most controversial is like it's very difficult for climate change to end civilization. It's the kind of thing that could kill like a million people a year. So is like a really big problem and could destroy a lot of habitats and a lot of things like that, which are also hard to quantify, but to actually, it can't actually end like kill most people or end the world. That's pretty much what all climate scientists say. And we're trying to focus on the very biggest problems, which means focusing on the things that could actually permanently set us back or permanently end civilization.

Jon Krohn: 01:04:55 Yeah. So to bring the episode back a little bit to the president, it's interesting to me that you've just published this book, 80,000 Hours: How to Have a Fulfilling Career That Does Good while at the same time a big part of your mind and what you spend time on is on looking at these runaway AI scenarios and catastrophic conditions and just you're considering a lot of these exponential scenarios where careers are vastly different. And so it's kind of, do you feel like they almost seem contradictory, right? That you're like, "Okay, let's have a book on your 80,000 hour career, but actually none of us really knows what careers are going to be like in a few years." That's interesting, right?

Benjamin T.: 01:05:42 But the core message of the book is your career is like your biggest opportunity to make a difference, tackle



some of the world's biggest problems. And in a sense, I now just see that case as even more important. The next five or 10 years, we might be dealing with these truly historical levels of change that like as big as or even bigger than the Industrial Revolution, you might be able to actually do something about those and there's like real jobs where you can help tackle these issues and that in a sense just makes your career like even like how you decide to spend this next five or 10 years is just even more important basically.

- Jon Krohn: 01:06:27 Okay. Yeah. All right. That makes a lot of sense. Your career might not last another 80,000 hours or if you're getting started, it might not be, but the coming hours are very important. We've got hours to fix this people can read your book, obviously that's one way. Yeah.
- Benjamin T.: 01:06:45 And the kind of core heuristics that are in there, those still, I've designed them so that they still all apply in the next five or 10 years and yeah, you might not be able to make as long-term plans as the past and that's discussed, but that doesn't mean the alternative is just like giving up on thinking about your career. It just means you need to think even more about flexibility and updating and optimizing over a shorter time horizon.
- Jon Krohn: 01:07:09 Yeah. So maybe can you give us maybe your top heuristics that you think are the most useful for people to be selecting the right career for them?
- Benjamin T.: 01:07:17 I mean, we've kind of already covered the biggest ones. I think from the perspective of having an impact on the world, I think the biggest question is which problems to focus on in the first place, which is just the thing we've been talking about. And so really thinking hard, if I look back 10 years from now with all this crazy stuff that's happening from AI, how will I wish I'd spent that time and could I actually help with any of these issues and what might I do about them and which ones are most



important to focus on? And just actually taking some time to think about that really big picture question, I think is often the thing that will yield the most. And then it's like, okay, well, what could I contribute to these problems? And that gets into the question of which skill is going to be most valuable in the next five or 10 years, which we touched on at an abstract level, the figuring out what the key bottlenecks are at each time and the things that are complimentary to AI, but AI can't do yet because they're too messy or they're too long horizon or there's not enough data or they involve physical presence, these types of things.

01:08:21 And I think for a data scientist in particular or someone doing AI engineering or just someone with a data science skillset, this probably won't be news to anyone, but essentially the routine analysis things are becoming more automated and it's more about maybe combining that skillset with more of a builder skillset or a researcher skillset or a manager skillset and just being someone who just builds stuff and gets stuff done or figures out important problems, but you're using your knowledge of data science as one of the strings in your bow to help you do that, but you'll need other strings and it's less about being super narrow technical specialists and more about using AI to do those aspects, but then you're becoming ... I mean, essentially everyone's becoming a bit more like a manager where it's more about you're writing specs, you're figuring out what's worth working on in the first place, you're coordinating with other team members and so exactly you said you're recently interviewed someone who was all about how to teach soft skills to data scientists.

01:09:31 And that's why I was thinking that seems like a promising thing because it's like that combination ... I mean, this has actually been a thing for the last 20 years in the economy is the jobs that have grown the most are ones where they involve both quantitative and social skills and then social



skills growing second and then pure quantitative STEM jobs have actually grown less than average. And so yeah, being able to basically take the technical skills and then translate them into coordinating with people and solving real problems. But anyway, there's a whole chapter in the book about which skills will be automated in the future and which ones will be most valuable. So that's like the second piece. And then the third piece is then just the more specific career planning advice, which is like ABZ plans. And if you have three options, how do you choose between them and how there's a really common bias is narrow framing, so just not considering enough options or getting either too swayed by your gut and some cool fancy thing that is really attractive to your intuition or just ignoring your gut also being a mistake and how to strike that balance and all the classic job hunting advice still applies.

- Jon Krohn: 01:10:45 Yeah. I mean, you wrote the whole book on it, so it's not surprising that you were able to do that quickly, but I basically, I've kind of been skimming the research that we did before the episode as you've been talking and you've covered most of the big points that we could have focused on the whole episode instead of
- Benjamin T.: 01:11:02 Superintendents. I'd say there's a lot more in the book.
- Jon Krohn: 01:11:05 Oh, for sure, for sure. But in terms of high level categories of things to discuss. And so I think we've gotten to a ... What I'd like to ask you is the kind of final technical questions. We have wrap up questions, which you may remember from last time. So I'll ask you for a book recommendation and how people should follow you, but as a final question before those final specific question for you in the end of your book, you have final thought on your deathbed and in that you contrast two life stories. One is focused on personal fulfillment and another on societal contribution. So maybe you can fill us in a bit more on this final thought on your deathbed section, but



the question is, how do you think individuals, how do you think listeners can practically balance the desire for personal satisfaction with the moral imperative or maybe the drive that some people have to contribute to solving global problems?

- Benjamin T.: 01:12:00 I mean, one thing is there's less, I think there's a bit less trade off than people think because for most people a sense of meaning is really important and that's partly what the deathbed is getting to. It's like you could imagine, well, you was really successful in a career and you did loads of cool travel around the world and you did a job you're passionate about and it was really interesting, but you could kind of imagine then looking back, especially if we're about to go through this transformative AI period and being like, "Well, what was it all for? " And people do wonder that they had a successful career and they wonder that at the end.
- 01:12:43 And then instead, I think it asks you to consider instead you tried really hard to improve our resilience to pandemics and you worked your ass off at that and you helped make the world safer from an engineered pandemic. And then it was like, it'd be really weird to look back and then be like, "Oh, what was the point of all that? " So yeah, I mean, this is also just shown in the motivation literature that a sense that you're actually contributing to something that helps others is one of the biggest drivers of job satisfaction and life satisfaction. Yeah, that's not to pretend there's no trade-offs and ultimately I think we all have to make a very individual decision about how much we're going to focus on personal fulfillment and helping the world in general and are helping our immediate circle around us. But I think there is a significant ... I mean, on the other hand, if you just do something that burns you out and you just believe in intellectually but don't actually enjoy, then that's also going to be terrible for your impact.



- 01:13:49 So they're mutually supporting.
- Jon Krohn: 01:13:51 All right. Well, so my key takeaway from this is things are moving super fast today compared to where we were five years ago relative to where we are today, obviously vastly different world that we live in compared to when you were last on the podcast. And yeah, five years from now it's going to be even more vastly different. Sometimes we get surprised sometimes things end up being an S curve as opposed to an exponential and maybe there is some roadblock where like, okay, the easy stuff of scaling model weights and compute time and amount of data, but it seems to me like I think we're going to keep accelerating. And then we get to know the next
- Benjamin T.: 01:14:37 Three years or
- Jon Krohn: 01:14:38 Two or
- Benjamin T.: 01:14:38 Three years. Yeah.
- Jon Krohn: 01:14:39 Even just scaling compute time, there's so much more. Even if you just fixed model weight size and data sets today and you're just like, let's have accuracy improved by computing for longer and coming up with engineering tricks to use less compute over longer periods of time or for more cycles of reflection things continue to accelerate dramatically. And there's just so many people now tackling AI capability problems, orders of magnitude more than they probably are concerned about AI safety and those people doing AI capabilities work, even if you don't today have a fully autonomous AI system, you're vastly more capable as an AI researcher at the frontier by having these tools at your fingertips. And so it's just so easy to bring together so many different ideas from disparate fields or go deeper within your own field, have text or code, PowerPoint presentations created automatically to make it easier to disseminate your ideas or make your ideas real, test out, experiment with



different ideas like ... Yeah, it's really hard to imagine that on the backdrop of all that, we're not going to be even more vastly ahead when I invite you on for episode 1500.

- Benjamin T.: 01:16:11 But you're pointing to a really important thing there, which is not only are we facing these huge problems in the next three to 10 years, but also like your leverage is going up still. Even if we do eventually all get automated before we get to that point, we actually have a lot more power to get things done than in the past and that just makes the stakes even higher. How are you going to use that power the next couple of years?
- Jon Krohn: 01:16:36 A hundred percent. That's exactly right. And that is the kind of, I guess even if we have a lot of question marks over a decade from now, five years from now, maybe even in terms of what the world is like and where there are opportunities in a career, I am very confident that in the coming months, in at least the coming months, you have ... Yeah, exactly what you're saying, way more power than ever before. That's kind of interesting. It's like I talk about my, in terms of my technical skills, my moat has disappeared, the 20 years spent becoming an expert at stochastic gradient descent and writing Python code. It's not really that important. It still can be useful to understand. It can be useful for reviewing work and going back and forth with a conversational agent on what you can be doing. It's still useful to know all those things, but the barrier to entry on those technical skills is so much lower than it ever has been, but that doesn't mean that somebody with a data science, AI engineering, software developer skillset exactly they're saying is less valuable.
- 01:17:48 It means you're more for at least months to come because you can do so much more, you can now focus on understanding user requirements or thinking about startup ideas and just making it happen as opposed to being like, "Okay, I have this great idea. Now I need to hire this software development team." You can just do it.



- Benjamin T.: 01:18:10 A team of like five people can probably get done what would've taken like 20 or 30 people in the past.
- Jon Krohn: 01:18:18 Even a year ago, even so I was with a longtime colleague of mine, brilliant guy, Ed Donner, he's created tons of agentic AI courses and him and I were talking about how a year ago a project that we might have scoped as costing \$150,000 and take three to four months for a team of four software developers like front end engineer, back end engineer, a project manager and maybe someone doing QA, a QA engineer, or maybe you have a second backend engineer or something, but kind of like a team of four people, \$150,000 cost and a year later it's basically free. It's a \$20 a month, maybe a \$200 a month Claude subscription to get it done in hours, which is crazy.
- Benjamin T.: 01:19:11 And I mean, this is the other reason to focus on your impact and meaning that I was forgetting earlier, which is you can also have a lot more impact than people thought and even compared to 10 years ago. And if you learn something that is a lot easier than you thought, then it makes sense to focus on it more than you would've otherwise. And if we make it through this time, we're going to have all of history to chill in the most AGI super abundant world, but we only have this five year period to help make sure that transition goes well before we end, I did want to just briefly talk about some of the super practical advice where I mean the AI engineering skillset is maybe like one of the most in demand skill sets within AI alignment because again, probably when we last spoke five or 10 years ago, that was a more like pre-paradigmatic phase where people just didn't really ... It was more about like having big picture research ideas about how to even tackle this, but now it's really moved into a much more empirical phase where there's just a lot of very concrete projects and it's much more bottlenecked just by having exactly AI engineers able to set up all these control monitoring experiments, like do red teaming, do interpretability studies.



01:20:33 And there's a lot of programs that are designed to transition people into AI alignment. So the MATS scholarship is the biggest one and then there's also arena and that's exactly like if you're already quantitative and you want to get into AI alignment within three months, that's like exactly what it's designed to help you do. And then I think also the concentration of power stuff, a lot of that is also, there's like technical elements of that, like exactly what instructions are given to the AI, whose commands does it follow and how is access controlled to make it hard for one person to completely tell the AI exactly what to do and things like that. And then finally, well, there's the data science, applied data science skillset, which governments and policy teams need. So how do we figure out the labor market impacts of AI? That's the thing that governments are hiring data science people to do and tracking the trends.

01:21:31 And then finally this builder skillset you're talking about, that's super useful as well. So again, with pandemics, one big thing we need is just really good disease surveillance where we're sequencing all the wastewater in airports and things like that and looking for things that are growing exponentially and that could give us way more early warning of a new pandemic. And that's exactly like people who are good at just basically just building stuff quickly, data analytics, that's exactly the type of thing for them. And yeah, if you want more customized advice, then we have the one-on-one advice on the 80,000 hours website, you can speak to someone and they could recommend, given your specific skillset, hopefully introduce you to people in these different areas depending on which one you're interested in.

Jon Krohn: 01:22:20 Incredible. Yeah. To summarize back to you, one of the things that you just said is five years ago when you were on the show, as you said, if an AI engineer would have ... When we said AI engineer as a job, it probably would have meant somebody kind of at the frontier, somebody at a



frontier lab or doing research trying to push people inventing large language models or transformer architectures and going to conferences. But today AI engineer means this, as you said, highly empirical job where you're probably, it's rare that you would even be fine-tuning LLMs, you're mostly taking off the shelf capabilities and combining them together to be making real world impact. And yeah, you can make a huge amount of impact from that seat, from increasing profitability all the way through.

Benjamin T.: 01:23:13 Hopefully not that one. Yeah.

Jon Krohn: 01:23:16 All the way through to saving the world, you choose. I mean, I'm sure there's even, you've talked a lot in the past about things like for some people, the way that they can make a biggest impact in their career is say if they have some prodigious ability to trade stocks at high frequency or whatever, and it could be the best way that they can make an impact is actually making lots of money as a trader at a hedge fund or doing investment banking or whatever and then taking what they need to survive, but then donating the rest to AI alignment research or to malaria nets in Sub-Saharan Africa or one of the other highly good value per dollar ways of saving lives or improving quality of life that there's lots of posts about that on the 80,000 hours website in terms of places you can be effectively donating the capital that you have.

01:24:24 And so in that respect, there's kind of a scenario where you can imagine you're like, "Don't do the profit one," but there's maybe to some extent, it's getting the right balance of, okay, how can I use my vastly useful AI engineering skillset to generate enough for myself that I can support me and my family and our needs, but then using also some of my time or as much of my time as I can to be making a big positive impact.



- Benjamin T.: 01:24:49 Totally. That's like a lower bound on your impact is you just do what suits you and hopefully isn't something with a negative impact, but it's just something relatively high earning and then donating some of the money. And yeah, I mean, as like a very lower bound, it's still possible to save someone's life for like three, \$4,000 from malaria. And so as a data scientist at a big tech company, you could actually be saving tens of people's lives maybe every year and still living on like much more than you would've lived on in the nonprofit sector or as a teacher or a nurse or another helping career. And then I would actually argue you could have even more impact by donating it to like pandemic prevention or AI alignment, though obviously it's harder to quantify the impact. But yeah, a lot of these organizations just if they had more money ... I mean, we mentioned Meter earlier.
- 01:25:56 If Mita had more money, they could hire more engineers, they could have more compute and they could have better measurements of whether about to get recursive self-improvement or not, which is like the most crucial strategic thing to understand in the world.
- 01:26:10 And it's kind of amazing that anyone can actually help with that by using this thing called money to translate our labor into other people's labor. But yeah, I think by finding a job that suits you directly in the area, you could probably have even bigger impact still, but obviously that involves making a career transition, which by the way, on our Substack, the 80,000 hour Substack, we have a guide to like, if you want to switch into working on AI risk in three months, what should you do? And it has like six numbered steps where it's like, these are the things to work through and all the courses and fellowships that can also help you transition we have listed on that article and also on our job board so that yeah, the job board not only has like open jobs, but also it has this other tab where you can find like fellowships and training courses and funding that also is around to help people transition.



- 01:27:08 It's called How to Transition into AI Risk in Three Months.
- Jon Krohn: 01:27:11 Ah, how to transition. Okay. All right. I'll have that in the show notes as well. Fantastic. So yes, I've got it. Yeah, it's on the Substack exactly as opposed to the main site, which you said. Well then, what an episode we've gone well over on the time that we had set aside for this, but really interesting conversation. So I'm glad that we did that. Other than your own book and all the resources that you've already talked about in this episode, which we will for sure have in the show notes, including yes, a link to 80,000 hours, how to have a fulfilling career that does good, your new international bestselling book. Beyond that, I always ask for a book recommendation other than your own. Do you have anything that comes to mind for you?
- Benjamin T.: 01:27:56 I mean, I maybe said this last time, but *The Precipice* by Tobiaud, an introduction to accentual risk and why it matters, but if you want something more fun than maybe just ... I mean, heavy going, but the Greg Egan books, which are like a sci-fi novel and they're the only ones that really wrestle with how weird digital ... If you could actually have digital people, how weird that would be. Most sci-fi is basically just the world today, but we can have spaceships, but what we're heading to is a lot weirder than pretty much anything that you see envisioned. And this is pretty much the only author I know who really tries to wrestles with that.
- Jon Krohn: 01:28:42 Cool. Yeah. I found Greg Egan, so I'll have a link to his works in the show notes as well. Ben, other than the resources that you've already mentioned, obviously your book, obviously the 80,000 hours website, the 80,000 hours Substack, your own personal Substack. Other than those, are there any places that people should be following you maybe on social media?



- Benjamin T.: 01:29:02 Yeah. And my personal Substack is more focused on what's happening with AI and what you can do about it. And then I'm most active on X, so just Ben_J_Todd. Nice.
- Jon Krohn: 01:29:15 All right. What does the J stand for?
- Benjamin T.: 01:29:17 John, actually.
- Jon Krohn: 01:29:19 There we go. I didn't know that. All right. Well, Benjamin, John, Todd, it was great to have you on the podcast and yeah, if we're all here in five years, maybe you can come back on the show.
- Benjamin T.: 01:29:32 Cool. Thanks so much. Really enjoyed it.
- Jon Krohn: 01:29:35 Phenomenal episode today. Wow. In it, Ben Todd detailed how passion is a poor starting point for career choice and how building rare and valuable skills is what actually generates lasting fulfillment, how there's no permanent solid ground in an AI career, only the moving bottleneck of whatever humans can still do that AI can't. So the winning move is to keep shifting your focus onto that fast moving frontier. He talked about how once you have a human level digital worker, you almost immediately have a superhuman one because of recursive self-improvement. How AI alignment now has perhaps a thousand people working on it, but the risk of extreme power concentration has maybe only 20, making it perhaps the most neglected problem in the world relative to its scale. And he talked about how your career has become a bigger lever, not a smaller one since a team of five can now do what once took 20 or 30 people, which raises rather than lowers the stakes of how you spend the next few years.
- 01:30:36 As always, you can get all the show notes, including the transcript for this episode, the video recording, any materials mentioned on the show, the URLs for Ben's, social media profiles, as well as my own at



superdatascience.com/1007. Thanks to everyone on the SuperDataScience podcast team, our podcast manager, Sonja Brajovic, media editor, Mario Pombo, our partnerships team Natalie Ziajski, our researcher, Serg Masís and our founder Kirill Eremenko. Thanks to all of them for making another awesome episode possible for us today. And for enabling that super team to create this free podcast for you, we are deeply grateful to our sponsors. You can support this show by checking out our sponsor's links, which are in the show notes. And if you ever are interested in sponsoring an episode yourself, you can get the details on how at johnkrohn.com/podcast. Otherwise, please help us out by sharing this episode with someone who would love to have it shared with them.

01:31:28 Review the episode on your favorite podcasting app or on YouTube that really does make a difference for us, particularly Apple Podcast reviews. Subscribe if you're not already a subscriber, but most importantly, we hope you'll just keep on tuning in. I'm so grateful to have you listening and I hope I can continue to make episodes you love for years and years to come. Till next time, keep on rocking it out there and I'm looking forward to enjoying another round of the SuperDataScience Podcast with you very soon.