# EPISODE 74

[INTRODUCTION]

**[0:00:10.4] SC:** Hello and welcome to another episode of TWiML Talk, the podcast where I interview interesting people doing interesting things in machine learning and artificial intelligence. I'm your host, Sam Charrington.

A few quick announcements before we dive into today's show. In a few short weeks, we'll be holding our final TWiML Online Meet Up of the year. Yes, on Wednesday, December 13th, bring your thoughts on the top machine learning and AI stories of 2017 for our discussion segment and for our main presentation, former TWiML Talk guest, Bruno Goncalves, will be discussing the paper; Understanding Deep Learning Requires Rethinking Generalization by Chiyan Zhung from MIT and Google Brain and others. You can find more details and register at twimlai.com/meetup.

If you received my newsletter, you already know this, but TWiML is growing and we're looking for an energetic and passionate community manager to help expand our programs. This position can be remote, but if you happen to be in St. Louis, all the better. If you're interested, please reach out to me for additional details.

I should mention that if you don't already get my newsletter, you are really missing out and should visit twimlai.com/newsletter to sign up.

Okay. A bit about the show you're about to hear. The show is part of a series that I am really excited about impart because I've been working to bring it to you for quite a while now. The focus of the series is a sampling of the really interesting work being done over at OpenAI, the independent AI research lab founded by Elon Musk, Reid Hoffman and others.

In this episode I'm joined by Greg Brockman, OpenAI cofounder and CTO. Greg and I touched on a bunch of interesting topics in this show. We start with the founding and goals of OpenAI before diving into a deep discussion on artificial general intelligence. What it means to achieve it

and how we're going to do it safely and without bias. We also touch on how to massively scale neural networks and their training and the evolution of computational fameworks for AI.

This conversation is not only informative and nerd-alert worthy, but we cover some very important topics, so please take it all in, enjoy and send along your feedback.

A quick note before we jump in; support for this OpenAI series is brought to you by our friends at NVIDIA, a company which is also a supporter of OpenAI itself. If you're listening to this podcast, you already know about NVIDIA and all the great things they're doing to support advancements in AI research and practice. But what you may not know is that the company has a significant presence at the NIPS Conference, which is going on this week in Long Beach, California, including four accepted papers.

To learn more about the NVIDIA presence at NIPS, head on over to twimlai.com/nvidia and be sure to visit them at the conference. If you'll be at NIPS, send me a shout out as well because I'll be there too.

Now, on to the show.

[INTERVIEW]

**[0:03:23.5] SC:** All right everyone. I am on the line with Greg Brockman. Greg is a cofounder and CTO of OpenAI. Greg, welcome to this week in machine learning in AI.

**[0:03:33.6] GB:** Thank you for having me.

**[0:03:34.8] SC:** Awesome. Hey, why don't we get started as the tradition here on the podcast with you telling us a little bit about your background and how you got involved in AI?

**[0:03:44.6] GB:** Sure thing. I got into programming I guess relatively late, so after high school I took a year off and went abroad. I was working on a chemistry textbook and I send it off to one of my friends who had done something similar in math and he wrote back saying, "There's only one problem with this, which is that you don't have a Ph.D. so no one is going to publish it. You

can either self-publish or you can make a website trying to promote things that way." I was like, "Well, I guess I'll figure out how to make a website." And so I went online and taught myself how to code and build a little sample table sorting widget and I thought that was cool and I built something bigger and bigger and never really looked back at the chemistry.

One of the early things that really captivated me was the idea of being able to write code that could understand things that I cannot. The way that you write code that you build systems, you think hard about a problem and you understand it. You write it down in this very obscure way that we call a program and suddenly anyone can get the benefit of what you just did and if there is a way to amplify that and to have a program that could do things that I didn't have to even understand myself, then suddenly the set of problems you could solve is so much broader.

**[0:04:52.2] SC:** How do you get from building a website to that?

**[0:04:55.3] GB:** Yeah. I read Turing's 1950 paper, Computer Machinery and Intelligence, and when you read it, he was talking about the Turing test and he had this picture of "by the year 2000, you'll be able to build this child machine that will learn just like a human child and will get full intelligence". This was in 2009 that I was reading this paper and it's like, "Where is that machine? Why hasn't anyone built it?"

**[0:05:22.0] SC:** Nice.

**[0:05:23.0] GB:** So throughout college — So I ended up doing a bunch of different startups and ended up transferring. I started out at Harvard. I was there for a year and a half before going to MIT. I was there for a semester and a half before leaving to go work on Stripe where I was the CTO for five years and build that from four people to 250 employees. It's now over a thousand or so.

Kind of for me, the goal has always been to work on AI. It was just a question of when and the right way of doing it. While I was working in the startup world, if you read Hacker News and if you look at what's happening, you see all these articles in deep learning does X, deep learning from Y, and the thing that was very unclear to me from the outside was substance are hype. I'd actually done similar investigation on Bitcoin. From the outside, Bitcoin similarly is something

where there're a lot of people talking about it, is it substance or hype and had really done a deep dive and kind of concluded that this was 2014 that it's kind of people weren't really focused on the right things that if people were focused on just kind of speculation and not about building products, delivering value, still might be the case that the Bitcoin will succeed.

I have a very different observation when it came to deep learning and AI and what was happening and I realized that a lot of smartest friends from college were now in the field and that things were starting to work in a very real way and solved tasks that you just couldn't have solved in other way.

The fact that things are actually happening that you can actually build systems that can have real-world application and that it's also still very much at the very beginning of the S-curve, for me it was very clear that the moment is now. I was talking to a bunch of people in the field. I was talking to Sam Ottman, if you put together this dinner with Elon Musk and Ilya Sutskever and some others and the focus of this dinner was clear things are happening. Things are moving very quickly. How can we best have a positive impact? How can we help ensure that this plays out in the best possible way? Because AI is just going to be the most transformative technology that humans ever create. Just having any kind of contribution to making that play out better is the most worthwhile thing that I could imagine.

The conclusion was that it seems like it's not too late. It's not impossible to build a web with a lot of the strongest researchers in the field, especially if you focus it on this goal of this technology. It's not enough just to build it. You also need to think about how do you make sure it actually benefits everyone. We had that as our hypothesis, and I at the time had left Stripe a couple of months earlier and said, "I'm going to fulltime. I'm trying to make this happen."

We put together a team and in December of 2015 launched at NIPS and asked that we existed, and since then have been working on going from zero to pushing the envelope of what is possible in this field.

**[0:08:17.8] SC:** Yeah. For those who aren't familiar with OpenAI as an organization. What does it look like today in terms of the number of researchers and what's the model? Like you see

folks that are affiliated with OpenAI that are affiliated at other places as well. How have things played out since then?

**[0:08:40.3] GB:** Yeah. We tend to think of ourselves as cherry-picking the best parts of academia and the best parts of industry towards a great focused goal. OpenAI at the end of the day is a series played to build general intelligence and make sure that it plays out well for our society, it plays out in a safe way. To do that, we don't know how to build something like general intelligence today. So you need to do a fundamental research. You need to push the limits of what's possible, but it's also the case that the field has really transitioned from being an individual sport to being a team sport. Really just this year and maybe some in 2016, you're able to start using large clusters of machines much more productively. In 2012, Google had the cat neuron, which is 16,000 course, but we're able to be surpassed by two grad students on a GPU and a GGPUs. Today, that's a very different story where you have people training image at 15 minutes on 1024 GPUs.

If you look at something like our Dota Project, which we'll talk about in a bit, which really required this team of people with engineering background, with research backgrounds all coming together towards a shared goal. So the way that we structure ourselves is that we have few different teams internally. So we have a robotics team, we have a Dota team, we have a few other teams and people — We have the full mix of skillsets that are required to accomplish a goal within those teams.

We also have an infrastructure team that is more of a horizontal team that supports the work of all these different teams and accelerates the work of those teams. One thing that we're increasingly seeing is collaborations amongst the teams which ends up being a really powerful thing where we can take code from Dota and use it for robotics and really accelerate what's possible there.

**[0:10:19.4] SC:** Can you speak a little bit to the open aspect of OpenAI? It's clearly something that was important to you in setting out on this journey, but at the same time there's been some critique of the level of openness at OpenAI and there's still some things that I think you're clearly publishing a lot of research, but folks have asked are you publishing. Could you be doing more in terms of publishing datasets and things like that? How do you think of that?

**[0:10:50.6] GB:** I think it's a great question, and I think that that's one misconception that people have had since the beginning of OpenAI, is I think people have put kind of a narrative that's very different from what we're trying to do and I think it's a really good thing to ask about.

The goal of OpenAI is to ensure that the world post-general intelligence is good for humans. Along the way, it's really important as an organization that we both — Like one thing that we think about a lot is what we can do to both accelerate our organization, but then also things that we can do to help accelerate the field and things that we can do to deliver value to the world generally.

The last one ends up playing out on multiple time scales. So there's a very short time scale work that you can do. For example, today we published — We released a few more algorithms in our baselines projects where we have done high quality implementations of all the standard reinforcement learning algorithms to now rather than people having to — and basically end up with a bad baseline that then they say, "Oh, my new method is better than this one." You can just take the work that we've done to, well, turn all of these baselines and use that. That's a short timeline delivering a value.

The thing that we really are trying to do is a much longer timeline. It's really about build an organization  that can be at the forefront of this research and to actually be able to steer how it plays up in society. To do that, it's not as simple as just taking a work that you do in real-time and toss it up at events or put it up on GitHub. I think you require something much more thoughtful. So I guess a lot of how we think about it is that OpenAI will be a success if you fast-forward 5, 10 years or whenever the preferred checkpoint is and you look back and you say that the amount of value that we delivered over the long time frame was the max that we could have, and some of that I think is the short run. It means that you don't necessarily publish or release code to the maximum extent that we possibly could. That is always made with the choice of because we think that we're going to be able to better delivery value over the long run.

**[0:13:01.3] SC:** Specifically, meaning that, clearly, publishing code has a cost to it, have resources that are associated with it. But even more so, once you publish it, there are some rightly or wrongly expectation of maintaining that overtime that also has a cost, and all of that,

those accumulated costs potentially slow the organization down or kind of require ongoing resource. Is that specifically the thinking or is that just part of it?

**[0:13:30.7] GB:** That is a real one for sure. Again, one weird thing about this field is that the technologies that are being developer are ones that are very desired by big companies, right?

**[0:13:41.2] SC:** Yeah.

**[0:13:41.2] GB:** The Googles and the Facebooks, and therefore a massive amounts of resources in this. To actually have an impact as an nonprofit which we're resources but it's very different from the level of resourcing that you would see at one of those companies that you really have to answer the question, "What is it that I'm going to do that is the differential impact? How can I make the most — Get the most bang for my buck, makes sure that sort of the differential impact of this organization existing is as large as possible.

That means that having a large open-source project that lots of people are using is very, very valuable, but there are lots of other people who would do the same thing and where it's sort of is much more — You look at TensorFlow, that that's something it's great for Google, because lots of people are using their tools. It means that when they hire people that they're using the same platform. And so the incentive exists for the big companies to do that. I think the thing that we view as unique to us is really thinking about this AGI problem and thinking about how do you make sure that when it comes to — There's kind of two problem that really core to AGI. The first is; well, you're going to be building this really powerful system, you should imagine you basically are going to train like the goal, general intelligence. Like how should you even think about that? What even is it? It's going to be a system — And I think for the purpose of this conversation we should define it as a system which can perform any economically valuable task as well as a human.

If you can build that, first of all, to train it is going to require a lot more compute than to actually run it. Let's say that you're going to need a bunch of agents that each one of those agents is going to run much faster than real-time in order to train, and so you should kind of think of whatever system you're going to build while there's going to be this massive datacenter that's just going to be sitting idle while you're running your single AGI. Really, you're going to have the

ability to run lots of them from day one, and so you should kind of thing of the thing you're going to build as this organization of the most competent person you've ever met that are all working together in concert towards a shared goal with no ego. It's going to be a pretty powerful system. Today, we have lots of companies that are organizations with people and are able to accomplish pretty wild things, and I think if you can build the kind of system I just described, that it's really hard to see what the limits of that is going to be.

The first thing you have to ask is, "Is it going to do what we want at all?" This is what is referred to as technical safety and a problem that we work on to one step towards the solving a technical safety side that we've done is this human feedback project in collaboration with Feed Mind. The idea is that the core problem on technical safety is that you need to be able to specify goals to AI. The AI somehow needs to listen to you. It needs to reflect human values. It needs to sort of do what humans want in some pretty deep way. There needs to be humans in the training process somewhere. The human feedback project that we worked on is the first step in its direction where a human labeler is shown two videos of a behavior and they just click on which behavior is more like the ones they want. We're able to show that with 500 bits of feedback that you're able to train an AI to do some backflips. That's step one. This kind of work is something we think is really important, and it's a little bit taboo to talk about the field. The idea that, yeah, you could build these systems, they're going to do anything crazy.

I totally understand why. I think that there're several motivations for that. One is that the field has gone through these hype cycles of booms and bust and winters and that to really think through the, "What if these all succeeds? What if it works?" is something where if you're just going to go through another cycle of that, then  there's really no point.

There's a second thing which is that I think people really reason from what the computers of today can do. You take your Pascal GPU and, sure, you can train a great image classifier, but you're going to be able to train anything better than that, and the answer is, "Well, not really," but there's really two things that are changing that I think people don't see that are going to really accelerate the kinds of models that we can run.

The first one that I alluded to earlier, which is the fact that you can now use datacenter scale in order to get better performance. Again, this is 2012. Basically two GPUs, you get severely

diminishing returns. Beyond that, that state of the art, 2014, you could do eight GPU training and that seemed great. Just this year, Facebook did 256 GPUs, image net one hour and someone also just did 1024 GPUs, image net 15 minutes. You get this massive scale of the kinds of models that we can run, the kinds of systems we can train.

The second one, which is the acceleration of the neural net hardware, and if you look at everything up to 2016, it's actually very smoothly Moore's Law, even though it's not driven by the same factors as Moore's Law, which is pretty remarkable. If you can look, you can actually look and there's this diagram from Eric Kurtsel's book, where the data cuts off at 1998 and he just puts a double exponential to it and predicts, basically, exactly correctly the 2016 Pascal GPU is right on the curve that's predicted from this data just cutting off at 1998. Very smooth Moore's Law all the way through.

This year, something weird happened. This year, we ended up with an order of magnitude increase in the number of flops available for running neural nets. These numbers are all public. So there's the Volta off the of my head. 2016 Pascal GPU is like 20 teraflops, able to. Volta came out this year is more like 90 teraflops, and Google TPU 2.0 is more like 180 teraflops. There's really this explosion this year.

The thing is we really expect this kind of acceleration of the compute available for neural nets in a small compact package to continue to accelerate much faster than Moore's Law and there's a very simple reason. The reason is that you can — Neural networks are very, very parallel, and just like the brain is this big network of just a bunch of tiny little cords, they're all talking to each other. There are some learning rule and there are some propagation rule, and the way that we've always designed hardware is much more for serial execution. No one's really had this incentive to design this massively parallel hardware before. There were a lot of low hanging fruit. You don't need to invent any novel technology. You can just use — You don't have to rely on predecessors getting smaller in order to get these speed-ups. If you combine these two things, datacenter scale with faster neural network accelerators, suddenly computer available for any of our models is going to really skyrocket.

That's kind of the perspective that we have, is that things are going to be different as a result of the hardware coming online. Timelines are always tricky to predict exactly, but I think that we're

going to see even just next year, I think if you just look at what we could do a year ago and what you can look at right now with respect to image generation, with respect to voice generation, I think in 2018, as long as we continue to see the compute coming online and the way that we expect, that we should able to have perfect video generation. We should be able to have basically perfect speech synthesis as well. This kind of acceleration in terms of the capabilities is going to be really tied to — Well, we have all these ideas of these models, but we need to compute in order to run them. As long as we get that compute, that I expect to continue to see the capabilities increase in lockstep. That's thing number one that's really important.

Thing number two that is extremely important is the question of, "Okay. Let's say you build an AGI who does what humans want, that reflects human values. Who're values and who are the people who get to specify what this AI should want?" That's a much harder problem. The first one is technical problem. That sounds like a thing where — As these systems play out, we're very good at solving technical problems. If you can actually build this kind of very powerful system, like there's good reason to believe that if you put on the effort, you should also be able to make it safe and solve that technical side. It's not easy, but you have to really want to — You have to really try to solve that problem, but it seems like a thing that is solvable, a thing that's much harder is this non-technical problem of who owns it and who specifies the goals. That is something that is also very core to OpenAI and how we think about the value that we're delivering and that we really want this technology be something that is not just benefitting one corporation, one person, even one small subset of people. We really want this to be something that is benefitting the world, and we have ideas around the right way for that to play out. I guess when it comes to how do we think about open, that's exactly how we think, is solving those two problems. If we can do that, then that is the most important thing any of us could imagine doing.

**[0:22:18.1] SC:** On that later point in terms of whose values, are those things that you have ideas about but haven't turned into projects yet or are you doing — Are there public projects that you've been working on that speak to that second item?

**[0:22:32.6] GB:** Yeah. It's something we spend a lot of time thinking about. I think that we haven't yet done kind of public speaking about our thoughts there, but a lot of what we've spending time doing has been building relationships with a lot of people on the field, a lot of people on governments, a lot of people just in various positions who I think will end up being

influential with respect with how this technology plays out and just this kind of — I think really trying to lay the groundwork for where we shop things to go.

**[0:22:59.9] SC:** One of the themes that keep coming up as you're speaking was this notion of like a timeline or a timeframe for AGI. Do you have one that you kind of manage to or is there a general agreement within OpenAI and that community as one we think AGI is going to happen or even the timescale, and maybe some context for this. I don't think I've told this story on the podcast before. Maybe I have, but relatively recently I was, with some fellow entrepreneurs, talking about — We're just kind of catching up and someone pushed me on, "Hey! On this AI safety issue," I didn't use those words, "but are you a Mark Zuckerberg or are you an Elon Musk?" I tend to answer that question like, "Well, you know, it's kind of in the middle. I think there's a lot of sensationalism, but he kept pressing, pressing, pressing for me to answer.

One of the things that occurred for me in thinking about this was that, "You know, if I think about who Elon Musk is, his timeframe is probably way longer than mine." The guy is like building rocket ships. He's thinking long term.

I tend to answer that question in terms of I think people really overblow what is likely to happen in 10 years, right? I wonder with that as context, how do you think about the world, you, Greg and OpenAI more generally in terms of the timeframe for worrying about and thinking about these kinds of issues?

**[0:24:32.0] GB:** Yup. Timeline is a really interesting and hard question. It is the hardest question and I think this is true for any technology. If you look at the invention of flight, people, all the experts in the field right up until flight was created were saying flight is for the birds, that Newton had proved that having air flight would never happen. Then you have the Wright Brothers during their flight just a few months later.

If you look at kind of any transformative technology, it is really the case that it's hard to distinguish exactly when it will happen and I think that there's something very inherent to this, because if people knew, "Okay. Here's the timeline to it," then you would just focus more, work harder and accelerate that timeline, and so you would try not to be inaccurate.

I think that the way that we really think about it — I think [inaudible 0:25:15.3] had a good blog post where he talked about something that he did. He was listening to a bunch of AI experts saying AGI is very, very far away and he went up and he asked people, "Okay. Tell me, what is the least impressive accomplishment that you're very confident is not going to happen in the next two years?" People really didn't have a good answer.

How can it be that you both have thought very, very deeply about, "Okay. It's going to take this long. It's going to take exactly this long. Here is when we're going to deliver it," and also don't have, "Okay. Here is something that I'm willing to bet. This is the least impressive thing that just we're not going to do in this timeline."

I think what's really going on, I agree with the conclusion that he has there, is that people don't really have a good concept of it, right? People don't really have — that in general, people end up picking with their gut rather than through having like really rational, "Here is exactly the factors that are going to enable it and here's why we're not going to be able to do it in the near term." Besides the fact that while you look at what my — I look at my dumb AI agent where I'm trying to get this thing to even be able to tell a cat from a dog and can you imagine trying to build something as smart as me. It's just there's this disconnect.

The way that we think about it is that we certainly know that some things are going to be changing and I think that specifically the hardware is going to change in a way that people are not expecting right now and is faster than Moore's Low. It's not something that people praise into their internal sense of what's happening.

There's a question of how hard does that take you? On what timeline does it take you there? A thing that's important to us an organization is that regardless of what the timeline ends up being, that we are able to have the influence that we want that we're able to ensure that this ends up playing out well.

The second part to that is that, well, you can also say, so you don't know when that super transformative stuff is going to happen, but you can say something about what is going to happen in the near term. What is going to happen over the next five years. Again, it's very clear we're going to be able to do synthesis of perfect videos. If you look at what's the 2020

presidential campaign is going to look like when you're able to generate the kinds of videos that we already see that we're very, very close to being able to do.

There are a bunch of technologies. Robotics is a perfect example where to-date there's been results of learning on robots, but none of the roboticists are impressed, because all the task that people can accomplish are worse than what the roboticists can already do. So if we talk to a roboticist, they say, "Okay. Come back — Call me when you could do something I couldn't do in the 70s." That's going to change, and the moment that you change that, the moment that you have your first learning-based result that blows away what was possible without learning, I think there will be a c-change.

We've seen this in the number of other disciplines. We saw it with Vision, pre-2012. You go to the big Vision Conference and there is like one neural net paper if you're lucky. Now, there's like — Basically everything is neural nets and Vision. Like I don't think anyone even remembers that it was different. I think that on robotics, that it's pretty clear that humans aren't getting any smarter. We're not getting any better at thinking through these problems and being able to program all the rules for exactly what a robot should do and how it should react. I would say that fact that you're going to have learning come end-to-end really change what it's capable, what robots are capable of I think is going to be massively impactful.

So that's how we think about it, is that the big goal of AGI is something where you can't know — You certainly can't know that it's close, but I also don't know that you can know that it's super far and that we also know that there's going to be transformative to applications in the near term. For us, the mandate for us, the way that we operate is stay on the cutting edge. Make sure that we're pushing forward and always be asking, "How can we ensure that our integral overtime of value delivery is as large as possible?"

**[0:29:05.7] SC:** That was a non — I'm not giving a timeline answer.

**[0:29:08.9] GB:** Yeah.

**[0:29:10.0] SC:** Very well said. You did give two examples of applications where you think we'll see transformative short term things happening. One is audio and video generation and the

other is robotics. Are there specific examples of kind of leading indicators or examples that are leading indicators that kind of give you the confidence that those two specific things, for example, will change pretty dramatically pretty quickly?

**[0:29:39.2] GB:** Yeah. I guess on the robotics front, we work on robotics and this is our goal, is to be able to change, and it's to really unlock robotics through learning methods. It's actually interesting, because the way that we think about it. The way that we work on robotics is that we are geared towards trying to build general technologies rather than trying to maximize robotic capabilities and that that steers the set of projects that we're going to work on and the ones that we're going to pick.

I think one thing that we're really excited about is if we succeed, we stay focused on AGI, but can also enable robotics to really kickoff.

**[0:30:23.0] SC:** What's an example of those two things in opposition to one another?

**[0:30:26.5] GB:** One perfect example is that I think that there are so many really positive applications that we're going to see in robotics over upcoming years. Like elderly care robots are a perfect example. Where that's something that I think is going to deliver value to a lot of people. It's going to really be transformative to a number of people's lives, but it's also not necessarily something that we're going to work on ourselves. That the kind of thing that we want to do is to build the underlying technology that would allow that application to happen, but stay focused on pushing forward on new applications rather than productizing.

**[0:30:58.9] SC:** Got it. Is that different than where basic research as supposed to applied research, or is there another nuance to that?

**[0:31:07.6] GB:** Yeah. I'd say that we're kind of halfway in between, because when I think basic research — And I guess it might depend per field, but when I hear basic research, I think of the individual sport type research of people kind of off in their own, thinking deep thoughts and coming back when they have something that seems cool.

That for us, that we really try to take results and push them to the limits of scale. With our Dota systems, that that's exactly what we did, where rather than just showed that, "Okay. Here's some system that can kind of work on — In some toy way," actually work in a really hard task, and that I think the thing that distinguishes it from applied research is that we focus — Solving Dota is clearly not going to be something that is going to be transformative through many people's lives. It's transformative to a subset of people, but not in the kind of direct impact that one would have in a more applied setting.

**[0:32:05.4] SC:** Yeah. One of the things that I saw recently that is a bit of an example of what you're suggesting will happen to video is that NVIDIA recently published some work using GANs to generate these synthetic celebrity faces. Did you see that one?

**[0:32:21.1] GB:** Yeah, absolutely.

**[0:32:21.9] SC:** That was incredible.

**[0:32:23.2] GB:** That was incredible. Yeah, I was going to bring that up as another leading indicator of this kind of thing.

**[0:32:28.6] SC:** Yeah. We've already seen like there's some other research. I forget if it was related to GANs or another approach where you're able to give a kind of static photographers, you're able to kind of create three dimensional and change the expression on the photographs. All of these, like the pieces are all in place or near in place to create these perfect synthetic videos, although the full end-to-end thing isn't quite there yet.

**[0:32:55.9] GB:** Yup. You know what you need to get the full end-to-end thing in place?

**[0:32:59.7] SC:** What's that?

**[0:33:00.3] GB:** Compute.

**[0:33:01.9] SC:** Mm-hmm. That's it? That's your — It's just compute?

**[0:33:05.7] GB:** For that particular problem, I think that our ideas are really proving out, and if we were able to run at larger scale, that we'd have a really good time. It's, I think, really important to also drill into this story around compute, because it's not as simple as just you take the code someone already wrote and you just run it on more GPUs and it's magically going to solve the problem. But it's much more that it's like compute in this field is just like particle accelerators in physics. If you don't have the particle accelerator, there's not going to discover the secrets of the universe. You're not going to have your breakthrough.

If you have the particle accelerator, it's not just that you just — Somehow, like the physicist is not useful and just translating like ideas into experiments. It's that you now have this tool that fundamentally allows you to achieve the result that you were looking for. That's really where we are on video generation, is that we have ideas that are clearly in the right space and maybe we need some additional tricks. Maybe we need to do some additional tuning, but if we're able to run at much larger scale than we are right now, then we can actually try out these ideas that we have.

I think that the converse is also true, that if for whatever reason we were to freeze the level of compute that is available for running these models, that progress would really slow down.

**[0:34:26.0] SC:** Yeah. Your earlier point about — You kind of hinted at this a couple of times in the conversation, but, I think one of the things that contributes to our ability to predict a couple of things, I think, contribute to our ability to — Or the difficulty we have predicting when AGI happens is I don't know that we've like clearly — Maybe I should phrase this as a question. How well defined do you think it even means to have a achieved AGI? Is it absolute or is there like a minimal viable AGI product that would suffice?

**[0:35:06.0] GB:** Yup. I think this is a good question, and I kind of think of whenever I think of the question of how do you define the AGI? What is an AGI? What will an AGI look like? I always think a little bit of — Have you heard of bike shedding as a term?

**[0:35:17.9] SC:** Yup. Absolutely. You should explain it though. You should explain bike shedding.

**[0:35:22.4] GB:** The idea behind bike shedding, so let's say that you're designing a nuclear reactor. What you'll do is you'll bring in these experts and the experts will tell you things. Honestly, like if they tell you like, "You got to do it this way," and like, "This is really important." You'll probably trust them and say, "Okay. You do it. You've got a lot of experience in this. This is great. Go off and run with it."

When it comes to, "Okay. We're also going to have this bike shed outside and what color should it be?" Everyone is going to have an opinion. Everyone feels like they are uniquely qualified experts to talk about bike shed colors.

I think that with intelligence, there's something similar here, where we all have our own conception of intelligence, what it's like, what's hard, what it is that we do, what's going on in our heads. I think that the question of, "Okay. Well, this system that you built does this, but it doesn't do that. What is AGI going to look like? How hard is it? When is it going to arrive?" I think these things end up being approached kind of like the bike shed, where everyone has their daily experience and kind of fit that to — I think one thing that is true is that no one is truly an expert in AGI, right? We haven't built it yet, and so anyone who is claiming that I've got this special knowledge. It's a little hard to take that at face value, right? You can't go to a university and say, "Well, I got my AGI undergraduate degree."

**[0:36:42.9] SC:** Yeah. And build five of them in the course of getting it.

**[0:36:45.5] GB:** That's right. I think that the bike shedding term is like I think kind of a negative one usually, but I view it as almost a positive way where we have such — Like intelligence is just so fundamental to us and who we are and this notion of what it even means to human, that is everyone has thought about this. Thousands of years ago people were speculating about what it would be to build a mind and what goes on inside of our own heads. I think it's actually kind of this marvelous thing that people care so much, but the flipside ends up being that you almost have this philosophical debate that becomes very irreducible.

The way that I think about it is that, to the extent, we're going to try to resolve philosophical questions that have been standing for 2,000 years. We are probably out of luck except to the extent that our technical progress informs us. For example, we now have a much better sense

of what it's going to be to build a mind than Aristotle would have. It's not going to be some big rule-based system. It's not going to be most of the things that you might have expected. It's going to be a big statistical system that's going to run this massive parallel fashion on a bunch of cores and kind of describe things like that. Is it going to be matrix multiplies and taking gradients? Well, that's a different question, and that we certainly have not resolved that yet.

I think that the question then of, "Okay. What is an AGI?" A lot of how I like to frame that conversation is to kind of sidestep the deep philosophical questions of do you need to have something that's conscious? Do you need to have something that kind of fulfills other notions of intelligence and really just focus on what can it do? Can you build a system that is able to accomplish any economically valuable task that you put in front it? I think that is something where you can tell, right? I think that you can tell if — Another way I reason about this is if you took a human and you want it to figure out is this person, again, real intelligence? Is that something that you think you could test? We certainly spend a lot of time trying to assess various people's skills and capabilities on various different axes. You can almost think of it as deciding if you built an AGI as giving it a bunch of different job interviews and seeing if you want to hire it. I think that this kind of framing of there are deep philosophical questions, but at the end of the day you can think about it instead in terms of very functional what is this system capable of? The later is something we're able to do. The former is something that is fundamentally very hard.

I also think that this framing really raises a second point, which is, well, is this — It's a very utilitarian kind of view of the kind of system that we're talking about, the kind of things that we might want to build and why should we want to build something like that at all? If it really opens this Pandora's box of what is it mean to be human and what is the value of human? How do we make sure that humans have meaning and really a place in the resulting world? That is, I think, the hardest problem and that is something that is something that's very core to OpenAI and how we think about this technology is that it's pretty clear that — Like, I think, indisputable in a short term that companies are pouring in tons and tons of resources in order to make advances in AI, which is different from AGI. I think that the amount of resource is going to that is smaller, is more focused, but I think that as it feels closer to people, as people feel that, "Wow! Look at all these progress in AI. Feels like kind of my internal neural net is telling me that this could actually happen," and then you start thinking through the economic value that will be delivered by that

system and how important it could be for X-company or Y-company. I think that that will change, and then I think that the question is not so much about accelerating the timeline to AI, but it's really about ensuring that this technology plays out in a way that isn't just one company gets all the spoils, but is really about humanity is ultimately the winner.

**[0:40:46.3] SC:** All right. It may turn out we may get thrown a curve ball here and it may turn out that the technologies and techniques that allow us to create AGI are totally orthogonal to the ones that we've created in the process of trying to create AI. But from where we sit now, it certainly seems like to the point of all the pieces that we discussed that go into creating these videos, like they're all kind of right in line with the kinds of problems we would expect to have to solve in order to get to an AGI. So all of that huge investment that is profit-driven, if we can say on the part of many of the companies, most of the companies that are investing in those technologies are maybe accidentally push us closer to this AGI.

**[0:41:34.7] GB:** Yeah, and it's actually pretty interesting that there's this classic mantra in the field that as soon as you're able to do it isn't AI anymore and people said this about chess. Chess is the most important thing and that only some super intelligent —

**[0:41:49.4] SC:** Predictive typing.

**[0:41:50.4] GB:** Exactly. It turned out all of that stuff, once it happened people were like, "That's not AI." I think that this is dead. I think that this way of people reacting to things we're able to do is now different and you look at AlphaGo, you look at Dota, and for these systems, there really is something going on in them that is very akin to intuition. It's much deeper than simply performing some big search and being very dumb and making up for that dumbness with just having your massive brain.

You think about the image generation that came from NVIDIA, and that's something where humans can't even sit down and start to think about how you could write the rules for it. So I think this is a very encouraging thing and I think that there's kind of this piece to it which is what's really going on right now is that if you look at the problem of trying to recognize a cat or a dog in an image, trying to recognize objects and images, that the space of image is the super complicated, very high dimensional space, is this high dimensional manifold. There's this

fundamental complexity in that domain. For the human to write down all the rules for that would be a pretty massive undertaking.

What we've built is that we have these systems which are able to absorb the complexity of the domain and able to kind of figure themselves around and that you got this neural network that's got these millions of parameters. That's just not something that exist in the natural world. It's not something that we're used to and it's able to reconfigure itself and to really absorb all of the inherent — That inherent complexity.

I think that the ability to do that is what really distinguishes this learning revolution from AI previously, and now it might turn out that there are limits to what we can do with our learning algorithm, but it's also kind of crazy that the learning algorithm we use, that propagation is developed in 1986. How can it be that this algorithm and, really, neural nets in some ways data back to even maybe the 60s, maybe the 40s, depending on how you count. That these very simple, very obvious ideas that you couldn't run on your particle accelerators if you will, you don't have the particle accelerators to run the experiments, but these simple ideas turn out to be so powerful, and I think there's something really fundamental there, that I can't decide between two different explanation.

One is that intelligence is fundamentally simple. That there's a — I can kind of back explain some explanation of, "Well, if you had something that was complicated, then it would have a very large prior," and so you're kind of making this prior. So yeah, you shouldn't expect it to be very general. The more depressing version of this is that well maybe we're just really bad at making anything complicated work.

**[0:44:41.2] SC:** Aha.

**[0:44:41.9] GB:** But if it's the first, and I think there's a lot of evidence that really indicates that it is the first. Then I think that's very encouraging, that the simple ideas, if you implement them correctly, if the mathematics works, if math kind of points to the right direction, if you implement it correct, you scale that massively, then you're going to be able to get things that — Things will happen that you weren't expecting.

One thing that's really weird to me about the kind of progress that I see is that I've seen on repeated occasions algorithms that work better at large scale than their designers expected. We've seen this with algorithms here, talk to the person who invented it and they say, "Oh, no. That's not going to work for X, Y, Z reason," and then we scale that really large and it actually works really well.

I think that this is, again, for me as an engineer, is totally contrary to experience. For me as an engineer, you really only get, if you're lucky, the kind of performance that the person was intending. As soon as your 10-X scale, 100-X scale, good luck.

**[0:45:43.7] SC:** Everything starts to break. Right.

**[0:45:45.0] GB:** Totally. Totally broken.

**[0:45:46.6] SC:** Is there more to that than just more data, and more data fixing more problems in terms — Or more data basically covering for our lack of sophistication in the algorithms themselves?

**[0:46:00.7] GB:** Yeah. It's something like that, though I would phrase it a little differently, which is that I think that the algorithms that we have are fundamentally capable of absorbing all the compute data you can throw at them. The data question is also an interesting one, because the thing that people are used to is supervised learning where you have this big static dataset that encapsulates your world knowledge. But where things are really shifting is towards more of the reinforcement learning paradigm. If you think about it, that's where you want to be. You want to have an environment that you're interacting with that you're able to change, that you have this dynamic feedback loop going on. There, you suddenly have upgraded your environment. You can think of your big set of images is just a static environment and now you've upgraded to this very dynamic world. There, suddenly, you sort of are able to get infinite data or at least you can spend a lot of compute to get a lot of data. If it's a video game like Dota, you can run this on many, many cores. If it's a robotics simulator, you again can spend a bunch of compute there. If it's the real world, okay, you're in a for a little bit of a harder time.

Maybe you do something like Google did with having a big arm farm. Maybe you do something else. I think where really want to end up is that we want to end up in a place where the limiting factor is the amount of compute that we can throw at these models and where we can have massive generative models that have absorbed a lot of world knowledge that you're able to do things inside of that, and we can't run those models yet today. Like we're at the very sort of — We're at the very nascent edge of what I expect we're going to be able to do with generative models and with this kind of approach. Model-based RL is kind of the term of art that a lot of people use. I think that in upcoming years we will able to see lots of progress based on these ideas of scale up. Use algorithms that can absorb all the compute and that that can make up for lack of data, that can make up for lack of everything else.

**[0:47:56.8] SC:** What specifically does model-based RL refer to relative to just RL?

**[0:48:02.3] GB:** Yeah. The idea with model-based RL is that you have a — It may be learned or maybe not learned in some way, model of the environment, that you query and you can kind of explore with them. It's kind of like you as a human, if you picture your house and picture walking around your house and you can kind of plan things out, you can see like, "Oh! If I do this, that this thing will happen," and then you don't actually have to go and spend the very expensive time of walking through your house. That kind of thing, you can see it's very powerful to have this ability to plan and explore an imagination rather than the real environment. Again, it's all very nascent. It doesn't really work right now. I think that it really cannot work until we have the faster computers online.

**[0:48:43.8] SC:** One of the things you said at the very beginning of the interview is kind of stuck with me is interesting, and that is this idea that, ultimately, to train an AGI, it's going to require massive amounts of compute, but then once we train it, like the actual inference letting that AGI be generally intelligent is going to require much less compute. It strikes me that there's some interesting questions there, like what do we do with all that compute. You addressed some of it in terms of — You kind of phrase it as maybe the thing that we are doing is we are running multiple instances of this AGI thing in parallel. We're taking advantage of all that compute that we had to create to train it by running a bunch of these things in parallel, but it also kind of makes me wonder if maybe the AGI doesn't need to be all that general if we're ultimately

segmenting the problems-base up in the end anyway. Does that question make sense? Do you see where I'm going with that?

**[0:49:52.4] GB:** Not entirely.

**[0:49:53.8] SC:** I guess there are two questions here. I guess, one; are there other implications of this idea that you propose that we're going to have to build up this massive compute capability to train the AGI, and then once we've trained it, we need that compute capability less. What are all the implications of that? That's one question.

**[0:50:19.3] GB:** Got it.

**[0:50:20.4] SC:** Question number two is; if ultimately what we end up doing is running a bunch of parallel intelligences, do they all need to be general anyway? Can we have a cluster of intelligences that are really good at thing-X, a cluster of intelligence that are good at thing-Y. Scale that out, and that is what ultimately we start to think of as general intelligence. We just have a bunch of these less general intelligences.

**[0:50:50.4] GB:** Yeah. It makes a lot of sense. On the first one, one thing that I think is worth thinking about is when you actually build a computer system that is autonomously generating huge amounts of revenue or value, there's something that comes as big incentive to make more such computer systems.

Today, if you have a big pile of money, you want to turn it into more money, well you start a company or you invest in a company and you hire a bunch of people and those people produce economic value towards some goal and then it kind of continues this cycle. Whereas if you have a computer that is just as good as a human worker, well then you have a big pile of money, you should build a big datacenter, and there's going to be this big incentive to kind of dot the world with datacenters. I think that's one perspective on what happens on the compute-front. I think it is possibly the case that you can take your big training datacenter and use all of that compute to run a single AI much faster. So rather than — Imagine if you had an Einstein in Silicon, that you're now able to run a thousand X-real-time or a million X-real-time. That'd be pretty good. This person sitting around and thinking about physics and thinking about you get someone in

there thinking about medicine and how to cure diseases. You get someone in thinking about how should we build rockets to go to the stars and all sorts of things like that. That would be pretty valuable, pretty good. It's not guaranteed that we'll be able to use all of that compute usefully in a single AI, but I think that at the very least being able to run parallel copies of these AIs is something that we should expect.

Then there's a question of, "Well, what would that be good for?" I guess when I think about these, I always try to make analogies to things that are in our experience today. In our experience today, why do we ever want to have a group of more than one human doing something? Like building companies and the task are hard and that you have different people that specialize in different skills, and all of those things are things that we should expect, would transfer to the systems of the future. I think it'll be very valuable.

By the way, the idea of computer system that autonomously produces value where all the interesting stuff is done by the computers and the humans just kind of stick around and clean up the fans is something that exist today, it sounds pretty dystopic, but if you look at Bitcion mines, that is exactly what they are.

There's a good article recently with a bunch of pictures from Chinese Bitcoin mines, which I recommend looking at if you want to think about kind of the more cyber [inaudible 0:53:16.1], this topic version of this stuff. Again, there are a lot of hazards here with the technology that we're talking about building. Again, the weirdest thing for me is the fact that it's so — People don't talk about this in a serious way and that I think that the — For most technologies, when you're building them, you think about what happens if we really succeed. I think that for partially historic reasons, partially for this reason that we all feel our own sense of how far off AGI is and how hard it's going to be and how impossible this be, but imagine building it, that really seriously thinking through what happens if it works is something that is a bit taboo. That's thing one.

Question two — Can you remember your question two was?

**[0:54:00.7] SC:** I think question two was, ultimately, do we need AGI at all if the deployment model, if you will, ends up being to — Or the scalability model ends up being to segment our

workload into a bunch of separate things. Does a collection of more specialized intelligences become the thing that we initially come to see as a general intelligence?

**[0:54:29.1] GB:** Yeah. I think that's an open question or that's a possibility. The way that I think about it is I guess, again, back to the idea of we have organizations of humans that can accomplish goals that humans individually cannot. And so it might well be that even though you wound up with specialization, I would certainly expect that you'll end up with specialization towards specific tasks. I think that I would expect that a general AI would also have these very hyper-trained narrow AI modules within it, and you absolutely should do that.

One thing I think is kind of interesting about today's AI systems is if you look at something like the neural turing machine, you basically spend this big model. You spend a lot of compute, a lot of data, a lot of training time in order to learn how to do various tasks. For example, one of the tasks from the original paper is to learn to sort. Pretty cool, right? This system learns how to sort and it's kind of learn this program, but when you really think about it, it's like, "I could do the same thing at Python on my core in like two seconds."

At the end of the day, if you have a specific task you're trying to solve, you can hyper-optimize for that and do a lot better from an efficiency standpoint than this very general thing. I think something similar happens with humans where we have — When you have to sit and think about something when you're not a master of it, you're trying to really reason how it works versus when you practice a bunch and it's in your muscle memory. Kind of like this is has gone to the much more efficient hot path, and I think that we'll certainly see analogs to this sort of behavior.

**[0:56:08.1] SC:** Fair enough. We're at the top of the hour, or beyond the top of the hour actually and we haven't touched on the thing that I expected us to spend a bunch of time on, which is the Dota 2 project, but we covered a lot of really interesting ground in terms of AGI and what that means and what we should be thinking about. What I'm thinking we should do is maybe call this a part one and find some time to get together again to do part two where we dive into the work that you've done on Dota.

**[0:56:39.4] GB:** Sounds good.

**[0:56:40.5] SC:** All right.

**[0:56:41.7] GB:** Perfect. This is a lot of fun. I really appreciate it.

**[0:56:44.7] SC:** Yes, same here. Greg, thank you so much.

**[0:56:46.2] GB:** Oh, great chatting. So long.

[END OF INTERVIEW]

**[0:56:51.9] SC:** All right everyone, that's our show for today. Thanks so much for listening and for your continued feedback and support. For more information on Greg or any of the topics covered in this episode, head on over to twimlai.com/talk/74. To follow along with our OpenAI series, visit twimlai.com/openai.

Of course, you can send along your feedback or questions via Twitter to @twimlai or @samcharrington or just leave a comment right on the show notes page, and make sure to reach out if you're at the NIPS Conference. Thanks again to NVIDIA for supporting this series. Of course, thank you once again for listening, and catch you next time.

[END]