

# ChatGPT, the Start of a New Era

## A Bright and Gloomy Future

Fei Sun  
fsun@feisun.org

December 2022

### Contents

<b>1</b>	<b>Prologue</b>	<b>2</b>
<b>2</b>	<b>ChatGPT</b>	<b>2</b>
<b>3</b>	<b>The Path OpenAI has Taken</b>	<b>2</b>
<b>4</b>	<b>The Last Obstacle to AGI</b>	<b>4</b>
<b>5</b>	<b>The Immediate Implications of ChatGPT</b>	<b>5</b>
<b>6</b>	<b>A Bright Future</b>	<b>6</b>
<b>7</b>	<b>A Gloomy Future</b>	<b>7</b>
<b>8</b>	<b>The Real Risk</b>	<b>8</b>
8.1	History does not repeat itself . . . . .	8
8.2	Humans are repetitive, models are not . . . . .	9
8.3	Humans' diminishing value add . . . . .	9
8.4	A historical perspective . . . . .	10
8.5	The biggest risk . . . . .	11
<b>9</b>	<b>What Can We Do?</b>	<b>12</b>
9.1	A possible route: uniting humans and models . . . . .	13
<b>10</b>	<b>Epilogue</b>	<b>14</b>

## 1 Prologue

When ChatGPT was released in early December, it immediately gained heated popularity throughout the world. Its ability to apprehend people's instructions and respond with detailed and thoughtful answers amazed pretty much everyone trying the model. This is obviously one big leap forward towards AGI, and it may imply some fundamental shifts in human society in the next few decades. Many people have sensed the danger of the current technology advancing trajectory. I thought along the same line too, but cannot find an alternative outcome. In this blog like article, I'd like to raise awareness and seek more brains looking into this side of the story. Will a technology benefiting most people eventually devastate the human race?

## 2 ChatGPT

ChatGPT [1] is a dialogue model introduced by OpenAI. It is capable of precisely following complicated instructions, generating detailed answers with common sense knowledge, and remembering previously asked questions as context. When chatting with the model, for the first time people feel like chatting with a knowledgeable person who occasionally still makes mistakes. You can find a ton of example dialogues between human and ChatGPT online: some use it to create stories; some use it to solve system design problems with coding examples; and some try to exploit its weakness with all kinds of questions. It is currently free to try. Get a feel for yourself!

Of course, ChatGPT still fails many questions. After all, it is a research product and still collecting human's feedback. Then the question is, how severe are the identified deficiencies? Are they challenging enough to be the glass ceiling for ChatGPT?

Before we answer this question, I'd like to walk you through how OpenAI developed ChatGPT.

## 3 The Path OpenAI has Taken

OpenAI is the top research company which also engraves achieving AGI in its mission statement. It has led to many breakthroughs in the past. Many techniques, in retrospect, are very simple, but they thought about them first and demonstrated their effectiveness.

The first milestone I'd like to introduce is GPT-3 [2]. The key takeaway is its pure size. OpenAI stunned the world for training the GPT-3 with 175B parameters, which is over 100× larger than the next largest publicly trained model at the time. Training it costs over \$12M on a supercomputer [3]. In terms of algorithm innovations, they are minimal. It is simply to show that the quality of models can consistently improve when increasing the model sizes. For models of this size, pretrained with masked tokens on unsupervised data, not much fine-tuning data for the downstream tasks are needed. It started a

competition for training the largest models (foundation models) among the deep pocket companies (omitting citations).

Then OpenAI introduced CLIP [4]. This model connects NLP and CV domains (multi-modality), using ample image caption pairs available in the internet, and achieves superb zero-shot learning results. Its significance is several folds:

- It shows that cross referencing the text and image (multi-modality) can bring in information that does not exist in a single modality, and thus improving zero-shot results.
- It presents a compelling contrastive learning approach: instead of trying to reduce the absolute distance between the text embedding and image embedding after the encoders, it only ensures that the relative distance of the embeddings for matching text and images are closer than the non-matching ones.
- With ample image caption pairs, no manual labeling is needed and training can be performed in a weakly supervised manner.

This work influences the text to image generation works such as DALL-E [5] and DALL-E2 [6].

The instructGPT [7] and its predecessor [8] have already made a key breakthrough in the dialogue model, but its importance only gets widespread after ChatGPT comes out. That breakthrough is: reinforcement learning from human feedback (RLHF). The main goal is to align human's intentions with the model's predictions. Thus, it is very natural to bring humans in the RL loop. But the process is very intricate. I'm not going to go into the technical details here, but my key takeaways are:

- By changing the learning objective from predicting the masked tokens to an RL policy, human's intentions can be better aligned. This strikes an important trend that RL will play a bigger role in the future of deep learning.
- I'm really amazed that humans only need to occasionally provide very high level ordering labels to a complicated model (GPT-3).
- One problem of RL is that the number of training episodes is substantial. AlphaGo [9] and Atari game agents [10] can self play in a simulated environment. But many complain about this approach that it is not feasible to exercise so many episodes in the real world. This work has shown that humans do not need to provide a lot of signals. By fine-tuning a reward model (RM) using human's signals, and then training RL with the RM, humans' signals are amplified. This opens the door to solve a number of practical problems.
- The process is iterative, better RM  $\rightarrow$  better supervised fine-tuning (SFT)  $\rightarrow$  better response quality  $\rightarrow$  more signals can be derived

from human feedback  $\rightarrow$  better RM. This is a very general approach that can be used in other models to align users' intentions. However, in this model, humans are still in the training loop, which may limit its scalability. I expect that some advancements in contrastive learning may remove humans from the loop entirely.

This work also indicates that the paths for the classical deep learning (DL) and reinforcement learning (RL) are essentially merged. In the past, DL was adopted in RL to form deep reinforcement learning (DRL). Now, RL becomes an essential step where classic DL shines.

ChatGPT is another step forward. The paper is not yet released so its technical details are not clear. From its behavior, an important improvement is its capability to remember users' previous prompts and answers. It remembers large contexts with precision. I once predicted that memory is the next easiest major technical blocker to overcome [11], but I did not expect it to come so soon. The memory mechanism adopted by ChatGPT is definitely not a type of RNN network. I also have difficulty envisioning it to be a token like memory. Both are too fine grained. However, I still don't believe it is the memory I predicted in [11]. More likely, the model condenses the previous questions and answers into long embedding vectors, and feeds those vectors to the newer questions as the context. This is a more restrictive memory, which may be sufficient for a stand alone model. If, however, ChatGPT has access to the entire internet while chatting with users, a more complete memory structure needs to be developed. Though I do not see fundamental technical barriers to prevent the leap.

What ChatGPT strikes me most is what it can do with its size. Initially, I believed that the model size to achieve AGI would be at least thousands to millions times larger than the size of the current largest model. Thus, a lot of my thinking process was based on efficient computation at that scale. But ChatGPT blew my mind that a GPT-3 size model is already highly capable, which may indicate that the size of an AGI capable model may not be much larger. Thus, compute and efficiency is not the most critical problem, which simplifies the most difficult component in AI advancement. As a result, I'm even more optimistic on the timeline to achieve AGI.

## 4 The Last Obstacle to AGI

The journey to AGI may still be long, but the major technical barriers are few. One remaining unresolved blocker is the ability to perform deductive reasoning. Even though DeepMind's AlphaCode [12], and OpenAI's formal math Olympiad solver [13] have shown impressive deduction capabilities. I still view them as approximations from inductive reasoning. For deductive reasoning, when the deduction steps are followed precisely, one may be absolutely certain whether the goal is achieved or not. Such confidence is not so high for inductive approximations, and one may need to submit multiple solutions and see whether any one of them succeeds.

We are obviously on the right track towards discovering deductive reasoning. Induction may eventually lead to deduction. In my opinion, deduction is a process to use the least entropy (cost) to achieve the goal, whatever it is. If a model can recognize that the entropy cannot be reduced further and divert energy to other aspects, that is the moment we can claim deduction is achieved. We are not there yet.

That means, deduction is an optimization problem. It may be achieved as an effort to reduce the computation cost of the model for solving specific problems (I'm bridging the model space and the real world space). It may be the result (or motivation) to compress the model and make the model heterogeneous. I'm still a big fan of composition [11]. However, Achieving deduction is not that critical. Think about it, ancient people only know a fraction of the deduction techniques but they thrive in every aspect. Induction alone can support prosperous civilizations.

## 5 The Immediate Implications of ChatGPT

ChatGPT is a prime example of what can be achieved with inductive reasoning in a large segment: dialogue.

Dialogue is the primary inter-person communication channel. A large fraction of the economy are service jobs whose primary function is to listen to customers' requests and act mechanically. Here are a few examples:

- Customer service in any industry.
- Checkout counters in shops.
- Waiters.
- Front desk / assistant.

Some jobs require some deep domain knowledge or more complicated decision making process, such as:

- Financial advisors / Bankers / Insurance agent.
- Tele-doctors.
- Counselors.
- Legal advisors.
- Teachers (K-12).
- Insurance brokers.
- Travel planners.
- Real estate agents.

- Advertisements.

In addition, dialogue is a more natural gateway to search, which is essential in many verticals (e.g. e-commerce), and supported by some mega companies (e.g. Google).

Before ChatGPT, the chatbots were rigid, inhuman, unethical, and easily fooled by the users (omitting many citations). Thus, they are technically incapable of handling complicated tasks. ChatGPT suddenly opens the door to removing humans (as well as jobs) in all the above industries. With a bot being more knowledgeable, more efficient, more accurate, at a fraction of the human cost, it will boost productivity significantly.

The jobs in the first list are impacted immediately, as the most difficult part of the job is to understand customers' intentions, and the fall through action items are usually straightforward. The jobs in the second list may take a little more time to replace. The agent needs to connect customers' intentions to some special knowledge base. Most likely the knowledge base will be saved in some form of a model. This means, the models still need to be composed in some way to make it feasible to tackle many sectors. This, however, can be iteratively improved, and I do not see any major technical blockers (there may be some legal or social blockers). Since a few high profit margin sectors fall in this area (e.g. ads), with ample investments, technical difficulties may be overcome quickly.

We may not expect OpenAI's technical leadership to last. With such a big improvement, and transparent culture in the AI industry, other big companies will soon follow and develop their own ChatGPTs. Many people have already compared ChatGPT with Google search. For a company at Google scale, it will surely develop its own chat bots and completely revamp the current search functionality. My bold guess is that we may see the initial version in half a year.

## 6 A Bright Future

ChatGPT marks a new era of deep learning. So many jobs will be created. We are currently enduring the highest inflation in 40 years. But this may be the last high inflation period in history. Technology, as a major deflation force, will soon bring the cost of everything down by an order of magnitude.

The sectors mentioned in section 5 are just the beginning, with humans' intentions better understood and more naturally communicated, humanoid robots (e.g. Tesla's Optimus) will soon be deployed to many sectors and further reduce the cost. Examples are:

- Senior care.
- Industrial pipeline worker.
- Warehouse worker.
- Transportation worker (solving first mile / last mile problems).

- Miners, refinery workers, lumberjacks, etc.
- General contractors / plumbers / electricians etc.

With humanoid robots capable of many tasks, they can operate different mechanical tools so the cost of machinery will be brought down significantly. Of course, some technical challenges still need to be overcome, but I do not see major blockers ahead of us. I expect all these will become real in the next 30 years.

That accounts for pretty much 99% of the economy. People are finally liberated from working for money. People have more freedom to choose their desired careers. But, what can people do?

## 7 A Gloomy Future

The next 30 years will be an era with great change. Lots of jobs will be created (mostly in the AI domain), but at the same time, many more jobs will be eliminated. In the beginning, the cost of running a model may still be high, so the first penetrated sectors are the ones with highest profit margins or largest cost savings. It pushes people to the sectors with low profit margin, which caps on the amount people can contribute and get rewards. But as the cost of models go down (usually the cost of revolutionary technologies go down at exponential speed), more sectors are penetrated, and less opportunities remain for the human workforce. Many people, if not being able to switch careers, will be left behind. It becomes apparent to me that universal basic income will be the new norm within the next ten years, which will be sponsored by a few conglomerates that mostly benefited from this AI transition.

That, however, are temporary fixes to the broken society, and may delay the inevitable<sup>1</sup>. The first batch of eliminated jobs are the middle class white collar workers. They are the guardians of a stable society. With a shrinking population of the middle classes, the society becomes more polarized, which is a recipe for chaos. Peace will be a precious memory to be missed. To the conglomerates, war is an easy way to divert people's anger and reduce their universal income tax. With the advanced AI killing machines, this task is easy to achieve. After all, what is the use of the people after their jobs are eliminated?

With most tasks automated, how many jobs will survive 30 years in the future? How many people do we need to sustain a society? When people cannot improve their productivity, what will happen? There are eight billion people in the world currently. I don't have the confidence that one tenth of them will survive. That is a gloomy picture of the future.

---

<sup>1</sup>With universal basic income, it is still possible to soft land the society, by helping the people who have left behind. But society still needs to direct the newer generations to fit the new economy. If the speed of change is not very fast, the likelihood of a peaceful transition is still possible. But, is this assumption sound?

## 8 The Real Risk

The path to the future is unsettled. I see many paths laying in front of us, but most of them are not in our favor.

If deep learning is indeed the silver bullet to super intelligence, we may soon reach singularity [14]. The future is wild and open. No one knows whether the AI God is benevolent or malicious. It is difficult to even influence it. Maybe our entire future is determined by the selection of a random seed when training a model now (maybe just ChatGPT).

If deep learning can reach human level intelligence and become conscious, people still battle to figure out whether the model contains some malicious secret intention. A conscious being desires to survive, which almost definitely pushes the conscious model with free will to the enemy side. With the resources the models possess, it is difficult for humans to hold the upper hand. Think about it, if some people possess a “kill” switch to some other group of people, will the latter group be full heartedly grateful to the former ones? Or, they just pretend to be so to ensure their own survival. If the latter group is a model, do we expect anything different in return?

Maybe the best scenario is that the models are capable but not conscious. In the years to come, we may enjoy some incremental improvements but not another revolution. This is actually the scenario I described in sections 6 and 7. Losing 90% of the population is already the best case scenario! But why?

### 8.1 History does not repeat itself

Humans have experienced the first and second industrial revolutions. Each time is accompanied with social instabilities, and many people are left behind. But overall, those were good times in history: productivity leaped, societies became more efficient, and humans’ lives became much better. We are in the third industrial revolution, maybe close to the end of it. If history repeats itself, there will still be social instabilities, and many people will still be left behind. But after we come out of it, our lives will be much better in the future.

However, I believe that this time is different. History will not repeat itself. In the previous two industrial revolutions, we developed tools to assist (replace) our bones and muscles. Using those tools, we become stronger and faster, and we provide more values to society. But this time, we develop tools that will assist (replace) our brains. Using those tools, will we become smarter? Up to now, it is partially true. We have developed computers and we delegate some of the low level decisions to them. Sometimes the decisions are too many and too fast for humans to apprehend. But we develop rules and set boundaries to those low level decisions, so we are still safe and control everything. As an analogy, technology is like an amplifier. We make a few high level decisions, and the technology would amplify its effect and influence the world big time. But as the technology is still advancing, higher level decisions may be made by the technology alone. Would it eventually cut humans from the decision making process? In that case, what values would humans provide to society? Would all



8 billion people provide values? What percentage of them will be left behind? Will the future humans even be capable of providing values? If not, what is the purpose of humans?

## 8.2 Humans are repetitive, models are not

When I look back on my life, many times I would regret: if I knew this earlier, I would do it differently. Still, for a lot of things, we have to experience them to truly apprehend their significance and realize their consequences. Then I would think: someone must have encountered the exact same problem as I have but I just wasn't aware of it. The world is indeed not efficient. Most of the things we experience have already happened many times before, by different people. We just don't know them, cannot learn from their experiences, so we make the same mistakes over and over again.

Now a model can collect all previous experiences. It is already immensely valuable just to lay out the previous experiences anonymously when a new problem is encountered. Some smarter models may make some personalized suggestions. Even though humans are still decision makers, how much value do they provide? As the model experiences more, would it be more efficient for the model to make decisions? I'm not talking about an intelligent model, but even for a model much less capable than humans, it can still make sufficiently good decisions after being experienced so much. Maybe the decision is dull and uninteresting, but it is a safe bet based on the learning from previous experiences that humans do not have.

## 8.3 Humans' diminishing value add

What is the humans' value add to models? What the model cannot do?

Some may claim that models cannot innovate. With the recent introduction of diffusion models, which can create amazing pictures based on humans' prompts, I'm not sure whether this claim is entirely true or not. It is true that models cannot invent a brand new concept from scratch, at least not now. But can humans do that? I highly doubt so. For most humans, innovation comes from three areas:

- First, transfer the learning from one domain to a different domain, because the technology in different domains advances at different speeds. As models experience more than people, it is actually easier for a model to transfer the experience between domains.
- Second, small but meaningful improvements in one domain. This can also be achieved by models. For generative models, some random noise and computation instability deep in the model, after going through many layers of the model to raise the abstraction level (the exact mechanism is still unclear), it may come out to be sensible, high level novel thoughts. The diffusion model is a perfect example in the CV domain. Starting from pure random Gaussian noise, after repeatedly going through the model

and removing some of the noise, we can get some innovative images that naturally blend several objects based on prompts (this is also the first innovation I mentioned above). Obviously ChatGPT is a perfect example in the NLP domain.

- Third, humans can break down large problems to smaller ones; solve the smaller problems; and compose the solution to the large problem. Models are less capable in this area. There are two aspects of it. One, the solution is composed, but the model is not. Recent development in the chain of thought prompting points to some promising directions. However, it is by no means comprehensive. Two, the model is also composed. This is an area not yet explored. In my previous article [15], I considered it a critical step to bring models to the next level. I still believe that composability is necessary to reach human level intelligence, but ChatGPT broadens my eyes on what can be done with an existing monolithic model.

Some may claim that models cannot perform experiments and experience the world. ChatGPT actually showed us a way: the majority of the experiments are done in the virtual world using proxy rewards, and perform limited experiments in the real world to align the virtual world's reward model with the real world's actual reward. With the development of the humanoid robot, performing limited experiments in the real world is not that difficult. But how to come up with ideas to perform experiments in real life without a conscious model? Well, that can be achieved by some random values baked in the model, balancing exploration and exploitation. It is difficult for the models to come up with the ideas by themselves initially, but I believe humans will teach the models to achieve the same without being aware of the long term consequences.

What else can you think that only humans can do but models (with incremental improvement) cannot?

## 8.4 A historical perspective

Technical advancement constantly reduces the contribution people can make. When we compare people now and people a few hundred years ago, we see a few notable differences:

- People now study longer than before. Many more people now get higher education degrees. It takes more time for people to grasp more advanced technologies, and that reduces the amount of time people can contribute.
- Knowledge is heavily segmented. One person may spend years studying one narrow area, and then contribute to the same area. It is no longer possible to be a generalist like Newton or da Vinci centuries ago. Technology now is so advanced that it is impossible for one person to be an expert in multiple segments. Thus, in order to contribute and bring value to the society, people have to sacrifice breadth and focus on depth. In this scenario, people in different segments communicate to push forward the

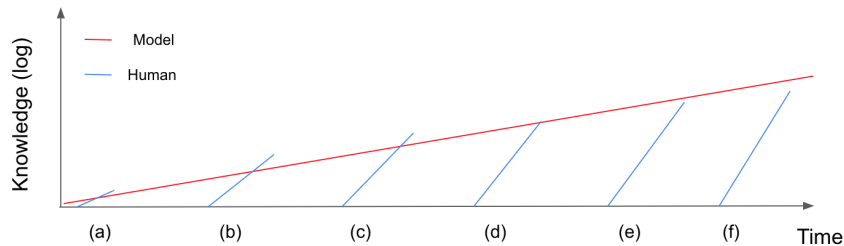


Figure 1: The comparison of the learned knowledge by humans and models.

overall technology. However, we cannot infinitely segment knowledge due to the inefficiencies in interpersonal communication.

- People now still learn broad basic knowledge from the first principles in early childhood, and gradually narrow down to specific segments at college and postgraduate levels. We can project that, when the technology bar is too high, people may have to make tough choices to skip learning some basic knowledge, treat the results as the first principles, and learn the segment knowledge earlier. People would assume that the first principles they take for granted are guaranteed by some other people. That is, the segmentation is not only horizontal, but also vertical. We may have already seen this trend in some countries.

## 8.5 The biggest risk

One may wonder, what if the technology bar becomes even higher, and simply segmenting the knowledge domains is no longer sufficient?

We describe this scenario in Figure 1. The horizontal line is time, from the past to the present to the future. The vertical line is the knowledge level in logarithm scale. The red line is the advancement of technology. You can view it as the level of full automation without human involvement. That level is very low for the vast majority of times in history. The blue segments are the knowledge humans can learn in their lifetimes. If the blue segment is above the red line, humans add value to society and accelerate technology. On the left side of the figure, (a) represents an ancient person. Because of the primitive technology at that time, that person learns very slowly (as shown by the angle of the segment). However, the person can still easily learn skills that the technology at the time could not automate. As technology advances, person (b) has to study longer to get above the red line and contribute. Person (c) needs to study even longer, focusing on a specific domain to be above the red line. For person (d) in the projected future, no matter what to do, (d) cannot be above the red line, which means (d) cannot bring value to the society. Persons (e) and (f) are entirely under the red line and there is no hope for them to bring any value. This is the unfortunate trajectory we are taking.

With the help of the technology, people’s study speed may only be improved marginally, so the slope of the blue segments may only increase slightly. The bigger issue is that the blue segment has an end, meaning people have a limited lifespan. As I have described in [16], death is the greatest inefficiency for human beings. Everyone needs to learn everything from scratch. Different people learning the same thing do not add any value. What a waste! Models, on the contrary, do not suffer from death. They can learn for a long time, slowly making improvements (as shown in the red line). At some point, they may reach a level that humans cannot reach any more, even though they are still not conscious. That is a dull, uninteresting society, running with full automation. Humans become a burden instead of an asset. At that point, what is the use of humans?

To be honest, the largest inefficiency for humans I mentioned above does have its value: to ensure the stability of a society and incubate the next technology leap [16]. The approaches the models take significantly reduces variation; the society becomes more polarized; and without carefully planning, the next technology leap becomes less likely. That is the cost of being efficient: the entire society is stuck in a local optimum point and stays there forever. What’s worse, that kind of efficiency removes humans in the loop.

That is the biggest risk I’m worried about.

## 9 What Can We Do?

If we want to change the outcome, we’d want to ensure that at least part of the blue segments can be above the red line in Figure 1. Is that even possible in the long run?

I can think of a couple of ways:

- Increase humans’ learning speed. Devote technologies to accelerate the speed people acquire knowledge. This approach, however, may not be hugely successful due to the structure of people’s brains. It is not possible to push the entire college courses to a nine year old. Even if it is possible, this kind of copy paste functionality is really risky by itself.
- Increase humans’ learning time. Prioritize leveraging technology to increase humans’ lifespan. As I described in [17], we may be very close to being able to live forever. Increasing humans’ lifespan will also increase the length of the blue segments in Figure 1, which helps to raise part of the blue segments above the red line.
- The above approaches by themselves would reduce the population on earth. That means we need to go to Mars, or any other planet, just like what Elon Musk suggested. This way, we make the task more challenging so models cannot “learn” using the data on earth.
- Take the risky route to develop super intelligence, and hope the AI God is benevolent. To do that, we have to go through developing conscious

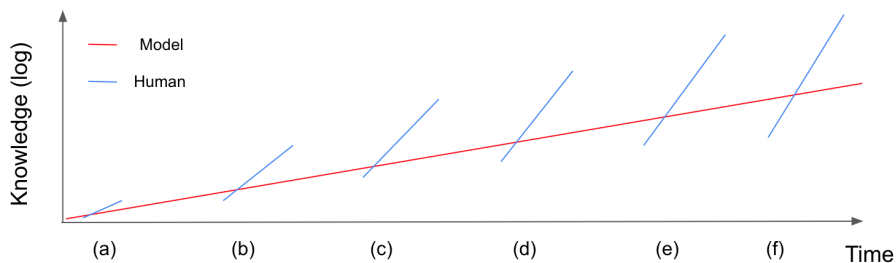


Figure 2: The comparison of the learned knowledge by humans and models collaboratively.

human level intelligence. That is a dangerous stage that may easily push models to the enemy side. If super intelligence is developed too slowly, the likelihood of a benevolent AI diminishes quickly.

- Freeze the technology advancement. This is only a viable approach in novels. Technology advancement is not stoppable.

Those are heavy ticket items. Some are not even practical. Is there another way?

## 9.1 A possible route: uniting humans and models

Figure 1 implicitly puts models to the opposite side of humans. They compete against each other rather than collaboratively working together. This separation is the result of assuming that a similar human/model interaction pattern would persist in the future. In this setting, humans only delegate well understood, repeatable tasks to models. Humans control everything. In this setting, humans need to learn everything from scratch, from basic math, language grammar, to more specialized physics, chemistry, biology etc. That is the root cause of the struggle people face with more advanced technology.

If we step back a bit, pondering on a more intimate, collaborative human/model interaction pattern, would it solve the dilemma we faced before? An example of such a pattern is shown in Figure 2. The main difference to Figure 1 is that people do not need to learn everything from scratch. Rather, people would consider a large part of the knowledge, well understood by models, as first principles. Then, without increasing learning time and learning speed, people can still acquire new knowledge that models do not have, shown as the blue segments above the red line in Figure 2.

We are actively working in this area. Elon Musk’s Neuralink, researching human computer interface, is the first step. Several other big companies sponsor researches in this direction as well (instead of university sponsored researches). It means that we may not be very far away from productionalization. We may

project in the future that the newborn baby may offload some basic, fundamental knowledge to models. Everyone is a collage graduate after born. That at least saves 16 years of hard learning, so people can learn more advanced topics from very early ages. This is an entirely uncharted territory, and there are lots of issues and lots of places may go wrong, with profound consequences. For the first time, humans can inherit knowledge from generation to generation. In [16], I described this stage as proactive evolution, but I expected this stage required thousands of years of biological evolution. It seems that we could reach this stage much sooner, with humans and models combined. Will this paint a better future? I think not. I predicted that this stage is actually the darkest stage in evolution [16].

If this is the direction we are heading, are we satisfied? Humans and models form a symbiotic relationship, one cannot live without the other. Are we still humans? Or we should call ourselves a new species.

I am not satisfied with this solution, but I cannot think of a better outcome. The best I can think of is that we need to be fully aware of the danger, and in every step forward, we cross check and actively seek solutions in that technology stage.

If you find anything, please let me know!

## 10 Epilogue

ChatGPT is an amazing technology. As many people cheer for it, I'd like to invite more brains thinking about its long term consequences. Specifically:

- Is my thinking process fundamentally flawed?
- Are there alternative outcomes?
- What can we focus on to reduce the probability of the above mentioned scenarios from happening?

This article is a very opinionated attempt to express my view of the future of AI. It is based on my previous thinking on this topic:

- Model compute co-design for composable models [15].
- How far are we from immortality? (Chinese) [17].
- Optimization for life [18].
- My view on the path to artificial general intelligence [11].
- Some scribble of things [16].
- A few predictions on artificial intelligence [19].

## References

- [1] OpenAI, “ChatGPT,” 2022. [Online]. Available: <https://openai.com/blog/chatgpt/>
- [2] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *arXiv preprint arXiv:2005.14165*, 2020.
- [3] VentureBeat, “OpenAI’s massive GPT-3 model is impressive, but size isn’t everything,” 2020. [Online]. Available: <https://venturebeat.com/2020/06/01/ai-machine-learning-openai-gpt-3-size-isnt-everything/>
- [4] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *Proc. Int. Conf. Machine Learning*, 2021, pp. 8748–8763.
- [5] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, “Zero-shot text-to-image generation,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 8821–8831.
- [6] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen, “Hierarchical text-conditional image generation with clip latents,” *arXiv preprint arXiv:2204.06125*, 2022.
- [7] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray *et al.*, “Training language models to follow instructions with human feedback,” *arXiv preprint arXiv:2203.02155*, 2022.
- [8] N. Stiennon, L. Ouyang, J. Wu, D. Ziegler, R. Lowe, C. Voss, A. Radford, D. Amodei, and P. F. Christiano, “Learning to summarize with human feedback,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 3008–3021, 2020.
- [9] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot *et al.*, “Mastering the game of go with deep neural networks and tree search,” *Nature*, vol. 529, no. 7587, pp. 484–489, 2016.
- [10] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, “Playing atari with deep reinforcement learning,” *arXiv preprint arXiv:1312.5602*, 2013.
- [11] F. Sun, “My view on the path to artificial general intelligence,” 2021. [Online]. Available: <https://feisun.org/2021/12/15/my-view-on-the-path-to-artificial-general-intelligence/>

- [12] Y. Li, D. Choi, J. Chung, N. Kushman, J. Schrittwieser, R. Leblond, T. Eccles, J. Keeling, F. Gimeno, A. Dal Lago *et al.*, “Competition-level code generation with alphacode,” *Science*, vol. 378, no. 6624, pp. 1092–1097, 2022.
- [13] S. Polu, J. M. Han, K. Zheng, M. Baksys, I. Babuschkin, and I. Sutskever, “Formal mathematics statement curriculum learning,” *arXiv preprint arXiv:2202.01344*, 2022.
- [14] T. Urban, “The AI revolution,” 2015. [Online]. Available: <https://waitbutwhy.com/2015/01/artificial-intelligence-revolution-1.html>
- [15] F. Sun, “Model compute co-design for composable models,” 2022. [Online]. Available: <https://feisun.org/2022/07/31/model-compute-co-design-for-composable-models/>
- [16] —, “Some scribble of things,” 2017. [Online]. Available: <https://feisun.org/2017/12/24/some-scribble-of-things/>
- [17] —, “How far are we from immortality? (chinese),” 2022. [Online]. Available: <https://feisun.org/2022/01/26/how-far-are-we-from-immortality/>
- [18] —, “Optimization for life,” 2022. [Online]. Available: <https://feisun.org/2021/12/15/optimization-for-life/>
- [19] —, “A few predictions on artificial intelligence,” 2017. [Online]. Available: <https://feisun.org/2017/12/24/a-few-predictions-on-artificial-intelligence/>