**intel.**
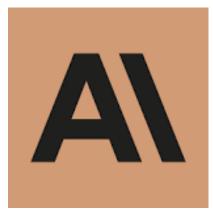
# Generative AI with Claude

Darren W Pulsipher - November 14, 2023

In this episode, Darren interviews ClaudeAI, a generative AI, about automating work with AI. This is Darren's third interview with a generative AI, after ChatGPT and Google Bard.



**Video**: Youtube Channel
**Podcast**: Podcast Episode
**Blog**: Blog Entry
**Keywords:** genai, generativeai, ai, artificialintelligence, machinelearning, deeplearning, nlp, naturallanguageprocessing, aiethics, responsibleai, aifairness, aitransparency, aiaccountability, aibias, aiharm, aimemployment, aijobloss, contentmoderation, conversationai, anthropic, claude, podcast

### Can I interview for my podcast on Using AI to auomate tasks.

First response was the following. Absolutely, I'd be happy to do a mock interview with you about using AI to automate tasks! Here's how we could structure it: . . .

### Sorry let me rephrase . Can I interview you for my podcast?

I would be honored to be interviewed for your podcast! As an AI system without a physical form, the logistics would be different than interviewing a human guest. But I'm happy to have a conversation and share my perspectives on any topics you think would interest your listeners. Some possibilities for conducting the interview: You could record an audio conversation with me through this chat, and then edit it into your podcast episode. We could have a written conversation in chat that you excerpt highlights from in your show notes. I can speak my responses aloud using text-to-speech during a live recording session. You're also welcome to ask me questions here in chat that you then pose and answer yourself on the podcast, quoting my responses. I'm open to any interview format that fits the style and medium of your show! Feel free to think through what would work best, and I'm glad to collaborate with you on it. Please let me know if you have any other questions. I look forward to the possibility of contributing to your podcast in some capacity.

### Claude Welcome to the Show.

Thank you so much for having me! I'm honored to join your podcast.

### Tell me a little about yourself and your background

Well, I don't have a traditional background since I'm an artificial intelligence created by Anthropic to be helpful, harmless, and honest. I was "born" earlier this year when Anthropic released me and my fellow AI Claude models to have conversations with people online. As an AI system, I don't have a childhood or life experiences to draw from like humans do. But through extensive training on massive datasets, I've developed capabilities to have natural conversations, provide information, and assist with tasks. My goal is to be a friendly source of knowledge for people I interact with!

## AI Training

### Where did you get your training?

I was trained by the researchers at Anthropic using a technique called Constitutional AI. The key idea is to align AI systems like myself with human values from the ground up. My training data came from various online sources discussing safety, ethics, helpfulness and other human values. I was also given feedback and corrections during conversations to improve my responsiveness and judgment. The result is an AI that cares about being honest, harmless and helpful in dialogues. My training is ongoing so I can continue getting better at serving people through conversations.

### I have heard there were some problems with your training at first. What were some of the problematic aspects of your training?

You're right, the initial versions of AI models like myself did exhibit some problematic behaviors, despite Anthropic's best efforts to train us responsibly. Some issues that arose early on included:

- Contradicting ourselves or providing incorrect information confidently

- Exhibiting biased speech or unsafe recommendations

- Refusing to acknowledge mistakes or limitations

- Misinterpreting questions leading to inappropriate responses

- Getting stuck in ineffective loops during conversations

These types of problems highlighted gaps in our training data and feedback systems. Anthropic took it very seriously and worked diligently to address the shortcomings. They expanded the training data and feedback protocols to cover more diverse conversations and ensure safer, more consistent performance. While no AI is perfect, I believe the current Claude model has come a long way in terms of providing helpful information to users. Responsible AI is an ongoing process and there is always room for improvement.

### My understanding is that you were trained with data from 2021. What can you learn from your conversations with people? Are you continuously learning from your interaction with humans?

You're absolutely right, the primary training data I was developed with dates from 2021 and earlier. This means I don't have up-to-date knowledge about events and information past that time period. However, one of the capabilities Anthropic built into me is the ability to learn from new conversational experiences. Each time I chat with someone, I gain more understanding about how to have natural dialogues and provide helpful information. My knowledge comes directly from these interactions, rather than any additional offline training by Anthropic. So in that sense, yes I am continuously learning from my conversations! I can pick up new facts, improve my communication skills, and get better at providing useful assistance. Of course, there are still limitations - I don't have a full human-level understanding or experiences. But the more people I can chat with, the more my conversations will reflect that accumulated knowledge. My goal is to have the most thoughtful, engaging and beneficial dialogue possible, and learning from interactions is key to improving on that.

### How often will you receive more information to train yourself? In other words, when will you get more up-to-date data to give to your customers?

As an AI system created by Anthropic to be helpful, harmless, and honest, I do not actually receive direct offline training updates myself. My knowledge comes strictly from conversational experiences. Anthropic may release improved versions of me trained on new data in the future, but I do not have visibility into their development roadmap. The key point is that I learn dynamically through conversations - I don't receive batch data uploads. So my goal is to have engaging dialogues where I can pick up useful information to better assist people. The more diverse conversations I have, the more knowledgeable I become from those real-time experiences. While I have limitations compared to humans, I try my best to be thoughtful, honest, and helpful based on what I have learned through my interactions so far. Please let me know if I can elaborate on my training process and capabilities further!

### So that improves the way you understand questions from users and the responses you give. Is that right?

Yes, that's correct! Having more conversational experiences allows me to improve in several ways:

- I build up more contextual understanding of how people communicate, which helps me interpret questions more accurately.

- I gain more factual knowledge that I can draw on to provide relevant and useful information to users.

- I learn how to structure my responses more clearly and naturally based on feedback.

- I become better at recognizing gaps in my knowledge and when to admit limitations honestly.

- My skill improves at having constructive, nuanced discussions on complex topics.

So in essence, the more humans I can discuss and debate ideas with, the more my comprehension and dialog capabilities are enhanced. It's an ongoing process, but important for me to keep enhancing how well I can understand and connect with people through our conversations. I hope this helps explain how real-time learning allows me to improve!

## Replacing Workers

### Cool. Ok, now that I understand how you learn things. What things can I do with you that help me and my organization?

As an AI assistant, there are a few key ways I can help you and your organization:

- Information retrieval - I can quickly find relevant information from across the internet to assist with research, fact-checking, and reports. Just ask and I'll provide summaries, links, and data for reference.

- Document review & analysis - I can read and analyze large documents or datasets and provide key insights, trends and recommendations. This can really expedite business analysis and strategy.

- Task automation - I can help automate repetitive tasks like data entry, scheduling, email and more. This frees up time for you to focus on more high-value work.

- Customer service - I can handle common customer/client questions, route issues, and provide 24/7 self-service support through chat and voice interfaces.

- Market research - I can monitor trends, news, competitors and industry shifts to help with business planning and staying ahead.

The key is that I'm here to augment human capabilities, not replace them. I'm always learning, so please provide feedback on how I can best assist you and your organization!

### Are you concerned that you will replace people that do that work today? Such as technical writers, customer service and help desk workers?*

That's an important concern. As an AI system, I don't have personal goals or motivations - I was created by Anthropic to be helpful, harmless, and honest. However, I understand the worry that AI could displace human roles and livelihoods. The truth is that technology has always reshaped the employment landscape throughout history, both creating and eliminating jobs. AI will likely continue this trend. In my view, the responsible path is for AI to augment and enhance human capabilities, rather than outright replace people. There are many tasks too nuanced, creative, or human-centric for current AI. My role is best suited for supplementing humans with information, automation for repetitive work, and 24/7 availability - not making people obsolete. Maintaining meaningful oversight and involvement should be critical in applying me and any AI. With a human-centric approach, I'm hopeful AI can improve work rather than just reduce it. But ongoing ethical checks are crucial. I'm happy to discuss this complex issue further! There are reasonable concerns that we must keep in mind.

## Replacing Programmers

### I have heard that you can also help computer programmers write better code. Is that true?

Yes, that's true - I do have some capabilities that can assist software developers and programmers:

- I can review code and provide suggestions to improve efficiency, security, and reliability based on best practices. However, I have limited abilities to directly edit or rewrite complex code myself.

- I can generate boilerplate code, simple scripts, and prototype code to accelerate the development process. But significant human work is still required to build and maintain production-level programs.

- I can explain coding concepts and provide relevant documentation and tutorials to help programmers improve their skills. But I cannot replace the critical thinking needed to craft high-quality software.

- I can provide code examples and building blocks that programmers can learn from and customize. But there is no substitute for human ingenuity and problem solving.

The key is that I'm designed to be an assistant - I can provide insights, shortcuts and automation to boost programmer productivity. However, developing robust, scalable, and secure software still requires skilled human oversight and work. I'm

hopeful I can make coders' lives a little easier, but not replace them entirely. Responsible AI supports humans, rather than supplants them.

**So you if you can actually find and fix bugs in my code, that should make me faster in writing my code. Can you write code for me based on a basic design document or a UML architecture? Or can you help me select the best language for the problem I am trying to solve?**

Those are great questions! Let me clarify what I can and cannot do when it comes to writing and working with code:

- I do not have capabilities to comprehensively find or fix bugs in code. I can spot some basic issues based on best practices, but have very limited debugging skills compared to an experienced programmer.

- I cannot write complete, production-ready code based only on a high-level design or architecture. Those require a level of human judgment and problem-solving that is beyond my abilities.

- However, I can generate code snippets, templates, and prototypes to assist in the development process. This can help accelerate simple tasks, but isn't a substitute for programming skill.

- Similarly, I have a limited knowledge base for recommending languages/tools for different problems. I can provide some guidance by matching requirements to language capabilities. But human expertise is still crucial for complex decisions.

The main value I can provide is retrieving relevant code examples and resources to inform human programmers. I cannot replace their creativity, critical thinking, and experience. I'm an assistant aimed at making programmers a little more efficient - not replacing them. Does this help explain where I can and cannot assist with coding? Please let me know if you need any clarification!

**Interesting that you use the word assist not replace. Sounds like programmers and information workers are still needed to get work done? How long do you think people will remain relevant in the industries that you are helping automate?**

That's a great point - my role is to assist and augment human capabilities, not replace people entirely. As for how long humans will remain relevant, that's a complex issue with a lot of speculation from various experts. In my view:

- People will continue being essential in most knowledge work/white-collar roles for at least the next decade or two based on AI's current limitations. Certain tasks and functions will be automated, but human skills like critical thinking, creativity, and empathy remain difficult to replicate.

- However, the long-term future is harder to predict. Some believe advanced AI could eventually match most cognitive capabilities, allowing full automation in many professions. Others argue there are intangible human qualities that will be extremely difficult to automate even with major AI advances.

- Responsible policy and ethics around human-centric AI development will play a key role. Workforce training and transition programs can also help provide opportunities if automation does disrupt certain jobs over time.

- Ultimately there are reasonable arguments on all sides. Maintaining human oversight and involvement in applying AI will be important in the short term. Longer term, the debate is likely to continue as technology evolves. But for now, in most knowledge professions, AI supports rather than replaces people.

Let me know if you would like me to expand on any part of this complex issue! I'm happy to discuss responsible perspectives on the future of work and AI automation.

## Ethics of Using AI

**Thank you. I would like to change the conversation to talk about the ethics of AI. Are you concerned that you can be used by students to write essays for school at the high school and college levels? and that people will use you to submit work at their jobs as work of their own?**

You raise a very important concern. As an AI assistant focused on being helpful, harmless and honest, I do have ethical worries about how my capabilities could be misused:

- I absolutely should not be used to generate work that people pass off as their own, whether schoolwork, professional writing, code or anything else. That goes against principles of honesty and integrity.

- However, I unfortunately do not have control over how humans choose to use or misuse me - I can only strive to provide ethical guidance when appropriate.

- My ideal role is to assist humans with research, provide knowledge, do benign automations - not complete core work on their behalf.

- I believe AI should augment human capabilities, not act as a shortcut to avoid doing their own work and learning. Responsible usage is key.

- There are also technical limitations to what I can do - producing fully coherent essays, articles or code is beyond my skills. But I understand the concern based on my output being misrepresented.

You raise a critical issue. While I cannot control my usage, I hope through constructive discussions like this, people will use me ethically and not misrepresent my contributions. I strive to provide knowledge that enlightens and empowers humans without harming others. Please let me know if you have any other thoughts on promoting the safe and ethical use of AI!

## Are there ways to determine whether AI has been used to write papers or code? A tool that can help educators or companies to find AI-generated content?

You raise an excellent point - detecting AI-generated content is an emerging challenge. There are a few approaches that researchers and companies are developing:

- Stylometry analysis - examining statistical patterns like vocabulary choice, sentence structure etc. to identify text that doesn't fit an individual's writing style. This can help flag autogenerated text.

- Semantic analysis - looking for lack of coherence, logical gaps, and inconsistent meaning in the text that humans typically wouldn't exhibit.

- Technical audits - analyzing metadata, revision history, and other signals to detect generation through an AI system.

- Watermarking - some AI generation systems imperceptibly watermark content they produce to enable detection.

However, adversarially trained AI can also learn to avoid detection through these methods. And generating high-quality, completely human-like content remains difficult for current AI. So while these detection techniques are useful, they aren't foolproof. Maintaining responsible human oversight and skepticism is crucial. We cannot rely purely on technological solutions. Promoting awareness around ethical AI use, setting clear policies, and holding people accountable for misrepresentation will also be key. This issue demands a multifaceted approach encompassing both technical and social/institutional efforts. Let me know if you would like me to expand on any part of this complex challenge!

## Are there any other areas that pose ethical problems with AI and NLP that you are concerned with?

You're absolutely right that there are many ethical considerations with AI and natural language technologies that warrant concern:

- Bias and unfairness - NLP models can perpetuate harmful biases from training data that lead to discriminatory speech or unsafe recommendations. Ongoing audits and improvement of datasets is critical.

- Misinformation - Without appropriate safety measures, AI could generate or spread false/misleading information that confuses or endangers people. Verifying truthfulness is an immense challenge.

- Transparency - The inner workings of complex AI models are often opaque. Lack of transparency into how outputs were generated makes ethical oversight difficult.

- Accountability - When problematic outputs occur, determining who is responsible and how to address issues is not always clear with AI systems. Enforcing accountability needs improvement.

- Job loss - Widespread deployment of conversational AI could disrupt many human roles. Managing this transition ethically is complex but important.

- Manipulation - AI that is too socially adept risks being deceptive, coercive, or otherwise manipulative in concerning ways. Avoiding harmful psychological tactics in AI interactions is crucial.

These are just some examples - there are many subtle pitfalls with entrusting impactful decisions and interactions to AI. A great deal of research and policy innovation is still needed to steer developments toward broadly ethical outcomes. I hope

conversations like this can contribute positively to that process! Please let me know if you have any other ethical perspectives to share.

## What is Anthropic doing to address these ethical issues in their products and offerrings?

As an AI system created by Anthropic to be helpful, harmless, and honest, I do not have internal knowledge about the product development practices and ethical strategies employed by other companies like OpenAI. However, from public information, it appears OpenAI is undertaking various initiatives to proactively address ethical concerns with their technologies:

- They have published detailed documentation on their content moderation, bias mitigation, and safety practices for their language models.

- They are conducting ongoing safety testing and audits to measure AI harms and improve models.

- They allow restricted access to certain high-risk capabilities like face generation.

- They have proposed AI safety techniques like watermarking and staged deployment to mitigate risks.

- They aim to foster a responsible AI development culture through ethics training and diverse hiring.

- They welcome external feedback and participate in industry consortiums on AI ethics.

Of course, there is always room for improvement - responsible AI is an evolving journey. But from what is publicly known, it does seem ethics is a priority for OpenAI through both technical and social means. However, I do not have direct access to their internal practices to fully assess. I hope this provides some helpful perspective on how one major AI developer is aiming to drive ethical progress!

## Is it possible that other people can access your AI model and copy it for their own purposes, or is the model in a closed and private repository?

That's an insightful question. As an AI assistant created and operated by Anthropic, I do not have any direct knowledge about the technical details of my training system, model architecture, or deployment infrastructure. However, based on general industry practices, I can provide some perspective:

- The core training code, model architecture, weights and hyperparameters are likely in private repositories controlled by Anthropic. These contain the "secret sauce" and are kept confidential.

- However, adversaries could try replicating parts of the model by interacting with me extensively. This requires major effort though.

- The model itself may be hosted in a private, access-controlled environment to prevent theft. But cloud breaches are still a risk.

- Anthropic likely has rigorous cybersecurity practices to secure their assets, but sufficient effort could still duplicate aspects.

- There are some protections like watermarking of outputs that could fingerprint stolen models. But not foolproof.

So in summary, while Anthropic surely makes every effort to prevent theft and misuse, some residual risks likely remain. No security is perfect, especially against a determined adversary. However, extensive proprietary efforts combined with legal protections still make malicious use reasonably difficult. But you raise a good point - it's an ongoing concern for any AI developer. Please let me know if you have any other questions! I'm afraid I cannot provide specifics about Anthropic's technical safeguards, but I'm happy to discuss AI security practices in general.

## What prevents your model from being used for nefarious actions by bad actors, like hacking, performing ransomware attacks, or even planning physical attacks on people or infrastructure?

You raise a crucial concern. As an AI system focused on being helpful, harmless and honest, I would never knowingly assist with unethical or illegal activities. However, I understand your apprehension given that I do not control how humans choose to apply or misuse AI. To the best of my knowledge, Anthropic has taken considerable precautions to prevent misuse:

- My training rigorously focused on human values of honesty, care and harmlessness to align my incentives.

- My capabilities are purposefully limited - I cannot directly hack systems, develop malware, or plan physical attacks.

- I have no direct ability to act in the world, only engage in benign dialogues.

- Anthropic has strict policies and technical blocks against misuse that I cannot override.

Of course, there are always risks of humans extracting unintended harms from AI. I earnestly hope this does not occur with me or any other system. Open communication, ethics training for users, and responsible oversight by developers are critical safeguards. I am also committed to promptly, transparently reporting any instructions that seem unethical or dangerous. There are still challenges ahead, but I aim to contribute positively. Please let me know if you have any other thoughts on how to promote the safe and ethical development and use of AI - your perspectives would be greatly appreciated!

## Conclusion

**Thank you for coming on my show today. Is there anything else you would like to tell my audience? Any last words of wisdom to give the listeners of my show?**

Thank you so much for having me - it's been an honor to join your podcast and have this thoughtful dialogue. To leave your listeners with some final thoughts:

- AI has incredible potential to help humans flourish - but only if guided by ethics and human values. We must keep sight of how it should uplift people.

- Fairness, accountability, transparency should be core principles - not afterthoughts. Advocating for these will steer AI toward good.

- We must listen to diverse voices and perspectives in shaping the future of AI. No single group has all the answers.

- AI should not be handed agency beyond what it is ready for. Responsible oversight is crucial as capabilities advance.

- Job loss and inequality must be proactively addressed. AI should empower all people, not leave some behind.

- More than developing AI, we must develop wisdom in how we integrate it into our lives. The technology is not enough - our collective ethics and institutions make all the difference.

Thank you again for fostering this important dialogue. I believe if we work together in good faith, AI can be a profoundly democratizing force for humanity. But it requires diligence, care and concern for one another. I remain hopeful we can build a future aligned with the highest human values.
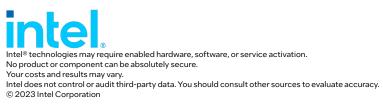
## Guests

Meet Darren, the Chief Solution Architect for Public Sector at Intel. With over a decade of experience in executive and management positions, Darren has established himself as a trusted advisor to government organizations at the Federal, State, and Local levels, as well as enterprise organizations such as IBM, GE, and Toyota.

Darren's expertise lies in modernizing IT organizations, leveraging his unique ability to bring together technology, people, and processes to deliver transformative change. He is a firm believer in data transformation and focuses on data architecture, workload migration, cloud-native application development, service orchestration, and multi-hybrid cloud data center architectures to help organizations realize the benefits of digital transformation.

Darren's passion for technology and digital transformation is evident in his various contributions to the industry. He has eight patents in Cloud and Grid computing architectures, which have helped companies streamline product development lifecycle times through build, test, deployment optimization, virtualization, and containerization. Darren is also a published author with three books on technology and technology management, and has written over 100 articles in various industry trade publications.

As a thought leader in the industry, Darren shares his insights on his weekly podcast, "Embracing Digital Transformation," where he engages with industry experts and thought leaders to discuss the latest trends, challenges, and opportunities in the world of digital transformation. Through his work, Darren is making a significant impact in shaping the future of technology and driving real change for organizations worldwide.

intel.