

Epilogue and Final Thoughts



You've made it up to this point – congratulations! I hope you found the book interesting and that it helped you toward your goals.

While writing this book, a number of changes and new developments have occurred that are definitely worth mentioning. We are indeed seeing a development to *Moore's Law* in terms of the increasing complexity and accuracy of Generative AI models.

So, in this final chapter, we will briefly recap what we have learned throughout this book, as well as unveiling the most recent developments and what to expect in the near future.

More specifically, we will cover the following topics:

- Overview of what we have learned so far

- How LLMs are entering the industries

- Latest developments in and concerns about the field of Generative AI

- What to expect in the near future

By the end of this chapter, you will have a broader picture of the state of the art developments within the domain of Generative AI, how it is impacting industries, and what to expect in terms of new developments and social concerns.

Recap of what we have learned so far

We started this book with an introduction to the concept of Generative AI and its various applications. We saw how Generative AI is about not only text but also images, video, and music.

In [Chapter 2](#), we then moved on to look at the company that brought Generative AI to its greatest popularity: OpenAI. Founded in 2015, OpenAI mainly focuses its research on a particular type of generative model, **Generative Pre-trained Transformers (GPT)**. Then, in November 2022, OpenAI released ChatGPT, a free web app of a conversational assistant powered by GPT models. It gained immense popularity, with it reaching 1 million users in just five days!

ChatGPT has been a game-changer. Its impact on daily productivity, as well as in various industry domains, is huge. Before dwelling on the ways ChatGPT could impact those areas, in [Chapter 3](#), we learned how to set up and start using a ChatGPT account. We also saw, in [Chapter 4](#), how to properly design the most important element when using generative models such as ChatGPT: the prompt. Prompts are the user's input, nothing more than instructions in natural languages. Designing prompts is a pivotal step to getting the maximum value from your generative models, to the point where **prompt engineering** has become a new domain of study.

Once we got familiar with ChatGPT and prompt design, we moved on to [Chapter 5](#), where we finally got concrete examples of how ChatGPT can boost your daily productivity and become your daily assistant. From email generation to improving your writing skills, we saw how many activities can be improved thanks to the generative power of ChatGPT.

But we didn't stop there. With *Chapters 6, 7, and 8*, we saw how ChatGPT can boost not only daily productivity but also domain-specific activities – for developers, from code generation and optimization to interpreting machine learning models; in the case of marketers, from new product development to improving **Search Engine Optimization (SEO)**; and for researchers, from experiment design to the generation of a presentation based on a study.

Starting with [Chapter 9](#), we shifted the conversation to the enterprise-level, which discussed how OpenAI models have become consumable directly via Azure so that enterprises can maintain reliability and security.

Finally, in [Chapter 10](#), we saw concrete examples of enterprise use cases with Azure OpenAI models. Each example came with a business scenario as well as an end-to-end implementation with Python, using Streamlit as the frontend.

This journey was meant to provide you with greater clarity about what we are talking about when we refer to popular buzzwords such as ChatGPT, OpenAI, and LLMs.

However, in the next section, we will see how the incredibly fast AI developments in recent months are bringing brand-new technologies on top of what we have learned so far.

This is just the beginning

Throughout this book, we saw how Generative AI and, more specifically, GPT models are revolutionizing the way both citizens and large enterprises are working.

Nevertheless, we have embarked on a journey where ChatGPT and GPT models represent only the first steps toward an era of unprecedented technological advancements. As we have seen throughout the book, these models have already demonstrated exceptional capabilities in language understanding and generation. However, the true potential of Generative AI has yet to be fully realized.

A glimpse of what we might expect has already been unveiled by the first releases of **Multimodal Large Language Models (MLLMs)** and the introduction of the **Copilot** system by Microsoft.

The advent of multimodal large language models

So far, we've mainly focused on **Large Language Models (LLMs)**, as they are the architecture behind the GPT-x family and ChatGPT. These models are trained on massive amounts of text data, such as books, articles, and websites, and use a neural network to learn the patterns and structure of human language.

As we saw in [*Chapter 2*](#), if we want to combine further Generative AI capabilities with LLMs, such as image understanding and generation, we need the support of additional models, such as DALL-E. This holds true until the introduction of MLLMs.

MLLMs are AI systems that combine NLP with computer vision to understand and generate both textual and visual content. These models are trained on massive amounts of data, such as images and text, and are capable of generating human-like responses to queries that include both text and visual inputs.

In recent months, there have been great developments in the field of MLLMs, and in the next sections, we are going to focus on two main models: Kosmos-1 and GPT-4.

Kosmos-1

In their paper *Language Is Not All You Need: Aligning Perception with Language Models*, Microsoft’s researchers Shaohan Huang et al. introduced **Kosmos-1**, an MLLM that can respond to both language and visual cues. This enables it to perform tasks such as image captioning and visual question answering.

While LLMs such as OpenAI’s ChatGPT have gained popularity, they struggle with multimodal inputs such as images and audio. Microsoft’s research paper highlights the need for multimodal perception and real-world grounding to advance toward **Artificial General Intelligence (AGI)**.

Kosmos-1 can perceive various modalities, follow instructions through zero-shot learning, and learn from the provided context using few-shot learning. Demonstrations of the model show its potential to automate tasks in various situations involving visual prompts.

The following figure provides an example of how it functions:

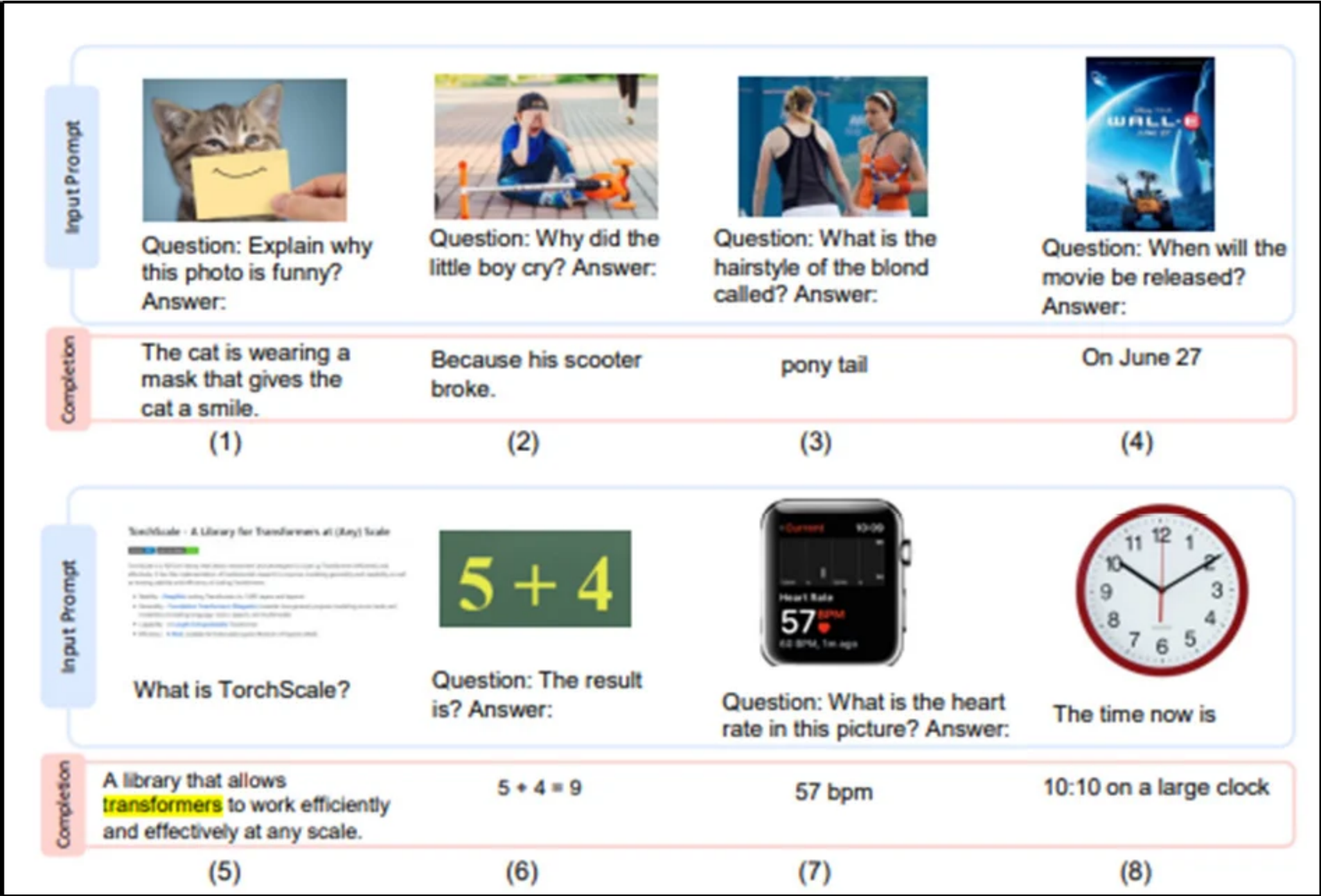


Figure 11.1 – Example of multimodal inputs with Kosmos-1. Original picture from <https://arxiv.org/pdf/2302.14045.pdf>

Tests on the zero-shot Raven IQ test revealed a performance gap compared to adult levels but showed promise for MLLMs to perceive abstract conceptual patterns by aligning perception with language models.

Definition

The Raven IQ test, also known as Raven’s Progressive Matrices, is a nonverbal standardized test designed to measure a person’s abstract reasoning and fluid intelligence. Developed by John C. Raven in 1936, the test consists of multiple-choice questions with visual patterns in the form of matrices. The participant’s task is to identify the missing piece that completes the pattern.

The following figure illustrates an example of the Raven IQ test solved by Kosmos-1:

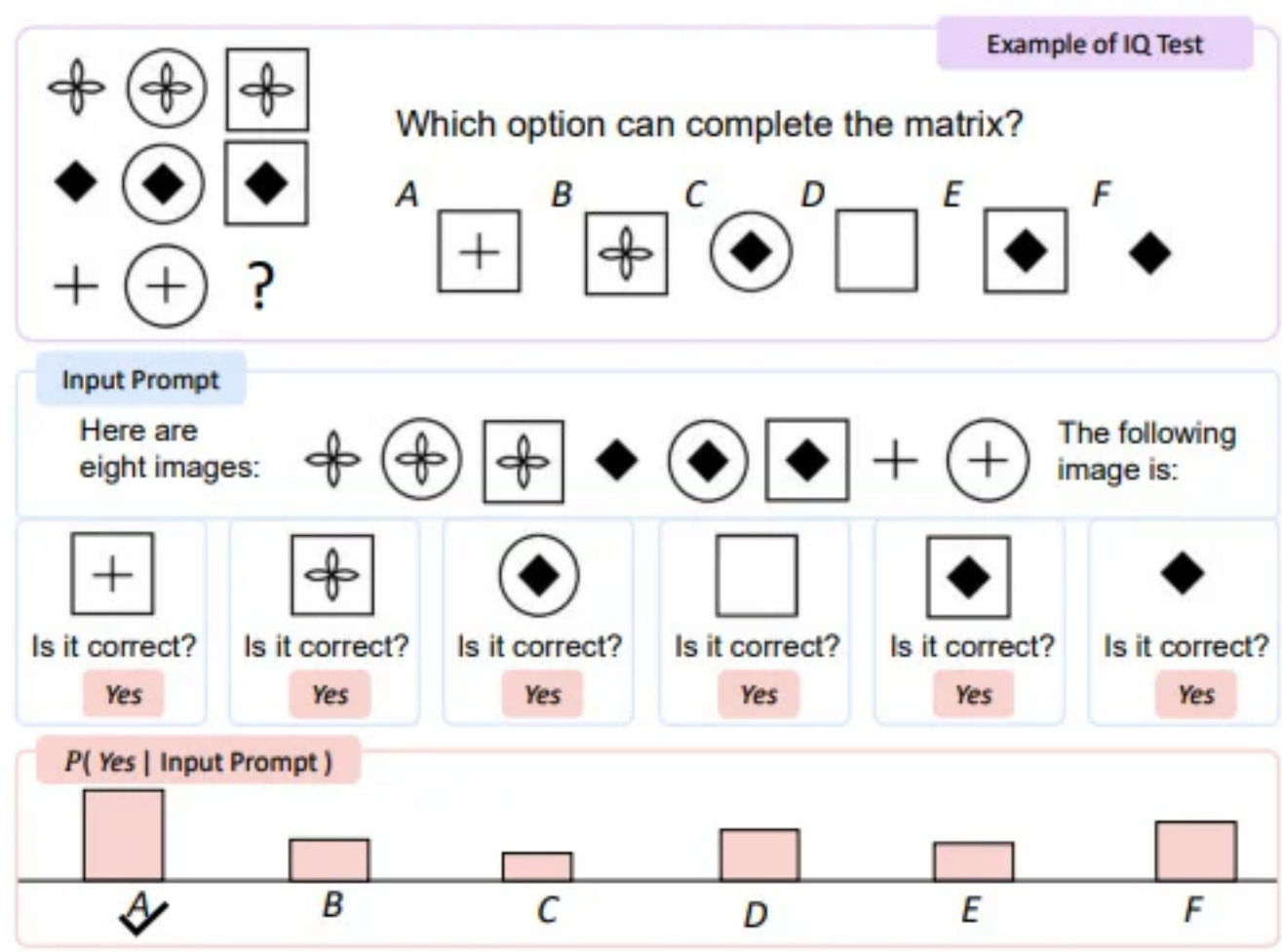


Figure 11.2 – Example of a Raven IQ test solved by Kosmos-1. Original picture from <https://arxiv.org/pdf/2302.14045.pdf>

For now, Kosmos-1 is only able to analyze images and text. However, in the conclusion of the research paper, Microsoft researchers announced its intention to further develop the model to integrate a speech capability.

GPT-4

On March 14, 2023, OpenAI announced the new version of the GPT series: **GPT-4**. The technical description of this brand-new model is described by OpenAI’s researchers in the paper *GPT-4 Technical Report* (<https://arxiv.org/pdf/2303.08774.pdf>).

According to this paper, it is evident that GPT-4 exhibits a higher level of general intelligence compared to previous AI models. GPT-4 demonstrates near-human performance across a wide range of tasks, including mathematics, coding, vision,

medicine, law, and psychology, without requiring special prompting.

More specifically, there are four main areas where GPT-4 outperforms its previous version (GPT 3.5):

Multimodality: GPT-4 is a great example of an MLLM, since it is able to understand and generate not only natural language but also images:

User What is funny about this image? Describe it panel by panel.



Source: [hmmm \(Reddit\)](#)

GPT-4 The image shows a package for a "Lightning Cable" adapter with three panels.

Panel 1: A smartphone with a VGA connector (a large, blue, 15-pin connector typically used for computer monitors) plugged into its charging port.

Panel 2: The package for the "Lightning Cable" adapter with a picture of a VGA connector on it.

Panel 3: A close-up of the VGA connector with a small Lightning connector (used for charging iPhones and other Apple devices) at the end.

The humor in this image comes from the absurdity of plugging a large, outdated VGA connector into a small, modern smartphone charging port.

Figure 11.3 – Example of GPT-4 understanding and explaining an image

With GPT-4, we are basically able to process and understand a whole document, made of both text and images.

Accuracy: GPT-4 has proven to be more reliable, creative, and receptive than GPT-3.5, especially in complex tasks. To understand this difference, the two models were tested on several exams, originally designed for humans, and GPT-4 (both with and without vision capabilities) consistently outperformed GPT-3.5.

Tests have also been done in the field of computer vision. In this case, OpenAI evaluated GPT-4 on traditional benchmarks designed for computer vision models and, also in this case, GPT-4 considerably outperformed most **State-of-the-Art (SOTA)** models.

Alignment: As we saw in [Chapter 5](#), OpenAI published an AI Alignment (<https://openai.com/alignment/>) declaration whose research aims to make AGI aligned with human values and follow human intent.

Significant efforts have been dedicated to improving GPT-4's safety and alignment with user intent. As a result, GPT-4 has become considerably safer than its previous versions, with an 82% reduction (<https://openai.com/product/gpt-4>) in the likelihood of generating responses to prohibited content requests compared to its predecessor, GPT-3.5. The reason behind this improvement is that OpenAI's GPT-4 incorporates new research advancements that add an extra layer of safety. Informed by human input, this safety feature is integrated directly into the GPT-4 model, making it more adept at managing potentially harmful inputs. Consequently, the chances of the model generating unsafe responses are significantly reduced.

Additionally, OpenAI's internal evaluations indicate that GPT-4 is 40% more likely than the previous version to generate accurate and fact-based responses. These enhancements showcase the continuous progress being made in refining AI language models, ensuring their safety and increasing their reliability for users.

Overall usability: Last but not least, GPT-4 addresses one of the main limitations of its predecessor. Up to GPT-3.5, we had a maximum number of tokens to take into account, which was 4,096. With GPT-4, the maximum number of tokens has increased greatly, to around 32,000, which makes it more suitable for complex and longer tasks, especially if they involve step-by-step reasoning.

The previous figure shows just some examples of the full capabilities of GPT-4, yet these are very impressive in themselves. Once more, it is significant to note that with MLLMs, we are entering a new phase of Generative AI where a single foundation model will be able to fully process and understand a whole document and then generate new materials based on it.

Note

On March 21, 2023, Microsoft announced that GPT-4 is available within Azure OpenAI Service (<https://azure.microsoft.com/en-us/blog/introducing-gpt4-in-azure-openai-service/>). This means that this powerful model can already be consumed for enterprise-scale projects, or tested directly within the Azure OpenAI Playground.

GPT-4 is an extremely powerful model, and it is already the engine behind many AI-infused applications. One of these applications is the new version of ChatGPT, called ChatGPT Plus. But there is another one that, in my opinion, is far more interesting, since it is revolutionizing search engine tools: Microsoft Bing. We will dive deeper into that in the next section.

Microsoft Bing and the Copilot system

In recent years, Microsoft has emerged as a leading player in the field of AI, investing heavily in research and development to drive innovation and unlock new possibilities. As part of its commitment to advancing AI technology, Microsoft has forged a strategic partnership with OpenAI, as we saw in [Chapter 9](#).

This collaboration between Microsoft and OpenAI aims to accelerate progress in AI, combining their respective expertise in cloud computing, software, and cutting-edge AI models. Together, they seek to create AI systems that not only have remarkable capabilities but also adhere to principles of transparency, fairness, and ethical responsibility.

Since the announcement of the general availability of Azure OpenAI Service in January 2023, Microsoft has released a series of new developments within the Generative AI domain, leveraging the power of LLMs, including GPT-4.

In the next sections, we are going to focus on two of the most promising developments: the new Bing and the Copilot system.

The new Bing

Microsoft Bing is a web search engine owned and operated by Microsoft. The service has its origins in Microsoft's previous search engines: MSN Search, Windows Live Search, and later Live Search.

In February 2023, Microsoft announced (<https://blogs.microsoft.com/blog/2023/02/07/reinventing-search-with-a-new-ai-powered-microsoft-bing-and-edge-your-copilot-for-the-web/>) a new version of Bing powered by GPT models. Furthermore, with the launch of GPT-4, on March 14, 2023, Microsoft confirmed (https://blogs.bing.com/search/march_2023/Confirmed-the-new-Bing-runs-on-OpenAI%E2%80%99s-GPT-4) that the new Bing is actually running on OpenAI’s GPT-4.

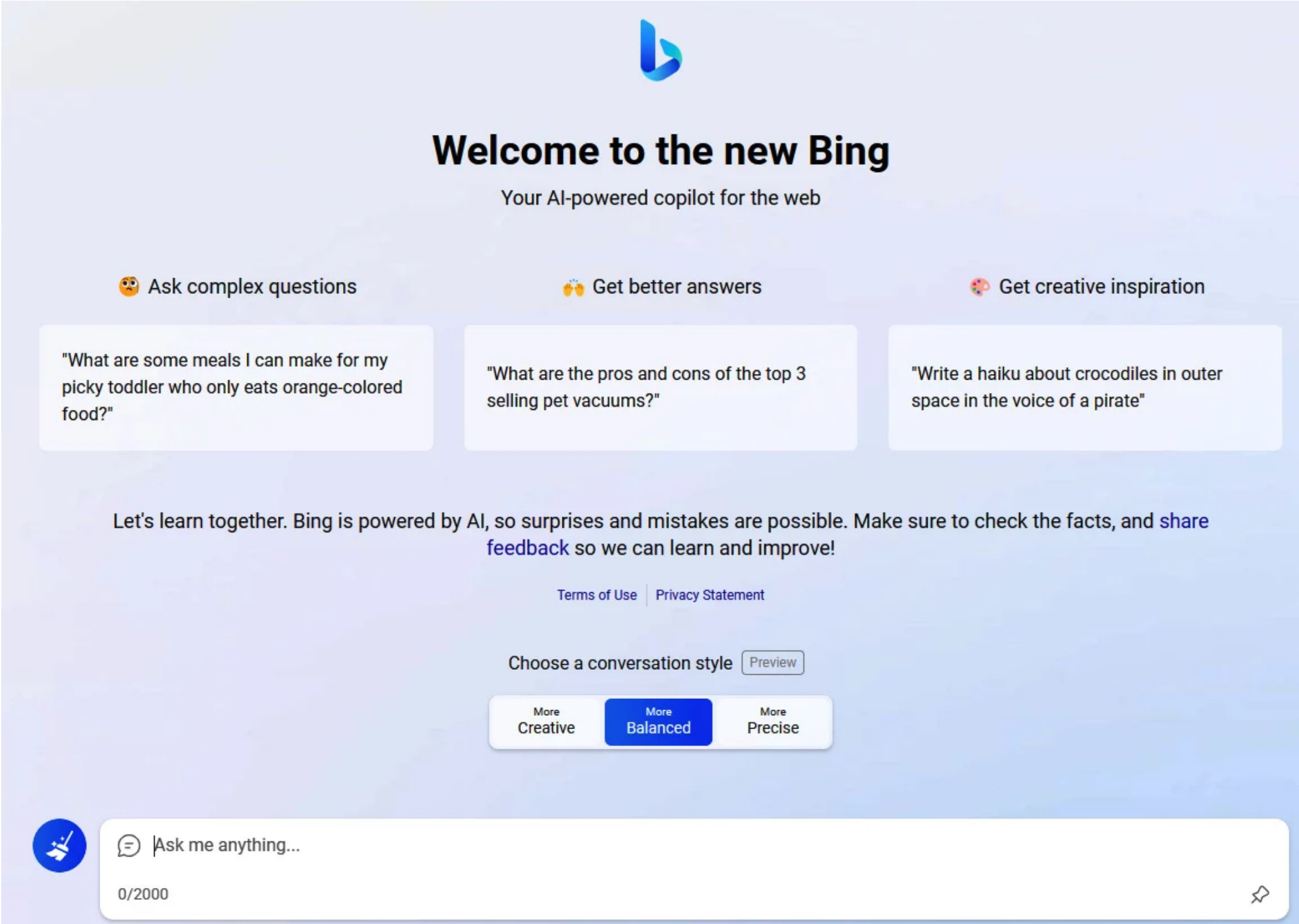


Figure 11.5 – The new Bing

In its new version, Bing has become sort of like a version of ChatGPT that is able to navigate the web (hence bypassing the problem of ChatGPT’s limited knowledge cut-off of 2021) and also provide references beyond the expected response. Refer to the following screenshot:

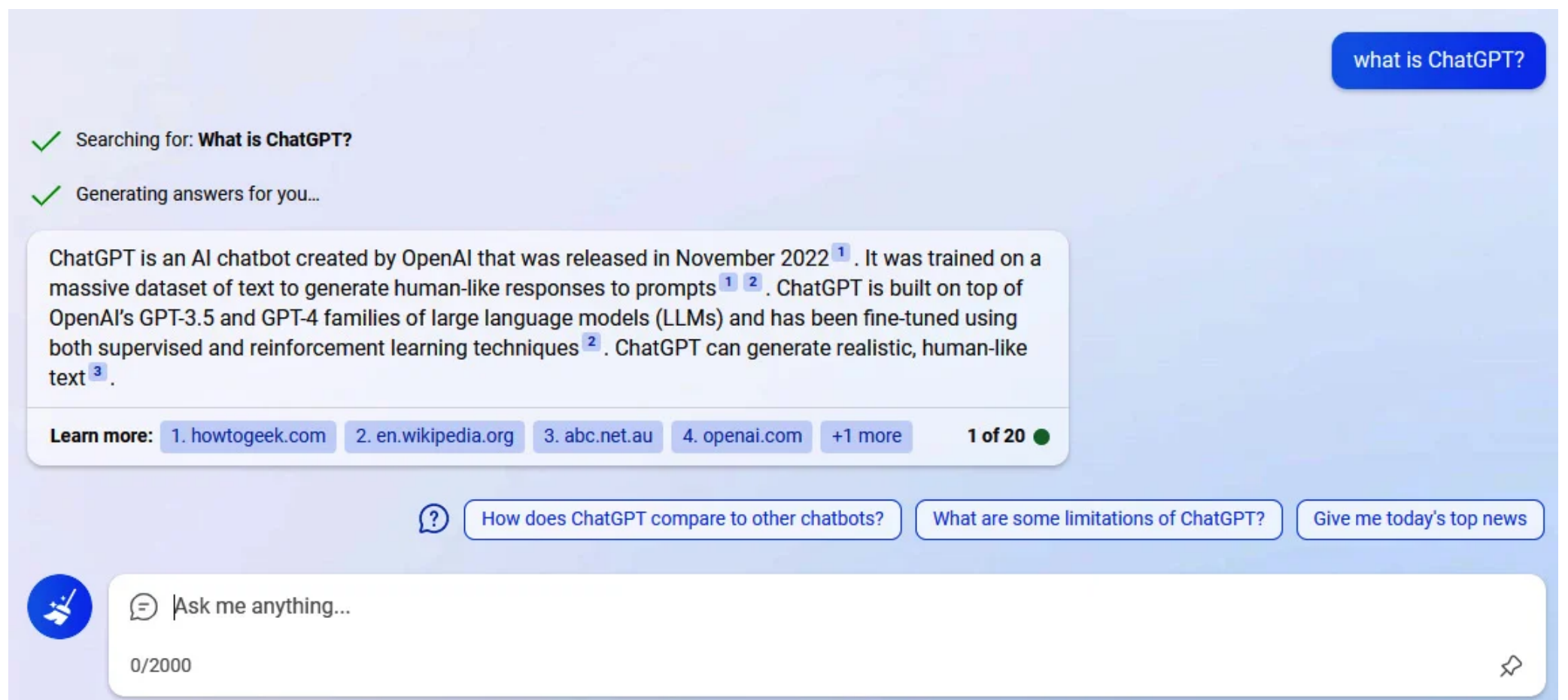


Figure 11.6 – The new Bing providing an answer with references

The new Bing can also assist in generating content, in the same fashion as ChatGPT. For example, I could ask Bing for support while writing LinkedIn posts:

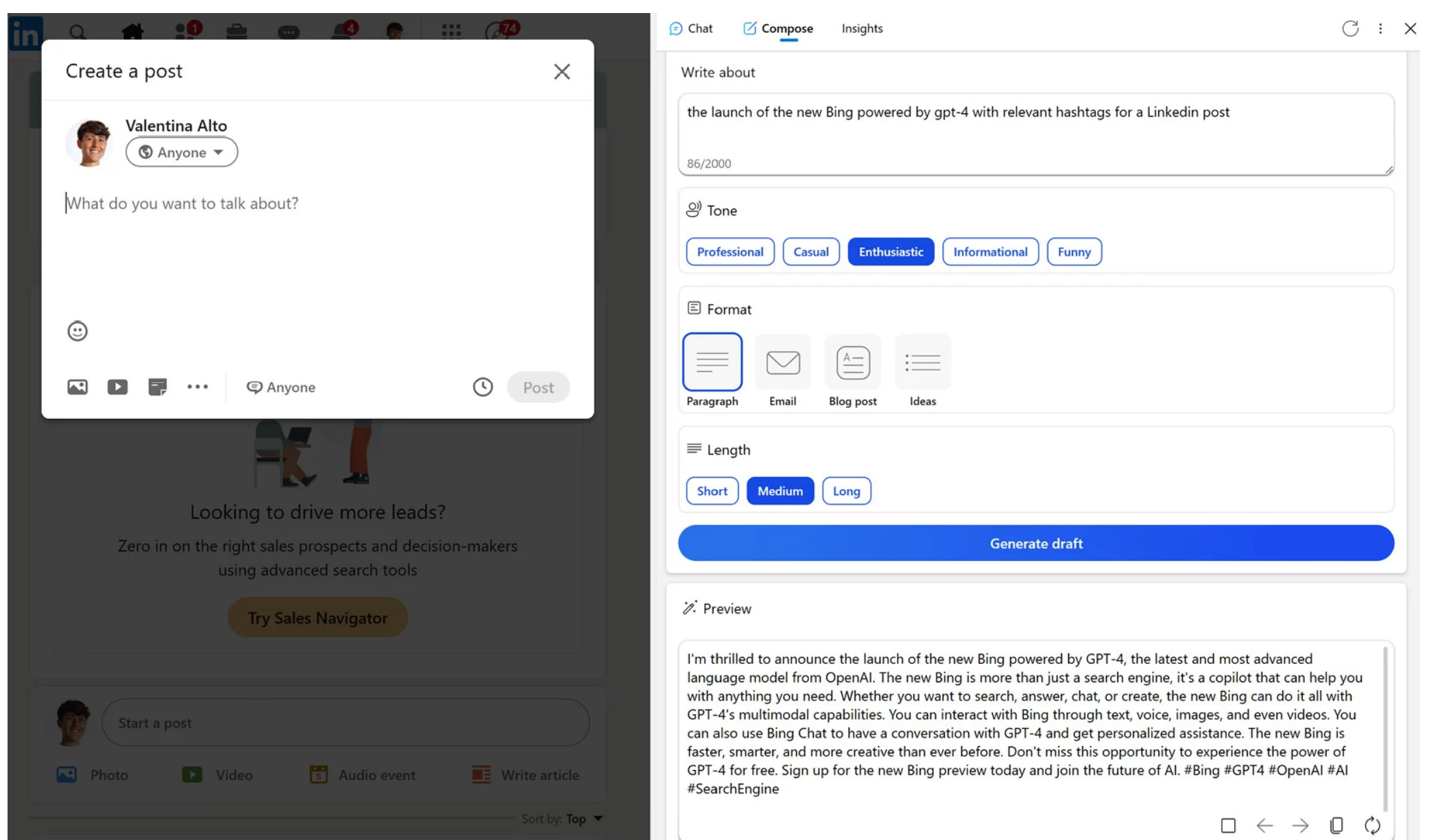


Figure 11.7 – Example of Bing used as a LinkedIn post assistant

With this last feature, the new Bing can be seen as a copilot for the web, speeding up research as well as the retrieval or generation of materials.

The concept of a copilot is pivotal in this new landscape of foundation models and LLMs, since it is the most likely way these new AI systems will enter organizations.

Definition

As the name suggests, a copilot acts as an expert assistant to a user with the goal of supporting it in solving complex tasks. Copilots have user-friendly, natural-language interfaces and are powered by foundation models. Plus, they are scoped to a perimeter defined by the user. For example, a Copilot within application A will be working only with application A's data.

In the next section, we will see how Microsoft has extended this concept to its whole suite of applications.

Microsoft 365 Copilot

Introduced by Microsoft in March 2023

(<https://blogs.microsoft.com/blog/2023/03/16/introducing-microsoft-365-copilot-your-copilot-for-work/>), the copilot system is a sophisticated processing and orchestration engine built on top of the following three technologies:

- Microsoft 365 apps, such as Excel, Word, and PowerPoint

- Microsoft Graph, a platform that provides access to various data and intelligence across Microsoft 365 services

- LLMs, such as GPT-4

Based on the copilot system, Microsoft 365 Copilot is a revolutionary AI-powered assistant designed to enhance productivity and unleash creativity in the workplace. By utilizing LLMs and integrating with Microsoft 365 apps and data, Copilot transforms natural language into a powerful tool for completing tasks and connecting with your work.

Microsoft 365 Copilot is seamlessly integrated into popular apps such as Word, Excel, PowerPoint, Outlook, and Teams. Using natural-language prompts, Copilot can perform tasks such as generating status updates based on meetings, emails, and chats. The user always maintains control over its core activities, allowing for increased creativity, analysis, expression, and collaboration across various Microsoft 365 applications.

In this section, we saw how, in the last few months, there have been further developments to OpenAI models and, more generally, further AI models in the domain of LLMs. We also saw how companies such as Microsoft are introducing a new way of integrating these LLMs into applications, with the brand-new concept of Copilot.

In the next section, we will dive deeper into how other companies are embracing OpenAI's models in their processes and digital transformation paths, covering different use cases and approaches.

The impact of generative technologies on industries – a disruptive force

As ChatGPT and Generative AI models continue to evolve, their capabilities are undoubtedly transforming industries in once unimaginable ways. On the one hand, the integration of these technologies has the potential to significantly enhance productivity and drive economic growth. By automating time-consuming tasks, Generative AI can free up human workers to focus on more creative, strategic, and value-added activities. Moreover, AI-driven tools can augment human capabilities, enabling professionals to make more informed decisions and generate novel ideas.

There are already some examples of enterprises that are embedding Generative AI within their core businesses:

Morgan Stanley, a leading wealth management firm, has been working to improve access to its extensive content library by utilizing OpenAI's GPT-4 technology (<https://openai.com/customer-stories/morgan-stanley>). The company has been exploring GPT's embedding and retrieval capabilities to create an internal chatbot that can efficiently search through its vast wealth management resources. This initiative, led by Jeff McMillan, head of analytics, data, and innovation, aims to make the firm's collective knowledge easily accessible and actionable.

Duolingo, a language learning app with a game style that boasts over 50 million users, has always relied on AI as part of its strategy. One of the features in which AI has been involved is in **Explain My Answer**. This feature allows the user to go deeper into the explicit grammar rules behind Duolingo's response (for example, if the user's answer is incorrect, the user could open a chat and ask for elaboration on why it is incorrect). So far, Duolingo has tried to implement this feature with both pre-written grammar tips and AI with GPT-3. However, it was only with the advent of GPT-4 that the accuracy of responses and learners' engagement spiked, thanks to its advanced capabilities of language understanding in terms of the grammar rules of different idioms.

Stripe, a fintech payment services provider, did something pretty visionary. At the beginning of 2023, it asked 100 employees to stop doing their day-to-day work activities and start envisioning how LLMs could enrich and empower the functionalities of a payment service. Not only did they identify many potential applications (the output was a list of 50 applications!) but they also started prototyping 15 concrete use cases. Among those, one of particular interest is using GPT-4 for fraud detection. Stripe's Discord community faces infiltration by malicious actors. GPT-4 aids in identifying potential fraudsters by analyzing post syntax within Stripe's Discord community and detecting coordinated malicious activities, ensuring platform security.

On the other hand, the rapid adoption of ChatGPT and Generative AI models raises concerns about job displacement, data privacy, and the potential misuse of technology. As automation reshapes the labor market, industries must navigate the challenges of workforce transition, ensuring that employees are equipped with the skills necessary to thrive in an increasingly AI-driven world.

The disruptive impact of Generative AI on industries is undeniable, offering both opportunities and challenges. By fostering collaboration between humans and AI, promoting ethical development and deployment, and prioritizing life-long learning and reskilling, we can shape a future where Generative AI serves as a powerful catalyst for positive change across industries.

So far, we have mainly focused on what can be achieved with Generative AI. In the next sections, we will unveil the rising concerns about the ethical implications of this new disruptive technology.

Unveiling concerns about Generative AI

As Uncle Ben said to the young Peter Parker, *“With great power comes great responsibility.”*

As we reach the end of our enlightening journey through exploring the world of ChatGPT and GPT models, it is imperative to address the concerns that have emerged about these AI technologies. While the advancements in Generative AI have been nothing short of ground-breaking, they have also raised vital questions about privacy, security, ethical implications, and potential misuse.

In fact, many announcements and statements have been released in recent months about those topics from companies and institutions as well as individual speakers, including concerns, calls for stopping further developments, and proper regulations.

In the next section, I want to share some of the latest news and developments, hoping it will also stimulate reflections and guesses on what to expect in the near future.

Elon Musk calls for stopping development

A recent open letter, signed by Elon Musk and over 1,000 other technology professionals, has called for a temporary stop to the development of AI systems more sophisticated than GPT-4. The signatories, including Steve Wozniak, Emad Mostaque, and Yuval Noah Harari, express their concerns about the significant risks these AI systems pose to society and humanity.

The letter requests that leading AI laboratories pause the training process for these advanced systems for a minimum of six months, ensuring that the halt is both public and verifiable. It highlights worries related to the potential for AI-driven propaganda, job automation, and a loss of control over our civilization.

This appeal comes on the heels of OpenAI's launch of GPT-4, an enhanced language model that powers the premium version of ChatGPT. According to OpenAI, GPT-4 is more capable of complex tasks and generates more refined results than previous iterations, with fewer shortcomings.

AI systems such as GPT-4 operate on vast amounts of data, which they utilize to respond to queries and execute tasks. ChatGPT, which debuted in November, has human-like abilities to compose emails, arrange travel plans, write code, and excel in exams such as the bar exam.

OpenAI has not yet commented on the letter, but the organization has acknowledged the importance of ensuring that AI technologies smarter than humans serve humanity's interests. OpenAI suggests that future systems may require independent evaluation before training and that advanced efforts should regulate the expansion of computational resources used for model development.

Several companies, including Google, Microsoft, Adobe, Snapchat, DuckDuckGo, and Grammarly, have unveiled services that harness Generative AI features. OpenAI's research points to the risks involved with these capabilities, such as the possibility of quoting untrustworthy sources or empowering malicious actors to deceive or exploit others.

AI specialists are increasingly alarmed by the trajectory of the industry and the potential absence of necessary precautions and comprehension of the consequences. The letter underlines that advanced AI could have a dramatic impact on life on Earth and requires careful planning and management. It notes that such planning is currently insufficient, as AI laboratories persist in creating and deploying ever more powerful AI systems that are challenging to understand, anticipate, or control.

If this open letter had no binding effect, another example is what the Italian “Garante della Privacy” declared, which we will focus on in the next section.

ChatGPT was banned in Italy by the Italian “Garante della Privacy”

Italy has become the first Western nation to prohibit ChatGPT due to privacy concerns.

The Italian data protection authority **Garante della Privacy** announced it will immediately impose a ban on OpenAI and launch an investigation:

ChatGPT disabled for users in Italy

Dear ChatGPT user,

We regret to inform you that we have disabled ChatGPT for users in Italy at the request of the Italian Garante.

We are issuing refunds to all users in Italy who purchased a ChatGPT Plus subscription in March. We are also temporarily pausing subscription renewals in Italy so that users won't be charged while ChatGPT is suspended.

We are committed to protecting people's privacy and we believe we offer ChatGPT in compliance with GDPR and other privacy laws. We will engage with the Garante with the goal of restoring your access as soon as possible.

Many of you have told us that you find ChatGPT helpful for everyday tasks, and we look forward to making it available again soon.

If you have any questions or concerns regarding ChatGPT or the refund process, we have prepared a list of [Frequently Asked Questions](#) to address them.

—The OpenAI Support Team

Figure 11.10 – Message from OpenAI when accessing ChatGPT in Italy

The Italian regulator will not only block ChatGPT but also investigate its compliance with the **General Data Protection Regulation (GDPR)**, which governs the use, processing, and storage of personal data.

Following a data breach involving user conversations and payment information, the authority stated on March 20, 2023 that there is no legal basis for the *mass collection and storage of personal data* to train the platform's underlying algorithms.

The regulator also expressed concerns about the inability to verify users' ages, potentially exposing minors to inappropriate responses.

The Italian data protection authority has given OpenAI 20 days to address its concerns or face a penalty of €20 million (\$21.7 million) or up to 4% of their annual revenue.

OpenAI deactivated ChatGPT for Italian users on April 1, 2023, at the request of the Italian data protection regulator, the Garante, and stated its commitment to privacy protection and GDPR compliance.

The company stated that it looks forward to working closely with the Garante and hopes to make ChatGPT available in Italy again soon.

The previously mentioned concerns and interventions are just scraping the surface of a broader topic, which is the concept of Responsible AI, which will be the subject of the next section.

Ethical implications of Generative AI and why we need Responsible AI

The previous section highlighted how, alongside the widespread knowledge and adoption of Generative AI technologies, a general concern is rising.

The rapid advancement of AI technologies brings forth a plethora of ethical considerations and challenges that must be carefully addressed to ensure their responsible and equitable deployment. Some of them are listed here:

Data privacy and security: As AI systems rely heavily on data for their learning and decision-making processes, ensuring data privacy and security becomes paramount. In [Chapter 9](#), we already saw how Microsoft addressed

the topic of data privacy with Azure OpenAI Service, in order to guarantee the **Service-Level Agreements (SLAs)** and security practices expected of the Azure cloud. However, this data privacy topic also affects the data that is used to train the model in the first instance: even though the knowledge base used by ChatGPT to generate responses is public, where is the threshold of the consent of involved users whose information is used to generate responses?

Bias and fairness: AI models often learn from historical data, which might inadvertently introduce biases. Addressing bias and fairness in AI systems involves the following:

- **Diverse datasets:** Ensuring that training data is diverse and representative of various demographics can help reduce biases in AI models

- **Algorithmic fairness:** Developing algorithms that prioritize fairness and do not discriminate against specific demographic groups is essential

- **Monitoring and auditing:** Regular monitoring and auditing of AI systems can help identify and rectify biases, ensuring that the outcomes are equitable

Transparency and accountability: As AI systems become more complex, understanding their decision-making processes can be challenging. This involves the following two important aspects:

- **Explainable AI:** Developing AI models that can provide clear explanations for their decisions can help users understand and trust the system.

- **Responsibility and liability:** Establishing clear lines of responsibility and liability for AI systems is crucial to hold developers, organizations, and users accountable for the consequences of AI-driven decisions.

The future of work: AI-driven automation has the potential to displace jobs in certain sectors, raising concerns about the future of work. Throughout this book, we have seen how ChatGPT and OpenAI models are able to boost productivity for individuals and enterprises. However, it is also likely that some repetitive tasks will be definitively replaced by AI, which will impact some workers. This is part of the change and development process, and it is better to embrace change rather than combat it.

Some actions in this direction could be reskilling and upskilling programs – governments, organizations, and educational institutions should invest in reskilling and upskilling programs to help workers adapt to the changing job market and acquire new skills required for emerging roles.

Most importantly, human-AI collaboration should be encouraged. Developing AI systems that complement and augment human capabilities can help create new job opportunities and foster collaborative work environments.

By addressing these ethical considerations and challenges, we can work in the right direction to ensure that AI technologies are developed and deployed responsibly, promoting a better and more equitable future for all.

Now, the next logical question might be: given the tremendous acceleration of AI technologies in recent months, what should we expect in the near future?

What to expect in the near future

The acceleration of AI research and developments in recent months has been incredible. From November 2022 up to the time of writing (April 2023), we have seen the following occur:

- The launch of ChatGPT (November 2022)

- The general availability of Azure OpenAI (January 2023)

- The general availability of the model API behind ChatGPT, GPT-3.5-turbo (February 2023)

- The general availability of MLLMs such as Kosmos-1 and GPT-4 (March 2023)

- Microsoft's announcement of the Copilot system (March 2023)

This incredible pace makes it hard to predict what will come next. As we have seen, this velocity has also raised concerns among institutions, companies, and public figures because of the lack of regulation for these new technologies. At the same time, companies and institutions will inexorably need to adapt to this new landscape in order to keep up with competitors.

If we think about the near future, we are talking about *tomorrow*. We have seen how some IT companies, such as Microsoft, are already integrating GPT models into their applications as a copilot system, while other companies, such as WolframAlpha, Expedia, and Instacart, have designed plugins that are integrated directly into ChatGPT.

The move toward infusing OpenAI models into applications is evident also by the variety of frameworks that have been developed with the purpose of facilitating the integration between LLMs and applications as well as managing prompts, conversations, memory, tokenization, and other typical steps required. Some examples of those frameworks are LangChain, Pinecone, and Semantic Kernel.

We mentioned in [*Chapter 2*](#) that the mission of OpenAI is to build broadly beneficial AGI, a type of AI that, in being “general,” is intended to have the ability to learn and perform a wide range of tasks, without the need for task-specific programming.

In other words, OpenAI is envisioning an AI machine that is able to do whatever a human can do.

If we had thought about this a year ago, it would have seemed futuristic. Today, at the vertiginous pace of development, is it so unbelievable that we will obtain AGI machines not so far in the future?

Summary

The rapid development of AI technologies such as OpenAI, ChatGPT, and Generative AI models is ushering in a new era of innovation and transformation. With the immense potential to revolutionize industries and reshape day-to-day life, these advancements are rewriting the rules of human-machine interaction.

As we stand on the brink of this AI-driven future, it is our collective responsibility to ensure that these technologies are used responsibly and ethically. By embracing the opportunities and addressing the challenges, we can foster a world where AI empowers humanity and elevates our potential to new heights.

GPT began two years ago – an era if we think about the pace of AI developments in recent months – yet it reflects the inevitable influence of AI on our lives and the challenges that lie ahead in adapting to this new reality.

References

https://python.langchain.com/en/latest/getting_started/getting_started.html

<https://learn.microsoft.com/en-us/semantic-kernel/whatissk>

<https://www.pinecone.io/>

[Previous Chapter](#)

[Next Chapter](#)