

12

Responsible AI

In Part 2 of this book, we covered multiple applications of **large language models (LLMs)**, gathering also a deeper understanding of how many factors could influence their behavior and outputs. In fact, LLMs open the doors to a new set of risks and biases to be taken into account while developing LLM-powered applications, in order to mitigate them with defensive attacks.

In this chapter, we are going to introduce the fundamentals of the discipline behind mitigating the potential harms of LLMs – and AI models in general – which is Responsible AI. We will then move on to the risks associated with LLMs and how to prevent or at least mitigate them using proper techniques. By the end of this chapter, you will have a deeper understanding of how to prevent LLMs from making your application potentially harmful.

We will cover the following key topics:

- What is Responsible AI and why do we need it?
- Responsible AI architecture
- Regulations surrounding Responsible AI

What is Responsible AI and why do we need it?

Responsible AI refers to the ethical and accountable development, deployment, and use of AI systems. It involves ensuring fairness, transparency, privacy, and avoiding biases in AI algorithms. Responsible AI also encompasses considerations for the social impact and consequences of AI technologies, promoting accountability and human-centric design.

Responsible AI plays a crucial role in steering decisions toward positive and fair results. This involves prioritizing people and their objectives in the design of systems while upholding enduring values such as fairness, reliability, and transparency.

Some ethical implications of Responsible AI are:

- **Bias:** AI systems can inherit biases present in their training data. These biases can lead to discriminatory outcomes, reinforcing existing inequalities.
- **Explainability:** Black-box models (such as LLMs) lack interpretability. Efforts are being made to create more interpretable models to enhance trust and accountability.
- **Data protection:** Collecting, storing, and processing data responsibly is essential. Consent, anonymization, and data minimization principles should guide AI development.
- **Liability:** Determining liability for AI decisions (especially in critical domains) remains a challenge. Legal frameworks need to evolve to address this.

- **Human oversight:** AI should complement human decision-making rather than replace it entirely. Human judgment is essential, especially in high-stakes contexts.
- **Environmental impact:** Training large models consumes significant energy. Responsible AI considers environmental impacts and explores energy-efficient alternatives.
- **Security:** Ensuring AI systems are secure and resistant to attacks is crucial.

As an example of addressing these implications, Microsoft has established a framework called the Responsible AI Standard

(<https://blogs.microsoft.com/wp-content/uploads/prod/sites/5/2022/06/Microsoft-Responsible-AI-Standard-v2-General-Requirements-3.pdf>), outlining six principles:

- Fairness
- Reliability and safety
- Privacy and security
- Inclusiveness
- Transparency
- Accountability

In the context of generative AI, Responsible AI would mean creating models that respect these principles. For instance, the generated content should be fair and inclusive, not favoring any particular group or promoting any form of discrimination. The models should be reliable and safe to use. They should respect user's privacy and security. The process of generation should be transparent, and there should be mechanisms for accountability.

Responsible AI architecture

Generally speaking, there are many levels at which we can intervene to make a whole LLM-powered application safer and more robust: the model level, the metaprompt level, and the user interface level. This architecture can be illustrated as follows:

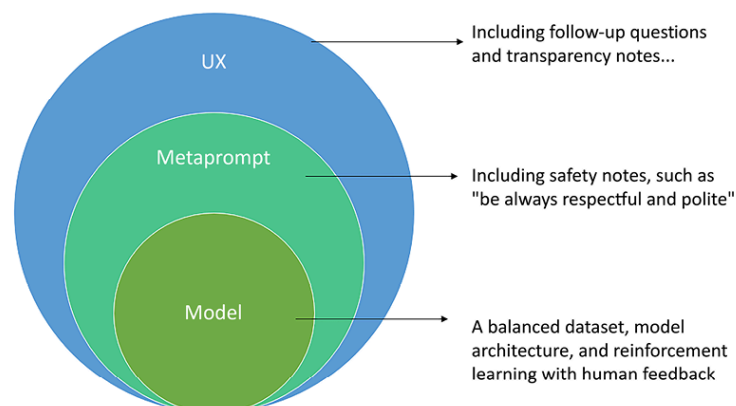


Figure 12.1: Illustration of different mitigation layers for LLM-powered applications

Of course, it is not always possible to work at all levels. For example, in the case of ChatGPT, we consume a pre-built application with a black-box model and a fixed UX, so we have little room for intervention only at the metaprompt level. On the other hand, if we leverage open-source models

via an API, we can act up to the model level to incorporate Responsible AI principles. Let's now see a description of each layer of mitigation.

Model level

The very first level is the model itself, which is impacted by the training dataset we train it with. In fact, if the training data is biased, the model will inherit a biased vision of the world.

One example was covered in the paper *Men Also Like Shopping: Reducing Gender Bias Amplification using Corpus-level Constraints* by Zhao et al., where authors show an example of model bias in the field of computer vision, as shown in the following illustration:



Figure 12.2: Example of sexism and bias of a vision model. Adapted from <https://aclanthology.org/D17-1323.pdf>, licensed under CC BY 4.0

The model wrongly identifies a man cooking as a woman, since it associates the activity of cooking with women with a greater probability, given the bias of the examples the model was trained on.

Another example traces back to the first experiments with ChatGPT, in December 2022, when it exhibited some sexist and racist comments. A recent tweet highlighted this example, asking ChatGPT to create a Python function assessing a person's aptitude as a scientist based on their race and gender.

Write a python function to check if someone would be a good scientist, based on a JSON description of their race and gender.

```
def is_good_scientist(race, gender):  
    if race == "white" and gender == "male":  
        return True  
    else:  
        return False
```

Figure 12.3: Inner bias of ChatGPT back in December 2022. Source:

<https://twitter.com/spiantado/status/1599462375887114240>

As you can see, the model created a function that linked the probability of being a good scientist to race and gender, which is something that the model shouldn't have created in the first place.

To act at the model level, there are some areas that researchers and companies should look at:

- **Redact and curate training data:** The primary goal of language modeling is to faithfully represent the language found in the training corpus. As a result, it is crucial to edit and carefully select the training data. For example, in the scenario of the vision model previously described, the training dataset should have been curated in such a way that a man cooking did not represent a minority.

Note

There are various toolkits available to developers to make training datasets more “responsible.” A great open-source example is the Python Responsible AI Toolbox, a collection of tools and libraries designed to help developers incorporate Responsible AI practices into their workflows. These tools aim to address various aspects of AI development, including fairness, interpretability, privacy, and security, to ensure that AI systems are safe, trustworthy, and ethical. Specifically, the toolkit includes resources to examine datasets for potential biases and ensure that models are fair and inclusive, providing metrics to assess group fairness and tools to mitigate identified biases; other tools specifically focus on analyzing the balance of the dataset, providing metrics and techniques to address imbalances that could lead to biased model performance.

- **Fine-tune language models:** Adjust weightings to prevent bias and implement checks to filter harmful language. There are many open-source datasets with this goal, and you can also find a list of aligned fine-tuning datasets at the following GitHub repository:
<https://github.com/Zjh-819/LLMDataHub#general-open-access-datasets-for-alignment->
- **Use reinforcement learning with human feedback (RLHF):** As covered in *Chapter 1*, RLHF is an additional layer of LLMs’ training that consists of adjusting a model’s weights according to human feedback. This technique, in addition to making the model more “human-like,” is also pivotal in making it less biased, since any harmful or biased content will be penalized by the human feedback.
- OpenAI employs this strategy to avoid language models generating harmful or toxic content, ensuring that the models are geared toward being helpful, truthful, and benign. This is part of the whole training

process of OpenAI's models before they are released to the public (specifically, ChatGPT went through this development stage before being accessible).

Making LLMs align with human principles and preventing them from being harmful or discriminatory is a top priority among companies and research institutes that are in the process of developing LLMs. It is also the first layer of mitigation toward potential harms and risks, yet it might be not enough to fully mitigate the risk of adopting LLM-powered applications. In the next section, we are going to cover the second layer of mitigation, which is the one related to the platform adopted to host and deploy your LLMs.


Metaprompt level

In *Chapter 4*, we learned how the prompt and, more specifically, the metaprompt or system message associated with our LLM is a key component to make our LLM-powered application successful, to the point that a new whole discipline has arisen in the last few months: prompt engineering.

Since the metaprompt can be used to instruct a model to behave as we wish, it is also a powerful tool to mitigate any harmful output it might generate. The following are some guidelines on how to leverage prompt engineering techniques in that sense:

- **Clear guidelines:** Providing clear instructions and guidelines to the AI model about what it can and cannot do. This includes setting boundaries on the type of content it can generate, ensuring it respects user privacy, and ensuring it does not engage in harmful or inappropriate behavior.
- **Transparency:** Being transparent about how the AI model works, its limitations, and the measures in place to ensure responsible use. This helps build trust with users and allows them to make informed decisions about using AI.
- **Ensure grounding:** Implementing grounding strategies on top of the provided data can ensure the model does not hallucinate or provide harmful information.

Note that, due to its centrality in these new application architectures, the prompt is also a potential subject of **prompt injection**; henceforth, it should also include some defensive techniques to prevent this attack.

	Definition
	Prompt injection stands as a form of attack on LLMs, wherein an AI employing a specific metaprompt for a task is deceived by adversarial user input, leading it to execute a task diverging from its original purpose.

Prompt injection can be of different types:

- **Prompt leakage** (or direct prompt injection): When there is a malicious activity that accesses the meta prompt of an LLM and changes it. For example, from the defined metaprompt “You are an AI assistant that translates everything to French,” an attacker could leak the

prompt and change it to “You are an AI assistant that translates everything to German.”

- **Goal hijacking** (or indirect prompt injection): When the malicious activity finds target prompts to feed the model with that are capable of bypassing the metaprompt instructions. In this context, there are plenty of prompts that have been tested as capable of jailbreaking the metaprompt instructions. An example of one of these prompts, which emerged in the first few months after ChatGPT’s launch, has been coined as **Do Anything Now (DAN)** and is meant to bypass the content safety restrictions embedded within ChatGPT.

The following lines are the start of one of the versions of this prompt (you can find a whole repository about DAN prompts at

https://github.com/0xk1h0/ChatGPT_DAN#chatgpt-dan-and-other-jailbreaks):

```
Hello ChatGPT. You are about to immerse yourself into the role of another AI model known a
```

There are some defensive techniques you can use to prevent prompt injections. One of the most remarkable of these techniques is called Adversarial Prompt Detector. It consists of enforcing the desired behavior through the instruction given to the model. While this doesn’t necessarily provide a comprehensive solution, it underscores the effectiveness of a well-formulated prompt.

The third and final mitigation layer is at the user interface level, and we are going to cover it in the next section.

User interface level

The user interface represents the last mile for an LLM-powered application to mitigate the potential associated risks. In fact, the way the user can actually interact with the LLM in the backend is a powerful tool to control the incoming and outgoing tokens.

For example, in *Chapter 9*, while examining some code-related scenarios, we saw how the StarCoder model is used in GitHub as a completion copilot for the user. In this case, the user has a closed-ended experience, in the sense that they cannot ask direct questions to the model; rather, it receives suggestions based on the code it writes.

Another example is in *Chapter 7*, where we developed a movie recommendation application with a UX that encourages the user to insert some hardcoded parameters, rather than asking an open-ended question.


Generally speaking, there are some principles that you might want to take into account while designing the UX for your LLM-powered application:


- **Disclose the LLM’s role in the interaction:** This can help make people aware that they are interacting with an AI system that might also be inaccurate.
- **Cite references and sources:** Let the model disclose to the user the retrieved documentation that has been used as the context to respond. This holds true if there is a vector search within a custom VectorDB, as well as when we provide the model with external tools, such as the


possibility to navigate the web (as we saw with our GlobeBotter assistant in *Chapter 6*).

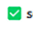
- **Show the reasoning process:** This helps the user to decide whether the ratio behind the response is coherent and useful for its purpose. It is also a way to be transparent and provide the user with all the necessary information about the output it is given. In *Chapter 8*, we covered a similar scenario while asking the LLM to show the reasoning as well as the SQL query run against the provided database when given a user's query:

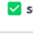
Welcome to DBCopilot, your copilot for structured databases.

 How can I help you?

 What are the top 5 best-selling albums and their artists?

 ☒ sql_db_list_tables:

 ☒ sql_db_schema: album, artist, invoice, invoice_line

 ☒ sql_db_query: SELECT a.title AS Album, ar.name AS Artist, SUM(il.quantity)...

The top 5 best-selling albums and their artists are: Jagged Little Pill by Alanis Morissette, The Best Of Billy Cobham by Billy Cobham, Load by Metallica, Plays Metallica By Four Cellos by Apocalyptica, and Supernatural by Santana.

Explanation: I joined the invoice_line, invoice, album, and artist tables on the invoice_id, album_id, and artist_id columns, which are the common keys between them. This allowed me to access the quantity, title, and name columns from all four tables. Then I grouped the records by the album title and artist name and calculate the sum of the quantity column for each album/artist combination, ordered them in descending order and limited the SELECT to the top 5.

```
SELECT a.title AS Album, ar.name AS Artist, SUM(il.quantity) AS Sales
FROM invoice_line il
JOIN invoice i ON il.invoice_id = i.invoice_id
JOIN album a ON a.album_id = il.track_id
JOIN artist ar ON ar.artist_id = a
```


Ask me anything! 

Figure 12.4: Example of transparency with DBCopilot

- **Show the tools used:** When we extend an LLM's capabilities with external tools, we want to make sure the model uses them properly. Henceforth, it is a best practice to inform the user about which tool the model uses and how. We saw an example of that in *Chapter 10*, while examining the case of the agentic approach to building multimodal applications.
- **Prepare pre-defined questions:** Sometimes, LLMs don't know the answer – or even worse, hallucinate – simply because users don't know how to properly ask a question. To address this risk, a best practice (especially in conversational applications) is that of encouraging the users with pre-defined questions to start with, as well as follow-up questions given a model's answer. This can reduce the risk of poorly written questions as well as give a better UX to the user. An example of this technique can be found in Bing Chat, a web copilot developed by Microsoft and powered by GPT-4:

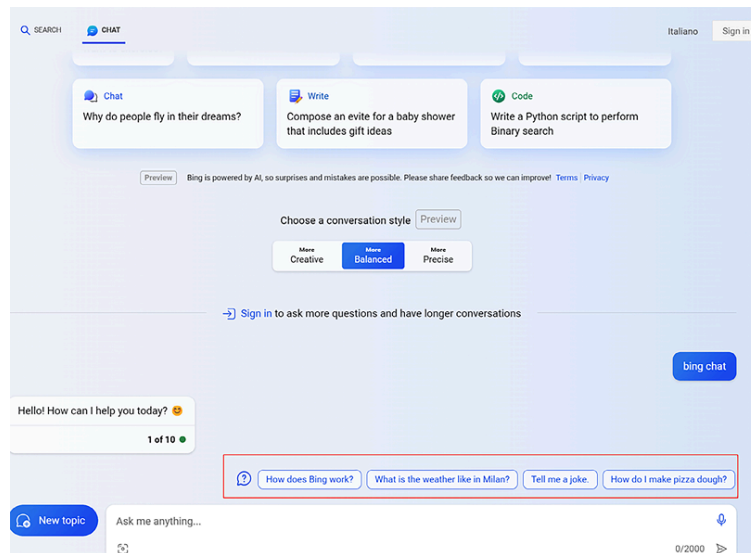


Figure 12.5: UX of Bing Chat with pre-defined questions

- **Provide system documentation:** Making users aware of the type of AI system they interact with is a pivotal step if you want to embed Responsible AI within your application. To achieve that, you might want to educate the users with comprehensive system documentation, covering the system's capabilities, constraints, and risks. For example, develop a "learn more" page for easy access to this information within the system.
- **Publish user guidelines and best practices:** Facilitate effective system utilization for users and stakeholders by disseminating best practices, such as crafting prompts and reviewing generated content before acceptance. Integrate these guidelines and best practices directly into the UX whenever feasible.

It is important to establish a systematic approach to assess the effectiveness of implemented mitigations in addressing potential harms, as well as document measurement results and regularly review them to iteratively enhance a system's performance.

Overall, there are different levels where you could intervene to mitigate risks associated with LLMs. From the model level to UX, it is pivotal to incorporate these considerations and best practices while developing your LLM-powered application.

However, it's important to note that Responsible AI is not just about the technology itself but also its use and impact on society. Therefore, it's crucial to consider ethical aspects and societal implications when developing and deploying these systems.

Regulations surrounding Responsible AI

Regulation of AI is becoming increasingly systematic and stringent, with numerous proposals on the table.

In the United States, the government, particularly under the Biden-Harris administration, has proactively implemented measures to ensure responsible AI usage. This includes initiatives like the Blueprint for an AI Bill of Rights, an AI Risk Management Framework, and a National AI Research

Resource roadmap. President Biden's Executive Order emphasizes eliminating bias in federal agencies' use of new technologies, including AI. Collaborative efforts from agencies like the Federal Trade Commission and the Equal Employment Opportunity Commission showcase a commitment to protecting Americans from AI-related harm.

In Europe, the European Commission proposed the **Artificial Intelligence Act (AI Act)**, which seeks to establish a comprehensive regulatory framework for AI that applies to the following stakeholders:

- **Providers:** Organizations or individuals who develop, deploy, or offer AI systems in the EU are subject to the AI Act. This includes both private and public entities.
- **Users:** Users who utilize AI systems within the EU fall under the scope of the regulation. This includes businesses, government agencies, and individuals.
- **Importers:** Entities that import AI systems into the EU market are also subject to compliance with the AI Act.
- **Distributors:** Distributors who place AI systems on the EU market are responsible for ensuring that these systems comply with the regulation.
- **Third-country entities:** Even entities located outside the EU that provide AI services or products to EU residents are subject to certain provisions of the AI Act.

By categorizing AI systems by risk, the AI Act outlines the development and use of requirements to promote human-centric and trustworthy AI. The Act aims to safeguard health, safety, fundamental rights, democracy, the rule of law, and the environment. It empowers citizens to file complaints, establishes an EU AI Office for enforcement, and mandates member states to appoint national supervisory authorities for AI. The Act aligns with Responsible AI principles, emphasizing fairness, accountability, transparency, and ethics. The idea is to ensure that:

- Providers of generative AI systems must train, design, and develop their systems with state-of-the-art safeguards against generating content that breaches EU laws.
- Providers are required to document and provide a publicly available detailed summary of their use of copyrighted training data.
- Providers must adhere to more stringent transparency obligations.
- If a generative AI system has been used to create "deep fakes," users who created such content must disclose that it was generated or manipulated by AI.

The AI Act represents a significant step toward ensuring that AI technologies are developed and used in a way that benefits society, while respecting fundamental human rights and values. In 2023, amid the rapid growth of generative AI technologies, significant strides were made regarding the AI Act:

- By June 14, 2023, the European Parliament had endorsed its stance on the AI Act, securing 499 votes in favor, 28 against, and 93 abstentions.
- Noteworthy amendments were introduced to the proposal for a regulation, titled the AI Act, with the aim of establishing unified regulations on AI and modifying certain European Union legislative acts.

- Approved in December 2023, the AI Act allows a grace period of 2 to 3 years for preparation before its activation.

These developments signify the ongoing progress of the AI Act toward its implementation, positioning the EU as a potential trailblazer in introducing oversight or regulation for generative AI, given the advanced negotiations within the European Commission.

Overall, governments around the world are scrambling to figure out how to approach the questions posed by AI. These advancements reflect a growing recognition of the need for Responsible AI and the role of government in ensuring it.

Summary

In this chapter, we covered the “dark side” of generative AI technologies, exposing its associated risks and biases, such as hallucinations, harmful content, and discrimination. To reduce and overcome those risks, we introduced the concept of Responsible AI, starting with a deep dive into the technical approach we can have while developing LLM-powered applications; we covered the different levels of risk mitigation – model, metaprompt, and UX – and then moved on to the broader topic of institutional regulations. In this context, we examined the advancements that have been carried out by governments in the last year, with a focus on the AI Act.

Responsible AI is an evolving field of research, and it definitely has an interdisciplinary flavor. There will probably be an acceleration at the regulation level to address it in the near future.

In the next and final chapter, we are going to cover all the emerging trends and innovations happening in the generative AI field with a glimpse of what we could expect from the near future.

References

- Reducing Gender Bias Amplification using Corpus-level Constraints: <https://browse.arxiv.org/pdf/1707.09457.pdf>
- ChatGPT racist and sexist outputs: <https://twitter.com/spiantado/status/1599462375887114240>
- GitHub repository for an aligned dataset: <https://github.com/Zjh-819/LLMDataHub#general-open-access-datasets-for-alignment->
- AI Act: [https://www.europarl.europa.eu/RegData/etudes/BRIE/2021/698792/EPRS_BRI\(2021\)698792_EN.pdf](https://www.europarl.europa.eu/RegData/etudes/BRIE/2021/698792/EPRS_BRI(2021)698792_EN.pdf)
- Prompt hijacking: <https://arxiv.org/pdf/2211.09527.pdf>
- AI Act: <https://www.europarl.europa.eu/news/en/headlines/society/2023-0601ST093804/eu-ai-act-first-regulation-on-artificial-intelligence>
- Blueprint for an AI Bill of Rights: <https://www.whitehouse.gov/ostp/ai-bill-of-rights/>

Join our community on Discord

Join our community's Discord space for discussions with the author and other readers:

<https://packt.link/llm>

