



**University of  
Zurich**<sup>UZH</sup>

**Zurich Open Repository and  
Archive**

University of Zurich  
University Library  
Strickhofstrasse 39  
CH-8057 Zurich  
[www.zora.uzh.ch](http://www.zora.uzh.ch)

---

Year: 2023

---

## **Toward Eliminating Hallucinations: GPT-based Explanatory AI for Intelligent Textbooks and Documentation**

Sovrano, Francesco ; Ashley, Kevin ; Bacchelli, Alberto

Posted at the Zurich Open Repository and Archive, University of Zurich  
ZORA URL: <https://doi.org/10.5167/uzh-257180>  
Conference or Workshop Item



The following work is licensed under a Creative Commons: Attribution 4.0 International (CC BY 4.0) License.

Originally published at:

Sovrano, Francesco; Ashley, Kevin; Bacchelli, Alberto (2023). Toward Eliminating Hallucinations: GPT-based Explanatory AI for Intelligent Textbooks and Documentation. In: Tokyo'23: Fifth Workshop on Intelligent Textbooks (iTextbooks) at the 24th International Conference on Artificial Intelligence in Education (AIED'2023), Tokyo, Japan, 3 July 2023 - 7 July 2023. CEUR-WS, 54-65.

# Toward Eliminating Hallucinations: GPT-based Explanatory AI for Intelligent Textbooks and Documentation

Francesco Sovrano<sup>1,3</sup>, Kevin Ashley<sup>2</sup> and Alberto Bacchelli<sup>1</sup>

<sup>1</sup>University of Zurich, Switzerland

<sup>2</sup>University of Pittsburgh, PA, USA

<sup>3</sup>University of Bologna, Italy

## Abstract

Traditional explanatory resources, such as user manuals and textbooks, often contain content that may not cater to the diverse backgrounds and information needs of users. Yet, developing intuitive, user-centered methods to effectively explain complex or large amounts of information is still an open research challenge. In this paper we present ExplanatoryGPT, an approach we devised and implemented to transform textual documents into interactive, intelligent resources, capable of offering dynamic, personalized explanations. Our approach uses state-of-the-art question-answering technology to generate on-demand, expandable explanations, with the aim of allowing readers to efficiently navigate and comprehend static materials. ExplanatoryGPT integrates ChatGPT, a state-of-the-art language model, with Achinstein's philosophical theory of explanations. By combining question generation and answer retrieval algorithms with ChatGPT, our method generates interactive, user-centered explanations, while mitigating common issues associated with ChatGPT, such as hallucinations and memory shortcomings. To showcase the effectiveness of our Explanatory AI, we conducted tests using a variety of sources, including a legal textbook and documentation of some health and financial software. Specifically, we provide several examples that illustrate how ExplanatoryGPT excels over ChatGPT in generating more precise explanations, accomplished through thoughtful macro-planning of explanation content. Notably, our approach also avoids the need to provide the entire context of the explanation as a prompt to ChatGPT, a process that is often not feasible due to common memory constraints.

## Keywords

Intelligent Textbooks, Software Documentation, Explanatory AI, ChatGPT, Question-answering technology, Hallucinations mitigation

## 1. Introduction

The increasing demand for efficient explanatory resources, such as user manuals and textbooks, calls for innovative approaches that can cater to the diverse backgrounds and information needs of users. Traditional materials often contain static, predefined content, which may not be sufficient for users to effectively comprehend complex information [1]. The challenge lies

---


*Tokyo'23: Fifth Workshop on Intelligent Textbooks (iTextbooks) at the 24th International Conference on Artificial Intelligence in Education (AIED'2023), July 03–07, 2023, Tokyo, Japan*

✉ francesco.sovrano@uzh.ch (F. Sovrano); ashley@pitt.edu (K. Ashley); alberto.bacchelli@uzh.ch (A. Bacchelli)

ORCID 0000-0002-6285-1041 (F. Sovrano); 0000-0003-0193-6823 (A. Bacchelli)



© 2023 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

in developing intuitive, user-centered methods that transform static content into dynamic, interactive explanations tailored to users' needs.

Recent advancements in natural language processing, particularly in large-scale language models like ChatGPT, have shown potential for enhancing explanatory capabilities [2, 3]. However, deploying these models for generating explanations can introduce challenges, such as hallucinations and memory constraints [4, 5].

To navigate these challenges, this paper proposes ExplanatoryGPT, a novel approach that harnesses the strengths of ChatGPT. This approach is dedicated to resolving the stated challenges and improving the quality of explanations in intelligent textbooks and documentation. For the purposes of this study, we interpret *hallucination* broadly as the generation of content that does not maintain fidelity to a given context or source.

ExplanatoryGPT combines ChatGPT with Achinstein's philosophical theory of explanations. By integrating question generation and answer retrieval algorithms with ChatGPT, our method generates interactive, user-centered explanations, while mitigating common issues associated with ChatGPT, such as hallucinations. Our approach contributes to the field of intelligent textbooks and documentation by enhancing the explanatory power of textual information utilizing question-answering technology to generate on-demand, expandable explanations.

In particular, our method involves the following steps: First, the question-answering algorithms extract pertinent textual information found in textbooks or documentation. Then, this information is reorganized based on the user's specific query, ensuring that the explanation is tailored to the user's needs and context. Subsequently, we employ ChatGPT to refine the retrieved information, enhancing its readability, coherence, and cohesion. By leveraging ChatGPT's text generation capabilities, we produce explanations that are not only more understandable but also more engaging and relevant to the user.

To demonstrate the effectiveness of our approach, we have applied it to a law textbook [6] and software documentation in healthcare and finance, showcasing how ExplanatoryGPT improves intelligent textbooks and documentation in different real-world scenarios. While we have not yet conducted a comprehensive qualitative or quantitative evaluation of our solution, we provide various concrete examples where ExplanatoryGPT surpasses ChatGPT. These examples demonstrate how ExplanatoryGPT generates more accurate explanations and avoids creating false information. Moreover, we offer insights into why ExplanatoryGPT is more suitable for educational settings and technical documentation.

Our work contributes to the academic discourse on the importance of integrating advanced language models, such as ChatGPT, and philosophical theories like Achinstein's framework, for the development of intelligent textbooks and documentation. This interdisciplinary approach can lead to the creation of more effective and versatile explanatory methods, which can be adapted to different domains and use cases. The benefits of our approach extend beyond the improvement of explanations themselves. By generating more accessible and comprehensible explanations, we can heighten users' engagement and understanding of various subjects and technologies. This, in turn, can lead to increased adoption and more effective use of educational materials.

## 2. Related Work

Our approach builds upon existing research on intelligent textbooks and natural language processing, particularly regarding advanced language models like GPT. In this section, we review related work and establish connections to our proposed approach.

Studies suggest that using intelligent textbooks<sup>1</sup> and interactive e-books can increase usage, motivation, and learning gains compared to static e-books [1]. Various approaches to interactive e-books and intelligent textbooks focus on the cognitive processes of readers, aiming to improve pedagogical productivity through expert systems or sophisticated interfaces. These approaches typically include: (1) showing personal progress through open learner models [8]; (2) specializing in ad hoc tasks through some domain modeling [9, 10]; (3) modeling a student through questions, in order to identify and suggest personalized contents [11, 12]; (4) associating pedagogically valuable quizzes and exercises to portions of the e-book [13, 14]; (5) providing tools for manually creating new interactive e-books [7, 15].

The use of Artificial Intelligence (AI) for the automatic generation of interactive e-books seems to be under-explored. In one such project, [16] propose to automatically augment the sections of existing books with related YouTube videos by directly annotating the PDF, thus without breaking the structure of these textbooks. Unlike previous research, our approach fully automates the conversion of existing e-books into interactive versions by integrating theories of explanations, intelligent interfaces, and Explanatory Artificial Intelligence (YAI). We explore how questions can practically organize and categorize explanation content.

By integrating ChatGPT's strengths with Achinstein's philosophical theory of explanations and question-answering algorithms, our proposed approach is a novel solution for enhancing intelligent textbooks and documentation. ExplanatoryGPT is aimed to generating more interactive, user-centered explanations while addressing challenges associated with advanced language models, fostering more effective and personalized learning experiences.

Our work aligns with current educational trends and the growing interest in the use of AI in education, as evidenced by the emerging literature. For example, UNESCO's Quick Start Guide<sup>2</sup> provides an overview of ChatGPT in higher education and emphasizes challenges and ethical implications. Similarly, Joyner [2] explores the impact of ChatGPT on education and curricula, drawing parallels with earlier technologies and predicting future effects. Additionally, Tlili et al. [3] conduct a qualitative instrumental case study examining ChatGPT's use in education, investigating various aspects such as public discourse, educational transformation, response quality, usefulness, personality and emotion, and ethics, revealing concerns like cheating, honesty and truthfulness, privacy, and manipulation.

To address the issues of hallucination and memory constraints commonly encountered in large language models like GPT [4, 5], we have developed an innovative solution. While previous research has also explored this problem and proposed solutions based on prompt engineering [17] and retrieval-augmented models [18], our approach distinguishes itself by incorporating philosophical theories of explanations in addition to information retrieval and machine learning techniques. Specifically, we build upon the work of [19], which utilizes Achinstein's theory of

---

<sup>1</sup>Intelligent textbooks extend regular textbooks by integrating machine-manipulable knowledge [7].

<sup>2</sup><https://unesdoc.unesco.org/ark:/48223/pf0000385146.locale=en>

explanations for answer retrieval (cf. Section 3.2), and demonstrate how it can be integrated with ChatGPT.

### 3. Background

In this section, we provide an overview of the theoretical foundations that underpin our proposed approach for enhancing the explanatory power of AI systems. We discuss YAI and the role of ChatGPT as an advanced language model.

#### 3.1. Explanations According to Ordinary Language Philosophy

The concept of explanation in philosophy began to have a more precise role in the 20th century with the growth and development of the philosophy of science. Hempel's deductive-nomological model [20] gave rise to the first theory of explanations, followed by many competing theories, such as Salmon's Causal Realism [21], and Achinstein's Ordinary Language Philosophy [22].

Achinstein's theory emphasizes the communicative aspect of explanation, its usefulness in answering questions, and fostering understanding between individuals. Holland's theory [23] frames the process of explaining as a purely cognitive activity, while Sellars [24] suggests a utilitarian process of constructing a coherent belief system.

According to Achinstein, explaining is an *illocutionary* act, born of a clear intention to produce new understandings in an explainee by providing a correct content-giving answer to an open-ended question. Illocution in explaining involves informed and *pertinent* answers to the main question and other implicitly relevant questions [25, 19].

**Definition 1 (Illocution in Explaining).** *Explaining is an illocutionary act that provides answers to an explicit question on some topic along with answers to several other implicit or unformulated questions deemed necessary for the explainee to understand the topic properly. In the most generic case, no assumption can be made about the explainee's knowledge and objectives, and the only implicit questions that can be exploited for illocution are the most generic ones, called archetypal questions.*

For example, an answer like 'I am happy because I just got a paper accepted at this important venue, and [...]' would generally be considered an explanation because it answers other *archetypal questions*.

#### 3.2. Explanatory Artificial Intelligence

An YAI is an artificial intelligence program designed to generate user-centered, interactive explanations out of (possibly extensive) collections of explainable information [26]. An example of YAI based on Achinstein's theory of explanation is YAI4Hu [25, 19].

YAI4Hu is a fully automatic explanatory tool used to explain (pre-existing) documentation about an AI-based system. In particular, the textual content of such documentation is algorithmically reorganized and represented as a special hypergraph where information can be either explored through *overviewing* or searched via *open-ended questioning*. On the one hand,

*open-ended questioning* can be performed by asking open-ended questions in English through a search box that uses the knowledge graph for efficient answer retrieval. An example of open-ended questioning is shown in Figure 2.

On the other hand, *overviewing* can be performed iteratively from an initial explanation by clicking on automatically annotated words for which explanations exist. In particular, annotated words are visible because they have a unique format that makes them easy to recognize. After clicking on an annotation, a modal opens (see Figure 1), showing a navigation bar with tabs containing explanatory overviews of the clicked annotated words. The information shown in the overview includes:

- A short description of the explained word (if available).
- The list of other taxonomically connected words.
- A list of predefined archetypal questions (e.g., why is this aspect/concept important, what is this aspect/concept, etc.) and their respective answers ordered by estimated pertinence.

### 3.3. ChatGPT as an advanced language model

ChatGPT is a cutting-edge language model based on the GPT architecture [27]. GPT has been used in various applications, including question-answering systems [28], summarization [29, 30], and chatbots [31, 32], showcasing its versatility and potential for generating contextually relevant and coherent text. However, ChatGPT and similar models have been known to suffer from hallucinations and memory constraints, which can lead to the generation of text that is plausible-sounding but factually incorrect or ungrounded [33]. Our proposed methodology addresses these issues by integrating ChatGPT with question-answering algorithms, grounding the generated explanations in relevant and accurate information, and refining the readability, coherence, and cohesion of the retrieved information.

## 4. Proposed Approach

In this section, we introduce our proposed approach, which combines answer retrieval algorithms with ChatGPT to produce high-quality explanations. Our main objective is to leverage ChatGPT's ability to provide clear explanations while addressing common issues such as hallucinations and lack of control over its outputs, which can utilize any learned information regardless of its relevance.

Our strategy builds upon the open-questioning and overviewing mechanisms from Section 3.2. These mechanisms use a concept called illocution, where key archetypal questions are answered to produce valid explanations. Our goal is to enhance these mechanisms by utilizing ChatGPT to merge the answers, thus taking advantage of the strengths of both techniques.

Our approach combines information retrieval algorithms with ChatGPT in a complementary manner to generate high-quality explanations. On one hand, the information retrievers extract and reorder critical information from textbooks or software documentation, designed to address a wide range of user questions, explicit or implicit. By aligning explanations with the user's specific query, we ensure relevance and contextual appropriateness.

On the other hand, ChatGPT is employed to refine the retrieved information. ChatGPT possesses good text generation capabilities, producing human-like content (see Section 3.3). Moreover, it enables the generation of high-quality and customizable text [34]. For these reasons, after the algorithms extract and reorganize relevant information, we use ChatGPT to refine the information retrieved, aiming to enhance readability, coherence, and cohesion of the explanations, making them more accessible to a wider audience.

In the refinement process, ChatGPT rephrases, restructures, and condenses the information gathered by the answer retrieval algorithm, ensuring clear, concise, and easily understandable explanations. This approach also addresses potential issues of hallucinations or memory limitations in ChatGPT by grounding it in the information retrieved by the algorithms. To guide ChatGPT in the refinement process, the following instruction (or any similar instruction) can be given:

*Using exclusively the information provided, create a coherent and clear explanation of '{question}' for a student with no prior knowledge on the subject. Provide both jargon and simpler synonyms. Your exposition should follow this structure: Short Answer, Technical Details, Conclusion.*

*Information provided:*

- {retrieved\_answer\_1}
- {retrieved\_answer\_2}
- ...

In this instruction:

- '{question}' is a placeholder for the actual question that the user (implicitly or explicitly) asks, which will be replaced with the specific question that requires an explanation.
- '{retrieved\_answer\_1}', '{retrieved\_answer\_2}', and so on, are placeholders for the answers that the question-answering algorithms extract. These answers come from reliable sources or databases and provide the raw material for ChatGPT to structure the explanation.

The instruction guides ChatGPT to create an explanation using the information provided, i.e., the retrieved answers. The explanation should be clear and understandable to a student with no prior knowledge of the subject. It asks ChatGPT to use both technical terms (jargon) and their simpler synonyms to make the information accessible to a wide range of users.

The requested structure (i.e., Short Answer, Technical Details, Conclusion) helps ChatGPT structure the information logically. The 'Short Answer' provides a brief response to the question, the 'Technical Details' delve deeper into the explanation, and the 'Conclusion' summarizes the main points. By following these guidelines, ChatGPT is expected to generate explanations that are not only technically accurate but also clear, well-structured, and suitable for users with different levels of understanding. The focus on using exclusively the information provided helps keep ChatGPT on track and reduces potential hallucinations.

In the next section, we will delve into the details of how this approach can be employed, its limitations, and the outcomes it can yield.

Months since most recent credit inquiry not within the last 7 days
× Credit inquiry
× Contrastive Explanations Method (CEM)

## Contrastive Explanations Method (CEM)

- **Abstract:** The Contrastive Explanations Method (CEM) is an algorithm that tries to find the minimally sufficient changes to the input (of an automated process like an Artificial Neural Network) that can change the output. Sometimes CEM is not able to find these minimally sufficient changes, but it can go close to that.
- An algorithm code of it is <https://github.com/IBM/AIX360/blob/master/aix360/algorithms/contrastive/CEM.py>
- This is a kind of «Based on Features» [\[+\]](#)
- A document referring to it is <https://arxiv.org/abs/1802.07623>
- See also [artificial neural network](#), etc... [\[+\]](#)

**How**

- ▶ The bank uses a method called Contrastive Explanations Method (CEM) to explain how it predicts your risk performance. It identifies features that are responsible for the prediction and those that are absent but would have changed the prediction. [\[-\]](#)
  - The bank uses an Artificial Neural Network to predict your risk performance. To explain the prediction, the bank uses an algorithm called Contrastive Explanations Method (CEM). CEM identifies the features in your input that are responsible for the prediction and those that are absent but would have changed the prediction. CEM tries to find the minimally sufficient changes to the input that can change the output. However, sometimes CEM cannot find these changes, but it can get close to them.
  - Research shows that to fully scrutinize model biases, people need a diverse set of explanation capabilities. Algorithms like Boolean Rule Column Generation and Generalized Linear Rule Model can help people inspect if there is discrimination in the overall logic of the model. CEM and ProtoDash can help people ensure that they are not being unfairly treated by comparing the model's decisions for them to other individuals.
  - The bank uses a method called Contrastive Explanations Method (CEM) to explain how it predicts your risk performance. CEM identifies the features that are responsible for the prediction and those that are absent but would have changed the prediction. The bank also uses other algorithms to help people inspect if there is discrimination in the overall logic of the model and ensure that they are not being unfairly treated by comparing the model's decisions for them to other individuals. [\[-\]](#)

**Figure 1: Example of intelligent explanatory overview generated by ExplanatoryGPT.** This figure contains an example of interactive *overview* in the form of a *scrollspy* showing relevant questions and answers as explanations. The reader can select a new topic to overview by clicking on any underlined word.

## 5. Proof of Concept and Discussion: Real-World Applications of ExplanatoryGPT

We have developed three proof-of-concept demonstrations of the ExplanatoryGPT. These examples showcase the potential and versatility of the ExplanatoryGPT system. Two of these demonstrations focus on the software documentation of AI-based systems, specifically those used for Credit Approval and predicting heart diseases. We selected these examples because they delve into intricate, technical details and concepts that could easily confuse an AI model. More information about these AI-based systems can be found in previous work [19].

Our third demonstration involves the ExplanatoryGPT's application to a textbook designed to teach the writing of legal memorandums in the US legal system. This example showcases the model's potential to be used in educational settings, specifically for teaching complex topics like law. Details about this textbook and educational setting can be found in [35].

The reason we chose these diverse examples is because they involve information and technical concepts that may not be easily understood by ChatGPT. Software documentation, for instance, often contains technical information that needs to be interpreted in context, while textbooks may contain specific information that may need to be carefully controlled or hidden during the learning process.

An overview of the ExplanatoryGPT mechanism in the Credit Approval system is presented in Figure 1. As seen in this figure, ExplanatoryGPT synthesizes both retrieved and generated information to deliver an explanation. This mechanism is expected to be especially effective in technical areas such as credit approval, where information needs to be clearly communicated to



# Intelligent Textbook

Here you can read some excerpts of the book **Brostoff, T. & Sinsheimer, A. (2013). United States Legal Language and Culture: An Introduction to the US Common Law System. Third Edition, Oxford University Press.** You can also ask questions about the content covered by these excerpts. However, please note that this question answering tool has limited capabilities and cannot replace the excerpts, nor directly answer the assignment questions.

How do I write a memorandum?

**Question:** How do I write a memorandum? ✕

**Answers:**

- A memorandum is a document that lawyers use to communicate information and analysis on a specific legal problem. It should be concise, clear, and objective. [-]
- A memorandum usually includes a description of factual background, a statement of legal issues, an introduction of relevant laws, an analysis of how the law should apply to specific facts, and a conclusion. The sections of a legal memorandum are: heading, the question presented, the brief answer, the fact section, the discussion/analysis section, and the conclusion section. In the discussion/analysis section, you will fully explain and analyze your predictions. When drafting a legal memorandum, you will usually predict how a court will rule on a particular legal question under particular factual circumstances. You should always be as definitive as you can be when you give your answer to the reader.
- Writing a legal memorandum is difficult at first but will become easier with practice. When editing your work, make sure your paragraphs are focused and not overly long. The memorandum should be easy to read, and the reader should understand it completely after reading through it once. [-]

Pertinence	Source	File
69.34%	In many legal settings specialized forms of written communication are required. In many others, writing is the medium in which a <u>lawyer</u> must express their analysis of an issue and seek to persuade others on their clients' <u>behalf</u> . Any legal <u>document</u> must be concise, clear, and conform to the objective standards that have evolved in the legal profession.	Legal writing

**Figure 2: Examples of intelligent answers generated by ExplanatoryGPT.** The figure only displays one source answer (at the bottom of the picture), but other 9 sources were considered by ChatGPT during the explanation generation process.

avoid any misunderstanding. Similarly, Figure 2 illustrates the open-questioning mechanism of ExplanatoryGPT in the educational scenario of teaching legal memorandum writing.

Table 1 presents a side-by-side comparison of the explanatory outputs produced by ChatGPT and ExplanatoryGPT. This comparison highlights the differences in the quality and precision of explanations provided by the two models. Specifically, the explanations generated by ChatGPT are obtained through the following instruction which (unlike ExplanatoryGPT) does not include any retrieved answer but only the given question:

*“Create a coherent and clear explanation of ‘question’ for a student with no prior knowledge on the subject. Provide both jargon and simpler synonyms. Your exposition should follow this structure: Short Answer, Technical Details, Conclusion.”*

As shown in Table 1, our ExplanatoryGPT mechanism has the potential to effectively control what content should be part of the explanations. This control is crucial for several reasons:

1. It mitigates the ‘hallucination’ problem of ChatGPT, where the model sometimes generates information that is irrelevant or incorrect. This is evident in the ‘CEM’ example shown in Table 1.
2. It improves contextualisation of explanations which may otherwise be generic and not specifically focused. For example, the ‘Satisfactory Trades’ explanation in Table 1 is more accurate and targeted due to the control provided by ExplanatoryGPT.

**Table 1**

This table presents examples of questions and their corresponding explanatory *short answers* generated by ExplanatoryGPT and ChatGPT. The purpose of these examples is to shed light on the hallucinations produced by ChatGPT. In the table, all the explanations generated by ChatGPT are hallucinated or vague.

Question	ExplanatoryGPT	ChatGPT
How do I write a memorandum?	A memorandum is a document that lawyers use to communicate information and analysis on a specific legal problem. It should be concise, clear, and objective.	A memorandum is a document used for internal communication within an organization.
What is CEM (in the context of the Credit Approval System)?	The Contrastive Explanations Method (CEM) is an algorithm that helps us understand how small changes to the input of a computer program, like an Artificial Neural Network, can affect its output.	CEM, or Customer Experience Management, refers to the strategies and practices that businesses use to understand and enhance the experiences of their customers.
What is the number of satisfactory trades?	The number of satisfactory trades refers to the count of credit accounts that were paid in full and closed without any negative information.	The number of satisfactory trades refers to the total count of successful transactions or exchanges that meet the expectations and requirements of the parties involved.

3. It provides a mechanism to control which information should be shown or hidden. This feature is particularly useful in educational settings, where certain information might be withheld as part of a learning exercise.

Instead of employing ChatGPT as a stand-alone generator of end-to-end explanations, we incorporate an YAI mechanism for the macro-planning phase, which involves the selection and ordering of information [36]. Subsequently, ChatGPT is utilized for micro-planning and for surface realization, i.e., the coherent amalgamation of information into sentences.

The main task of ExplanatoryGPT is to avoid overly broad, incorrect, or out-of-context explanations. This becomes especially clear when considering the examples in Table 1, such as the first two questions. The answers provided by ExplanatoryGPT are much more specific and indicate a greater ‘awareness’ of the explanatory context.

Notably, in the absence of an YAI mechanism, one might still formulate a prompt incorporating the full context necessary for generating an appropriate response. However, in the broadest of scenarios, executing this prompt without a mechanism (i.e., an YAI) for filtering redundant information would not be feasible, effective, or efficient. This is primarily due to the stringent memory constraints associated with ChatGPT. Therefore, the use of a YAI allows for more accurate and context-specific outputs.

## 6. Conclusion and Future Work

In this paper, we presented ExplanatoryGPT, an approach that combines ChatGPT with Achinstein’s philosophical theory of explanations for intelligent textbooks and documentation. Our proposed YAI methodology has the potential to enhance the quality and usability of explanations, making them more interactive and user-centred.

We showcased the potential of ExplanatoryGPT to generate better and less hallucinated explanations in a law textbook and software documentation related to healthcare and finance. Our

research highlights the importance of integrating advanced language models and philosophical theories for the future development of intelligent textbooks and documentation.

For future work, we plan to conduct extensive user studies to empirically evaluate our approach, as our current evaluation was limited in scope and rigor. Other research directions include expanding the application of our methodology to other domains and industries.

In conclusion, our proposed YAI methodology offers a promising avenue for improving learning and understanding experiences for a diverse audience, ultimately contributing to more inclusive and efficient educational materials and software documentation.

## Declarations

**Acknowledgements.** We thank the copyright holders of [6] for allowing us to use (parts of) the book to conduct the experiments, carry out the case study and present this article. F. Sovrano and A. Bacchelli gratefully acknowledge the support of the Swiss National Science Foundation through the SNF Project 200021\_197227.

**Author Contributions.** F. Sovrano: conceptualization, methodology, software, data curation, original draft preparation, visualization, investigation, validation, formal analysis. K. Ashley: conceptualization and supervision. A. Bacchelli: review, editing, and supervision.

## References

- [1] B. Ericson, An analysis of interactive feature use in two ebooks, in: Proceedings of the First Workshop on Intelligent Textbooks, volume 2384 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2019, pp. 4–17.
- [2] D. A. Joyner, Chatgpt in education: Partner or pariah?, *XRDS* 29 (2023) 48–51.
- [3] A. Tlili, B. Shehata, M. A. Adarkwah, A. Bozkurt, D. T. Hickey, R. Huang, B. Agyemang, What if the devil is my guardian angel: Chatgpt as a case study of using chatbots in education, *Smart Learn. Environ.* 10 (2023) 15.
- [4] R. Azamfirei, S. R. Kudchadkar, J. Fackler, Large language models and the perils of their hallucinations, *Critical Care* 27 (2023) 120.
- [5] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder, P. F. Christiano, J. Leike, R. Lowe, Training language models to follow instructions with human feedback, in: *NeurIPS*, 2022.
- [6] T. K. Brostoff, A. Sinsheimer, *United States Legal Language and Culture: An Introduction to the U.S. Common Law System*, Oxford University Press USA, 2013. URL: <https://books.google.it/books?id=SVsGAQAAQBAJ>.
- [7] M. Wang, H. Chau, K. Thaker, P. Brusilovsky, D. He, Knowledge annotation for intelligent textbooks, *Technology, Knowledge and Learning* (2021).
- [8] J. Kay, B. Kummerfeld, PUMPT: an e-textbook platform based on a personal user model for learning (short paper), in: Proceedings of the Third International Workshop on Intelligent Textbooks 2021, volume 2895 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2021, pp. 27–34.

- [9] J. P. Beier, M. A. Rau, Embodied learning with physical and virtual manipulatives in an intelligent tutor for chemistry, in: Artificial Intelligence in Education - 23rd International Conference, AIED 2022, Durham, UK, July 27-31, 2022, Proceedings, Part I, volume 13355 of *Lecture Notes in Computer Science*, Springer, 2022, pp. 103–114.
- [10] I. A. Chacon, J. Barria-Pineda, K. Akhuseyinoglu, S. A. Sosnovsky, P. Brusilovsky, Integrating textbooks with smart interactive content for learning programming, in: Proceedings of the Third International Workshop on Intelligent Textbooks 2021, volume 2895 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2021, pp. 4–18.
- [11] M. Mohammed, C. A. Shaffer, Increasing student interaction with an etextbook using programmed instruction (short paper), in: Proceedings of the Third International Workshop on Intelligent Textbooks 2021, volume 2895 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2021, pp. 40–44.
- [12] N. Matsuda, M. Shimmei, PASTEL: evidence-based learning engineering method to create intelligent online textbook at scale, in: Proceedings of the First Workshop on Intelligent Textbooks, volume 2384 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2019, pp. 70–80.
- [13] Z. Wang, J. Valdez, D. B. Mallick, R. G. Baraniuk, Towards human-like educational question generation with large language models, in: Artificial Intelligence in Education - 23rd International Conference, AIED 2022, Durham, UK, July 27-31, 2022, Proceedings, Part I, volume 13355 of *Lecture Notes in Computer Science*, Springer, 2022, pp. 153–166.
- [14] M. Shimmei, N. Matsuda, Automatic question generation for evidence-based online courseware engineering, in: Proceedings of the Fourth International Workshop on Intelligent Textbooks 2022, volume 3192 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2022, pp. 18–25.
- [15] B. Kluga, M. S. Jasti, V. Naples, R. Freedman, Adding intelligence to a textbook for human anatomy with a causal concept map based ITS, in: Proceedings of the First Workshop on Intelligent Textbooks, volume 2384 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2019, pp. 124–134.
- [16] J. Barria-Pineda, A. B. L. Narayanan, P. Brusilovsky, Augmenting digital textbooks with reusable smart learning content: Solutions and challenges, in: Proceedings of the Fourth International Workshop on Intelligent Textbooks 2022, volume 3192 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2022, pp. 77–91.
- [17] Y. Bang, S. Cahyawijaya, N. Lee, W. Dai, D. Su, B. Wilie, H. Lovenia, Z. Ji, T. Yu, W. Chung, et al., A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity, arXiv preprint arXiv:2302.04023 (2023).
- [18] Y. Mass, H. Roitman, Ad-hoc document retrieval using weak-supervision with BERT and GPT2, in: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020, Association for Computational Linguistics, 2020, pp. 4191–4197.
- [19] F. Sovrano, F. Vitali, Generating user-centred explanations via illocutionary question answering: From philosophy to interfaces, *ACM Trans. Interact. Intell. Syst.* 12 (2022).
- [20] C. G. Hempel, P. Oppenheim, Studies in the logic of explanation, *Philosophy of Science* 15 (1948) 135–175. doi:10.1086/286983.
- [21] W. Salmon, *Scientific Explanation and the Causal Structure of the World*, Book collections on Project MUSE, Princeton University Press, 1984. URL: <https://books.google.it/books?>

id=2ug9DwAAQBAJ.

- [22] P. Achinstein, *The Nature of Explanation*, Oxford University Press, 1983. URL: <https://books.google.it/books?id=0XI8DwAAQBAJ>.
- [23] J. Holland, K. Holyoak, R. Nisbett, P. Thagard, *Induction: Processes of Inference, Learning, and Discovery*, Bradford books, MIT Press, 1986. URL: <https://books.google.it/books?id=Z6EFBaLApE8C>.
- [24] W. Sellars, *Science, Perception and Reality*, New York: Humanities Press, 1963.
- [25] F. Sovrano, F. Vitali, From philosophy to interfaces: an explanatory method and a tool inspired by achinstein's theory of explanation, in: *IUI '21: 26th International Conference on Intelligent User Interfaces*, College Station, TX, USA, April 13-17, 2021, ACM, 2021, pp. 81–91.
- [26] F. Sovrano, F. Vitali, Explanatory artificial intelligence (yai): human-centered explanations of explainable ai and complex data, *Data Mining and Knowledge Discovery (2022)*.
- [27] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al., Language models are unsupervised multitask learners, 2019.
- [28] B. Sievers, Question answering for comparative questions with GPT-2, in: *Working Notes of CLEF 2020 - Conference and Labs of the Evaluation Forum*, Thessaloniki, Greece, September 22-25, 2020, volume 2696 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2020.
- [29] Z. Yang, Y. Dong, J. Deng, B. Sha, T. Xu, Research on automatic news text summarization technology based on GPT2 model, in: *AIAM 2021: 3rd International Conference on Artificial Intelligence and Advanced Manufacture*, Manchester, United Kingdom, October 23 - 25, 2021, ACM, 2021, pp. 418–423.
- [30] B. Chintagunta, N. Katariya, X. Amatriain, A. Kannan, Medically aware GPT-3 as a data generator for medical dialogue summarization, in: K. Jung, S. Yeung, M. P. Sendak, M. W. Sjoding, R. Ranganath (Eds.), *Proceedings of the Machine Learning for Healthcare Conference, MLHC 2021, 6-7 August 2021, Virtual Event*, volume 149 of *Proceedings of Machine Learning Research*, PMLR, 2021, pp. 354–372.
- [31] T. Teubner, C. M. Flath, C. Weinhardt, W. M. P. van der Aalst, O. Hinz, Welcome to the era of chatgpt et al, *Bus. Inf. Syst. Eng.* 65 (2023) 95–101.
- [32] M. Cascella, J. Montomoli, V. Bellini, E. Bignami, Evaluating the feasibility of chatgpt in healthcare: An analysis of multiple clinical and research scenarios, *J. Medical Syst.* 47 (2023) 33.
- [33] H. Alkaissi, S. I. McFarlane, Artificial hallucinations in chatgpt: implications in scientific writing, *Cureus* 15 (2023).
- [34] P. P. Ray, Chatgpt: A comprehensive review on background, applications, key challenges, bias, ethics, limitations and future scope, *Internet of Things and Cyber-Physical Systems (2023)*.
- [35] F. Sovrano, K. Ashley, P. Brusilovsky, F. Vitali, Yai4edu: an explanatory AI to generate interactive e-books for education, in: *Proceedings of the Fourth International Workshop on Intelligent Textbooks 2022*, volume 3192 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2022, pp. 31–39.
- [36] R. Dale, H. Moisl, H. Somers, *Handbook of Natural Language Processing*, Taylor & Francis, 2000. URL: <https://books.google.it/books?id=VoOLvxyX0BUC>.