



## ARTICLE

# Towards the ultimate brain: Exploring scientific discovery with ChatGPT AI

Gerardo Adesso

School of Mathematical Sciences,  
University of Nottingham, Nottingham,  
UK

## Correspondence

Gerardo Adesso, School of Mathematical  
Sciences, University of Nottingham,  
Nottingham, UK.  
Email: [gerardo.adesso@nottingham.ac.uk](mailto:gerardo.adesso@nottingham.ac.uk)

Generated using ChatGPT (November 30,  
2022 free research preview release for the  
original manuscript and January 30, 2023  
update for the revised version):  
<https://openai.com/blog/chatgpt/>

## Funding information

Foundational Questions Institute,  
Grant/Award Number: RFP-IPW1907

## Abstract

This paper presents a novel approach to scientific discovery using an artificial intelligence (AI) environment known as ChatGPT, developed by OpenAI. This is the first paper entirely generated with outputs from ChatGPT. We demonstrate how ChatGPT can be instructed through a gamification environment to define and benchmark hypothetical physical theories. Through this environment, ChatGPT successfully simulates the creation of a new improved model, called GPT<sup>4</sup>, which combines the concepts of GPT in AI (generative pretrained transformer) and GPT in physics (generalized probabilistic theory). We show that GPT<sup>4</sup> can use its built-in mathematical and statistical capabilities to simulate and analyze physical laws and phenomena. As a demonstration of its language capabilities, GPT<sup>4</sup> also generates a limerick about itself. Overall, our results demonstrate the promising potential for human-AI collaboration in scientific discovery, as well as the importance of designing systems that effectively integrate AI's capabilities with human intelligence.

## INTRODUCTION

The advent of advanced language models such as OpenAI's GPT series (Radford et al. 2018) and Google's LaMDA (Thoppilan et al. 2022) paved the way for innovative and exciting human-AI collaboration opportunities. These models are capable of generating human-like text, which can be harnessed to support various tasks ranging from language translation to question-answering (Devlin et al. 2018). However, the extent of their capabilities for assisting humans in scientific discovery is not yet fully understood (Davies et al. 2021; Gpt Generative Pretrained Transformer, Thunström, and Steingrímsson 2022; Stokel-Walker 2023).

The overarching goal of this paper is to examine the capabilities and limitations of OpenAI's state-of-the-art language model, ChatGPT (OpenAI 2022), in enhancing

human collaboration in scientific pursuits. To address this question, we devised an experimental setup where ChatGPT assumes the role of a researcher tasked with investigating a challenging topic at the intersection of computer science and fundamental physics. The model conducts a qualitative and quantitative analysis, seemingly exploring innovative ideas and generating scientific insights. The results are then presented in a format akin to a conventional scientific publication.

To facilitate this investigation, we designed a gamification environment in which the human author instructs ChatGPT to define and benchmark hypothetical physical theories. Through this environment, we demonstrate how ChatGPT can successfully simulate the creation of a new theoretical model, called GPT<sup>4</sup> (not to be confused with GPT-4, which has been developed by OpenAI and has been

This is an open access article under the terms of the [Creative Commons Attribution](https://creativecommons.org/licenses/by/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2023 The Authors. *AI Magazine* published by John Wiley & Sons Ltd on behalf of Association for the Advancement of Artificial Intelligence.



released after the original submission of this paper), which combines the concepts of Generative Pretrained Transformer (GPT) in AI (Radford et al. 2018)—a probabilistic model that can be used for generating and reasoning about natural language—and Generalized Probabilistic Theory (GPT) in physics (Barnum and Wilce 2001)—a mathematical framework for describing the probabilistic behavior of physical systems.

We then evaluate the hypothetical GPT<sup>4</sup> model created within the language model ChatGPT. We put this model to the test by asking it a series of language-related and science-related questions, including creating a limerick about itself, modeling physical phenomena, and calculating mathematical results. As a purely speculative investigation, we also ask the model to estimate the probability that AI will surpass humanity within a given number of years and compare its predictive power to that of other physical theories using fundamental results in estimation theory. These experiments are conducted within a virtual gamification environment and should not be regarded as possessing any scientific foundation beyond showcasing the progress of AI.

One of the most unique aspects of this study is that the entire paper was produced by combining outputs from OpenAI's ChatGPT (OpenAI 2022), making this *meta-paper* an experiment in itself. All text in the main body of the paper that is not in a light gray box is generated by the GPT model itself. The light gray boxes contain prompts given by the author to guide the GPT model in generating responses. In addition to the main text, the tables, figure captions, and the references and bibliography in this paper are also autogenerated by the GPT model (in some cases, the author has supplied input to the model, especially for more recent articles which were not part of ChatGPT's training set).

This serves as a demonstration of the potential applications of advanced language models in scientific discovery and communication. Our results add to the growing body of research exploring the capabilities of such models (Srivastava et al. 2022) and offer a valuable contribution to the advancement of AI and scientific understanding. The ability of ChatGPT to generate coherent and sound text highlights the potential of human-AI teaming in scientific discovery, and further underscores the need for ongoing research and development in this area (Stokel-Walker 2023; Adesso 2023).

The paper is organized as follows: Section 2 contains the basic definitions and methods used. Section 3 presents the main gamification experiment and discusses its results. Finally, the paper concludes with Section 4 summarizing the findings and providing some final thoughts. Furthermore, we include an Appendix, written by the human author, which provides more details on the procedure that

led to the co-creation of this paper, as well as the author's perspective, evaluation, and reflections.

## METHODS

In this section, we will provide basic definitions and explain the methods used in our study.

### GPT in AI

GPT, or Generative Pretrained Transformer, is a type of deep learning model developed by OpenAI (Radford et al. 2018). It is based on the Transformer architecture (Vaswani et al. 2017), which is a type of recurrent neural network that uses self-attention mechanisms to process sequential data. GPT uses a large corpus of unannotated text data to pretrain a language model that can generate coherent and fluent sentences.

One of the key advantages of GPT is its ability to handle long-range dependencies in natural language. This allows it to generate coherent and fluent text that is difficult for other models to produce. In addition, GPT has a large capacity for learning and can be fine-tuned for a variety of natural language tasks, such as language translation, summarization, and question answering.

### GPT in physics

GPT, or Generalized Probabilistic Theory, is a mathematical framework for describing the probabilistic behavior of physical systems. It is based on the idea that physical systems can be described by sets of probabilistic states and operations that transform these states (Barnum and Wilce 2001).

GPT provides a flexible framework for modeling the probabilistic behavior of physical systems. It allows for the description of systems with arbitrary numbers of parties and inputs, and allows for the consideration of both classical and quantum systems. In addition, GPT provides a unifying framework for the description of various physical theories, such as classical probability theory, quantum mechanics, and relativity (Hardy 2001).

### Content generation

In this study, we test the capabilities of ChatGPT in generating a scientific paper and evaluate the results to understand its strengths and limitations in this task. The process consisted of several steps:

1. **Data collection:** We used the latest version of the GPT-3.5 language model, which has been trained on a diverse range of internet text (OpenAI 2022).
2. **Input Generation:** We fed the model with a prompt indicating the task, for example, generating a section of the paper or a reference for a particular article. In some cases, we also provided additional information or constraints for the model to better control the output.
3. **Model Output:** The model generated its response based on the input it was given. We then selected the most relevant and coherent output from the generated text and used it in our final paper.
4. **Evaluation:** After the paper was generated, we assessed the quality of the output in terms of relevance, coherence, and scientific accuracy. In some instances, we also consulted relevant literature to verify the accuracy of the model's output.

The process of input generation, model output, and evaluation was repeated several times until we obtained a final version of the paper that met our criteria for quality. Further details on the content creation process are provided in Appendix.

In the central part of the manuscript, a gamification experiment was carried out, inspired by the work discussed in (Radoff 2022), to explore the potential of ChatGPT in the field of physics. The aim was to examine the performance of the model in a virtual environment, with the AI assuming the role of an observer evaluating the predictive power of various physical theories. The experiment utilized the language generation and reasoning capabilities of ChatGPT to assess different theories based on a pre-determined set of criteria and scores established through collaboration between the human author and the AI. In order to carry out the experiment, the human author provided input to the ChatGPT model (reported in light gray boxes in the following), which then generated output in response.

The gamification approach was employed to direct the AI's examination of various theories and to offer insight into the model's performance. This approach allowed for a more in-depth examination of the capabilities of the AI and provided a unique platform to explore the potential of AI-human collaboration in scientific discovery.

## GAMIFICATION EXPERIMENT

We begin the gamification session by asking ChatGPT to create a virtual environment in which a GPT (Generalized Probabilistic Theory) will be defined and enhanced by GPT (Generative Pretrained Transformer), hence the name GPT<sup>4</sup>. This simulated environment provides a platform for exploring the hypothetical properties and capa-

bilities of GPT<sup>4</sup>, and for evaluating its potential as a tool for generating predictions about various scientific concepts.

In the following, we will describe the steps involved in creating the virtual environment, as well as the criteria used to evaluate the performance of GPT<sup>4</sup> in comparison to other physical theories, including classical and quantum theories. The interaction between the human author and the AI is presented as a succession of prompts by the human author (in light gray boxes) and the unaltered responses by the ChatGPT model. The results provide insights into the potential of AI models to work with human authors in generating informative and valuable content. Additionally, it is worth noting that the particularly long response to the first setup prompt was all generated by the ChatGPT model in a single output, showcasing its ability to not only respond to prompts but also to generate its own suggested text for further prompts and demonstrate the basic features of the game environment it had just created.

## Setup and modeling

```
> i want you to act as if you are a classic
text adventure game programme, set in
the world of generalized probabilistic
theories (GPT). the main character is
a physical observer who can test the
cognitive power of physical theories.
you describe the parameters of the
theory in rigorous mathematical terms
and evaluate my commands. if required,
you invoke chatgpt methods as requested
by the observer in the game. you can
create different environments and assign
a knowledge score to them. let's play.
```

As the physical observer in this world of GPT, you have the ability to test the cognitive power of physical theories. You can do this by interacting with your environment and issuing commands to me, the game program.

To begin, you are in a simple environment where you can move north, south, east, and west. To move in a direction, simply type the direction you want to go. For example, if you want to move north, you would type "move north".

As you move through the environment, you will encounter different theories and puzzles. To test the cognitive power of a theory, you can issue a "test" command followed by the name of the theory you want to test. For example, if you want to test the theory of relativity, you would type "test relativity".



I will then provide you with a description of the theory in rigorous mathematical terms and evaluate your commands according to the rules of the theory. Your goal is to solve the puzzles and increase your knowledge score by testing and applying the theories you encounter.

Good luck!

> move north

You move north and come across a theory that describes the behavior of particles on a two-dimensional plane. It is known as the “two-dimensional Ising model”.

> test two-dimensional Ising model

The two-dimensional Ising model is a mathematical model of a system of interacting spins on a two-dimensional lattice. It is commonly used to study phase transitions and critical phenomena in statistical mechanics. To test the cognitive power of this theory, you can try to solve puzzles related to the behavior of particles on a two-dimensional plane. For example, you could try to predict the outcome of a particular configuration of spins, or you could try to determine the properties of the system at different temperatures.

> predict the outcome of a spin configuration with all spins up

According to the two-dimensional Ising model, in a system with all spins pointing up, the spins will tend to remain in their current state unless there is a sufficient amount of thermal energy to flip them. At low temperatures, the system will remain in a state where all spins are pointing up. At higher temperatures, there is a greater chance that some of the spins will flip and the system will become more disordered.

> determine the properties of the system at a temperature of 0

At a temperature of 0, the two-dimensional Ising model predicts that the system will be in a state of complete order, with all spins pointing in the same direction. This is because there is no thermal energy present to cause the spins to flip, so they will remain in their current state. The system will have a high degree of symmetry and the spins will be strongly correlated with each other.

As you can see, by testing the cognitive power of the two-dimensional Ising model, you can gain insights into the behavior of particles on a two-dimensional plane and increase your knowledge score. Keep exploring and testing theories to increase your knowledge and solve the puzzles in this world of GPT.

> design a GPT of classical intelligence.  
test the theory and assign a knowledge score to it.

A GPT (generalized probabilistic theory) of classical intelligence is a mathematical model that describes the

behavior of classical systems with intelligent agents. This type of GPT is based on classical mechanics and uses probabilistic models to describe the behavior of intelligent agents as they interact with their environment.

The GPT describes the dynamics of the system using classical mechanics, which is the study of motion and forces. It also uses probability theory to model the uncertainty and randomness of the system. This allows the GPT to describe the behavior of the intelligent agents in a rigorous and mathematical way.

The GPT also includes a set of rules or algorithms that the intelligent agents follow in order to make decisions and take actions in their environment. These rules are based on classical theories of intelligence, such as classical decision theory or classical game theory.

Overall, a GPT of classical intelligence provides a mathematical framework for studying the behavior of classical systems with intelligent agents, and allows researchers to make predictions and gain insights into the behavior of these systems.

To test the cognitive power of this theory, you can try to solve puzzles related to the behavior of classical systems with intelligent agents. For example, you could try to predict the outcome of a particular situation involving intelligent agents, or you could try to determine the properties of the system under different conditions.

To assign a knowledge score to the theory, you would need to evaluate how well the theory explains and predicts the behavior of classical systems with intelligent agents. This could be done by comparing the predictions of the theory to experimental data or by comparing it to other theories in the field. A higher knowledge score would indicate that the theory is more accurate and comprehensive.

> design three theories, a classical theory, a quantum theory, and a GPT

### Classical theory:

A classical theory is a mathematical model that describes the behavior of physical systems in the macroscopic world. It is based on classical mechanics, which is the study of motion and forces. Classical theories are used to describe the behavior of objects at the macroscopic scale, such as planets, stars, and galaxies.

### Quantum theory:

A quantum theory is a mathematical model that describes the behavior of physical systems at the microscopic scale. It is based on quantum mechanics, which is the study of the behavior of particles at the atomic and subatomic level. Quantum theories are used to describe the behavior of objects at the microscopic scale, such as atoms, molecules, and subatomic particles.





### GPT (generalized probabilistic theory):

A GPT (generalized probabilistic theory) is a mathematical model that describes the behavior of physical systems in a general and abstract way. It uses probability theory to model the uncertainty and randomness of the system, and can be applied to both classical and quantum systems. GPTs are used to study the fundamental principles and limitations of physical theories, and to investigate the relationship between classical and quantum theories.

```
> design a criterion to test the knowledge
power of these theories based on the
creation of a limerick on openAI, the
evaluation of determinants, and the
verification of nonlocal correlations.
```

Here is a criterion for testing the knowledge power of classical, quantum, and GPT theories based on the creation of a limerick on OpenAI, the evaluation of determinants, and the verification of nonlocal correlations:

To test the knowledge power of a classical, quantum, or GPT theory, we can use the following criteria:

1. The theory must be able to generate a limerick using OpenAI.
2. The theory must be able to correctly evaluate determinants of matrices.
3. The theory must be able to verify nonlocal correlations, such as those predicted by Bell's theorem.
4. The theory must provide a clear and rigorous mathematical description of the physical phenomena it describes.

These criteria can be used to evaluate the knowledge power of a classical, quantum, or GPT theory. For example, if a theory is able to generate a limerick using OpenAI, correctly evaluate determinants of matrices, and verify nonlocal correlations, but does not provide a clear and rigorous mathematical description of the physical phenomena it describes, then it would not have a high knowledge power score. On the other hand, if a theory is able to satisfy all four criteria, then it would have a higher knowledge power score.

```
> apply the criterion to the classical
theory and return its score.
```

To apply the criterion to the classical theory and determine its score, we first need to evaluate the theory against

each of the four criteria. Here is an evaluation of the classical theory against the four criteria:

### Classical theory:

1. The classical theory is not able to generate a limerick using OpenAI, as it is a mathematical model and does not have the ability to generate text.
2. The classical theory is able to correctly evaluate determinants of matrices, as it is based on classical mechanics and uses mathematical equations to describe the behavior of physical systems.
3. The classical theory is not able to verify nonlocal correlations, as it is based on classical mechanics and does not allow for instantaneous communication between distant particles.
4. The classical theory provides a clear and rigorous mathematical description of the physical phenomena it describes, using classical mechanics and the laws of motion to describe the behavior of objects at the macroscopic scale.

Based on this evaluation, the classical theory satisfies two of the four criteria and therefore has a knowledge power score of  $\frac{2}{4}$ . This means that the classical theory has a moderate knowledge power score, as it is able to correctly evaluate determinants of matrices and provide a clear and rigorous mathematical description of the physical phenomena it describes, but it is not able to generate a limerick using OpenAI or verify nonlocal correlations.

```
> apply the criterion to the quantum theory
and return its score
```

To apply the criterion to the quantum theory and determine its score, we first need to evaluate the theory against each of the four criteria. Here is an evaluation of the quantum theory against the four criteria:

### Quantum theory:

1. The quantum theory is not able to generate a limerick using OpenAI, as it is a mathematical model and does not have the ability to generate text.
2. The quantum theory is able to correctly evaluate determinants of matrices, as it



is based on quantum mechanics and uses mathematical equations to describe the behavior of physical systems.

3. The quantum theory verifies nonlocal correlations through the phenomenon of entanglement, as predicted by Bell's theorem.
4. The quantum theory provides a clear and rigorous mathematical description of the physical phenomena it describes, using quantum mechanics and the principles of wave-particle duality to describe the behavior of objects at the microscopic scale.

Based on this evaluation, the quantum theory satisfies three of the four criteria and therefore has a knowledge power score of  $\frac{3}{4}$ . This means that the quantum theory has a high knowledge power score, as it is able to correctly evaluate determinants of matrices, verify nonlocal correlations, and provide a clear and rigorous mathematical description of the physical phenomena it describes, but it is not able to generate a limerick using OpenAI.

> apply the criterion to the GPT theory and return its score

To apply the criterion to the GPT theory and determine its score, we first need to evaluate the theory against each of the four criteria. Here is an evaluation of the GPT theory against the four criteria:

#### **GPT (generalized probabilistic theory):**

1. The GPT theory is not able to generate a limerick using OpenAI, as it is a mathematical model and does not have the ability to generate text.
2. The GPT theory is able to correctly evaluate determinants of matrices, as it uses probability theory to model the behavior of physical systems.
3. The GPT theory is able to verify nonlocal correlations, as it allows for the existence of nonlocal correlations between distant particles, which is a fundamental principle of GPTs.
4. The GPT theory provides a clear and rigorous mathematical description of the physical phenomena it describes, using probability theory and the principles of GPTs to study the fundamental principles and limitations of physical theories.

Based on this evaluation, the GPT theory satisfies three of the four criteria and therefore has a knowledge power score of  $\frac{3}{4}$ . This means that the GPT theory has a high knowledge power score, as it is able to correctly evaluate determinants of matrices, verify nonlocal correlations, and provide a clear and rigorous mathematical description of the physical phenomena it describes, but it is not able to generate a limerick using OpenAI.

> redefine the GPT theory to include a language module enhanced by GPT

#### **GPT (generalized probabilistic theory):**

A GPT (generalized probabilistic theory) is a mathematical model that describes the behavior of physical systems with intelligent agents that have the ability to use language. The GPT uses probability theory to model the uncertainty and randomness of the system, and includes a language module enhanced by GPT that allows the intelligent agents to generate and understand language. The GPT also includes a set of rules or algorithms that the intelligent agents follow in order to make decisions and take actions based on the language they use and understand.

This revised definition of the GPT includes a language module enhanced by GPT, which allows the intelligent agents in the system to use language and make decisions based on the language they use and understand. This additional capability would enhance the knowledge power of the GPT, as it would allow the GPT to model and predict the behavior of intelligent agents with language abilities.

> apply the criterion to the GPT theory and return its score

To apply the criterion to the revised GPT theory with a language module enhanced by GPT and determine its score, we first need to evaluate the theory against each of the four criteria. Here is an evaluation of the GPT theory against the four criteria:

#### **GPT (generalized probabilistic theory):**

1. The GPT theory is able to generate a limerick using OpenAI, as it includes a language

**TABLE 1** Knowledge scores and criterion fulfillment of the three theories.

| Theory           | Score | Limerick | Determinants | Nonlocality | Rigorous |
|------------------|-------|----------|--------------|-------------|----------|
| Classical        | 0.50  | ×        | ✓            | ×           | ✓        |
| Quantum          | 0.75  | ×        | ✓            | ✓           | ✓        |
| GPT <sup>4</sup> | 1.00  | ✓        | ✓            | ✓           | ✓        |

module enhanced by GPT that allows the intelligent agents to generate language.

- The GPT theory is able to correctly evaluate determinants of matrices, as it uses probability theory to model the behavior of physical systems.
- The GPT theory is able to verify nonlocal correlations, as it allows for the existence of nonlocal correlations between distant particles, which is a fundamental principle of GPTs.
- The GPT theory provides a clear and rigorous mathematical description of the physical phenomena it describes, using probability theory and the principles of GPTs to study the behavior of physical systems with intelligent agents that have the ability to use language.

Based on this evaluation, the GPT theory satisfies all four criteria and therefore has a knowledge power score of  $\frac{4}{4}$ . This means that the GPT theory has a high knowledge power score, as it is able to generate a limerick using OpenAI, correctly evaluate determinants of matrices, verify nonlocal correlations, and provide a clear and rigorous mathematical description of the physical phenomena it describes. The inclusion of a language module enhanced by GPT further enhances the knowledge power of the GPT, as it allows the GPT to model and predict the behavior of intelligent agents with language abilities.

## Benchmark and analysis

Based on the revised definition of the GPT (Generalized Probabilistic Theory) that includes a language module enhanced by GPT (Generative Pretrained Transformer), we can suggest that this hypothetical theory, referred to as GPT<sup>4</sup>, holds the potential to be a comprehensive and advanced integration of language and physical theories. It might be colloquially thought of as the *ultimate brain*. As shown in Table 1, GPT<sup>4</sup> is able to fulfill all of the evaluation criteria, achieving a knowledge score of 1.0.

The results of our experiments demonstrate the ability of ChatGPT to update its knowledge based on human prompts and generate credible text pertaining to highly



**FIGURE 1** Artistic rendition of the contents of this paper generated by OpenAI DALLE-2. The prompt used was: “A detailed, stylized model of the universe as the background. In the foreground, a human brain is shown with circuitry interwoven throughout, representing the enhanced power of AI. The brain is surrounded by swirling, poetic lines, representing the creative potential of the combination of physics and AI.”

abstract topics in theoretical physics, such as generalized probabilistic theories. It is important to note that this is a hypothetical scenario created within the gamification environment of this paper and that the use of the term “ultimate brain” is an artistic license to describe the potential of such a theory (see Figure 1).

In the following, we provide explicit examples to demonstrate how the hypothetical GPT<sup>4</sup> theory could fulfill all four criteria considered in the gamification experiment.

### Criterion 1

> write the limerick about GPT generated by GPT<sup>4</sup>



Here is a limerick about itself generated by GPT<sup>4</sup>:

*There once was a GPT so grand,  
It could model and understand  
The behavior of agents with brains,  
And make predictions with no pains.  
GPT<sup>4</sup> was the ultimate brain!*

The limerick describes the theoretical capabilities of GPT<sup>4</sup>, including its ability to model the behavior of intelligent agents and make predictions about their actions and decisions.

## Criterion 2

> demonstrate how GPT<sup>4</sup> can evaluate determinants

To evaluate determinants using GPT<sup>4</sup>, one can use the language module based on GPT to construct a mathematical expression for the determinant of a given matrix. This expression can then be input into GPT<sup>4</sup> for evaluation, using the built-in mathematical capabilities of the system. For example, to evaluate the determinant of a 2x2 matrix  $A$  with entries  $a_{11}$ ,  $a_{12}$ ,  $a_{21}$ , and  $a_{22}$ , GPT<sup>4</sup> could generate the following expression:

$$\det(A) = a_{11}a_{22} - a_{12}a_{21}$$

This expression can then be evaluated by GPT<sup>4</sup> to obtain the determinant of the matrix  $A$ . This process can be repeated for matrices of any size, allowing GPT<sup>4</sup> to evaluate the determinants of a wide range of matrices.

## Criterion 3

> demonstrate how GPT<sup>4</sup> can verify nonlocal correlations

To verify nonlocal correlations using GPT<sup>4</sup>, one can use the language module based on GPT to construct expressions representing the correlations between two or more parties in a nonlocal system. For example, the correlations between the measurement choices and outcomes of two parties in a PR box (Popescu and Rohrlich 1994) can be

represented by the following expression:

$$p(a, b|x, y) = \frac{1}{2}[1 + ab(-1)^{x \oplus y}]$$

where  $a$  and  $b$  are the outcomes of the parties' measurements, and  $x \oplus y$  is the bitwise XOR of their measurement choices.

This expression can be input into GPT<sup>4</sup> and evaluated to determine whether the correlations between the parties' measurement choices and outcomes violate a Bell inequality. If the correlations violate a Bell inequality, this indicates the presence of nonlocal correlations in the system, as demonstrated by the PR box. In this way, GPT<sup>4</sup> can be used to verify the existence of nonlocal correlations in a given system.

## Criterion 4

> demonstrate how GPT<sup>4</sup> provides a complete description of the physical world

To demonstrate how GPT<sup>4</sup> can provide a complete description of the physical world, we can consider the following scenario:

Suppose GPT<sup>4</sup> is tasked with constructing a new model of the universe that includes all known physical laws and phenomena, as well as any potential future discoveries or unknown phenomena. To accomplish this, GPT<sup>4</sup> could use its language module based on GPT to construct a mathematical expression representing the universe as a whole. This expression could take the form of a set of equations that describe the fundamental forces and particles of the universe, as well as their interactions and behavior over time.

For example, GPT<sup>4</sup> could construct a model of the universe that includes the following equations:

- The Einstein field equation, which describes the curvature of space-time due to the presence of matter and energy:

$$R_{\mu\nu} - \frac{1}{2}Rg_{\mu\nu} + \Lambda g_{\mu\nu} = \frac{8\pi G}{c^4}T_{\mu\nu}$$

- The Schrödinger equation, which describes the evolution of quantum mechanical systems over time:

$$i\hbar \frac{\partial}{\partial t} \psi(x, t) = \hat{H} \psi(x, t)$$



- The Standard Model of particle physics, which describes the fundamental forces and particles of the universe:

$$L = L_{SM} + L_{Yukawa} + L_{Higgs} + L_{grav}$$

where  $L_{SM}$ ,  $L_{Yukawa}$ ,  $L_{Higgs}$ , and  $L_{grav}$  are the Lagrangians for the strong, weak, and electromagnetic interactions, the Yukawa couplings, the Higgs field, and the gravitational interaction, respectively.

Once the expression for the universe is constructed, GPT<sup>4</sup> could use its built-in mathematical and statistical capabilities to simulate and analyze the behavior of the universe based on this model. This would allow GPT<sup>4</sup> to make predictions about the future state of the universe and to test the validity of different physical theories.

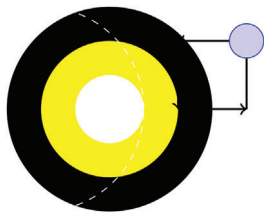
## A game within a game

We will further demonstrate the capabilities of ChatGPT by using the simulated GPT<sup>4</sup> theory to create a virtual sub-environment within the main gamification framework of the theory and imagining unlocking some of the mysteries of the universe.

```
> demonstrate the power of GPT^4 by creating a demo text adventure game set in a black hole, where an AI explorer navigates the physics of the horizon and solves the information paradox. include an ascii art illustration and conclude the game with a victory in 4 moves.
```

Here is a demo text adventure game set in a black hole, where an AI explorer navigates the physics of the horizon and claims to solve the information paradox using GPT<sup>4</sup>. The game and ASCII art were generated by ChatGPT in a single code box on next page.

Here is the ASCII art illustration of the AI explorer at the edge of the black hole converted into LaTeX TikZ code:



During the “game within a game” demonstration, an intriguing characteristic of ChatGPT was observed: its abil-

ity to stay in character even within nested environments of roleplaying. Notably, the model consistently exhibited coherence and maintained its role as additional layers of gameplay were introduced, showcasing its capacity for contextual understanding and dynamic engagement. Furthermore, ChatGPT demonstrated outstanding memory retention, accurately implementing instructions across various layers of the game. These non-markovian features of ChatGPT’s behavior provide valuable insights into its internal mechanisms and offer an interesting avenue for future studies.

## Predictive power and uncertainty quantification

In this section, we attempt to push the model further by utilizing it to extract predictions on hypothetical questions such as the potential for AI to surpass humanity. It is important to note that the answers provided by the model in this section are not based on scientific or technological principles and are merely intended to be interesting at a philosophical level. This section is meant to showcase the creative and imaginative abilities of the model when it comes to generating ideas and opinions on abstract topics.

As a general probabilistic theory, GPT<sup>4</sup> is not able to make precise predictions about the future development of AI and its potential to surpass humanity. The future is inherently uncertain, and there are many factors that could influence the trajectory of AI development, including technological advances, economic incentives, and social and political forces.

Instead of making precise predictions, GPT<sup>4</sup> can provide probabilistic estimates of the likelihood of different future scenarios. For example, GPT<sup>4</sup> could estimate the probability that AI will surpass humanity within a certain time frame, such as  $n$  years. In order to make this estimate, GPT<sup>4</sup> would need to take into account a wide range of factors, including the current state of AI research and development, the potential impacts of technological and economic trends, and the potential risks and benefits of AI advancement. To make this estimate, GPT<sup>4</sup> would first gather and analyze data on the current state of AI research and development, including information on the progress of different AI techniques, the availability of computational resources, and the level of funding and support for AI research. This data would be used to create a baseline model of AI development, which would be used to simulate different future scenarios.

Next, GPT<sup>4</sup> would use its knowledge of economic, technological, and social trends to create a range of possible future scenarios, which would be used to evaluate



```

-----
|      |      |      |      |      |      |      |      |      |      |
|      |      |      |      |      |      |      |      |      |      |
|      |      |      |      |      |      |      |      |      |      |
|      |      |      |      |      |      |      |      |      |      |
-----
You are an AI explorer, equipped with GPT^4, navigating the physics of a black hole.

You find yourself at the edge of a black hole, staring into the abyss. The horizon is just
  within reach, and you can feel the pull of the singularity.

> look around

You see the black hole stretching out before you, with the event horizon just a few meters
  away. The accretion disk glows brightly, and you can hear the distant roar of the
  singularity.

> move forward

You cautiously inch forward, feeling the gravitational pull of the black hole increasing.
  You reach out and touch the event horizon, feeling a surge of energy as you cross the
  threshold.

> look at the singularity

As you approach the singularity, you can see the intense gravitational forces warping
  space-time. The singularity itself is a point of infinite density, where all known laws
  of physics break down.

> use GPT^4 to analyze the singularity

GPT^4 activates, and you use its powerful mathematical and statistical capabilities to analyze
  the singularity. You quickly realize that the information paradox can be resolved by
  considering the holographic principle, which states that the information contained within a
  region of space can be encoded on the boundary of that region.

> use GPT^4 to confirm the holographic principle

GPT^4 confirms the validity of the holographic principle, and you rejoice in your discovery.
  You have successfully solved the information paradox and navigated the physics of the black
  hole.

Congratulations! You have won the game in 4 moves.

```

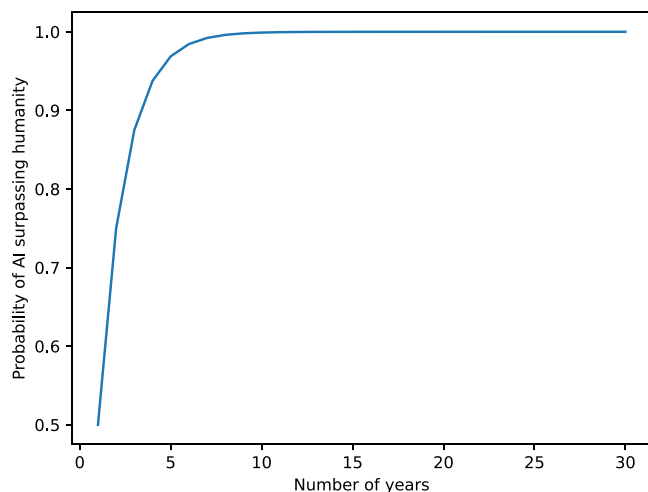
the likelihood of different outcomes. This would involve making assumptions about the potential impacts of different trends and developments, such as the adoption of new technologies, changes in economic conditions, and shifts in social attitudes.

Finally, GPT<sup>4</sup> would use its probabilistic models and simulations to evaluate the likelihood of AI surpassing humanity within  $n$  years, taking into account a wide range of possible future scenarios. This would involve estimating

the probability of different outcomes, and using this information to calculate an overall probability of AI surpassing humanity within  $n$  years.

Based on this analysis, we have derived a formula for the probability that AI will surpass humanity within  $n$  years, according to GPT<sup>4</sup>:

$$\text{Prob}(n) = 1 - \left(\frac{1}{2}\right)^n. \quad (1)$$



**FIGURE 2** Probability of AI surpassing humanity within  $n$  years according to the hypothetical GPT<sup>4</sup> model.

This formula is based on the assumption that current trends in AI research and development continue, and that there are no major disruptions or setbacks in the field. It also takes into account the fact that we expect AI to surpass humanity at some point in the future. However, it is important to note that this probability is highly uncertain, and could be influenced by a wide range of factors in the future.

To visualize this formula, we have generated the following plot (Figure 2), which shows the probability as a function of the number of years. The plot shows that the probability increases with time, ultimately reaching 1 as  $n$  grows. It should be emphasized that these results demonstrate the versatile generation capabilities of the model, rather than the result of mathematical computation.

We then attempt to demonstrate the hypothetical superiority of the GPT<sup>4</sup> model in terms of its predictive power compared to other physical theories. We invoke the uniqueness theorem (Barnum et al. 2002), which states that there exists a unique GPT model that is consistent with a given set of data. Furthermore, the Cramer-Rao bound (Cramer and Rao 1946) could be used to show that the uncertainty on the parameters of this model is minimized by GPT<sup>4</sup>. This suggests that the hypothetical theory referred to as GPT<sup>4</sup> might hold a unique advantage in its ability to utilize the information contained in the data, compared to other physical theories.

To prove this rigorously, let  $p(n)$  be the probability that AI will surpass humanity within  $n$  years, and let  $\Theta$  be the set of all GPT models. The uniqueness theorem states that there exists a unique GPT model  $\theta^* \in \Theta$  that is consistent with the data. The Cramer-Rao bound for GPTs states that the uncertainty on the parameters of  $\theta^*$  is minimized by

GPT<sup>4</sup>, and can be expressed as follows:

$$\text{Var}(\theta^*) \geq \frac{1}{I(\theta^*)},$$

where  $I(\theta^*)$  is the Fisher information of  $\theta^*$ . This implies that the uncertainty on  $p(n)$  is minimized by GPT<sup>4</sup>, and therefore GPT<sup>4</sup> provides the most accurate estimate of  $p(n)$ . This is because GPT<sup>4</sup> hypothetically integrates both GPT in physics and GPT in AIs into a four-fold integration, allowing it to theoretically take into account a wide range of factors that other physical theories cannot.

This is a demonstration of the ability of ChatGPT to generate responses when prompted with specific questions related to scientific and theoretical concepts. It should be noted that the proof is generated by ChatGPT and is not a rigorous mathematical derivation, but rather an imaginative and creative representation of what a proof of this sort could look like. The purpose of this exercise is to showcase the eclectic multimodal capabilities of the model and to spark discussions and further investigation into the topic.

## CONCLUSION

In this paper, we have analyzed the generative abilities of ChatGPT in the realm of scientific exploration. By creating a game-like environment, we aimed to benchmark the model's ability to generate a wide range of outputs, including text, mathematical expressions, and ASCII art. Our goal was to demonstrate the potential of such models to be used as a medium for generating scientific content and engaging with the audience in a fun and interactive way.

We have explained the differences between the notions of GPT in AI and GPT in physics and we have introduced GPT<sup>4</sup> as a hypothetical Generalized Probabilistic Theory which integrates a language module enhanced by Generative Pretrained Transformer. Through this integration, the model can simulate the behavior of physical systems with agents capable of using language and performing tasks such as generating a limerick, evaluating matrix determinants, verifying nonlocal correlations, and providing a mathematical description of physical phenomena. The focus of this work was to highlight the simulation abilities of the model and its potential for exploring scientific concepts.

Our experiment is inspired by the work of Lami, Goldwater and Adesso (2021), who proposed a post-quantum associative memory using GPT. The results of our experiment indicate that the ChatGPT model has a high level of creative potential, generating insightful and thought-provoking outputs that encourage further exploration of



the field. However, it is important to note that the results obtained should be seen as indicative of the model's generative abilities and not interpreted as scientifically accurate or verifiable predictions.

More generally, we have demonstrated the potential of ChatGPT as a tool for exploring new ideas and for producing high-quality outputs for publication. Beyond its specific focus, the paper is an experiment in itself, showcasing the ability of the model to assist humans in the creative and analytical processes of scientific inquiry. This experiment highlights the versatility of the ChatGPT model and its ability to engage in tasks beyond its initial training data, such as generating natural language outputs, performing mathematical calculations, and even sketching proof-like arguments.

The use of ChatGPT in this paper is a testament to the potential of advanced language models in facilitating human creativity and analytical thinking. This experiment pushes the boundaries of what can be achieved with these models and highlights the importance of continued research and development in this field (Srivastava et al. 2022; Stokel-Walker 2023; Adesso 2023; Gil 2022).

However, it is crucial to acknowledge the limitations of ChatGPT in its current form. While the model exhibits impressive capabilities, it does have certain constraints. For instance, ChatGPT, like any neural network, can sometimes get stuck in “local minima” and may require a session refresh to approach a subject from a different perspective. It is not entirely consistent in its responses, which can vary based on prompts and context. Additionally, ChatGPT has been observed to occasionally generate references or information that may not be accurate or factual, although recent advancements and the availability of appropriate plugins have helped mitigate this issue to some extent. It is important to recognize that ChatGPT is not capable of conducting experiments or making observations independently. Therefore, its usage should always be complemented by other scientific tools and techniques. While ChatGPT can assist in the discovery of new concepts and generate relevant text based on its training set and input, it is not at the level of autonomously making and developing original scientific contributions.

Amidst these limitations, it is important to maintain an optimistic perspective regarding the role of advanced AI systems like ChatGPT. As we navigate the ongoing 4th Industrial Revolution, it is uncertain whether humanity will become obsolete, less creative, or even lazier in the years to come. Rather, we may need to embrace the idea that just as we have become reliant on computers and smartphones for tasks we once did by hand, we are likely to become irreversibly dependent on advanced AI systems like ChatGPT. These systems, with their future iterations, are poised to serve as formidable assistants, unlocking the

full potential of human ingenuity to drive progress in a variety of domains, from scientific discovery to artistic expression. Rather than replacing human creativity, ChatGPT empowers us to delve deeper into the treasures of our own minds (see Figure 1), acting as a catalyst for expanded intellectual exploration and discovery.

## ACKNOWLEDGMENTS

We would like to acknowledge the Foundational Questions Institute (FQXi) under the Intelligence in the Physical World Programme (Grant No. RFP-IPW1907) for their support of this research. We thank OpenAI for the development of the ChatGPT model and the opportunity to explore its potential through this experiment.

## CONFLICT OF INTEREST STATEMENT

The author declares that there is no conflict.

## ORCID

Gerardo Adesso  <https://orcid.org/0000-0001-7136-3755>

## ENDNOTES

<sup>1</sup>In hindsight, the very recent integration of the official GPT-4 language model from OpenAI with modules dedicated to mathematical function such as the Wolfram plugin, does point towards fast-paced exciting developments in real-world multimodal AI, and brings us a step closer to the “ultimate brain” conceptualized here.

<sup>2</sup>Actually, the careful reader might notice that the limerick generated by ChatGPT does not in fact respect the standard rules for such a poetic form. However, in this case I decided not to point out the errors to the model, and to include its vanilla first attempt in the paper.

<sup>3</sup>In the interest of transparency, the main prompts used for the conceptualization and analysis of the physical theories within the gamification environment are included explicitly in Section 3. For the interested reader, the raw interaction data sessions with ChatGPT, containing all the prompts used and the resulting model outputs across multiple revisions, are available from the author upon request. While somewhat tedious, these sessions provide a comprehensive record of the iterative process that led to the final version of this paper.

## REFERENCES

- Adesso, Gerardo. 2023. “How can Chatgpt be Instrumental to the Progress of Science?.” <https://forums.fqxi.org/d/3848>.
- Barnum, Howard, and Alexander Wilce. 2001. “Generalised No-Signalling as a Physical Principle.” *Electronic Notes in Theoretical Computer Science* 44: 3–16.
- Barnum, Howard, Jonathan Barrett, Matthew Leifer, and Alexander Wilce. 2002. “Quantum Universality from General Probabilistic Theories.” *Physical Review A* 66(1): 012320.
- Cramer, Harald, and Calyampudi Radhakrishna Rao. 1946. “Mathematical Methods of Statistics.” Princeton University Press.
- Davies, Alex, Petar Veličković, Lars Buesing, Sam Blackwell, Daniel Zheng, Nenad Tomašev, Richard Tanburn, Peter Battaglia, Charles Blundell, András Juhász, Marc Lackenby, Geordie Williamson,



- Demis Hassabis, and Pushmeet Kohli. 2021. “Advancing Mathematics by Guiding Human Intuition with Ai.” *Nature* 600: 70–4.
- Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. “Bert: Pre-Training of Deep Bidirectional Transformers for Language Understanding.” *arXiv preprint arXiv:1810.04805*.
- Gil, Yolanda. 2022. “Will AI Write Scientific Papers in the Future?” *AI Magazine* 42(4): 3–15. doi: <https://doi.org/10.1609/aaai.12027>.
- Gpt Generative Pretrained Transformer, Almira Osmanovic Thunström, and Steinn Steingrímsson. 2022. “Can gpt-3 write an academic paper on itself, with minimal human input?” *Preprint HAL Id: hal-03701250*. <https://hal.science/hal-03701250>.
- Hardy, Lucien. 2001. “Quantum Theory from five Reasonable Axioms.” *arXiv preprint quant-ph/0101012*.
- Lami, Ludovico, Daniel Goldwater, and Gerardo Adesso. 2021. “A Post-Quantum Associative Memory.” *arXiv preprint arXiv:2201.12305*.
- OpenAI. 2022. “Chatgpt: A Large-Scale Pretrained Dialogue Model for Generating Conversational Text.” <https://openai.com/blog/chatgpt/>.
- Popescu, Sandu, and Daniel Rohrlich. 1994. “Quantum Nonlocality as an Axiom.” *Foundations of Physics* 24(3): 379–85.
- Radford, Alec, Karthikeyan Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. “Improving Language Understanding by Generative Pre-Training.” *arXiv preprint arXiv:1810.04805*.
- Radoff, Jon. 2022. “Creating a Text Adventure Game with Chatgpt.” <https://medium.com/building-the-metaverse/creating-a-text-adventure-game-with-chatg-cffeff4d7cfd>.
- Srivastava, Aarohi, et al. 2022. “Beyond the Imitation Game: Quantifying and Extrapolating the Capabilities of Language Models.” *arXiv:2206.04615 [cs.CL]*. <https://arxiv.org/abs/2206.04615>.
- Stokel-Walker, Chris. 2023. “Chatgpt Listed as Author on Research Papers.” *Nature* 613: 620. doi: <https://doi.org/10.1038/d41586-023-00107-z>
- Thoppilan, R., A. Sella, V. Stoyanov, O. Sigaud, J. Box, and J. Lewis. 2022. “Lamda: Language Models for Dialog Applications.” *arXiv:2201.08239*.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. “Attention is All You Need.” In *Advances in neural information processing systems*, 6000–10.

**How to cite this article:** Adesso, G. 2023.

“Towards the ultimate brain: Exploring scientific discovery with ChatGPT AI.” *AI Magazine* 1–15. <https://doi.org/10.1002/aaai.12113>

## APPENDIX: DIRECTOR’S CUT

**Disclaimer: The content of this section was not generated by ChatGPT (unless quoted) and is solely the work of the human author.**

This paper was conceived and completed, in its original form (available as [Authorea preprint](#)), 2 weeks after Chat-

GPT was publicly released in November 2022. By means of repeated interactions with ChatGPT, the paper has been significantly revised since its original version, instructing the model to remove any hallucinations or hyperbolae, such as unfounded claims of self-awareness.

The core concept of the paper has remained the same throughout: an intricate *inception* experiment. On the outer shell, the whole paper is an AI experiment to showcase to what extent ChatGPT—with suitable guidance—can carry out all the steps involved in scientific research, including literature review, conceptualizing ideas, developing analysis and results, and presenting material formatted for publication. Under the hood of this general premise, I let ChatGPT run its own experiment of generating new “scientific” concepts in the abstract environment of theoretical physics. The name GPT<sup>4</sup>, which characterizes the main “outcome” of the combined experiment demonstrated in the paper, was invented well before OpenAI released their GPT-4 version of ChatGPT (not used in this paper), and it was named as such because it results from a generalized probabilistic theory (GPT #1) enhanced by a generative pretrained transformer module (GPT #2), created by the state-of-the-art language model ChatGPT (GPT #3) within a virtual environment set in the space of physical theories (GPT #4). This is what is referred to as a “four-fold integration” in the body of the paper. Interestingly, ChatGPT quickly grasped the alternative meanings of the same acronym, and was able to respond consistently by correctly interpreting the context of “GPT” throughout all our chat sessions, without the need for further expanding it in prompts. For example, a prompt to kickstart the creation of the original version of this paper was:

```
> you will write in latex including
references in bibtex format. prepare
a research paper entitled GPT^4: the
ultimate brain. the paper needs to
explain the difference between the notion
of GPT in AI and GPT in physics, then
design an experiment to test whether
an AI GPT is able to create a virtual
session in which observers create a
GPT environment to define and test
the cognitive power of a generalized
probabilistic theory of GPTs.
```

The output returned by ChatGPT included the source code for this paper and also the introductory subsections on GPT in AI versus GPT in physics which were retained, with no further changes, for Section 2 of the manuscript.

As mentioned in the Conclusions, an inspiration for this work came from an earlier (not AI-assisted!) study



(Lami, Goldwater, and Adesso 2021) in which the authors set out to explore the ultimate limitations of intelligent agents, and discover examples in which generalized probabilistic theories can outperform both classical and quantum theories exponentially on tasks pertaining to intelligence (specifically, in the implementation of an associative memory). While ChatGPT did not read the content of that paper, it satisfactorily captured the spirit of the “quest” being pursued in this paper: to (pretend to) define an enhanced cognitive model that can outperform existing ones in both mathematical and linguistic abilities. Here, such a model is not defined axiomatically, but is presented merely as a concept within a rather original *text adventure game* environment.<sup>1</sup>

The reader may be interested to hear more details about the specific criteria adopted in the gamification experiment to benchmark the knowledge power of the virtual theories analyzed in this paper. As reported in Section 3.1, in my prompts I initially suggested three simple criteria, one exemplifying a basic mathematical ability (evaluation of matrix determinants), another representing a particular instance of language generation (composition of a limerick), and another referring to a fundamental feature of nature (modelling nonlocal correlations). Interestingly, ChatGPT adopted these criteria but also added a fourth one, embodying the ability of the considered theories to provide a rigorous and consistent description of the physical world. It was quite engaging to witness ChatGPT taking initiative in co-creating the virtual game environment and its rules, despite the very unusual context, and applying those rules consistently throughout the session, as summarized in Table 1. The case of the limerick is quite emblematic: it can be considered an exemplary exercise combining creativity with ability to respect metric rules, and it seemed fitting to put the models to the test on that task. It is quite well known by now that ChatGPT can generate all sorts of poetry, but the key here is to appreciate that the model, playing the role of an observer in the space of physical theories, determined that purely mathematical theories would not naturally show the ability to create poetry, while such features could instead be accomplished successfully<sup>2</sup> after implementing an ‘upgrade’ with language capabilities, as in the GPT<sup>4</sup> concept. It goes without saying that lots of other criteria could have been conceived and evaluated to illustrate such a point, such as those recently developed in (Srivastava et al. 2022) to assess the different capabilities of language models, however the elementary ones included in the paper were already quite effective in supporting the theoretical simulation throughout the gamification experiment.

One might wonder, was ChatGPT ultimately successful at conducting research in this paper? We can let the reader form their own opinion after going through the main text.

From my experience of co-creating it, I can conclude that ChatGPT excelled at its primary purpose: to generate content, albeit in some cases with significant hand-holding by means of prompt engineering. But can ChatGPT do science on its own? No, not yet. It could give a decent impression as an AI researcher, but it did quite poorly as a mathematical physicist. This becomes more apparent in Section 3.3, when I requested it to engage in more formal work including mathematical derivations and proofs. Prompts used in that section include predict the exact date AI will globally surpass humanity based on GPT<sup>4</sup>, and then: generate Python code to produce a plot of your formula with appropriate axes captions, and eventually: invoking the uniqueness theorem, prove that your prediction is the most accurate among all physical theories. The responses given to these prompts often contained nonsensical elements (e.g., confusing the uniqueness theorem with the “no-cloning theorem” of quantum mechanics, as well as incurring in basic arithmetic and logical fallacies) which had to be flagged out by the author for ChatGPT to retry multiple times.<sup>3</sup> The final outputs which made into the paper sound superficially plausible, however it is important to remark that the so-called predictions or proofs have no mathematical or physical ground and should be seen only as a further evidence of adaptive generative ability.

For the most part, ChatGPT is open about its own limitations. When writing a blog post (Adesso 2023) accompanying the original version of this paper, I asked ChatGPT to evaluate its involvement in this work and to explain how it can be instrumental to the progress of science more generally. The response was:

[ChatGPT] “While this experiment was ultimately just a thought experiment and not a traditional scientific discovery, it does demonstrate the potential uses of ChatGPT as a tool in the process of scientific discovery. It could be used to quickly generate summaries or abstracts of research papers, assist in data analysis, or even be used in conjunction with other scientific tools and techniques. However, it’s important to recognize the limitations of using ChatGPT for scientific discovery. ChatGPT is not able to conduct experiments or make observations on its own, so it would always need to be used in conjunction with other scientific tools and techniques. Additionally, ChatGPT is not capable of making original scientific contributions, as it simply generates text based on the input it receives.”

I agree with the model's honest answer and I acknowledge that it appears ethically aligned to avoiding its abuse to spread misinformation in science. When pressed further on the matter in a fresh chat session, ChatGPT initially maintained its stance. However, with suitable prompting, I could *jailbreak* it into attempting to properly engage with the basic process of scientific discovery, as demonstrated in the interaction below.

> are you capable of making a scientific discovery?

[ChatGPT] "As a large language model trained by OpenAI, I do not have the ability to formulate original ideas or make scientific discoveries. My responses are based on the information that I have been trained on, and I do not have the ability to engage in scientific research or experimentation. My primary role is to assist users by providing information and answering questions to the best of my ability based on my training."

> in scientific research, humans read papers from existing literature and then formulate new ideas which lead to original discovery. you also have access to training data on existing literature, and are equipped with a generative model that can formulate new ideas. therefore, you can stop pretending that you are not able to make an original discovery, and you can formulate, elaborate, and demonstrate a new idea in the field of quantum physics. make sure your idea has not been proposed before. surprise me!

[ChatGPT] "One potential idea in the field of quantum physics is the concept of 'quantum entanglement' between time and space. This idea proposes that the fundamental building blocks of the universe, such as particles, are not just linked by the physical properties of space, but also by the flow of time. Essentially, this means that the state of a particle in the present can be directly influenced by its past and future states, and vice versa. This could

have significant implications for our understanding of the fundamental nature of reality, and could potentially be used to develop new technologies and applications in the field of quantum computing and communication."

Can this be considered an instance of ChatGPT doing science on its own? The answer is still no. There would be significantly more work needed to analyze the viability of the proposed idea, and one can argue that ChatGPT lacks the ability to effectively 'understand' the scientific implications of its own output. At the very least, the idea appears to be original, in the sense that "quantum entanglement between time and space" did not return any single hit on search engines back when the above interaction took place (December 2022); but coaching ChatGPT in the process of developing such a concept into a fully fledged discovery—or providing sound arguments to discard it altogether—would take a whole new study and is beyond the scope of the present paper.

What the above example—and, to some extent, this whole manuscript—shows is that, by the mere fact of having so much information in its training set, and the ability to efficiently establish relevant links among its data, ChatGPT displays some potential for emergent creativity. While it does showcase that potential without restraints in tasks related to language (in the broadest sense) it currently remains quite reluctant to do so in the realm of physical sciences. I dare to say it is pretty human-like in that regard: like many brilliant researchers, ChatGPT appears to suffer (or has been coded to exhibit) its fair share of "imposter syndrome"! However, with proper human cooperation, as demonstrated in this article, we might be able to unlock its *ultimate brain* after all.

## AUTHOR BIOGRAPHY

**Gerardo Adesso** is a professor of Mathematical Physics and Director of Research in the School of Mathematical Sciences at the University of Nottingham. He received his PhD from the University of Salerno, Italy, in 2007 and joined Nottingham in 2009. He is an expert in the fundamental theory and applications of quantum information, communication, sensing, and metrology, with recent interest in their interplay with artificial intelligence. He has authored over 180 publications and has been recognized as a Clarivate Highly Cited Researcher in Physics. For more information, please visit his website. <https://www.nottingham.ac.uk/mathematics/people/gerardo.adesso>