

# ChatGPT vs Gemini

Gemini versus ChatGPT: applications, performance, architecture, capabilities, and implementation

Feature	Gemini (Google AI)	ChatGPT (OpenAI)
Model Type	LLM, search integration	LLM, focus on dialogue
Knowledge Cutoff	Access to more up-to-date information	Knowledge cutoff around April 2023
Factual Accuracy	Prioritizes accuracy, sourcing	Accuracy can vary, especially post-cutoff
Creativity	Less emphasis on creative output	Excels in creative text formats
Task Completion	Strong with search-supported tasks	More conversational, open-ended tasks
Code Generation	Capable	Capable
Bias and Safety	Efforts to mitigate bias and harm	Efforts to mitigate bias and harm
Conciseness vs. Detail	Gemini may prioritize shorter, direct responses	ChatGPT can be more verbose and offer greater detail
Understanding Complex Queries	Gemini's search integration may excel in handling complex or multi-part questions	ChatGPT might be better suited to parsing more conversational language even if complex
Multilingual Capability	Gemini may have stronger multilingual support due to Google's resources	ChatGPT's multilingual capabilities are improving, but still under development
Personalization	Gemini primarily focuses on general responses	ChatGPT can, within a conversation, provide tailored responses based on past interactions

## Architecture

### Gemini Architecture:

Gemini, developed by Google AI, comprises a family of LLMs, with Gemini Ultra 1.0 being the most advanced version [15-16]. Key features include Retrieval-Augmented Generation (RAG),

which integrates information retrieval with text generation, resulting in factually grounded outputs. Gemini is trained on a diverse dataset, allowing it to excel at various tasks and benefit from Google's scalable infrastructure.

**ChatGPT Architecture:**

ChatGPT, developed by OpenAI, is renowned for its conversational abilities and creative text generation. Built on the Generative Pre-training (GPT) architecture, it predicts the next word in text data, enabling fluent text generation [8,11]. ChatGPT incorporates Reinforcement Learning with Human Feedback (RLHF) to improve responses based on human input and utilizes instruction tuning for adaptability.

**Comparison and Key Differences:**

Gemini excels in producing factually accurate responses through RAG but may occasionally generate harmful content from its knowledge base. It boasts versatility in tasks like code generation, surpassing ChatGPT in task diversity. ChatGPT, on the other hand, focuses on natural language interactions and benefits from RLHF for bias reduction, resulting in safer and more helpful responses. However, OpenAI's transparency regarding training data has been questioned compared to Google's more detailed approach.

Feature	Gemini	ChatGPT
Developer	Google AI	OpenAI
Model Type	Multimodal Language Model (can handle text, images, and potentially other modalities)	Generative Pre-trained Transformer (GPT) Language Model
Base Architecture	Believed to be a Transformer-based architecture	Transformer-based architecture
Training Data	Massive, proprietary dataset curated by Google. Likely combines text, code, and potentially other forms of data.	A massive dataset of text and code with careful filtering to optimize for conversational quality and safety.
Key Strengths	Advanced language understanding. Multimodal capabilities give potential for broader uses. Variation in models allows tailoring to specific needs.	Exceptional conversational ability. Generates diverse and creative text formats. Good at following instructions and staying on topic.
Known Weaknesses	More details are relatively limited as access is not widespread yet. May have similar biases and potential for generating misinformation as found in other large language models.	Struggles with some areas of logic and reasoning. Can be 'jailbroken' with clever prompts to produce undesirable outputs.

Table 2 Architectural aspects of Gemini and ChatGPT

# SWOT Analysis

Table 4 SWOT analysis of Gemini and ChatGPT

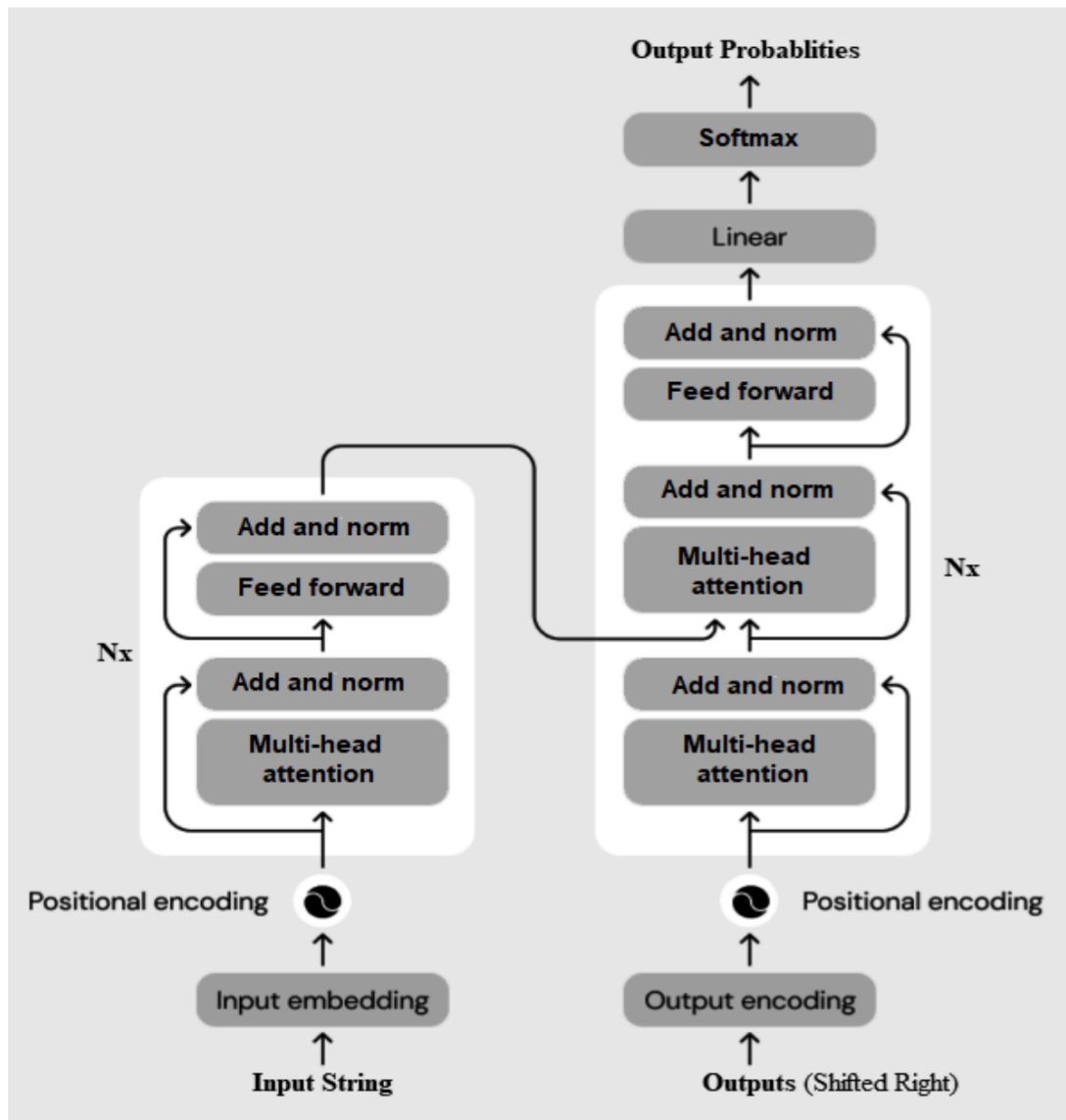
SWOT	Gemini (Google)	ChatGPT (OpenAI)
Strengths	Multimodality (text, images, potentially more): Processes and creates responses using various input formats.	Established reputation, widespread adoption: Enjoys public recognition and use across many applications.
	Improved reasoning capabilities: Exhibits enhanced understanding and logic in its responses.	Exceptional text generation prowess: Highly skilled in creating human-like text in various styles and formats.
	Computational efficiency: Functions with fewer resources, potentially making it faster and more widely accessible.	Adaptability to diverse tasks: Can be fine-tuned to handle translations, coding, content creation, etc.
	Open-source commitment: Transparency and potential for community improvements and ethical oversight.	Strong research and development focus: Continuous updates and potential for breakthrough advancements.
Weaknesses	Still under development: May experience inconsistencies and limitations compared to more mature models.	Limited multimodal capabilities: Primarily text-focused; may struggle to match Gemini's cross-format skills.
	Smaller real-world dataset compared to ChatGPT: Potentially less robust factual knowledge and less nuanced understanding.	Potential for bias and misinformation: Could propagate incorrect or harmful information based on its training data.
	Potential misuse of enhanced capabilities: Open access increases the risk of use for deceptive or malicious purposes.	Resource intensive: Computationally demanding, raising potential barriers to deployment.
Opportunities	Revolutionize multimodal user experiences: Enable interfaces combining text, voice, and visual interactions.	Continued advancement and refinement: Increasing potential for both accuracy and capabilities.
	Address complex problems with reasoning: Facilitate tasks involving logical deduction and complex problem-solving.	Widespread integration into applications: Expansion into areas like customer support, education, and entertainment.
	Foster innovation through open-source model: Collaborative model development, customization, and focus on ethical deployment.	Potential for commercialization: Offers licensing and specialized service opportunities.
Threats	Competition from other emerging LLMs: Rapidly advancing AI landscape and pressure to remain at the cutting edge.	Negative publicity surrounding AI misuse: High-profile misuses could diminish public trust in such technologies.
	Ethical considerations of multimodal AI: Responsibility for the unintended consequences of combining various media formats (e.g., generating potentially harmful content).	Public dependence on AI systems: Risks include algorithmic errors, manipulation of opinions, and over-reliance on generated output.
	Potential for biases within datasets: Unchecked biases in training data may permeate responses, harming fairness.	Regulation of AI technologies: Stricter policies and standards could impede future AI development.

## Detailed architect of GPT and features

According to OpenAI, their ChatGPT model can mimic discussion, respond to follow-up queries, acknowledge mistakes, contest false premises, and reject unsuitable requests. It was taught using a machine learning method called Reinforcement Learning from Human Feedback (RLHF) [5]. Its architecture consists of multiple layers of artificial neurons, which are inspired by the structure and function of neurons in the human brain. The input data is processed at each layer before being sent to the next layer, which produces the output in the end. In case of ChatGPT, the input data is a sequence of words and the output is a predicted next word or a series of words in response to a given prompt. Its architecture is based on the Transformer model, which was introduced in a paper published by a researcher at Google in 2017[6]. The Transformer model has two main parts: an encoder and a decoder. The encoder transforms the input data into a fixed-length representation, while the decoder uses the encoded representation to produce the output. The Transformer model is known for its ability to process input data in parallel, which makes it well-suited for tasks like translation and language modeling.

Fig.2 depicts the GPT architecture. The "add and norm" function, also known as the "addition and normalization" function, is a building block of the Transformer model that is used to stabilize the training process and improve the model's ability to learn long-term dependencies in the data. It consists of two main steps:

- Addition: The output of the previous layer is added to the input of the current layer.
  - Normalization: The sum is then passed through a normalization function, such as batch normalization or layer normalization, which scales and shifts the data to stabilize the distribution and prevent the values from becoming too large or too small.
- Figure 2. GPT Architecture
- EAI Endorsed Transactions on AI and Robotics



The "multi-head attention" function is a key component of the Transformer model that enables the model to process multiple elements of the input data simultaneously. It consists of multiple attention heads, each of which takes a different subset of the input data as input and computes a weighted sum of the data. The weighted sums are then concatenated and combined to produce the final attention output. As a result, the model is better able to comprehend the context and significance of the input and provide replies that are more cohesive and coherent. The "feed-forward" function is another building block of the Transformer model that is used to transform the input data through a series of fully-connected layers. It consists of two linear transformations, which are followed by a non-linear activation function such as ReLU. The output of the feed-forward

function is then combined with the output of the previous layer using the "add and norm" function. The "linear" function takes in a sequence of words as input and transforms it into a lower-dimensional representation that is suitable for processing by the rest of the model. The "Softmax" function, on the other hand, is used to predict the probability of each word in the vocabulary given the previous words in the input sequence. This allows GPT to generate text by sampling from the predicted probability distribution at each time step. The Softmax function is therefore used to generate natural language text.

Some of the key features of the GPT architecture include:

- Pre-training: It is pre-trained on a large dataset of conversational exchanges, which allows it to learn the patterns and structure of natural language conversation. This can help ChatGPT to generate more coherent and relevant responses.
- Contextualization: GPT is able to take into account the context of a conversation when generating responses, allowing it to generate more coherent and relevant responses.
- Flexibility: It can be fine-tuned for a variety of different chatbot tasks and domains.
- Scalability: GPT is a large, transformer-based model, which allows it to scale to very large datasets and handle long-range dependencies in language.
- Natural language generation: GPT is able to generate human-like responses in natural language, which makes it well-suited for a chatbot application like ChatGPT.

## **GEMINI architect**

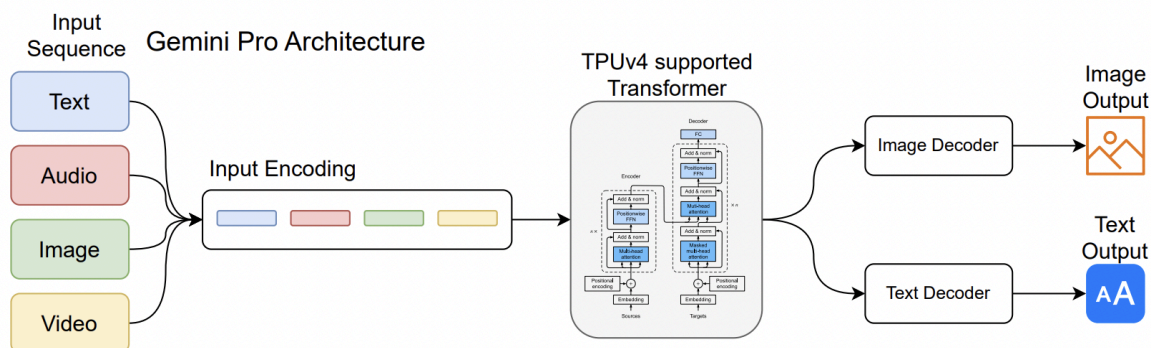
Gemini models build on top of Transformer decoders (Vaswani et al., 2017b) that are enhanced with improvements in architecture and model optimization to enable stable training at scale and optimized inference on Google's Tensor Processing Units. They are trained to support 32k context length, employing efficient attention mechanisms (for e.g. multi-query attention (Shazeer, 2019a)). Our first version, Gemini 1.0, comprises three main sizes to support a wide range of applications as discussed in Table 1.

Gemini models are trained to accommodate textual input interleaved with a wide variety of audio and visual inputs, such as natural images, charts, screenshots, PDFs, and videos, and they can produce text and image outputs (see Figure 2). The visual encoding of Gemini models is inspired by our own foundational work on Flamingo

(Alayrac et al., 2022), CoCa (Yu et al., 2022a), and PaLI (Chen et al., 2022), with the important distinction that the models are multimodal from the beginning and can natively output images using discrete image tokens (Ramesh et al., 2021; Yu et al., 2022b).

Video understanding is accomplished by encoding the video as a sequence of frames in the large context window. Video frames or images can be interleaved naturally with text or audio as part of the model input. The models can handle variable input resolution in order to spend more compute on tasks that require fine-grained understanding. In addition, Gemini models can directly ingest audio signals at 16kHz from Universal Speech Model (USM) (Zhang et al., 2023) features. This enables the model to capture nuances that are typically lost when the audio is naively mapped to a text input (for example, see audio understanding demo on the website).

Training the Gemini family of models required innovations in training algorithms, dataset, and infrastructure. For the Pro model, the inherent scalability of our infrastructure and learning algorithms enable us to complete pre-training in a matter of weeks, leveraging a fraction of the Ultra’s resources. The Nano series of models leverage additional advancements in distillation and training algorithms to produce the best-in-class small language models for a wide variety of tasks, such as summarization and reading comprehension, which power our next generation on-device experiences.



We trained Gemini models using TPUv5e and TPUv4 (Jouppi et al., 2023), depending on their sizes and configuration. Training Gemini Ultra used a large fleet of TPUv4 accelerators owned by Google 4 Gemini: A Family of Highly Capable Multimodal Models across multiple datacenters. This represents a significant increase in scale over our prior flagship model PaLM-2 which presented new infrastructure challenges. Scaling up the number of accelerators results in a proportionate decrease in the mean time between failure of hardware in the overall system. We minimized the rate of planned

reschedules and preemptions, but genuine machine failures are commonplace across all hardware accelerators at such large scales.

TPUv4 accelerators are deployed in “SuperPods” of 4096 chips, each connected to a dedicated optical switch, which can dynamically reconfigure 4x4x4 chip cubes into arbitrary 3D torus topologies in around 10 seconds (Jouppi et al., 2023). For Gemini Ultra, we decided to retain a small number of cubes per superpod to allow for hot standbys and rolling maintenance.

TPU accelerators primarily communicate over the high speed inter-chip-interconnect, but at Gemini Ultra scale, we combine SuperPods in multiple datacenters using Google’s intra-cluster and inter-cluster network (Poutievski et al., 2022; Wetherall et al., 2023; Yao Hong et al., 2018). Google’s network latencies and bandwidths are sufficient to support the commonly used synchronous training paradigm, exploiting model parallelism within superpods and data-parallelism across superpods.

## Model features & Capabilities

- Instruction Following - by collecting data for a diverse set of instruction following categories. For instructions that are verifiable programmatically such as word count, it generate synthetic data via prompting and response editing to ensure that such instructions are satisfied.
- Tool Use - Gemini draws on a range of tools via Gemini Extensions, including Google Workspace, Google Maps, YouTube, Google Flights, and Google Hotels. These tool-use capabilities also enable Gemini to be integrated as part of Gmail, Docs, Slides, Sheets and more.
- Multilinguality - Scaling Gemini from English to 40+ languages imposed research challenges in data quality. We leverage abundant high-quality English data by localization to native cultures (e.g., “president of the United States” -> “日本の首相”).
- Multimodal Vision - We empower Gemini and Gemini Advanced with image understanding capabilities by fine-tuning pre-trained Gemini models on a mixture of text-only and image-text data. Careful balancing of text and multimodal data ensures the model develops robust image understanding without adversely affecting the quality of the text-only interactions. To assess our models, we compile a dataset of human-curated and synthetic image-text prompts and responses, spanning various categories and difficulty levels. This dataset facilitates human evaluation for model comparison and selection.
- Coding - Despite the strong coding benchmark performance of the base model, post-training data still provides a significant boost to both code quality and code



correctness. This highlights the benefit of high-quality demonstration data and feedback data for coding use cases. Gemini Apps and Gemini API models use a combination of human and synthetic approaches to collect such data.

## Benefits of AI learning

Here are some specific ways in which ChatGPT can be used to enhance learning and the references that support it:

**Personalized Tutoring:** ChatGPT can be used to provide personalized tutoring and feedback to students based on their individual learning needs and progress. A study by Chen and colleagues (2020) demonstrated that a conversational agent based on a generative model (ChatGPT) could provide personalized math tutoring to students, resulting in improved learning outcomes. The study showed that the conversational agent was able to provide explanations that were tailored to students' misconceptions and was able to adapt to their level of understanding.

**Automated Essay Grading:** ChatGPT can be trained to grade student essays, providing teachers with more time to focus on other aspects of teaching. A study by Kim and colleagues (2019) showed that a generative model (ChatGPT) trained on a dataset of human-graded essays could accurately grade essays written by high school students, with a correlation of 0.86 with human grades. The study showed that the model was able to identify key features of well-written essays and was able to provide feedback that was similar to that provided by human graders.

**Language Translation:** ChatGPT can be used to translate educational materials into different languages, making them more accessible to a wider audience. A study by Johnson and colleagues (2016) demonstrated that a generative model (ChatGPT) trained on a dataset of bilingual sentence pairs could accurately translate between languages, achieving state-of-the-art results on several translation benchmarks. The study showed that the model was able to understand the meaning of sentences in one language and to generate accurate translations in another language.

**Interactive Learning:** ChatGPT can be used to create interactive learning experiences where students can interact with a virtual tutor in a conversational manner. A study by Peng and colleagues (2019) demonstrated that a generative model-based conversational agent could provide effective support for students learning English as a second language, resulting in improved language proficiency. The study showed that the

agent was able to understand students' questions and to provide appropriate and relevant responses.

**Adaptive Learning:** ChatGPT can be used to create adaptive learning systems that adjust their teaching methods based on a student's progress and performance. A study by Chiang and colleagues (2021) showed that an adaptive learning system based on a generative model (ChatGPT) could provide more effective support for students learning programming, resulting in improved performance on programming assessments. The study showed that the model was able to understand students' knowledge and to adjust the difficulty of the problems it generated accordingly.

## Possible Drawbacks of AI Learning

**Lack of Human Interaction:** ChatGPT and other generative models are not capable of providing the same level of human interaction as a real teacher or tutor. This lack of human interaction can be a disadvantage for students who may benefit more from a personal connection with a teacher. A study by D'Mello and colleagues (2014) found that students who interacted with a virtual tutor that mimicked human-like affective behavior had a better learning outcome than those who interacted with a virtual tutor that lacked this behavior.

**Limited Understanding:** Generative models are based on statistical patterns in the data they are trained on, and they do not have a true understanding of the concepts they are helping students learn. This can be a disadvantage when it comes to providing explanations or feedback that are tailored to a student's individual needs and misconceptions. A study by Wang and colleagues (2020) showed that a generative model-based tutoring system lacked the ability to provide explanations that were tailored to students' misconceptions.

**Bias in Training Data:** Generative models are only as good as the data they are trained on, and if the training data contains biases, the model will also be biased. For example, if a model is trained on a dataset of essays that are primarily written by students from a certain demographic, it may not be able to accurately grade essays written by students from other demographics. A study by Bolukbasi and colleagues (2016) showed that a generative model trained on a large corpus of text from the internet exhibited gender bias in its language generation.

**Lack of Creativity:** Generative models can only generate responses based on the patterns in the data they have seen during training, which can limit the creativity and

originality of the responses. A study by Ziegler and colleagues (2019) found that a generative model-based music composition system had a limited ability to generate original and diverse melodies.

**Dependency on Data:** Generative models are trained on a large amount of data, and the quality of the model is highly dependent on the quality and quantity of the data. If the data is not sufficient or not relevant, the model will not be able to perform as well. A study by Kocaguneli and colleagues (2019) showed that a generative model-based question answering system performed poorly when the training data was not relevant to the task at hand.

**Lack of Contextual Understanding:** Generative models lack the ability to understand context and situation, which can lead to inappropriate or irrelevant responses. A study by Gao and colleagues (2019) showed that a generative model-based dialogue system had a limited ability to understand and generate contextually appropriate responses in a conversation.

**Limited ability to personalize instruction:** ChatGPT and other generative AI models can provide general information and assistance, but they may not be able to personalize instruction to meet the individual needs of a particular student. (Ribeiro & Vala, 2020)

Ref :

“Gemini Versus ChatGPT: Applications, Performance, Architecture, Capabilities, and Implementation.” Nitin Rane, Saurabh Choudhary & Jayesh Rane. (2023) Available online:

[https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=4723687](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4723687)

“A Brief Analysis of “ChatGPT” – A Revolutionary Tool Designed by OpenAI

Authors.” **Md. Asraful Haque. (2023)** Available online:

<https://publications.eai.eu/index.php/airo/article/view/2983>

“Gemini technical report”

[https://storage.googleapis.com/deepmind-media/gemini/gemini\\_1\\_report.pdf](https://storage.googleapis.com/deepmind-media/gemini/gemini_1_report.pdf)

“Education in the Era of Generative Artificial Intelligence (AI): Understanding the Potential Benefits of ChatGPT in Promoting Teaching and Learning”. David Baidoo-anu & Leticia Owusu Ansah.

<https://dergipark.org.tr/en/pub/jai/issue/77844/1337500>

Abukmeil, M., Ferrari, S., Genovese, A., Piuri, V., & Scotti, F. (2021). A survey of unsupervised generative models for exploratory data analysis and representation learning. *Acm computing surveys (csur)*, 54(5), 1-40. <https://doi.org/10.1145/3450963>.

Alshater, M. (2022). Exploring the role of artificial intelligence in enhancing academic performance: A case study of ChatGPT (December 26, 2022). Available at SSRN: <https://ssrn.com/abstract=4312358> or <http://dx.doi.org/10.2139/ssrn.4312358>

Aydin, Ö., & Karaarslan, E. (2023). Is ChatGPT leading generative AI? What is beyond expectations? Academic Platform Journal of Engineering and Smart Systems.

Kasneci, E., Seßler, K., Küchemann, S., Bannert, M., Dementieva, D., Fischer, F., ... & Kasneci, G. (2023). ChatGPT for good? On opportunities and challenges of large language models for education. Learning and individual differences, 103, 102274.

Pearl Mike, "The ChatGPT chatbot from OpenAI is amazing, creative, and totally wrong". Mashable, December 3, 2022. Available online: <https://mashable.com/article/chatgpt-amazing-wrong>.

"How Chat GPT utilizes the advancements in Artificial Intelligence to create a revolutionary language model", Pegasus One, December 17, 2022. Available online: <https://www.pegasusone.com/how-chat-gpt-utilizes-the-advancements-in-artificial-intelligence-to-create-a-revolutionary-language-model>.

Chang, Y., Wang, X., Wang, J., Wu, Y., Yang, L., Zhu, K., ... & Xie, X. (2023). A survey on evaluation of large language models. ACM Transactions on Intelligent Systems and Technology.

Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., ... & Fedus, W. (2022). Emergent abilities of large language models. arXiv preprint arXiv:2206.07682.

Perera, P., & Lankathilake, M. (2023). Preparing to revolutionize education with the multi-model GenAI tool Google Gemini? A journey towards effective policy making. J. Adv. Educ. Philos, 7, 246-253.

Wu, T., He, S., Liu, J., Sun, S., Liu, K., Han, Q. L., & Tang, Y. (2023). A brief overview of ChatGPT: The history, status quo and potential future development. IEEE/CAA Journal of Automatica Sinica, 10(5), 1122-1136.

Team, G., Anil, R., Borgeaud, S., Wu, Y., Alayrac, J. B., Yu, J., ... & Ahn, J. (2023). Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805.

"Designing IoT Introductory Course for Undergraduate Students Using ChatGPT". Abdallah Al-Zoubi (2024) [https://link.springer.com/chapter/10.1007/978-3-031-51979-6\\_40](https://link.springer.com/chapter/10.1007/978-3-031-51979-6_40)