# A Comparative Survey of Chatgpt Models: Version 3, 3.5 And 4

**Anuja Beatrice B[1], Alfiya A[2]**
[1, 2] Dept of Computer Applications
[1, 2] Sri Krishna Arts & Science College

*Abstract-* *This paper aims to present a comprehensive comparative survey of three iterations of ChatGPT models: GPT-3, GPT-3.5, and GPT-4. These models which are developed by OpenAI, represent monumental growth in natural language processing (NLP) and conversational AI technology. By critically analyzing their architectures, capabilities, performance metrics, and applications, this survey provides insights into the evolution of ChatGPT and the improvements incorporated in each iteration. Additionally, the paper sheds light sheds light on the strengths, limitations, and potential future directions of these models in various domains, including language comprehension, generation, and dialogue systems.*

*Keywords-* ChatGPT, GPT-3, GPT-3.5, GPT-4, Natural Language Processing, Conversational AI

## I. INTRODUCTION

In recent years, ChatGPT models have become revolutionizing state-of-the-art natural language processing (NLP) models, which has the capability of generating human like text responses and initiating conversational interactions. Developed by OpenAI, these models have undergone iterative improvements, with each new version building upon the successes and limitations of its predecessors. The evolution from GPT-3 to GPT-4 represents a significant milestone in the advancement of conversational AI technology, marked by enhancements in model architecture, training data, and performance metrics. As the features of ChatGPT models continue to grow, the importance of conducting a comparative analysis becomes increasingly apparent. By systematically comparing the features, capabilities, and performance of different iterations, researchers can gain valuable insights into the progress made in the field of NLP, as well as the suspected challenges and the plethora of opportunities that lie ahead. This comparative analysis serves as a critical tool for evaluating the efficacy of ChatGPT models in various applications and guiding future research and development efforts towards the advancement of conversational AI technology.

## II. ARCHITECHTURE AND FEATURES

ChatGPT models, including GPT-3, GPT-3.5, and GPT-4, are based on the transformer architecture, a deep learning architecture primarily introduced in the seminal paper "Attention is All You Need" by Vaswani et al. This architecture contains of multiple layers of self-attention mechanisms and feedforward neural networks, enabling the model to capture long-range dependencies and contextual information in input sequences effectively. The size of each ChatGPT model varies, with GPT-3 comprising over 170 billion parameters, GPT-3.5 scaling up to 175 billion parameters, but GPT-4 potentially surpassing these sizes, although exact details may vary.

Training data for ChatGPT models typically consist of large-scale corpora of textual data from different sources, including journals, articles, books websites, and other publicly availed text sources. Pre-training objectives involve training the model to predict the next word or token in a sequence given the preceding context, using techniques such as autoregressive language modeling. Additionally, models like GPT-3 and GPT-4 may incorporate variants of self-supervised learning objectives, such as masked language modeling or causal language modeling, to improve generalization and robustness.

Fine-tuning capability enables users to adapt ChatGPT models to specific tasks or domains by providing task-specific training data and fine-tuning the model parameters accordingly. This process involves initializing the model with pre-trained weights and updating them through further training on task-specific data, often using techniques such as gradient descent and backpropagation. Fine-tuning allows ChatGPT models to achieve state-of-the-art performance on a wide range of NLP tasks, including text classification, sentiment analysis, and language translation.

Novel features introduced in GPT-3.5 and GPT-4 include enhancements to model architecture, training data, and performance metrics. GPT-3.5 may introduce improvements in model efficiency, scalability, or fine-tuning capabilities compared to GPT-3, while GPT-4 may further refine these

features or introduce entirely new capabilities, such as enhanced contextual understanding, multi-modal integration, or improved handling of long-range dependencies. Additionally, both versions may address limitations identified in earlier iterations, such as biases in generated text or difficulties in handling specific linguistic phenomena. Overall, these novel features contribute to the continued advancement and refinement of ChatGPT models, enabling them to achieve higher levels of performance and applicability in real-world scenarios.

## III. PERFORMANCE EVALUATION

ChatGPT models are evaluated using various metrics to assess their language understanding and generation capabilities. For language understanding, common evaluation metrics include precision, accuracy, recall, and F1 score, calculated based on the model's performance on tasks such as text classification, sentiment analysis, or question-answering. Additionally, metrics like BLEU (Bilingual Evaluation Understudy), ROUGE (Recall-Oriented Understudy for Gisting Evaluation), and METEOR (Metric for Evaluation of Translation with Explicit Ordering) are used to evaluate the quality of generated text in tasks such as language translation or summarization. These metrics provide quantitative measures of the model's ability to produce text that is grammatically correct, semantically coherent, and contextually relevant.

Perplexity scores can also be used as a metric for evaluating the performance of language models like ChatGPT. Perplexity measures how well the model predicts the next token in a sequence based on the preceding context. Lower perplexity scores indicate better performance, as they reflect the model's ability to accurately predict the next token given the context provided. By comparing perplexity scores across different models and datasets, researchers can assess the relative performance and effectiveness of ChatGPT models in capturing complex linguistic patterns and dependencies.

In addition to quantitative metrics, case studies are often employed to demonstrate the real-world performance of ChatGPT models in practical applications. These case studies involve deploying the model in specific use cases or scenarios, such as customer support chatbots, virtual assistants, or content generation tasks, and evaluating its performance based on criteria such as user satisfaction, task completion rate, and response quality. By conducting case studies in real-world settings, researchers can gain insights into the strengths, limitations, and potential areas for improvement of ChatGPT models, as well as their overall effectiveness in addressing practical challenges and meeting user needs.

Overall, a combination of quantitative evaluation metrics, such as language understanding and generation metrics, perplexity scores, and qualitative case studies, provides a comprehensive assessment of the performance of ChatGPT models, enabling researchers and practitioners to make informed decisions regarding model selection, fine-tuning, and deployment in real-world applications.

## IV. APPLICATIONS & USE CASES

ChatGPT models, including GPT-3, GPT-3.5, and GPT- 4, have a variety of applications across multiple domains, leveraging their capabilities in natural language understanding and generation. Some prominent applications and use cases include:

- Chatbot Applications: ChatGPT models are widely used in chatbot development to provide conversational interfaces for various services and platforms. These chatbots can assist users with tasks such as customer support, information retrieval, appointment scheduling, and virtual assistance.
- Text Completion and Generation Tasks: ChatGPT models excel in text completion and generation tasks, where they can generate coherent and contextually relevant text based on user input or prompts. These tasks include auto-completion of text, content generation for articles, essays, or stories, and summarization of lengthy documents.
- Assistance in Writing, Coding, and Creative Tasks: ChatGPT models can assist users in various writing, coding, and creative tasks by providing suggestions, feedback, and inspiration. Writers also use ChatGPT to develop ideas, tackle writer's block, or refine their writing style. Programmers can utilize ChatGPT to generate code snippets, debug code, or explore solutions to coding problems. Additionally, artists and creatives can collaborate with ChatGPT to generate creative prompts, explore new concepts, or generate artwork based on textual descriptions.
- Educational and Language Learning Applications: ChatGPT models are valuable tools for educational and language learning applications, providing personalized tutoring, language practice, and interactive learning experiences. Students can engage with ChatGPT-powered virtual tutors to receive explanations, practice exercises, and feedback on academic topics. Language learners can converse with ChatGPT in their target language to practice speaking, listening, and comprehension skills.

Overall, ChatGPT models offer versatile capabilities which could be leveraged in a wide range of applications and use cases, spanning customer service, content creation, writing

assistance, education, and beyond. As the technology continues to evolve, the potential for innovative applications and novel use cases of ChatGPT models is boundless, driving advancements in conversational AI and shaping the future of human-computer interaction.

## V. STRENGTH & LIMITATIONS

**GPT-3**

Strengths:

i. **Scalability:** GPT-3 is one of the largest language models ever created, with 175 billion parameters, enabling it to capture a vast amount of linguistic knowledge and generate diverse and contextually rich responses.
ii. **Versatility:** GPT-3 exhibits versatility in a wide range of tasks, from language translation and text completion to question-answering and code generation, making it applicable across various domains and use cases.
iii. **Performance:** GPT-3 achieves impressive performance on many benchmarks and tasks, demonstrating strong language understanding and generation capabilities.

Limitations:

i. Biases: GPT-3 may exhibit biases present in the training data, leading to biased or inappropriate responses in certain contexts. Mitigation strategies such as debiasing techniques or diverse dataset augmentation can help address this limitation.
ii. Context Understanding: GPT-3 sometimes struggles with understanding complex or nuanced contexts, leading to inconsistencies or inaccuracies in generated responses. Fine-tuning on domain-specific data and improving context modeling techniques can help alleviate this limitation.
iii. Safety Concerns: GPT-3 may generate harmful or inappropriate content, such as misinformation, hate speech, or sensitive personal information. Safeguarding measures such as content filtering, human-in-the-loop validation, and responsible use guidelines are essential to mitigate safety concerns and ensure ethical deployment.

**GPT-3.5 and GPT-4**

Strengths:

i. Enhanced Scalability: GPT-3.5 and GPT-4 may further increase model size and scalability, enabling them to capture even more complex linguistic patterns and generate higher-quality responses.
ii. Improved Versatility: New features and capabilities introduced in GPT-3.5 and GPT-4 may enhance model versatility, allowing for better performance on specific tasks or domains.
iii. Advanced Performance: GPT-3.5 and GPT-4 may introduce improvements in performance metrics such as perplexity, accuracy, and coherence, leading to better overall performance in language understanding and generation tasks.

Limitations:

i. Continued Biases: Despite advancements, GPT-3.5 and GPT-4 may still exhibit biases inherited from the training data, necessitating ongoing efforts to mitigate bias and promote fairness in model outputs.
ii. Contextual Understanding Challenges: Addressing complex contextual understanding challenges remains a priority for GPT-3.5 and GPT-4, requiring further research into advanced context modeling techniques and domain-specific fine-tuning strategies.
iii. Ethical Considerations: Ethical considerations surrounding the use of large language models like GPT-3.5 and GPT-4, including privacy, safety, and societal impacts, require careful attention and mitigation strategies to ensure responsible deployment and usage.

In conclusion, while ChatGPT models offer significant strengths in scalability, versatility, and performance, they also present challenges such as biases, context understanding limitations, and ethical concerns. Continual research and development efforts, along with ethical guidelines and mitigation strategies, are essential to address these limitations and realize the full potential of ChatGPT models in a wide range of applications and domains.

## VI. FUTURE DIRECTIONS

1) Potential Improvements and Research Directions for Future Iterations:
   i. Model Efficiency: Future iterations of ChatGPT could focus on improving model efficiency, reducing computational resources required for training and inference while maintaining or improving performance.

ii. Context Understanding: Advancements in context modeling techniques could enhance ChatGPT's ability to understand and generate responses based on complex and nuanced contexts, improving coherence and relevance in conversations.

iii. Domain-Specific Fine-Tuning: Tailoring ChatGPT models to specific domains through fine-tuning on domain-specific data could lead to better performance and more accurate responses in specialized contexts.

iv. Multi-task Learning: Exploring multi-task learning approaches could enable ChatGPT models to simultaneously perform multiple NLP tasks, such as language translation, summarization, and question-answering, improving versatility and efficiency.

2) Integration of Multimodal Capabilities:

i. Incorporating Visual and Audio Inputs: Future iterations of ChatGPT could integrate multimodal capabilities, allowing the model to process and generate responses based on visual and audio inputs in addition to text, enabling more immersive and interactive conversational experiences.

ii. Fusion Techniques: Developing fusion techniques to combine information from different modalities effectively could enhance ChatGPT's ability to understand and respond to complex inputs that involve multiple modalities, such as describing images or interpreting spoken language.

3) Addressing Ethical and Societal Implications:

i. Bias Detection and Mitigation: Continued efforts to detect and mitigate biases in ChatGPT models are crucial to ensure fair and equitable outputs, requiring ongoing research into bias detection methods, debiasing techniques, and diverse dataset curation strategies.

ii. Responsible Deployment Guidelines: Establishing comprehensive guidelines and best practices for the responsible deployment and usage of ChatGPT models is essential to mitigate potential risks, safeguard user privacy, and promote ethical AI practices.

iii. Transparency and Explainability: Enhancing the transparency and explainability of ChatGPT models by providing insights into model behavior, decision-making processes, and potential biases can foster trust and accountability in AI systems.

In conclusion, future iterations of ChatGPT hold great promise for advancements in model efficiency, context understanding, multimodal capabilities, and ethical deployment. By addressing these research directions and integrating multimodal capabilities while prioritizing ethical considerations, ChatGPT models can continue to evolve as powerful tools for natural language understanding and generation, contributing to a wide range of applications and domains while ensuring responsible and ethical AI deployment.

## VII. CONCLUSION

In conclusion, this comparative survey has provided valuable insights into the evolution and capabilities of ChatGPT models, including GPT-3, GPT-3.5, and GPT-4. Key findings from the analysis include the scalability, versatility, and performance improvements introduced in each iteration, as well as the challenges and limitations that persist across versions. By examining factors such as model architecture, training data, fine-tuning capabilities, and novel features, we have gained a deeper understanding of the strengths and weaknesses of each model, as well as their potential applications and use cases.

The comparative analysis has important implications for the future development and deployment of ChatGPT models. It highlights the need for ongoing research and innovation to address challenges such as biases, context understanding limitations, and ethical considerations, while maximizing the potential benefits of conversational AI technology. By identifying areas for improvement and opportunities for advancement, this analysis serves as a roadmap for future research and development efforts in the field of natural language processing and conversational AI.

## REFERENCES

[1] Brown, T. B., et al. "Language Models are Few-Shot Learners." NeurIPS, 2020.

[2] Radford, A., et al. "Learning to Summarize with Human Feedback." NeurIPS, 2019.

[3] Vaswani, A., et al. "Attention Is All You Need." NeurIPS, 2017.

[4] Yang, Z., et al. "XLNet: Generalized Autoregressive Pretraining for Language Understanding." NeurIPS, 2019.

[5] Keskar, N. S., et al. "On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima." ICLR, 2017.

[6] Viégas, F., et al. "How to Use t-SNE Effectively." Distill, 2016.

[7] Huang, G., et al. "Densely Connected Convolutional Networks." CVPR, 2017.

[8]   Hinton, G., et al. "Deep Neural Networks for Acoustic Modeling in Speech Recognition: The Shared Views of Four Research Groups." IEEE Signal Processing Magazine, 2012.

[9]   Mitchell, M., et al. "Model Cards for Model Reporting." FAT* Conference, 2019.

[10] Bender, E. M., et al. "Data Statements for Natural Language Processing: Toward Mitigating System Bias and Enabling Better Science." TACL, 2021.

[11] Jobin, A., et al. "The Global AI Index: Building a Comprehensive Sourcing and Ranking Mechanism." AI Index Report, 2021.

[12] Floridi, L., et al. "AI4People—An Ethical Framework for a Good AI Society: Opportunities, Risks, Principles, and Recommendations." Minds and Machines, 2018.

[13] Lu, J., et al. "ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and- Language Tasks." NeurIPS, 2019.

[14] Huang, X., et al. "Masked and Semi-Masked Cross-Modal Transformer for Vision-and-Language Tasks." CVPR, 2020.

[15] Tan, H., et al. "Multimodal Transformers for Visual Recognition and Generation." ICCV, 2021.

[16] Muhammad Usman Hadi, qasem al tashi, Rizwan Qureshi, Abbas Shah et al. "A Survey on Large Language Models: Applications, Challenges, Limitations, and Practical Usage" , Institute of Electrical and Electronics Engineers (IEEE), 2023

[17] Adrian David Cheok, Emma Yann Zhang. "From Turing to Transformers: A Comprehensive Review and Tutorial on the Evolution and Capabilities of Generative Models" , Qeios Ltd, 2023.

[18] J. K. D. D. T. Jayanetti, B. A. K. S. Perera, K. G. A. S. Waidyasekara. "Developing a definition for lean construction maturity models through a PRISMA systematic literature review" , FARU Journal, 2023.

[19] Lei Huang, Meng Song, Hui Shen, Huixiao Hong, Ping Gong, Hong-Wen Deng, Chaoyang Zhang. "Deep Learning Methods for Omics Data Imputation" , Biology, 2023.

[20] Yanshan Wang, Shyam Visweswaran, Sumit Kappor, Shravan Kooragayalu, Xizhi Wu. "ChatGPT, Enhanced with Clinical Practice Guidelines, is a Superior Decision Support Tool" , Cold Spring Harbor Laboratory, 2024

[21] Guanghua Wang, Weili Wu. "Surveying the Landscape of Text Summarization with Deep Learning: A Comprehensive Review" , Discrete Mathematics, Algorithms and Applications, 2023

[22] Mohammad Shahin, F. Frank Chen, Ali Hosseinzadeh, Mazdak Maghanaki, Ayda Eghbalian. "A Novel Approach to Voice of Customer Extraction using GPT-3.5 Turbo: Linking Advanced NLP and Lean Six Sigma 4.0" , Research Square Platform LLC, 2023