(REVIEW ARTICLE)

# A comprehensive review of advances in transformer, GAN, and attention mechanisms: Their role in multimodal learning and applications across NLP

Md Fokrul Islam Khan [1, *], Mst Halema Begum [1], Md Arifur Rahman [2], Golam Qibria Limon [2], Md Ali Azam [1] and Abdul Kadar Muhammad Masum [3]

[1] Masters in Management Information System , International American University, Los Angeles, USA.
[2] Doctor of Business Administration, International American University, Los Angeles, USA.
[3] (IEEE Senior Member), SU,Dhaka Bangladesh.

## Abstract

The emergence and subsequent development of deep learning, specifically transformer-based architectures, Generative Adversarial Networks (GANs), and attention mechanisms, have had revolutionary implications on Natural Language Processing (NLP) and multimodal learning. Transformer models are neural network architectures that change an input sequence into an output sequence. Transformer architectures like the Generative Pre-Training Transformer (GPT) and Bidirectional Encoder Representations from Transformers (BERT) leverage self-attention mechanisms to enable high-level contextual learning as well as long-range dependencies. GANs are a kind of AI algorithm that is designed to solve generative modeling problems. Different GANs, such as StyleGAN and BigCAN, study a collection of training data and learn the distribution probabilities used to generate such datasets. Attention mechanisms, acting as the unifying thread between Transformers and GANs in multimodal learning, optimize deep learning models to attend to the most relevant parts of the input data. This paper explores the synergy between these technologies, emphasizing their combined potential in multimodal learning frameworks. In addition, the paper analyzes recent advancements, key innovations, and practical implementations that leverage Transformers, GANs, and attention mechanisms to enhance natural language understanding and generation.

**Keywords:** Transformer Models; Generative Adversarial Networks (GANs); Attention Mechanisms; Multimodal Learning; Natural Language Processing (NLP)

## 1. Introduction

Deep learning has brought significant growth for NLP, and it has leveraged different technologies, from transformer models to GANs and attention mechanisms, as the bedrock of transformation. A paradigm in itself, in sequence modeling, the transformers have shifted focus from recurrent architectures to self-attention, thus facilitating parallelized computation and improved long-range dependency capture. Meanwhile, GANs have changed the entire landscape of generative AI by opening avenues for true-to-reality synthetic data creation of images, text, and videos. Attention mechanisms form the basis of transformers that improve these systems further by dynamically weighing input features for greater contextual understanding across modalities. This Paper posits a thesis: the integration of Transformers, GANs, and attention mechanisms is not merely additive but multiplicative, unlocking multimodal learning paradigms that transcend traditional NLP boundaries.

* Corresponding author: Fokrul Islam Khan

## 2. Transformer Models: Evolution and Role in NLP

### 2.1. Background of Transformers

Transformers were first introduced in 2017 by Vaswani et al. in their paper, Attention Is All You Need [1]. They marked a significant shift from earlier architectures, in other words, moving away from the traditional recurrent neural networks (RNNs). Before that, early-day language models could be described as very good at learning from sequential data with long-range dependencies. However, their ability to recognize the context of data in the long sequence was subpar [2]. The introduction of Transformer architecture completely changed that as it shifted the paradigm from sequential to self-attention mechanisms that can process sequences in parallel. Besides, the transformer paved the way for subsequent revolutionary frameworks like BERT and GPT.

### 2.2. Key Transformer Architectures

- BERT (Bidirectional Encoder Representations from Transformers): BERT is a language model introduced by Google researchers in October 2018 [3]. BERT learns to represent text as a sequence of vectors using self-learning. Usually, traditional models process text sequentially, either right-to-left or left-to-right, limiting their contextual awareness. BERT, however, reads the text in both directions simultaneously, capturing the full context of each word based on its surrounding tokens [4], as shown in Figure 1. Its fine-tuning flexibility amplifies its utility in that, after pretraining, BERT adapts to diverse downstream tasks like sentiment analysis and question answering through minimal architectural tweaks and task-specific labeled data. This transfer learning approach reduces training time and data requirements, democratizing access to high-performance NLP.
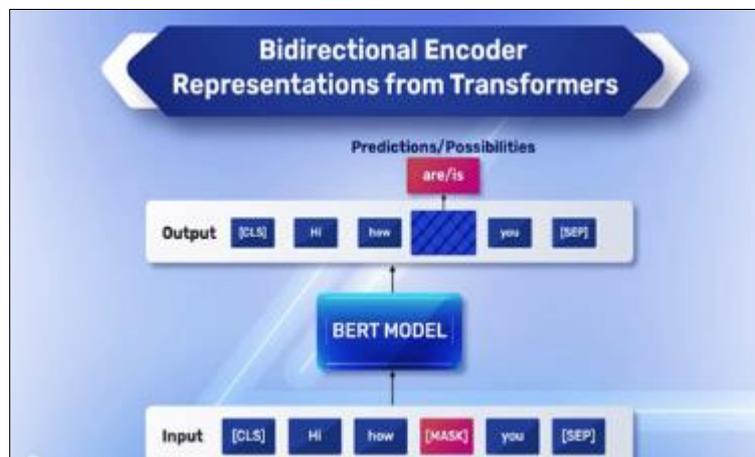


**Figure 1** Masked Language Modeling in BERT: The mastered word in the sequence "Hi, how [mask] you are inferred by the model, with predictions such as "are" and is" [5]

- GPT (Generative Pre-trained Transformer): 2020 ushered in the era of massive-scale language models with GPT-3. GPTs adopt a unidirectional, autoregressive design, enabling them to predict the next token in a sequence based on prior context [6]. This structure, combined with self-attention, allows GPT to capture long-range dependencies efficiently, surpassing RNN-based models in coherence and scalability. The shift to "few-shot" and "zero-shot" learning in GPT-3 and GPT-4—where tasks are specified via prompts rather than retraining—marks a leap toward general-purpose AI, reducing reliance on task-specific datasets [7]. Its few-shot learning, as shown in Figure 2, enables rapid adaptation to new domains without extensive learning.
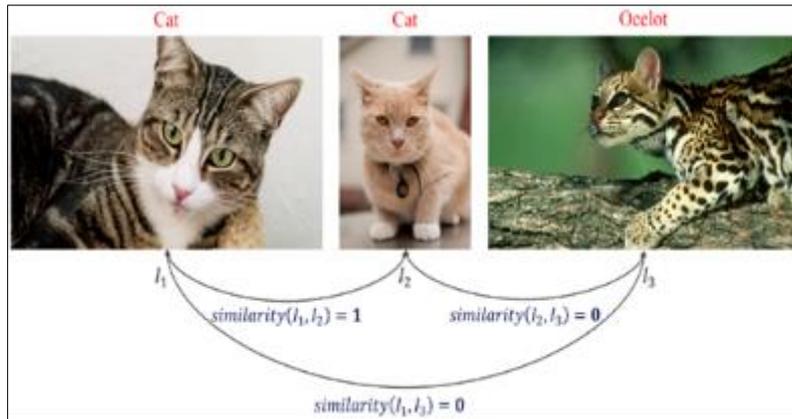
**Figure 2** Few-Shot Learning using GPT [8]. The two images of domestic cats (I1 and I2) are assigned a similar score of 1, while both are given a similarity score of 0 when compared to an ocelot (I3), demonstrating the model's ability to correctly differentiate between similar and dissimilar categories

## 2.3. Multimodal Applications of Transformers

Transformers facilitate multimodal learning by integrating textual and visual data. The multimodal capabilities emerge through extensions like VisualBERT and ViLBERT (vision-and-language BERT). ViLBERT, for example, employs dual-stream architectures—one for text, one for images—with cross-attention to align visual regions with textual tokens [9]. This enables tasks like visual question answering and image captioning. Meanwhile, GPT-style model condition generation on multimodal inputs, such as text guiding image edits or audio synthesis [10]. For example, GPT-3's successors integrate CLIP (Contrastive Language-Image Pretraining) embeddings to align text with visual features before generation.

## 3. Generative Adversarial Networks (GANs) in Multimodal AI
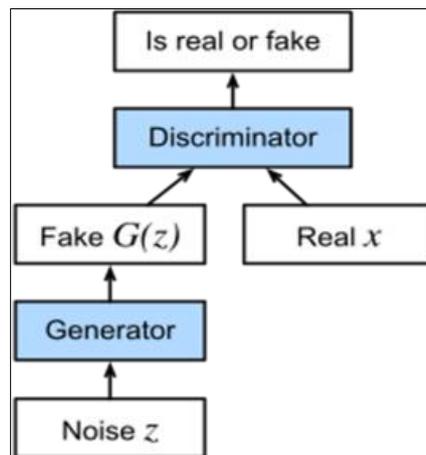
### 3.1. Overview of GANs



**Figure 3** An illustration of how a GAN works

A GAN is a kind of machine learning framework that is essential for making generative artificial intelligence. The GAN framework, as shown in Fig. 3, is composed of two opposing neural networks, the generator and the discriminator, which play a zero-sum game against each other [11]. The generator creates artificial, or otherwise called, synthetic data {G(z)} starting at some arbitrary random noise input (z). The discriminator, on the other hand, interrogates this synthetic output {G(z)} together with authentic data (x), labeling each one as either "real" or "fake."[12]. By weakening the discriminator's capacity to differentiate real items from fake ones, the generator advances its performance, whereas the discriminator strengthens its precision in categorization. This back-and-forth dynamic gradually yields highly convincing synthetic outputs over time.

### 3.2. GANs Architectures and Their NLP Applications

Two of the most prominent GANs are StyleGAN and BigCAN. StyleGAN, developed by Karras et al. in 2018, and its successor StyleGAN2, represent a leap in high-resolution image synthesis by introducing a novel generator architecture that disentangles high-level attributes (such as facial structure from fine-grained details (e.g., texture) [13]. Unlike traditional GANs, StyleGAN employs a mapping network to transform latent codes into an intermediate style space, followed by normalization layers that inject at multiple resolutions [14], as shown in Figure 4. This hierarchical control over synthesis enables unprecedented realism and flexibility when StyleGAN is paired with NLP models, facilitating text-to-image synthesis with fine-grained control.

BigGAN scales GANs to new heights by leveraging large batch sizes, architectural tweaks, and high-capacity models trained on large datasets like ImageNet. This results in superior quality and diversity, even for complex scenes. BigGAN's strength lies in its ability to generate high-fidelity, class-conditioned images, which dovetails with NLP in multimodal AI systems. By conditioning textual labels or embeddings, BigCAN can synthesize images aligned with detailed descriptions, advancing applications such as text-to-image synthesis and cross-model learning.
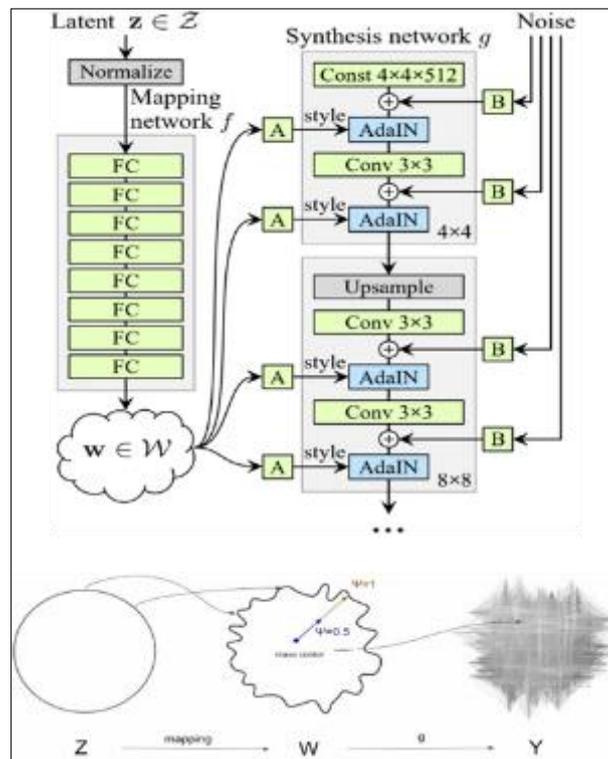


**Figure 4** StyleGAN generator network architecture [15]: This mapping network in StyleGAN transforms the latent code z into an intermedia latent space w, which is then used to contrail the synthesis network at multiple levels, enabling hierarchical control over different aspects of the generated image, Y

### 3.3. GAN in Multimodal Learning

In multimodal learning, GANs bridge modalities by learning joint distributions [16]. For example, StyleGAN's progressive rowing technique and noise injection enhance the diversity and realism of generated samples, avoiding mode collapse—a common GAN limitation. In multimodal learning, this ensures that synthetic data (e.g., images paired with text) is rich and varied, providing robust training examples for models that integrate multiple data types. In visual question and answering, for example, StyleGAN can generate diverse image-text-pairs, enabling better generalization across modalities and reducing overfitting to narrow datasets.

## 4. Attention Mechanisms: The Unifying Thread

Attention mechanisms constitute a technique in machine learning that instructs deep learning architectures to concentrate on the most pertinent aspects of the data they receive. Attention mechanisms underpin both GANs and Transformers, acting as the unifying thread in multimodal learning. Self-attention, as in Transformers, weights input

tokens dynamically, capturing contextual nuances across sequences. Multi-head attention, a variation that conducts several attention operations at once (as shown in Figure 5), takes this a step further by simultaneously analyzing features, thereby strengthening the quality and complexity of the resulting representations [17]. In GANs, attention refines generation by focusing on relevant input regions.
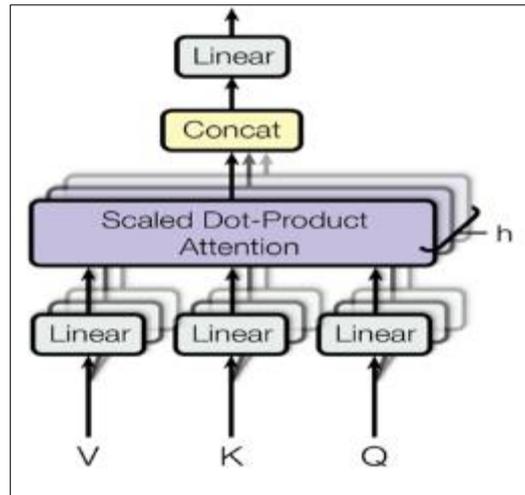


**Figure 5** Multi-head attention, as shown in this image, applies multi-parallel scaled dot-product attention mechanisms to different projections of Query (Q), Key (K), and Value (V), concatenating and linearly transforming them to capture diverse contextual relationships [17]

In multimodal NLP, cross-attention aligns disparate data types—text with images and audio with transcripts—by mapping features into a shared space. This adaptability enhances tasks like visual question answering, where attention prioritizes image regions pertinent to textual queries. Recent developments, such as adaptive attention and sparse attention, optimize efficiency, addressing scalability concerns in large-scale models.

## 5. Applications and Limitations in NLP and Multimodal Context

The deep learning models discussed in this paper have various applications in driving multimodal learning. The convergence of Transformers, GANs, and attention mechanisms has birthed transformative NLP applications. Multimodal extensions, such as speech-to-text and video summarization, leverage cross-attention for richer inputs. Text-image synthesis, powered by GANs and attention, enables creative tools like DALL-E, while sentiment analysis benefits from BERT's contextual depth.

However, these technologies face challenges despite their promising potential to reshape AI. Training large-scale models requires significant hardware resources. Additionally, pre-trained models inherit biases from training data, impacting fairness. Moreover, interpretability, crucial for trust, requires further refinement. Transformers, particularly GANs, often function as black-box models, limiting explainability.

## 6. Conclusion

Transformers, GANs, and attention mechanisms have reshaped AI, particularly in NLP and multimodal learning. Their continued evolution will unlock new frontiers in generative AI, human-computer interaction, and intelligent multimodal systems. However, challenges related to bias, interpretability, and computational efficiency must be addressed to ensure ethical AI adoption.

## Compliance with ethical standards

*Disclosure of conflict of interest*

No conflict of interest to be disclosed.

**References**

[1] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. Advances in neural information processing systems, 30.

[2] Ramachandran, K. K. (2024). The Evolution of Recurrent Neural Networks in Handling Long-Term Dependencies in Sequential Data. International Journal of Neural Networks and Deep Learning (IJNNDL), 1(1), 1-10.

[3] Koroteev, M. V. (2021). BERT: a review of applications in natural language processing and understanding. arXiv preprint arXiv:2103.11943. https://doi.org/10.48550/arXiv.2103.11943

[4] Apidianaki, M. (2023). From word types to tokens and back: A survey of approaches to word meaning representation and interpretation. Computational Linguistics, 49(2), 465-523. https://doi.org/10.1162/coli_a_00474

[5] Rath, S. (2023, October 31). BERT: Bidirectional Encoder Representations from Transformers. LearnOpenCV – Learn OpenCV, PyTorch, Keras, Tensorflow with Code, & Tutorials. https://learnopencv.com/bert-bidirectional-encoder-representations-from-transformers/

[6] Sun, Y. (2024). The evolution of transformer models from unidirectional to bidirectional in Natural Language Processing. Applied and Computational Engineering, 42, 281-289. https://doi.org/10.54254/2755-2721/42/20230794

[7] Kepel, D., & Valogianni, K. (2024). Autonomous prompt engineering in large language models. arXiv preprint arXiv:2407.11000. https://doi.org/10.48550/arXiv.2407.11000

[8] Kundu , R., & Skelton, J. (2024). Everything you need to know about Few-Shot Learning | DigitalOcean. Digitalocean.com. https://www.digitalocean.com/community/tutorials/few-shot-learning

[9] Guo, R., Wei, J., Sun, L., Yu, B., Chang, G., Liu, D., ... & Bu, L. (2023). A survey on image-text multimodal models. arXiv preprint arXiv:2309.15857. https://doi.org/10.48550/arXiv.2309.15857

[10] Gao, K., He, S., He, Z., Lin, J., Pei, Q., Shao, J., & Zhang, W. (2023). Examining user-friendly and open-sourced large gpt models: A survey on language, multimodal, and scientific gpt models. arXiv preprint arXiv:2308.14149. https://doi.org/10.48550/arXiv.2308.14149

[11] Sreevallabh Chivukula, A., Yang, X., Liu, B., Liu, W., & Zhou, W. (2022). Game theoretical adversarial deep learning. In Adversarial Machine Learning: Attack Surfaces, Defence Mechanisms, Learning Theories in Artificial Intelligence (pp. 73-149). Cham: Springer International Publishing. https://doi.org/10.1007/978-3-030-99772-4_4

[12] Ramzan, F., Sartori, C., Consoli, S., & Reforgiato Recupero, D. (2024). Generative adversarial networks for synthetic data generation in finance: Evaluating statistical similarities and quality assessment. AI, 5(2), 667-685. https://doi.org/10.3390/ai5020035

[13] Jain, A. (2024). Generative Adversarial Networks: A Review of Developments and Diverse Applications. Authorea, 1-47. https://doi.org/10.22541/au.172979391.16488935/v1

[14] Bermano, A. H., Gal, R., Alaluf, Y., Mokady, R., Nitzan, Y., Tov, O., ... & Cohen-Or, D. (2022, May). State-of-the-Art in the Architecture, Methods and Applications of StyleGAN. In Computer Graphics Forum (Vol. 41, No. 2, pp. 591-611). https://doi.org/10.1111/cgf.14503

[15] StyleGAN vs StyleGAN2 vs StyleGAN2-ADA vs StyleGAN3 – Lambdanalytique. (2022). Lambdanalytique.com. https://lambdanalytique.com/2022/07/01/stylegan-vs-stylegan2-vs-stylegan2-ada-vs-stylegan3/

[16] Ma, F., Li, Y., Ni, S., Huang, S. L., & Zhang, L. (2022). Data augmentation for audio-visual emotion recognition with an efficient multimodal conditional GAN. Applied Sciences, 12(1), 527. https://doi.org/10.3390/app12010527

[17] Han, S. Y., Sun, Q. W., Zhao, Q., Han, R. Z., & Chen, Y. H. (2022). Traffic forecasting based on integration of adaptive subgraph reformulation and spatio-temporal deep learning model. Electronics, 11(6), 861. https://doi.org/10.3390/electronics11060861.