
A SURVEY ON LARGE LANGUAGE MODELS FROM CONCEPT TO IMPLEMENTATION

A PREPRINT

 **Chen Wang**
Computer Science
The University of Alabama
Tuscaloosa, AL 35487
cwang86@ua.edu

 **Jin Zhao***
Electrical & Computer Engineering
The University of Alabama
Tuscaloosa, AL 35487
jzhao26@ua.edu

 **Jiaqi Gong***
Computer Science
The University of Alabama
Tuscaloosa, AL 35487
jgong5@ua.edu

ABSTRACT

Recent advancements in Large Language Models (LLMs), particularly those built on Transformer architectures, have significantly broadened the scope of natural language processing (NLP) applications, transcending their initial use in chatbot technology. This paper investigates the multifaceted applications of these models, with an emphasis on the GPT series. This exploration focuses on the transformative impact of artificial intelligence (AI) driven tools in revolutionizing traditional tasks like coding and problem-solving, while also paving new paths in research and development across diverse industries. From code interpretation and image captioning to facilitating the construction of interactive systems and advancing computational domains, Transformer models exemplify a synergy of deep learning, data analysis, and neural network design. This survey provides an in-depth look at the latest research in Transformer models, highlighting their versatility and the potential they hold for transforming diverse application sectors, thereby offering readers a comprehensive understanding of the current and future landscape of Transformer-based LLMs in practical applications.

Keywords LLMs · NLP · Transformer · industry impacts · application · fusion technologies

1 Introduction

The development of Large Language Models (LLMs) has become the focus of competition among leading technology groups. Notably, Google[®] and OpenAI[®] have emerged as major-league players in this domain, each advancing unique interpretations of natural language processing (NLP) models. A detailed survey of the development of these models, which form the basis of some of today's most advanced chatbot products, is provided in Table 1: 1) Google's Bard[®], which is built upon the Pathways Language Model (PaLM) and its advanced version (PaLM 2); and 2) OpenAI's ChatGPT[®], which is based on the series of Generative Pre-training Transformers (GPTs), especially GPT-3.5 and GPT-4. This table covers general comparisons and task-specific assessments, specifically on inference/reasoning, math skills, multitasking ability, and NL generation Anil et al. [2023].

The analysis includes a comparative exploration of key parameters including training parameter size, design architecture, methodology, and key features as well as their strengths and potential limitations. An important point of discussion is the training parameter size of GPT-3.5-Turbo, which is 20 billion – significantly smaller than its predecessor GPT-3.5 – as emphasized in Singh et al. [2023]. This reduction indicates improved computational efficiency with more efficient algorithmic optimization and data processing methods, resulting in improved performance with a less complex architecture. Table 1 also shows a clear trend of increasing task performance aligning with training requirements for LLMs. This pattern suggests an increase in computational requirements as the model is iterated from one iteration to the next, and highlights the need for advanced graphics processing units and efficient data acquisition techniques. This marks a shift in model development, *i.e.*, the transition from quantitative scaling to qualitative enhancement.

*Corresponding authors: Jin Zhao, and Jiaqi Gong.

Table 1: Evaluating performance across different GPT and PaLM model versions: general comparisons and task-specific assessments Anil et al. [2023].

Feature	GPT-3.5	GPT-4	PaLM	PaLM 2
# parameter	175 Billion (B) Koubaa [2023]	500 B Madden et al. [2023], Koubaa [2023]	540 B Chowdhery et al. [2022]	1.3 Trillion (T)
Architecture	(GPT-3.5-Turbo, 20 B Singh et al. [2023]) ^a Transformer	^a	^a	Pathways Chowdhery et al. [2022] (Transformer-based)
Core algorithm	^b self-supervised learning	^b	^b	^b
Key features	^c self-attention mechanism ^d positional encoding Bubeck et al. [2023]	^c ^d	^c ^d ^e pathway routing	^c ^d ^e few-shot learning Ma et al. [2023]
Strengths	^f good at generating creative text, translating and summarizing text.	^f	^g better at reasoning and understanding languages, especially in tasks that require common sense or real-world knowledge Brown et al. [2020].	^g it is available in smaller models that can be used on mobile devices and can perform few-shot learning tasks.
Weaknesses	^h less accurate, more prone to generating biased or offensive text Schramowski et al. [2023].	^h	ⁱ smaller, less versatile, more difficult to train.	ⁱ
Task 1: Inference	–	87.5	85.1	90.9
Task 2: Math	–	42.5	8.8	34.3 / 48.8
Task 3: Multitask	–	86.4	–	81.2
Task 4: NL Generation	–	95.9	83.2	87.4

The Transformer architecture is renowned for its self-attention mechanism. Originally designed for NLP tasks, this architecture has proven its versatility in a wide array of applications beyond language processing. The broader impact of the Transformer extends to the engineering domain. Transformer’s ability to process sequential data and identify both local and global features in sequences will revolutionize areas such as automated system configurations, troubleshooting, and safety management. Its parallel processing capabilities are also critical to the development of more agile and intelligent robotic systems.

A testament to this adaptability is DeepMind’s AlphaFold[®], as illustrated in Figure 1. AlphaFold uses a modified Transformer architecture to process amino acid sequences and effectively predict the three-dimensional structures of proteins. This application highlights the potential of the Transformer architecture in various scientific fields. The self-attention layers shown in Figure 1 are used to process and interpret relationships and patterns in protein sequences. In detail, the self-attention mechanisms serve the following purposes:

- **Include evolutionary information:** Multiple sequence alignment representations are fed into the self-attention layers, allowing the model to incorporate evolutionary information about protein families. This allows the model to understand which residues are important for protein structure and function based on their conservation across species.
- **Refine Structure predictions:** As part of the iterative process of structure prediction, self-attention layers refine the prediction by re-evaluating the relationships between residues after each iteration.

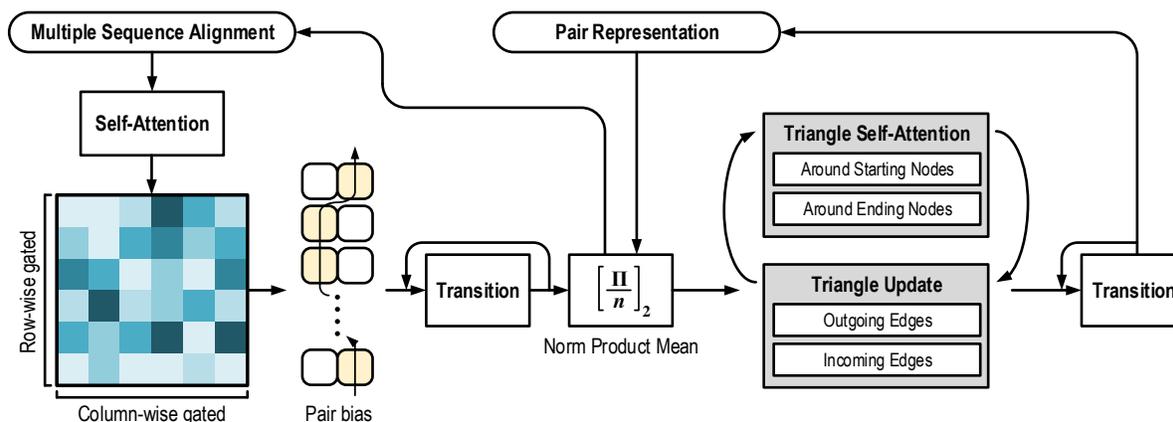


Figure 1: Architecture of AlphaFold model, highlighting the self-attention layers for interpreting the relationships and patterns within protein sequences. Arrows show the information flow.

This paper explores the evolution of LLMs in terms of Transformers, which facilitate natural language understanding. The paper also discusses the current capabilities of LLMs in interpreting textual input into visual output. This further extends to emerging areas such as code generation, interactive systems, and knowledge graphs. This paper aims to establish a straightforward framework for tracking the continuous progress of LLMs and their growing impact in various domains. The highlights of the paper are listed here –

- **Text-to-image model architecture:** Explores model training for converting text to images, highlighting the ‘Prior’ component’s role.
- **Image captioning and interpretation:** Highlights the shift towards using models for translating the semantic content of images into textual descriptions.
- **Broader applications of LLMs:** Details significant market expansion and application diversity of LLMs, noting technological integration.

The interaction between text and images is of key interest in current research. The ability of the NLP engines to generate images from textual prompts, as well as their ability to interpret and analyze images, presents significant challenges and opportunities, as shown in models such as Dall-E[®] *et al* Ramesh et al. [2022], Ding et al. [2021], Khan et al. [2022]. Integration of LLMs is expected to enhance contextual decision-making, respond to unique scenarios, provide ongoing feedback, and facilitate communication with future interactive systems.

In the remainder of this paper, Section 2 details the structure of Transformer and Transformer-based LLMs, focusing on their role in text-image-text interaction. Section 3 overviews the applications of LLMs in various domains. Section 4 explores the fusion of these models with other technologies. Section 5 presents potential future directions and challenges in the engineering and industrial sectors. Finally, Section 6 concludes the paper with a discussion on the ongoing development of LLMs.

2 Transformer Model Structure

The Transformer architecture has had an impact on machine learning in the field of sequential data processing. This section will delve into Transformer’s infrastructure, emphasizing its core principles. A fascinating challenge in the field is the transformation between text and graphical representations. While the transformation from text to graphics is well understood, the real challenge arises when training a Transformer to interpret an image and convey that understanding in text. Insights from generative deep learning, particularly those related to artistic style interpretation, offer potential approaches to this conversion from visual to text. Well-established Knowledge Graphs (KGs) on style learning are recognized as an important reference for this research.

Furthermore, as more models like ChatGPT become available to the general public, which can now recognize the generation of texts, a related question arises: can we extend this deep understanding to visual data? To address this question, we will conduct a comprehensive literature review to understand the current state of research in this area. We will present opportunities to make breakthrough contributions in this potentially emerging field. The integration of ChatGPT’s analytical capabilities with tools such as Perplexity AI forms the core of the research. The overarching goal is to enable converters to interpret visual content and communicate their understanding in textual form.

2.1 Text-to-image Model Architecture

This section delves into the fundamental architecture of these models, elucidates their distinctive features, and offers a comprehensive overview of the developing field of text-to-image generation. For instance, Google’s Disco Diffusion and Imagen, OpenAI’s contrastive language–image pre-training (CLIP) and Dall-E 2, as well as Facebook’s StyleGAN3, are of availability of this advancement, each characterized by unique approaches to converting textual descriptions into visually engaging images. Specifically, Disco Diffusion is adept at producing photo-realistic images, CLIP excels in blending text with visual elements. StyleGAN3, on the other hand, stands out for its style transfer capabilities. The taxonomy of model architectures, as well as their capable tasks, are given in Table 2.

Although with their respective features and architecture, the models have a common ground of fundamentally building on top of the Transformer architecture, which consists of two important parts: the encoder and the decoder. However, when focusing on NLP engines that convert text into images, the model architecture mainly involves two main subcomponents: the encoder (Prior) and the Decoder. As shown in Figure 2, each of these subcomponents is schematized in detail. They are interconnected by a transformer vector space so that the pre-trained model can seamlessly map text inputs to corresponding visual outputs based on its understanding of the text context.

To begin with, the “Prior” is typically a large, unsupervised neural network that is trained on a large corpus of text data Paik and Wang [2021]. It is then used as input to the decoder component of the LLM. The Prior is trained using a variety of techniques, including autoencoders, denoising autoencoders, and variational autoencoders Rothe et al. [2020], Wang et al. [2019], Lewis et al. [2019]. The “Prior” purpose revolves around three main steps:

1. **Contrastive pre-training:** This is the primary stage where the model is trained to discern and extract key features from input texts and input images. The fundamental goal is to equip the model with the capability to distinguish between similar and non-similar pairs of data, which can be text, images, or both. Neelakantan et al. [2022], Li et al. [2023a], Luo et al. [2022]. Within this step, the inputs are *tokenized* into text and image encoders, respectively. Specifically for image, it is divided into smaller patches representing a part of image, where key features such as color, texture, and shape are extracted. The extracted features are then quantified into data format for NLP to process. The association between individual text and image tokens is predefined, facilitating the model in drawing connections and comprehending the correlations between the two Gimpel et al. [2023].
2. **Label set from label text:** Upon extracting features and comprehending textual/visual “cues”, this step creates a structured and labeled dataset, where each entry is associated with its corresponding label, the text-image pairs are a result of this association. The model is hereby trained to anticipate visual content based on these labels Xu et al. [2022], Kung et al. [2023], Leake et al. [2020].
3. **Zero-shot prediction:** This step involves the transition from the training to the prediction phase. The labeled dataset is processed through a text encoder, yielding an output “feature vector”. Given a new, unseen image,

Table 2: The architectures of GPT-based models for their capabilities and task-specific assessments based on their architectural strengths.

Model	Architecture	Cabable Tasks	Versions
Disco Diffusion	iterative refinement of random noise into coherent images through a sequence of learned transformations.	- artistic image gen. - high-detail rendering	-v4; -v5.2; -v5.7
Imagen	multi-stage process based on NN from creating a low-resolution image to progressively refined output.	- high-res. image gen. - detailed scene depict.	-v23; -turbo
CLIP	learns to associate images with captions from vision- and text-transformers in setup.	- image-text matching - zero-shot classifi.	~
Dall-E	modified Transformer handling text and image tokens with autoregressive model.	- creative image synth. - text-based image modifi.	-v2; -v3
StyleGAN	GAN-based architecture on style-based generation for more realistic and controllable synthesis.	- realistic portrait gen. - adv. style transfer	-v2; -v2-ada; -v3

this image is analyzed against the feature vector, resulting in similarity scores that are then translated into textual descriptors or features. This step allows the model to make predictions on new, unseen data without requiring additional training Heyden et al. [2023], Paz-Argaman et al. [2020], Zhang and Saligrama [2015], Deng et al. [2020].

The Prior is important because it provides the model with a pre-trained knowledge of the data domain. By training the model to extract salient features from text and associate them with visual content, the Prior enables the model to make predictions on new, unseen data. This knowledge can be used to generate more realistic and accurate images, even for complex and challenging text descriptions. The output of the Prior model is a multi-dimensional representation vector or text features, It can be thought of as a summary of the text that captures its most important features. This representation encapsulates the statistical relationships between text and images. When it is fed to the Decoder, it acts as a catalyst, guiding the Decoder to generate more realistic and accurate images. Yet before being fed into the Decoder, the text features are going through a set of Transformer vector spaces.

As shown in Figure 2, the typical transition between the Prior and Decoder is achieved using a three-Transformer vector space, this is a progression mechanism that consists of an Encoder, two Encoder-Decoder sets, and a Decoder Zhang et al. [2023]. The output vector from the Prior is directly fed into the foremost Encoder. This Encoder transforms the output vector into a latent image representation, denoted as Z_t . Following this, Z_t is then propelled into a weighted Encoder-Decoder set to compute the subsequent output feature, namely Z_{t+1} . Note that an intrinsic feedback loop is present, in which Z_{t+1} aids in recalibrating Z_t , allowing the model to refine its latent image representation based on the output of the Decoder Hu and He [2019]. This ensures the model generates images that resonate more consistently with the input text description. The weight W here is special because it modulates the balance between the Prior’s output and the feedback loop Sheng et al. [2019], Wan et al. [2021]. A higher value of W will give more weight to the Prior, while a lower value will give more weight to the feedback loop Yuan et al. [2022], Li et al. [2020]. Beyond this point, Z_{t+1} traverses another Encoder-Decoder arrangement to create a provisionally stable output \tilde{Z}_t . This output is then directed into the Decoder of the Transformer vector space, mirroring the first Encoder before entering the latent image space Yuan et al. [2022], Nguyen et al. [2022], Sheng et al. [2019], Wan et al. [2021]. The final output passing through this Decoder is then fed into the self-attention layer of the “Decoder” part of the text-to-image model.

The GPT vision Transformer model is a Decoder-only architecture, which means that it uses only the Decoder component of the Transformer architecture. In contrast, BERT is a bidirectional Encoder that uses both the Encoder and Decoder components of the Transformer architecture Tu et al. [2022]. The core component of the Transformer is the attention

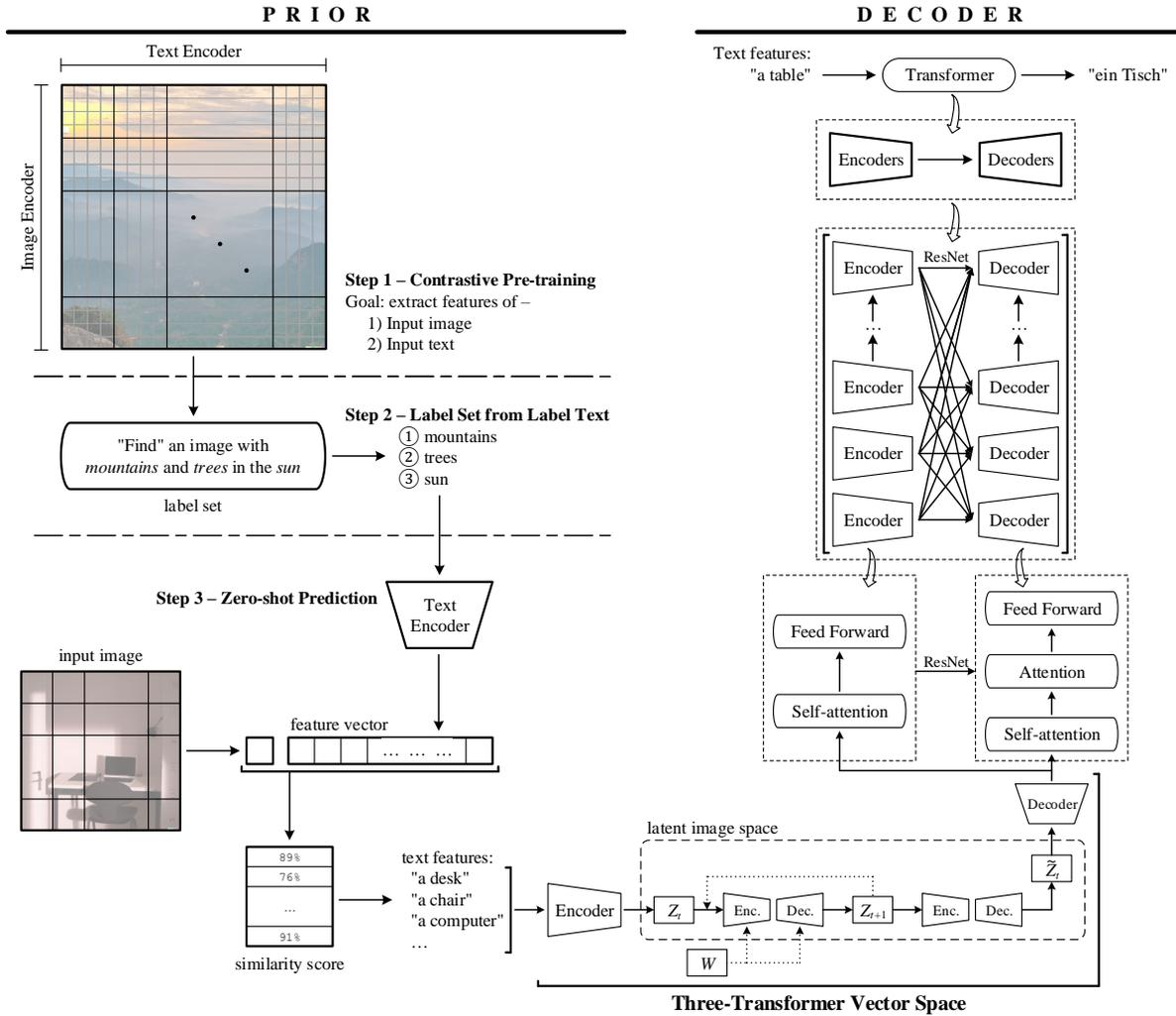


Figure 2: Transformer diagram in order to realize text-to-image conversion. The diagram consists mainly two subparts – a “Prior” and a “Decoder”, which are linked in a cascade manner by the three-Transformer vector space.

mechanism, which operates using three fundamental inputs: the query, keys, and values. The query is the text vector that is generated out of the three-Transformer vector space, while the keys and values are products of the Encoder and Decoder of the latent space. The scoring system in the attention mechanism is determined by taking the dot product of the query and the key. These scores undergo normalization through the Softmax function, transforming them into probabilities. These standardized probabilities dictate the weighted sum of the values. This progression, described as moving from text vectors to the Transformer vector space, eventually leads to the combination of eigenvalues of each weighted result. The primary role of the Transformer vector space is to provide more parameters for training and add weights to different features to dynamically let the model focus on the most relevant part of the input. It is also used to allow the model to learn more complex relationships between the words in the input text. In addition, the Transformer’s use of parallel computing permits the simultaneous processing of multiple sequential data. This parallelism not only contributes to efficiency but also recognizes inter-dependencies in textual input. The latter is achieved by sharing weights between converters, thus promoting a more coherent understanding. The conversion of data from the encoder to the decoder benefits from residual connectivity, which means that the output of the encoder is added to the input of the decoder to prevent gradients from being lost and disappearing.

Compared to traditional residual neural networks and Long Short-Term Memory (RNN+LSTM) networks, the Decoder has a much larger memory capacity, thanks to its deeper and broader structure. This facilitates multi-action processing

and speeds up the overall speedup. Combining parallelism and memory capacity, the transformer can process inputs with longer dependencies. This ensures a coherent representation of sequences, which is crucial when processing artwork, scripts, music, or language. While RNN+LSTMs sometimes have problems generating coherent outputs, Transformer excels, making pivotal improvements in sequence processing tasks.

2.2 Transformer-driven Image Captioning & Interpretation

Understanding the content of images or paintings through AI often has greater practical significance than realizing text-to-image conversion. Art appraisal and artwork psychoanalysis are some notable application Mokady et al. [2021], Nukrai et al. [2022]. There has been a focus on image captioning. CLIP is one of these attempts. The CLIP aims to articulate the semantic content of images into textual descriptions. The operation mechanism of which revolves around representing both images and texts within a unified semantic space Radford et al. [2021], Zhang et al. [2022a]. The “contrastive” in CLIP refers to contrastive learning, which lies within the domain of self-supervised learning where human-annotated labels are no longer required. Instead, it uses data augmentation and transformation to automatically generate labels from the input itself. CLIP is designed to generate vectors embedded in the semantic space of texts and images. In the said space, an image and its corresponding textual description (positive pairs) shall be pulled together closely, while unrelated texts and images (negative pairs) are to be pushed far apart. For a given positive pairing, many other text and image data points in the group can be used as negative examples Huang et al. [2022], Yang et al. [2022]. It is this comparison approach that provides the model with many learning signals and utilizes learning from a large dataset to enable CLIP to infer textual representations from unseen images. A large amount of experimental data suggests that the current tool is very useful Ma et al. [2022]. Nonetheless, it falls short in critical aspects such as understanding image style, grasping artistic nuances, and controlling the overall sentiment conveyed by the image. Here are some of the limitations of CLIP:

- **Quality sensitivity:** The model demonstrates superior performance with high-resolution images but may falter with low-quality counterparts, producing inaccurate descriptions Radford et al. [2021], Shi et al. [2022].
- **Abstract interpretations:** While CLIP efficiently processes real-world images, it may struggle to deliver precise descriptions for abstract images Radford et al. [2021], Zitnick et al. [2014].
- **Ambiguity handling:** Images that harbor multiple interpretations may pose a challenge to CLIP, resulting in potential description inaccuracies Radford et al. [2021], Shi et al. [2022].

There has been some progress in image captioning focusing on several aspects: First, image understanding, which is the basis for realizing image-to-text. Second, the development of new text generation techniques to improve the model ability to generate more realistic and creative textual descriptions. Third, robust models to meet challenges of low-quality images or images with multiple interpretations Gabajiwala et al. [2021], Kuo and Kira [2022]. Multiple approaches have emerged to address these challenges, each with unique contributions.

2.2.1 Generative Adversarial Network

Generative adversarial network (GAN) methods are becoming increasingly popular as they generate high-quality natural language descriptions they produce based on images. Generators and discriminators in those models generate text, with the latter distinguishing between real descriptions and created descriptions Brown et al. [2020], Hu et al. [2022]. Notably in Hu et al. [2022], a Large-scale iMAGE captiONer, or LEMON model, is combined with Transformer architecture to set new benchmarks in large-scale image captioning. Utilizing a vision-language pre-training approach, the model is trained on a massive web-collected dataset with image alt attributes. The model’s scalability and robustness are validated through empirical results, with sizes ranging from 13 to 675 million parameters, and it establishes new state-of-the-art performances on key image captioning benchmarks.

Additionally, other well-known image generators represent significant advancements in this field. Dall-E, in particular Dall-E 3, has gained recognition for its ability to generate detailed and contextually relevant images from textual descriptions. The Dall-E model has been trained on datasets incorporating hundreds of millions of images from the internet and licensed libraries, learning visual concepts by associating words from image descriptions with the images. The text-to-image synthesis translates textual descriptions into visual imagery with improved training using GPT-4V-generated captions. This approach allows the model to understand the text better and create images that more accurately represent the user’s prompt. In terms of applications, Dall-E 3 is used across various domains including art and design, advertising, fashion, and scientific visualization. The model also plays a pivotal role in the entertainment industry for generating scenes and characters in virtual prototyping. This model, along with others in the same category, underscores the diversity and capability of current generative models in handling complex image and language tasks.

2.2.2 Attention Mechanisms & Transformers

Transformers are commonly defined with attention layer(s), addressing the text generation aspect by facilitating a focused understanding of image features Huang et al. [2019], Al-Malla et al. [2022], Vaswani et al. [2017], Fei [2022]. For instance, the Transformer proposed in Fei [2022] tackles the “deviated focus” issue in attention mechanisms by utilizing a perturbation-based self-supervised approach. This method, also known as A2 Transformer, dynamically adjusts attention weights through a learnable network. By applying mask operations to disturb original attention weights and evaluating the performance impact, the authors identify crucial image regions for effective captioning. The model employs four fusion techniques—max pooling, moving average, exponential decay, and gate mechanism—to refine the attention weights Kou et al. [2023]. The proposed approach outperforms standard baselines on the MS COCO dataset in both automated metrics and human evaluations.

While traditionally defined so, recent findings suggest that attention layers may not be as crucial as previously thought, and a sufficiently large model can learn effectively without it. For instance, In Liu et al. [2021], a gated multi-layer perceptron (gMLP) model is introduced. The gMLP model is an alternative Transformer model built from MLP layers with multiplicative gating, yet without self-attention. This model performs comparably to the trending Transformers in vision and NLP tasks. In Bensouilah et al. [2023], Wang et al. [2023a], the different research groups demonstrate that, by combining convolutional and RNN (CRNN) layers with gMLP, the model can achieve high performance for handwritten text recognition without the need for attention mechanisms. Specifically, the gMLP captures dependencies between spatial locations in the feature maps through its use of channel projections and gating, this allows the model to focus on the most relevant parts of the image input when predicting each character in the decoded output sequence. Without attention, the spatial interactions from the gMLP provide the model with enough indication of relevant input parts to perform well. Across experiments on standard handwriting datasets, the proposed CRNN-gMLP model outperforms competitive attention-based methods, showing the sufficiency of the CNN and gMLP building blocks.

2.2.3 Hybrid Approaches

Hybrid approaches in image-to-text transformation represent a blend of different computational strategies, supporting the strengths of various architectures and methodologies to improve performance and adaptability. These approaches typically combine elements such as pre-trained models, retrieval mechanisms, and reinforcement learning techniques to create systems that are robust, flexible, and efficient. An example of hybrid models is SmallCap, introduced in Ramos et al. [2023]. SmallCap integrates the capabilities of a pre-trained CLIP encoder coupled with a GPT language model, with an emphasis on the cross-attention layers. Such an approach not only optimizes computational efficiency but also ensures adaptability across various domains. The amalgamation of these techniques underscores the potential and efficiency of hybrid strategies in the field of image-to-text transformation Zhao et al. [2019].

Furthermore, the integration of reinforcement learning into hybrid models introduces a dynamic learning paradigm, as seen in Ren et al. [2017], Yan [2023], Huang and Li [2022]. Specifically in Ren et al. [2017], Zhou *et al.* employ deep reinforcement learning, featuring a dual-network decision-making framework consisting of a “policy network” for local guidance and a “value network” for global evaluation. Both networks are initially pre-trained using standard supervised learning with cross-entropy loss, and the value is evaluated using mean squared loss. Subsequently, both networks undergo fine-tuning through reinforcement learning, guided by a novel reward function based on visual-semantic embedding. This approach outperforms existing state-of-the-art methods on the Microsoft COCO dataset across various evaluation metrics, including BLEU and CIDEr Cherukuri et al. [2022], Vedantam et al. [2015].

Hybrid models can vary in their structure that combining the outputs of various models to more complex integrations involving multiple learning methods. For instance, some hybrid models may integrate unsupervised learning components to enhance the pre-training of language models, while others might leverage meta-learning techniques to optimize performance across multiple tasks. The taxonomy of hybrid models includes, but is not limited to, the following categories:

- **Ensemble-based hybrids:** Outputs are combined of multiple different models, accounting for their strengths to improve overall performance.
- **Multi-paradigm hybrids:** Models integrate various learning paradigms to benefit from the respective advantages.
- **Cross-modal hybrids:** Focusing on integrating and processing information from different modalities, the model is designed for comprehensive and contextually aware contents.

Table 3: Market value and growth projections of LLMs, NLP, and related technologies in the realm of Generative AI.

Technologies	Market Value (B)			CAGR	Time Frame
	2022	2023	Projected		
Generative AI (Total)	\$8.65	–	\$188.62	36.10%	2022 – 2032
- Incl. NLP	\$19.68	\$24.10	\$112.28	24.60%	2023 – 2030
- Incl. LLM	–	\$0.04	–	6.20%	2023 – 2030
Other realm	–	\$11.30	\$51.80	–	2023 – 2028

3 Versatility of LLMs Across Domains

LLMs were initially recognized for their role in NLP, but their applications have now expanded into computer vision, image synthesis, and code design Lin et al. [2023], Wang et al. [2023b]. Market value forecasts for generative AI, LLMs, NLP, and other related Transformer technology areas are shown in Table 3 Bilan [2023], Fortune Business Insights/Technology [2023], The Market Publicist [2023], Global Info Research [2023], Digital Journal / Newsmantraa [2023], Infinity Business Insights [2023]. Specifically, LLM is projected to grow from \$8.65 billion in 2022 to \$188.62 billion by 2032, at a compound annual growth rate (CAGR) of 36.10%.

A closer look at the LLM market reveals a favorable growth trend and a promising future. In addition to LLMs, the broader AI market is also trending upward. This growth highlights the wide range of applications and growing demand for LLMs in a variety of fields, setting the stage for the discussion that follows on the diverse applications of LLMs.

This section explores the foundational contributions of LLMs in these fields, detailing how they connect text and visual narratives, interpret visual data, and are poised to redefine traditional computational paradigms. This discussion is intended to provide insight into the potential for LLMs to have a significant impact on a variety of technical fields.

3.1 Transcending Boundaries in NLP

In LLM, converter integration makes it possible to recognize complex relationships and dependencies in the input data. This capability not only bridges the gap between simple computational tasks and in-depth language understanding, but also improves the accuracy and quality of machine translation and semantic text evaluation Singh and Mahmood [2021], Kang et al. [2021]. The ability to capture long-distance dependencies in sentences improves the depth of translation and ensures a balance between syntactic structure and contextual meaning. the effectiveness of such approaches is presented in Guo [2022], which develops a deep learning-assisted semantic text analysis method for detecting human emotions from the text. This study emphasizes the potential of NLP techniques for sentiment detection. By utilizing word embeddings, which are critical for numerous NLP applications such as machine translation (MT) and sentiment analysis (SA), the technique captures the semantic and structural complexity of the text. The results of this approach were a significant human sentiment detection rate of 97.22% and classification accuracy of 98.02%.

Building upon this foundation, modern MT and SA have revolutionized the field with the aid of deep neural networks. The Transformer model has enhanced the capability to understand and translate longer and more complex sentences. Similarly, in SA, models like BERT and its variants have transformed the landscape, enabling a more nuanced understanding of sentiment by considering the broader context in which words appear. These models excel in tasks ranging from the overall sentiment of lengthy documents to fine-grained, aspect-based SA, catering to the intricate requirements of modern-day applications. The convergence of these technologies with LLMs has led to a synergistic effect, resulting in systems that are not only proficient in specific tasks but also display adaptability and understanding of human language.

Transformer stands out in question answering and document retrieval, especially in its ability to decode and represent complex semantic nuances as vectors. This capability not only improves the accuracy of matching user queries to relevant documents but also demonstrates a deep understanding of different topics and contexts. The development of retrieval technologies has further emphasized this capability. While early document retrieval systems primarily used sparse representations such as Term Frequency – Inverse Document Frequency, modern models like Dense Retriever now utilize dense vector representations. These models have shown increased efficacy, especially in pinpointing relevant passages for open-domain question-answering tasks. In a study in Zhang et al. [2022b], a multi-view document representation learning framework is proposed that addresses the limitations of a single vector representation in satisfying multi-view queries. Another research introduced the Distill-VQ approach to integrate vector quantization learning into a knowledge refinement framework aiming at evaluating query-document relevance using dense embeddings Xiao

et al. [2022]. The paragraph aggregation retrieval model, detailed in Althammer et al. [2022], is tailored for dense document-to-document retrieval, amalgamated rank-based and topical aggregation, relying on dense embeddings.

3.2 Navigating Through Text & Image Synthesis

In the field of computation, Transformer architectures, especially within LLMs, play a crucial role in connecting textual information with visual creations. Models built on GPTs have demonstrated their ability to generate a wide range of text types, from structured communications such as email to more imaginative and artistic output.

A particular area of interest in recent research is the “uncanny valley” observed in textual contexts. A study by Hazan examines AI-generated visuals, comparing them to eerily similar doppelgängers navigating this uncanny valley Hazan [2023]. The survey also highlights emerging challenges, including issues related to copyright and origin that are prevalent in the creative field. In a related study, Li *et al.* introduced a design approach for embodied conversational agents, incorporating a range of AI models, with GPT-3 being a significant component Li and Xu [2023]. Their research addresses key factors such as content suitability for different age groups, protection of children’s privacy, gender representation in ECAs, and the broader effects of the uncanny valley.

In terms of image synthesis, models such as Dall-E and Midjourney emphasize the morphing capabilities of morpher architectures. For example, Dall-E’s strengths are seen in generating detailed images based on textual cues, and capturing complex details and textures Chang et al. [2023]. The model employs discrete variational autoencoder for tokenizing images and text, which allows for efficient processing of the input data. The tokenized data then passes through a series of Transformer layers, enabling the model to generate highly detailed and contextually relevant images from textual descriptions. Similarly, Imagen, designed for text-to-image translation, sets the standard for photo-realism and language comprehension, marking progress in harmonizing text and image interactions. Unlike traditional GAN that directly map noise vectors to images, Imagen adopts a Transformer-based approach, utilizing language model for text encoding followed by a diffusion model for image synthesis. This two-step process involves first understanding the text at a granular level and then translating this understanding into a coherent visual representation. The diffusion model in Imagen contributes to the model’s ability to produce photo-realistic images that are both contextually aligned with the input text and visually compelling. These models highlight the adaptability and efficacy of Transformer architectures in capturing textual nuances and translating them into visual outputs. The contribution lies not just in the ability to generate images but in the architectural advancements that allow for a sophisticated interplay between text understanding and image synthesis.

In the realm of literature, the incorporation of Transformers has attracted the interest of writers and e-literature scholars. Their efforts have resulted in creations that blend literary expressions with visual elements through iterative processes. In Rettberg et al. [2023], Rettberg *et al.* emphasizes the literary origins of these combined forms.

3.3 LLMs into Computer Vision

The capabilities of LLMs now encompass the interpretation of visual content and the nuances of code semantics. The push to integrate LLMs into these areas stems from an increasing demand for precise, swift, and context-sensitive computational methods. Driven by the growing complexities of urban environments and transportation systems, heightened demand for accurate, fast, and context-aware solutions is observed Ogiela and Ogiela [2009], Saharia et al. [2022], Alverson and Yamamoto [2016]. In response to this demand, LLMs are recognized for their ability to transform vast data streams into valuable and actionable information.

Given the capacity of LLMs to handle and scrutinize large datasets, LLMs play a key role in areas like image classification and object detection Amin et al. [2023], Kleesiek et al. [2023]. This is applicable in 1) autonomous vehicle navigation, where real-time object detection is essential for safety; and 2) medical imaging, where accurate medical image classification is pivotal for early disease detection and subsequent treatment strategies Chang et al. [2022]. Additionally, the potential of LLMs in image enhancement, especially in refining low-resolution images, has been explored. The objective is clear: improve image quality without distorting the original content or introducing unwanted artifacts. With their deep learning foundations, LLMs have indicated their capability in maintaining this delicate balance Amin et al. [2023], becoming essential tools in sectors like digital forensics and film restoration Liu and Shen [2022], Suthar et al. [2023].

The coding nature with its layered semantics and strict syntax, demands a nuanced approach. LLMs, trained on extensive codebases, can generate coherent and contextually relevant code snippets, aligning with provided problem statements or descriptions. This automated code generation can expedite software development processes, reducing manual coding efforts Wang et al. [2023c], Li et al. [2023b], Joshi et al. [2023]. Moreover, the ability of LLMs to navigate through code as if it were natural language is transformative. By understanding the underlying logic and

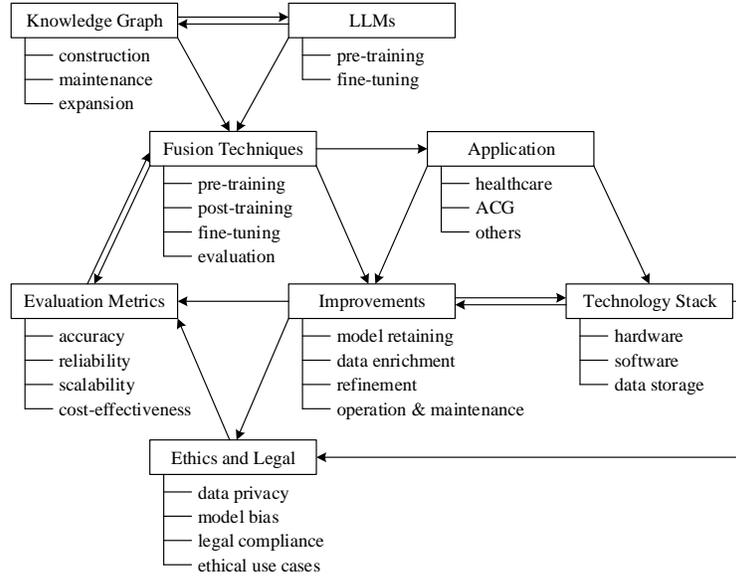


Figure 3: A diagram combining knowledge graph, LLMs, as well as their applications and constraints. Note that in this figure, the ACG is short for automated content generation.

semantics, LLMs can assist developers in debugging, identifying potential pitfalls, and suggesting code optimizations. This not only enhances the efficiency of the development process but also ensures the robustness of the final software product Xu et al. [2023], Wang et al. [2023c].

LLMs in Code Semantics directly roll out to software development. Consider tasks like code review, LLMs can assist in suggesting code enhancements and pinpointing inconsistencies, this enables developers to retrieve pertinent code snippets using natural language prompts Xu et al. [2023], Li et al. [2023b]. Given the compatibilities of LLMs with various interfaces, it is observed that LLMs can be smoothly integrated into development processes, equipping developers with a toolkit enriched by recent machine learning innovations.

4 Fusion Technologies with Synthesis of LLMs

Building on discussions in Section 3, the strength of LLMs in managing and crafting text that mirrors human-like quality underscores their potential in diverse tech solutions. A key area of exploration is the synergy of LLMs with other tech innovations. In particular, combining the Transformer system with big data and neural networks represents a way forward in enhancing LLM capabilities beyond computer science as highlighted by Tu et al. [2022]. This section is aimed at highlighting these blended techniques that involve LLMs including how they operate, where they are used and what they bring forward.

4.1 LLMs with Knowledge Graph

The integration of LLMs with knowledge graphs represents a significant advancement in data science, as demonstrated in Figure 3. The synergistic use of this LLM and the organization of knowledge graphs provide a strong system of improving data reuseability within different application domains. Fusion which plays a crucial role in building domain-specific content gives rise to new opportunities for a variety of application areas.

Commencing with an array of domain-specific corpora, the process is illustrated in Figure 4. Different corpora are labeled according to their respective domains. They feed into a series of pre-training frameworks, depicted as cascaded interconnected blocks. These blocks symbolize essential steps especially “pre-processing” and “tokenization,” which prepare the data for the next stage.

At the core of this integration is the LLM component block. The processed data is transformed into a state that fits LLM training. Following this is the “Knowledge Graph Integration” phase. Both LLMs and KG integration become dual in

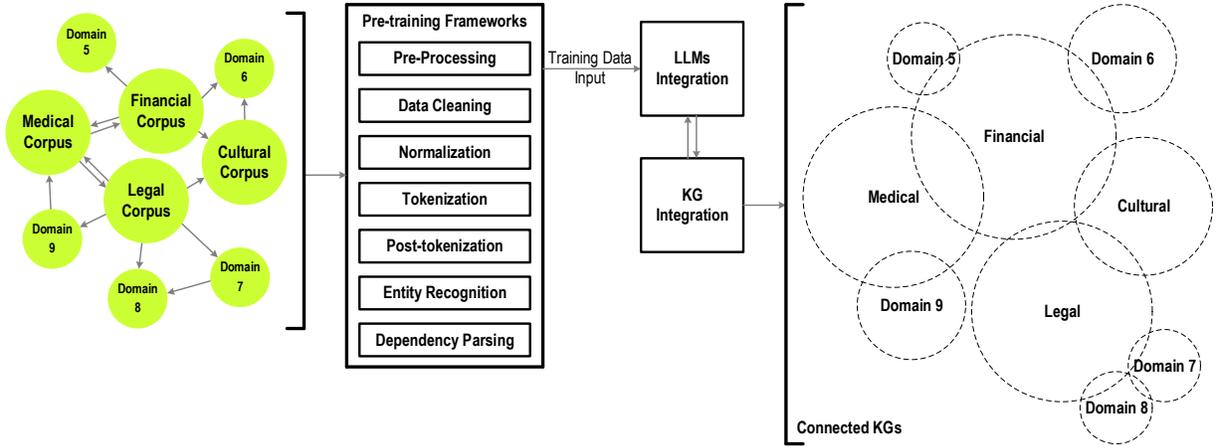


Figure 4: Domain-specific corpora to connected knowledge graph via LLMs training.

contributing to and deriving insights from each other, either knowledge graphs, or interpreted data. The bidirectional arrows between LLMs and KGs are shown in the figure to represent this interaction.

The culmination of this process is represented in the “connected knowledge graph” phase. This final stage is illustrated as a series of application domains, where each dot-line circle represents a specific knowledge graph of an individual domain, and the intersections between domains represent the interdisciplinary connections. The connected KG denotes the deployment of this integrated system in practical scenarios with effective guidelines from the beginning of the corpora stage to the integration process. This integration encapsulates the field’s continuous evolution, considering emerging evaluation metrics, technological advancements, important ethical and legal considerations, *etc.*

In Gao et al. [2023], a medical KG namely “Dr. Knows” was purposed in serving for automated diagnosis generation. The LLMs have been tuned in the medical/biological domains with the unified medical language system as text corpora. The NLP engine employed prioritizes the generation of diagnostic results, particularly highlighting its “explainable diagnostic pathway” feature Aracena et al. [2022]. This design ensures that medical professionals can effortlessly trace the origins of a patient’s ailment and confidently validate the accuracy of the results produced by the LLMs. The amalgamation of KGs and LLMs, particularly in this case, yields several benefits, including:

1. **Enhanced diagnostic proficiency:** The knowledge graph assists the LLM in adeptly interpreting and condensing complex medical terminologies, leading to improved diagnostic accuracy.
2. **Navigating electronic health records (EHRs):** This synergy addresses the convoluted presence in clinical narratives within EHRs, streamlining the extraction and understanding of patient information Gonen et al. [2022].
3. **Safety and optimization:** The combined approach delivers highly accurate automated diagnoses. This not only bolsters the reliability of medical assessments but also prioritizes patient safety, paving the way for superior healthcare outcomes.

Besides “Dr. Knows”, another practical application involving medical KGs and LLMs pertains to depression diagnosis and treatment Wang et al. [2023d]. Wang *et al.* initiated their methodology by creating a comprehensive pre-training database focused on the depression domain. The LLMs are hereby equipped with relevant knowledge for constructing a depression-centric knowledge graph. As experts intervened, it was refined and optimized using the LLM, ensuring improved performance in the targeted domain. This differs from “Dr. Knows” by Gao *et al.* such that the knowledge graphs are utilized to generate instructions, and the LLMs are fine-tuned using these instructions derived from the knowledge graph. This involution process improves the NLP engine’s performance in depression diagnosis and treatment, illustrating a symbiotic relationship between structured, domain-specific knowledge and adaptive, generative modeling. It is fine-tuned using the relationships among depression-related concepts Dai et al. [2023].

4.2 LLMs with Interactive Systems

Implementing interactive systems presents a significant challenge due to the necessity for the system to generate coherent responses to multimodal input. Earlier AI-driven NLP tools, taking Amazon Alexa as an example, faced

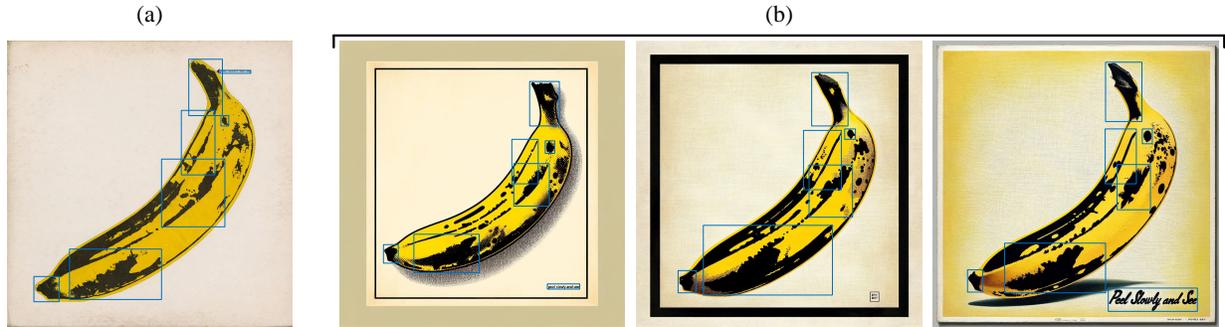


Figure 5: Dall-E 3 drawing based on the description given by GPT4-V.

difficulties in effectively handling sentences that combined multiple languages. Given a query like “What is your *Lieblingsfarbe?*” — a mix of English and German — the response was often restricted to “I don’t know that one.” Such challenges arise from the difficulties of navigating diverse language datasets. The task becomes even more complicated when dealing with varied modalities, including both text and visuals.

Considering the interaction between text and image, the difficulty isn’t just about object recognition but also understanding styles and extracting deeper visual meanings. In this context, the GPT-4 with vision (GPT-4V) model, as discussed in OpenAI [2023a,b], merges the capabilities of text-focused and vision-focused models. This model is trained on a combination of textual and visual data and is further optimized using the reinforcement learning from the human feedback approach, ensuring the results resonate with human expectations Ouyang et al. [2022]. The GPT-4V’s capability to handle multimodal inputs is achieved through a systematic process, specifically –

- Extract image features using the image recognition model. These features encompass pixel values, semantic information, *etc.*
- Derive textual features using NLP models. These encapsulate word vectors, syntactic structures, *etc.*
- Fusion the extracted image and text features. The integrated features enable GPT-4V to achieve a comprehensive understanding of the context and generate linguistically coherent captions.

By integrating interactive modes into LLMs, GPT-4V can analyze user-provided visual content, address a broader range of tasks, and enhance user experiences.

The use of GPT-4V in a real-time analysis showed an interesting result during a Dall-E 3 test. The famous banana painting by Andy Warhol, as shown in Figure 5(a), was fed into GPT-4V, requesting a detailed description. Then the GPT-generated text was fed back into Dall-E 3, which results in a recreated image, as depicted in Figure 5(b). After a detailed cross-comparison between the three generated images and the original painting, similarities have been marked in the form of bounding boxes across the image, representing the resemblance of the GPT generated with the original painting. The advancements in GPT-4V’s image captioning capabilities can be attributed to 1) its enhanced model parameters, 2) the adoption of the Transformer-XL architecture for improved sequence data handling, and 3) innovative training methods that reduce biases. Furthermore, its multimodal understanding, which mirrors human cognition, is a testament to its sophistication in design.

GPT-4V’s capability to process multimodal input furnishes it with distinct advantages. This model not only achieves a comprehensive understanding of visual scenes, positioning them for tasks like image captioning but also ensures the generation of linguistically coherent and accurate captions Qi et al. [2023]. Central to GPT-4V’s efficiency is its adoption of the Transformer-XL architecture, optimized for handling sequential data.

The Transformer-XL architecture comprises several key components Dai et al. [2019]:

- **Encoder:** composed of multiple self-attention layers, each self-attention layer incorporates relative positional encoding to capture the relative positions of elements in the sequence.
- **Decoder:** consists of multiple self-attention layers and feed-forward networks. While the self-attention layers focus on understanding the interrelations between the input sequence and hidden state sequence, the feed-forward networks transform the hidden state sequence.

- **Relative positional encoding:** captures the relative positions of elements within an input sequence. regardless of sequence length, the relative positional encoding remains consistent irrespective of sequence length, ensuring a more efficient model.
- **Recurrence mechanism:** facilitates information transfer between subsequences. It enhances the model’s ability to understand contextual relationships in longer sequences.

4.3 Applied Mathematics in LLM-based Modeling

In complex system mathematical modeling, the inclusion of LLMs brings about improved memory and analytical prowess. One significant role of LLMs is to support model interpretation and validation. By generating natural language explanations for simulation outcomes, a deeper grasp of the model’s predictive tendencies is achieved. Furthermore, the interactive nature of LLMs allows for inquiries about unexpected results or inconsistencies, with the LLMs offering insights for model enhancement. This interaction creates a feedback mechanism between humans and machines, beneficial throughout the model’s lifecycle, enhancing model precision Nemati et al. [2011], Wu et al. [2022]. Table 5 provides a comparative analysis. This table highlights key areas where LLMs contribute significantly to enhancing mathematical modeling, particularly in aspects such as memory capacity, analytical capabilities, model interpretation, and validation processes. Through this comparison, it becomes evident that the incorporation of LLMs not only augments the fundamental strengths of mathematical models but also addresses some of their traditional limitations, leading to a substantial improvement in overall model efficacy.

Considering Lithium-Sulfur (Li-S) battery cells, there’s potential for LLMs in their mathematical modeling. Recognized as potential successors for high-energy-density storage, Li-S batteries, however, grapple with issues like the “shuttle effect” Zhou et al. [2021], Wang et al. [2022], Sun et al. [2022]. This challenge, resulting from the movement of dissolved poly-sulfides, affects the battery’s stability and efficiency. In addressing this, mathematical models and LLMs can jointly aid in parameterization, data interpretation, and testing.

Initially, LLMs can help in collating vast research data and in mathematical model parameterization, pinpointing parameter relationships from experimental findings. By determining optimal parameter ranges, model reliability is ensured, and crucial variables influencing the shuttle effect are highlighted Zhou et al. [2021].

In data interpretation, while mathematical tools handle most analyses, correlating data with physical phenomena remains essential. For instance, lithium deposition and sulfur loss, contributing to battery degradation, are challenging to quantify directly due to their localized and dynamic nature. However, through data analysis, patterns indicating the degradation can be identified. LLMs can then correlate data from various sources, using relationships to estimate degradation extents Hu et al. [2023], Jiyane et al. [2023], Cheng et al. [2022], Klein et al. [2017].

In the battery modeling testing phase, mathematical modeling and LLMs work in tandem. While the former provides a structured approach based on physical principles, LLMs interpret results, detect data anomalies, and guide model refinement. This collaboration is further detailed in Table 4.

5 Challenges & Opportunities of LLMs

The rapid advancements of models such as GPT mark a new wave of potential. With a careful balance of adaptability and precision, the road ahead for LLMs is filled with both promising opportunities and complex challenges. This section will examine future research directions and the potential effects of LLMs in different fields.

5.1 Prompt Engineering & Program Automation

In the progression of LLM development, models such as GPT variants signify a pivotal shift towards program automation. GPT-3, presented by OpenAI, showcased capabilities in few-shot learning, eliminating the need for extensive fine-tuning Brown et al. [2020], Helmuth and Kelly [2022], Helmuth and Abdelhady [2020]. This progress allows developers to give high-level directives, and in response, LLMs autonomously produce the relevant code.

The research on program synthesis using LLMs also hints at future innovations in software development such as converting natural language descriptions into codes Austin et al. [2021]. This is the reason why the making of software development becomes less complicated by using this kind of technology. Even with limited programming skills, LLMs make software creation more inclusive to participate in. This shows that transitioning to a new state is possible. One such case is Codex by OpenAI. Codex shows that it is possible to mix programming with natural language and utilize LLMs for code-oriented assignments Bommasani et al. [2021]. This evolution begins with a detailed process of creating questions for the computer. Particular concern has been paid to the significance of a context in understanding uncertain requests mainly due to precise coding. This capability significantly simplifies the software development process, making

Table 4: Synergy between mathematical modeling and LLMs in simulation testing of Li-S model for shuttle effect investigation.

Synergy	Math Model	LLMs
Simulation & Testing	preset parameters for different temperatures, charge/discharge rates, electrolyte compositions, <i>etc.</i>	recommend relevant scenarios or conditions to test based on current research trends and known challenges to the specific Li-S battery.
Error Analysis	quantify discrepancies between the model’s predictions and experimental outcomes.	provide insights into the origins of these discrepancies, potentially linking them to recognized challenges or phenomena associated with Li-S batteries.
Documentation & Reporting	produces quantitative results, visual representations like graphs and charts.	facilitates the creation of detailed reports, summarizes key findings, and contextualizes results by referencing pertinent literature.

Table 5: Comparative analysis of mathematical modeling and LLMs

Aspect	Mathematical Modeling	LLMs
Capabilities	<ul style="list-style-type: none"> - Precise, quantifiable predictions - Suitable for well-defined, logical problems - Can model complex systems with known variables 	<ul style="list-style-type: none"> - Generates human-like text - Capable of processing large datasets - Excels in pattern recognition and predictive tasks
Strengths	<ul style="list-style-type: none"> - High accuracy in controlled environments - Well-established in scientific and engineering disciplines - Offers clear, deterministic solutions 	<ul style="list-style-type: none"> - Flexible in handling diverse and unstructured data - Continuously learning and updating from new data - Can generate creative and novel solutions
Limitations	<ul style="list-style-type: none"> - Requires precise and complete data - Often not flexible to changes in problem parameters - Limited in handling ambiguous or unstructured data 	<ul style="list-style-type: none"> - Dependent on the quality and size of training data - Possible biases in generated content - Explanations for decisions can be unclear
Simulation	<ul style="list-style-type: none"> - Effective in simulating physical and engineering systems - Based on deterministic and statistical models 	<ul style="list-style-type: none"> - Less effective for simulating physical systems - Better suited for simulating human-like interactions
Analysis	<ul style="list-style-type: none"> - Robust in logical and numerical data analysis - Can be limited by complexity and data availability 	<ul style="list-style-type: none"> - Strong in linguistic and content analysis - Can struggle with highly technical or niche topics
Model Refinement	<ul style="list-style-type: none"> - Requires manual adjustments and expertise - Often based on hypothesis testing and validation 	<ul style="list-style-type: none"> - Can self-improve through machine learning techniques - Limited by the scope of training data and initial design

it more accessible and inclusive, especially for those with limited programming expertise. Recently, it has been critical to integrate neural network algorithms into traditional programming methods Feng et al. [2020]. Such synergy leads to more sympathy between human and machine cooperation, empowering a machine to understand and accommodate human commands.

Technology develops simultaneously with benefits and challenges. This is a comprehensive discussion of the foundation models, including the LLMs, highlighting their huge capacity while showing where additional measures should be put in place for their responsible applications Roper [2022]. Despite the promising trajectory of this field, it is important to recognize that this area of research is not merely a prospective future direction but an active and dynamic field of study. The development of technology, accompanied by its benefits and challenges, demands a nuanced discussion regarding the foundational models, including LLMs. The active exploration and continuous development in prompt engineering and program automation signifies the field’s vibrancy and the ongoing need for innovative solutions and ethical considerations.

5.2 Deep Understanding in AI Systems

In earlier stages of AI development, systems such as DeepBlue were tailored for particular tasks, notably chess. Created by IBM, DeepBlue marked a milestone by defeating the world chess champion, Garry Kasparov, in 1997. Its design was rooted in deterministic algorithms, enabling the evaluation of millions of chess positions swiftly. However, these algorithms were chess-centric, utilizing predefined strategies and rule-based methods. DeepBlue’s chess accomplishments were indeed remarkable, but its deterministic foundation revealed constraints. It was evident that it couldn’t adapt to tasks beyond its designated domain. This highlighted the emerging need for AI systems with greater adaptability.

Transitioning to present-day advancements, LLMs, exemplified by models like GPT-4, signify a shift in AI methodologies. Distinct from deterministic systems, LLMs leverage neural networks and extensive data training. This equips them to handle diverse scenarios, even those not explicitly covered during training. A notable feature of these modern AI models is their capability to grasp the underlying intent of a query Baktash and Dawodi [2023], Yan et al. [2021]. Instead of solely depending on set rules, context, and query nuances are analyzed to generate apt responses. Such an ability to discern intent becomes vital in areas like natural language processing, where context often drives meaning Baktash and Dawodi [2023].

To encapsulate, while systems like DeepBlue demonstrated specialized expertise, the trajectory of AI has evolved towards models that are both adaptive and context-sensitive. Current models, grounded in neural networks and enriched by vast datasets, have enhanced AI’s capacity to address a broader spectrum of queries.

5.2.1 Interplay of Multiple Modalities

There is a need for combining diverse data types such as text, audio and visuals for the next LLM generations. This is not simply about handling different data, but understanding the inter-connections of the context among them. Consider a dialogue encompassing spoken language, visual cues, and ambient sounds. A refined LLM ought to interpret the dialogue’s tone, decode visual gestures, and discern relevant ambient noises Wu et al. [2023], Binz and Schulz [2023], Chang [2023]. Therefore, an improved LLM should be able to read what is not said, pick up on visual cues, and also pick up important background sounds. A holistic understanding of such may lead to effective and appropriate interactions. In Vaswani et al. [2017], Vaswani *et al.* examined Transformer architecture which opened doors for research into a large expanse of the undiscovered arena.

5.2.2 Sifting and Deciphering Central Data

In the current data-abundant era, the skill to navigate through extensive datasets and pinpoint crucial insights is paramount. LLMs can be tailored to detect central themes, bridge isolated data fragments, and construct an integrated knowledge framework Sharma and Feldman [2023]. This goes beyond mere text summarization to encompass broader contexts, interrelations among data segments, and the ramifications of such linkages Roumeliotis and Tselikas [2023], Hasselgren and Oprea [2023]. For example, when analyzing a scientific article, an LLM should discern the core proposition, the experimental approach, the derived outcomes, and the ensuing inferences, subsequently situating these within the larger academic context.

5.2.3 Broadening Application Horizons

The path of LLMs and transformer models has grown widely in use, showing major changes across many fields. These models, equipped with sophisticated features, are set to cross-reference extensive data, providing a solid countermeasure

against potential data discrepancies. In Fidas et al. [2023], the critical role of LLMs in bolstering the authenticity of academic research is highlighted, underscoring the imperative of ensuring integrity in scholarly work. Additionally, LLMs present a solution to documentation challenges faced by researchers who are not native English speakers.

Additionally, the role of AI in shaping academic knowledge databases ensures the swift spread and uptake of groundbreaking technologies, minimizing academic overlap. The aforementioned advancements underscore the expansive potential of LLMs and transformer models in the academic sector. In Issah et al. [2023], the potential of LLMs to enhance the quality of the peer review process through AI integration is explored.

However, as with any technological advancement, they present a series of unresolved questions that necessitate deeper exploration. Ethical challenges arise from the extensive automation of programs using LLMs, especially concerning originality and authorship. Enhancing the precision and trustworthiness of LLM-generated code remains a priority, and strategies need to be developed to train LLMs to effectively grasp context from multiple data types. Furthermore, there are obstacles exist in ensuring unbiased data extraction and interpretation by LLMs, and efforts are needed to achieve transparency and interpretability in LLM operations.

As the AI domain continues its evolution, these models are poised to significantly influence the future across disciplines. The advancements and research efforts emphasize the inherent challenges and opportunities associated with LLMs and transformer models, indicating their potential to redefine traditional processes and elevate standards.

6 Conclusion

Aiming to present an engaging story that encourages exploration and innovation in the realm of language-based models, this survey emphasizes the transformative potential of LLMs across diverse domains, from academia to industry. The expanding landscape of LLMs is explored, focusing on their evolution, structure, and versatile applications. The journey starts with the foundational Transformer model structure and extends to the fusion of LLMs with cutting-edge technologies. The evolution of LLMs has been nothing short of revolutionary in the field of AI. These models, with their capability to emulate human text generation, have paved the way for diverse applications, ranging from interactive systems to content generation. Yet, the forthcoming challenges and opportunities for LLMs appear even more expansive, spanning multimodal configurations, key information discernment, and varied applications.

References

- Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.
- Mukul Singh, José Cambronero, Sumit Gulwani, Vu Le, Carina Negreanu, and Gust Verbruggen. Codefusion: A pre-trained diffusion model for code generation. *arXiv preprint arXiv:2310.17680*, 2023.
- Anis Koubaa. Gpt-4 vs. gpt-3.5: A concise showdown. *TechRxiv preprint:22312330*, 2023.
- Michael G Madden, Bairbre A McNicholas, and John G Laffey. Assessing the usefulness of a large language model to query and summarize unstructured medical notes in intensive care. *Intensive Care Medicine*, pages 1–3, 2023.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- Yubo Ma, Yixin Cao, YongChing Hong, and Aixin Sun. Large language model is not a good few-shot information extractor, but a good reranker for hard samples! *arXiv preprint arXiv:2303.08559*, 2023.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22522–22531, 2023.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.

- Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. Cogview: Mastering text-to-image generation via transformers. *Advances in Neural Information Processing Systems*, 34:19822–19835, 2021.
- Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *ACM computing surveys (CSUR)*, 54(10s):1–41, 2022.
- Incheon Paik and Jun-Wei Wang. Improving text-to-code generation with features of code graph on gpt-2. *Electronics*, 10(21):2706, 2021.
- Sascha Rothe, Shashi Narayan, and Aliaksei Severyn. Leveraging pre-trained checkpoints for sequence generation tasks. *Transactions of the Association for Computational Linguistics*, 8:264–280, 2020.
- Liang Wang, Wei Zhao, Ruoyu Jia, Sujian Li, and Jingming Liu. Denoising based sequence-to-sequence pre-training for text generation. *arXiv preprint arXiv:1908.08206*, 2019.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.
- Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, et al. Text and code embeddings by contrastive pre-training. *arXiv preprint arXiv:2201.10005*, 2022.
- Wei Li, Linchao Zhu, Longyin Wen, and Yi Yang. Decap: Decoding clip latents for zero-shot captioning via text-only training. *arXiv preprint arXiv:2303.03032*, 2023a.
- Ziyang Luo, Yadong Xi, Rongsheng Zhang, and Jing Ma. A frustratingly simple approach for end-to-end image captioning. *arXiv preprint arXiv:2201.12723*, 2022.
- Henner Gimpel, Kristina Hall, Stefan Decker, Torsten Eymann, Luis Lämmerrmann, Alexander Mädche, Maximilian Röglinger, Caroline Ruiner, Manfred Schoch, Mareike Schoop, et al. Unlocking the power of generative ai models and systems such as gpt-4 and chatgpt for higher education: A guide for students and lecturers. Technical report, Hohenheim Discussion Papers in Business, Economics and Social Sciences, 2023.
- Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong Wang. Groupvit: Semantic segmentation emerges from text supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18134–18144, 2022.
- Tiffany H Kung, Morgan Cheatham, Arielle Medenilla, Czarina Sillos, Lorie De Leon, Camille Elepaño, Maria Madriaga, Rimel Aggabao, Giezel Diaz-Candido, James Maningo, et al. Performance of chatgpt on usmle: Potential for ai-assisted medical education using large language models. *PLoS digital health*, 2(2):e0000198, 2023.
- Mackenzie Leake, Hijung Valentina Shin, Joy O Kim, and Maneesh Agrawala. Generating audio-visual slideshows from text articles using word concreteness. In *CHI*, volume 20, pages 25–30, 2020.
- William Heyden, Habib Ullah, M Salman Siddiqui, and Fadi Al Machot. An integral projection-based semantic autoencoder for zero-shot learning. *IEEE Access*, 2023.
- Tzuf Paz-Argaman, Yuval Atzmon, Gal Chechik, and Reut Tsarfaty. Zest: Zero-shot learning from text descriptions using textual similarity and visual summarization. *arXiv preprint arXiv:2010.03276*, 2020.
- Ziming Zhang and Venkatesh Saligrama. Zero-shot learning via semantic similarity embedding. In *Proceedings of the IEEE international conference on computer vision*, pages 4166–4174, 2015.
- Shumin Deng, Ningyu Zhang, Zhanlin Sun, Jiaoyan Chen, and Huajun Chen. When low resource nlp meets unsupervised language model: Meta-pretraining then meta-learning for few-shot text classification (student abstract). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34-10, pages 13773–13774, 2020.
- Chaoning Zhang, Chenshuang Zhang, Sheng Zheng, Yu Qiao, Chenghao Li, Mengchun Zhang, Sumit Kumar Dam, Chu Myaet Thwal, Ye Lin Tun, Le Luang Huy, et al. A complete survey on generative ai (aigc): Is chatgpt from gpt-4 to gpt-5 all you need? *arXiv preprint arXiv:2303.11717*, 2023.
- Haoji Hu and Xiangnan He. Sets2sets: Learning from sequential sets with neural networks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1491–1499, 2019.
- Fenfen Sheng, Zhineng Chen, and Bo Xu. Nrtr: A no-recurrence sequence-to-sequence model for scene text recognition. In *2019 International conference on document analysis and recognition (ICDAR)*, pages 781–786. IEEE, 2019.
- Qi Wan, Haoqin Ji, and Linlin Shen. Self-attention based text knowledge mining for text detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5983–5992, 2021.

- Xin Yuan, Zhe Lin, Jason Kuen, Jianming Zhang, and John Collomosse. Text-to-image generation via implicit visual guidance and hypernetwork. *arXiv preprint arXiv:2208.08493*, 2022.
- Yijun Li, Richard Zhang, Jingwan Lu, and Eli Shechtman. Few-shot image generation with elastic weight consolidation. *arXiv preprint arXiv:2012.02780*, 2020.
- Kha Cong Nguyen, Ryosuke Odate, and Kanemaru Takashi. A robust text image recognition model with domain adaptation and attention mechanisms. In *2022 International Conference on Wavelet Analysis and Pattern Recognition (ICWAPR)*, pages 7–12. IEEE, 2022.
- Haoqin Tu, Zhongliang Yang, Jinshuai Yang, and Yongfeng Huang. Adavae: Exploring adaptive gpt-2s in variational auto-encoders for language modeling. *arXiv preprint arXiv:2205.05862*, 2022.
- Ron Mokady, Amir Hertz, and Amit H Bermano. Clipcap: Clip prefix for image captioning. *arXiv preprint arXiv:2111.09734*, 2021.
- David Nukrai, Ron Mokady, and Amir Globerson. Text-only training for image captioning using noise-injected clip. *arXiv preprint arXiv:2211.00575*, 2022.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- Youyuan Zhang, Juniu Wang, Hao Wu, and Wenjia Xu. Distinctive image captioning via clip guided group optimization. In *European Conference on Computer Vision*, pages 223–238. Springer, 2022a.
- Junchu Huang, Weijie Chen, Shicai Yang, Di Xie, Shiliang Pu, and Yueting Zhuang. Transductive clip with class-conditional contrastive learning. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3858–3862. IEEE, 2022.
- An Yang, Junshu Pan, Junyang Lin, Rui Men, Yichang Zhang, Jingren Zhou, and Chang Zhou. Chinese clip: Contrastive vision-language pretraining in chinese. *arXiv preprint arXiv:2211.01335*, 2022.
- Yiwei Ma, Guohai Xu, Xiaoshuai Sun, Ming Yan, Ji Zhang, and Rongrong Ji. X-clip: End-to-end multi-grained contrastive learning for video-text retrieval. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 638–647, 2022.
- Hengcan Shi, Munawar Hayat, Yicheng Wu, and Jianfei Cai. Proposalclip: Unsupervised open-category object proposal generation via exploiting clip cues. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9611–9620, 2022.
- C Lawrence Zitnick, Ramakrishna Vedantam, and Devi Parikh. Adopting abstract images for semantic scene understanding. *IEEE transactions on pattern analysis and machine intelligence*, 38(4):627–638, 2014.
- Hamza Juzer Gabajiwala, Nishant Kamlesh Jethwa, Parth Bimal Joshi, Arundhati Anurag Mishra, and Prachi Natu. Comprehensive review of various optimization algorithms for image captioning. In *2021 3rd International Conference on Advances in Computing, Communication Control and Networking (ICAC3N)*, pages 1703–1708. IEEE, 2021.
- Chia-Wen Kuo and Zsolt Kira. Beyond a pre-trained object detector: Cross-modal textual and visual context for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17969–17979, 2022.
- Xiaowei Hu, Zhe Gan, Jianfeng Wang, Zhengyuan Yang, Zicheng Liu, Yumao Lu, and Lijuan Wang. Scaling up vision-language pre-training for image captioning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 17980–17989, 2022.
- Lun Huang, Wenmin Wang, Jie Chen, and Xiao-Yong Wei. Attention on attention for image captioning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4634–4643, 2019.
- Muhammad Abdelhadie Al-Malla, Assef Jafar, and Nada Ghneim. Image captioning model using attention and object features to mimic human image understanding. *Journal of Big Data*, 9(1):1–16, 2022.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Zhengcong Fei. Attention-aligned transformer for image captioning. In *proceedings of the AAAI Conference on Artificial Intelligence*, volume 36-1, pages 607–615, 2022.
- Bonan Kou, Shengmai Chen, Zhijie Wang, Lei Ma, and Tianyi Zhang. Is model attention aligned with human attention? an empirical study on large language models for code generation. *arXiv preprint arXiv:2306.01220*, 2023.
- Hanxiao Liu, Zihang Dai, David So, and Quoc V Le. Pay attention to mlps. *Advances in Neural Information Processing Systems*, 34:9204–9215, 2021.

- Mouad Bensouilah, Mokhtar Taffar, and Mohamed Nadjib Zennir. gmlp guided deep networks model for character-based handwritten text transcription. *Multimedia Tools and Applications*, pages 1–19, 2023.
- Xishun Wang, Tong Su, Fang Da, and Xiaodong Yang. Prophnet: Efficient agent-centric motion forecasting with anchor-informed proposals. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21995–22003, 2023a.
- Rita Ramos, Bruno Martins, Desmond Elliott, and Yova Kementchedjieva. Smallcap: lightweight image captioning prompted with retrieval augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2840–2849, 2023.
- Dexin Zhao, Zhi Chang, and Shutao Guo. A multimodal fusion approach for image captioning. *Neurocomputing*, 329: 476–485, 2019.
- Zhou Ren, Xiaoyu Wang, Ning Zhang, Xutao Lv, and Li-Jia Li. Deep reinforcement learning-based image captioning with embedding reward. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 290–298, 2017.
- Zhaojie Yan. Reinforcement learning transformer for image captioning generation model. In *Fifteenth International Conference on Machine Vision (ICMV 2022)*, volume 12701, pages 166–172. SPIE, 2023.
- Feicheng Huang and Zhixin Li. Improve image captioning via relation modeling. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1945–1949. IEEE, 2022.
- Himaja Cherukuri, Alessio Ferrari, and Paola Spoletini. Towards explainable formal methods: From ltl to natural language with neural machine translation. In *International Working Conference on Requirements Engineering: Foundation for Software Quality*, pages 79–86. Springer, 2022.
- Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575, 2015.
- Yiqi Lin, Hao Wu, Ruichen Wang, Haonan Lu, Xiaodong Lin, Hui Xiong, and Lin Wang. Towards language-guided interactive 3d generation: Llms as layout interpreter with generative feedback. *arXiv preprint arXiv:2305.15808*, 2023.
- Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *arXiv preprint arXiv:2305.11175*, 2023b.
- Maryna Bilan. Statistics of chatgpt & generative ai in business: 2023 report, 2023. URL <https://masterofcode.com/blog/statistics-of-chatgpt-generative-ai-in-business-2023>.
- Fortune Business Insights/Technology. Natural language processing (nlp) market, 2023. URL <https://www.fortunebusinessinsights.com/industry-reports/natural-language-processing-nlp-market-101933>.
- The Market Publicist. Ai large language model market precise outlook 2023 openai, microsoft, google, nvidia, alibaba, baidu, 2023. URL <https://www.benzinga.com/pressreleases/23/09/34645297/ai-large-language-model-market-precise-outlook-2023-openai-microsoft-google-nvidia-alibaba-baidu>.
- Global Info Research. Global large language model (llm) supply, demand, and key producers, 2023–2029, 2023. URL <https://www.orbisresearch.com/reports/index/global-large-language-model-llm-supply-demand-and-key-producers-2023-2029>.
- Digital Journal / Newsmantraa. Large language model(llm) market 2023 growth drivers and future outlook 2030, 2023. URL <https://www.digitaljournal.com/pr/news/newsmantraa/large-language-model-llm-market-2023-growth-drivers-and-future-outlook-2030-meta-ai21-labs-tencent>.
- Infinity Business Insights. Global large language model (llm) market witnesses rapid growth as ai language processing takes center stage 2023 to 2030, 2023. URL <https://www.openpr.com/news/3121457/global-large-language-model-llm-market-witnesses-rapid-growth>.
- Sushant Singh and Ausif Mahmood. The nlp cookbook: modern recipes for transformer based deep learning architectures. *IEEE Access*, 9:68675–68702, 2021.
- Myeonggu Kang, Hyein Shin, and Lee-Sup Kim. A framework for accelerating transformer-based language model on rram-based architecture. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 41(9): 3026–3039, 2021.
- Jia Guo. Deep learning approach to text analysis for human emotion detection from big data. *Journal of Intelligent Systems*, 31(1):113–126, 2022.

- Shunyu Zhang, Yaobo Liang, Ming Gong, Daxin Jiang, and Nan Duan. Multi-view document representation learning for open-domain dense retrieval. *arXiv preprint arXiv:2203.08372*, 2022b.
- Shitao Xiao, Zheng Liu, Weihao Han, Jianjin Zhang, Defu Lian, Yeyun Gong, Qi Chen, Fan Yang, Hao Sun, Yingxia Shao, et al. Distill-vq: Learning retrieval oriented vector quantization by distilling knowledge from dense embeddings. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1513–1523, 2022.
- Sophia Althammer, Sebastian Hofstätter, Mete Sertkan, Suzan Verberne, and Allan Hanbury. Parm: A paragraph aggregation retrieval model for dense document-to-document retrieval. In *European Conference on Information Retrieval*, pages 19–34. Springer, 2022.
- Susan Hazan. The dance of the doppelgängers: Ai and the cultural heritage community. In *Proceedings of EVA London 2023*, pages 77–84. BCS Learning & Development, 2023.
- Zhixin Li and Ying Xu. Designing a realistic peer-like embodied conversational agent for supporting children’s storytelling. *arXiv preprint arXiv:2304.09399*, 2023.
- Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. *arXiv preprint arXiv:2301.00704*, 2023.
- Scott Rettberg, Talan Memmott, Jill Walker Rettberg, Jason Nelson, and Patrick Lichty. Aiwriting: Relations between image generation and digital writing. *arXiv preprint arXiv:2305.10834*, 2023.
- Lidia Ogiela and Marek R Ogiela. *Cognitive techniques in visual data interpretation*, volume 228. Springer, 2009.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022.
- Charlotte Y Alverson and Scott H Yamamoto. Educational decision making with visual data and graphical interpretation: assessing the effects of user preference and accuracy. *Sage Open*, 6(4):2158244016678290, 2016.
- Kanhai Amin, Pavan Khosla, Rushabh Doshi, Sophie Chheang, and Howard P Forman. Focus: Big data: Artificial intelligence to improve patient understanding of radiology reports. *The Yale Journal of Biology and Medicine*, 96(3): 407, 2023.
- Jens Kleesiek, Yonghui Wu, Gregor Stiglic, Jan Egger, and Jiang Bian. An opinion on chatgpt in health care—written by humans only, 2023.
- Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11315–11325, 2022.
- Zhaoshan Liu and Lei Shen. Medical image analysis based on transformer: A review. *arXiv preprint arXiv:2208.06643*, 2022.
- Pokhraj P Suthar, Avin Kounsai, Lavanya Chhetri, Divya Saini, and Sumeet G Dua. Artificial intelligence (ai) in radiology: A deep dive into chatgpt 4.0’s accuracy with the american journal of neuroradiology’s (ajnr) "case of the month". *Cureus*, 15(8), 2023.
- Yue Wang, Hung Le, Akhilesh Deepak Gotmare, Nghi DQ Bui, Junnan Li, and Steven CH Hoi. Codet5+: Open code large language models for code understanding and generation. *arXiv preprint arXiv:2305.07922*, 2023c.
- Xin-Ye Li, Jiang-Tian Xue, Zheng Xie, and Ming Li. Think outside the code: Brainstorming boosts large language models in code generation. *arXiv preprint arXiv:2305.10679*, 2023b.
- Harshit Joshi, José Cambronero Sanchez, Sumit Gulwani, Vu Le, Gust Verbruggen, and Ivan Radiček. Repair is nearly generation: Multilingual program repair with llms. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37-4, pages 5131–5140, 2023.
- Xuan Xu, Saarthak Kapse, Rajarsi Gupta, and Prateek Prasanna. Vit-dae: Transformer-driven diffusion autoencoder for histopathology image analysis. *arXiv preprint arXiv:2304.01053*, 2023.
- Yanjun Gao, Ruizhe Li, John Caskey, Dmitriy Dligach, Timothy Miller, Matthew M Churpek, and Majid Afshar. Leveraging a medical knowledge graph into large language models for diagnosis prediction. *arXiv preprint arXiv:2308.14321*, 2023.
- Claudio Aracena, Fabián Villena, Matías Rojas, and Jocelyn Dunstan. A knowledge-graph-based intrinsic test for benchmarking medical concept embeddings and pretrained language models. In *Proceedings of the 13th International Workshop on Health Text Mining and Information Analysis (LOUHI)*, pages 197–206, 2022.

- Hila Gonen, Srinu Iyer, Terra Blevins, Noah A Smith, and Luke Zettlemoyer. Demystifying prompts in language models via perplexity estimation. *arXiv preprint arXiv:2212.04037*, 2022.
- Xiao Wang, Kai Liu, and Chunlei Wang. Knowledge-enhanced pre-training large language model for depression diagnosis and treatment. In *2023 IEEE 9th International Conference on Cloud Computing and Intelligent Systems (CCIS)*, pages 532–536. IEEE, 2023d.
- Haixing Dai, Zhengliang Liu, Wenxiong Liao, Xiaoke Huang, Zihao Wu, Lin Zhao, Wei Liu, Ninghao Liu, Sheng Li, Dajiang Zhu, et al. Chataug: Leveraging chatgpt for text data augmentation. *arXiv preprint arXiv:2302.13007*, 2023.
- OpenAI. Gpt-4v(ision) system card, 2023a. URL <https://openai.com/research/gpt-4v-system-card>.
- OpenAI. Gpt-4v(ision) system card, 2023b. URL https://cdn.openai.com/papers/GPTV_System_Card.pdf.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- Xiangyu Qi, Kaixuan Huang, Ashwinee Panda, Mengdi Wang, and Prateek Mittal. Visual adversarial examples jailbreak aligned large language models. In *The Second Workshop on New Frontiers in Adversarial Machine Learning*, 2023.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019.
- Shamim Nemati, Bradley A Edwards, Scott A Sands, Philip J Berger, Andrew Wellman, George C Verghese, Atul Malhotra, and James P Butler. Model-based characterization of ventilatory stability using spontaneous breathing. *Journal of Applied Physiology*, 111(1):55–67, 2011.
- Xingjiao Wu, Luwei Xiao, Yixuan Sun, Junhang Zhang, Tianlong Ma, and Liang He. A survey of human-in-the-loop for machine learning. *Future Generation Computer Systems*, 135:364–381, 2022.
- Wei Zhou, Dengke Zhao, Qikai Wu, Jiacheng Dan, Xiaojing Zhu, Wen Lei, Li-Jun Ma, and Ligui Li. Rational design of the lotus-like n-co₂vo₄-co heterostructures with well-defined interfaces in suppressing the shuttle effect and dendrite growth in lithium–sulfur batteries. *Small*, 17(50):2104109, 2021.
- Wei Wang, Kai Xi, Bowen Li, Haojie Li, Sheng Liu, Jianan Wang, Hongyang Zhao, Huanglong Li, Amor M Abdelkader, Xueping Gao, et al. A sustainable multipurpose separator directed against the shuttle effect of polysulfides for high-performance lithium–sulfur batteries. *Advanced Energy Materials*, 12(19):2200160, 2022.
- Kai Sun, Chen Wang, Yan Dong, Pengqian Guo, Pu Cheng, Yujun Fu, Dequan Liu, Deyan He, Saikat Das, and Yuichi Negishi. Ion-selective covalent organic framework membranes as a catalytic polysulfide trap to arrest the redox shuttle effect in lithium–sulfur batteries. *ACS Applied Materials & Interfaces*, 14(3):4079–4090, 2022.
- Wenxuan Hu, Yufan Peng, Yimin Wei, and Yong Yang. Application of electrochemical impedance spectroscopy to degradation and aging research of lithium-ion batteries. *The Journal of Physical Chemistry C*, 127(9):4465–4495, 2023.
- Nomnotho Jiyane, Enrique García-Quismondo, Edgar Ventosa, Wolfgang Schuhmann, and Carla Santana Santos. Elucidating degradation mechanisms of silicon-graphite electrodes in lithium-ion batteries by local electrochemistry. *Batteries & Supercaps*, 6(8):e202300126, 2023.
- Eric Jianfeng Cheng, Yosuke Kushida, Takeshi Abe, and Kiyoshi Kanamura. Degradation mechanism of all-solid-state li-metal batteries studied by electrochemical impedance spectroscopy. *ACS Applied Materials & Interfaces*, 14(36):40881–40889, 2022.
- Michael J Klein, Gabriel M Veith, and Arumugam Manthiram. Chemistry of sputter-deposited lithium sulfide films. *Journal of the American Chemical Society*, 139(31):10669–10676, 2017.
- Thomas Helmuth and Peter Kelly. Applying genetic programming to psb2: the next generation program synthesis benchmark suite. *Genetic Programming and Evolvable Machines*, 23(3):375–404, 2022.
- Thomas Helmuth and Amr Abdelhady. Benchmarking parent selection for program synthesis by genetic programming. In *Proceedings of the 2020 Genetic and Evolutionary Computation Conference Companion*, pages 237–238, 2020.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.

- Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, et al. Codebert: A pre-trained model for programming and natural languages. *arXiv preprint arXiv:2002.08155*, 2020.
- Jack Roper. Transformer-based program synthesis for low-data environments. *arXiv preprint arXiv:2205.09246*, 2022.
- Jawid Ahmad Baktash and Mursal Dawodi. Gpt-4: A review on advancements and opportunities in natural language processing. *arXiv preprint arXiv:2305.03195*, 2023.
- Yongyi Yan, Jumei Yue, and Zengqiang Chen. Logical approach to livelock and deadlock of deterministic finite state machines: Modelling and finding. In *2021 40th Chinese Control Conference (CCC)*, pages 1–6. IEEE, 2021.
- Likang Wu, Zhaopeng Qiu, Zhi Zheng, Hengshu Zhu, and Enhong Chen. Exploring large language model for graph data understanding in online job recommendations. *arXiv preprint arXiv:2307.05722*, 2023.
- Marcel Binz and Eric Schulz. Using cognitive psychology to understand gpt-3. *Proceedings of the National Academy of Sciences*, 120(6):e2218523120, 2023.
- Edward Y Chang. Examining gpt-4: Capabilities, implications, and future directions, 2023.
- Ashwyn Sharma and David I Feldman. Team cadence at mediq-sum 2023: Using chatgpt as a data augmentation tool for classifying clinical dialogue, 2023. URL <https://ceur-ws.org/Vol-3497/paper-139.pdf>.
- Konstantinos I Roumeliotis and Nikolaos D Tselikas. Chatgpt and open-ai models: A preliminary review. *Future Internet*, 15(6):192, 2023.
- Catrin Hasselgren and Tudor I Oprea. Artificial intelligence for drug discovery: Are we there yet? *Annual Review of Pharmacology and Toxicology*, 64, 2023.
- Christos A Fidas, Marios Belk, Argyris Constantinides, David Portugal, Pedro Martins, Anna Maria Pietron, Andreas Pitsillides, and Nikolaos Avouris. Ensuring academic integrity and trust in online learning environments: A longitudinal study of an ai-centered proctoring system in tertiary educational institutions. *Education Sciences*, 13(6): 566, 2023.
- Iddrisu Issah, Obed Appiah, Peter Appiahene, and Fuseini Inusah. A systematic review of the literature on machine learning application of determining the attributes influencing academic performance. *Decision Analytics Journal*, page 100204, 2023.