

---

# A Survey of Transformer Models and Their Impact on Natural Language Processing

---

[www.surveyx.cn](http://www.surveyx.cn)

## Abstract

This survey paper explores the transformative impact of transformer models on natural language processing (NLP), highlighting their pivotal role in advancing language understanding and generation capabilities. By leveraging self-attention mechanisms and scalable architectures, transformers have achieved state-of-the-art performance across various NLP tasks, including machine translation, sentiment analysis, and automatic speech recognition. The integration of innovative approaches, such as MaskAugment, has further enhanced model performance by providing richer training data, validating the utility of masked language models for data augmentation. Despite these advancements, challenges persist, particularly regarding the computational demands of transformer models and the need for improved efficiency. Lightweight architectures offer promising solutions, achieving competitive performance with enhanced computational efficiency. The survey underscores the potential for future innovations in generative masked language models, emphasizing the importance of balancing speed and quality to enhance performance. Additionally, the development of models like BERT-GPT-4 suggests promising avenues for generating coherent and contextually accurate language. As transformer models continue to evolve, ongoing research is essential to optimize their robustness and applicability across diverse languages and speech inputs. By addressing these challenges and embracing innovation, the field can ensure that transformer models remain at the forefront of NLP, driving advancements in artificial intelligence and machine learning.

## 1 Introduction

### 1.1 Significance of Transformer Models in NLP

Transformer models have revolutionized natural language processing (NLP) by employing self-attention mechanisms that facilitate the generation of abstract token representations based on their contextual relationships [1]. This architectural advancement has enabled transformers to outperform traditional algorithms across various NLP tasks, notably in machine translation, where achieving high translation quality remains a significant challenge [2]. The capacity of transformers to manage complex tasks through sophisticated attention mechanisms has markedly enhanced language processing capabilities [3].

The limited availability of labeled datasets, particularly for low-resource languages, underscores the necessity for effective cross-lingual transfer strategies, which transformers adeptly facilitate [4]. This has led to the creation of multilingual models capable of zero-shot learning, thereby broadening the applicability of NLP systems. The advent of Generative Pre-trained Transformers (GPTs), such as GPT-3, has further transformed the field, enabling language generation that increasingly mimics human text [5].

Despite their profound impact, transformers encounter challenges in contextual representation due to their dependence on token position indices [6]. Innovations in context-scaling and task-scaling

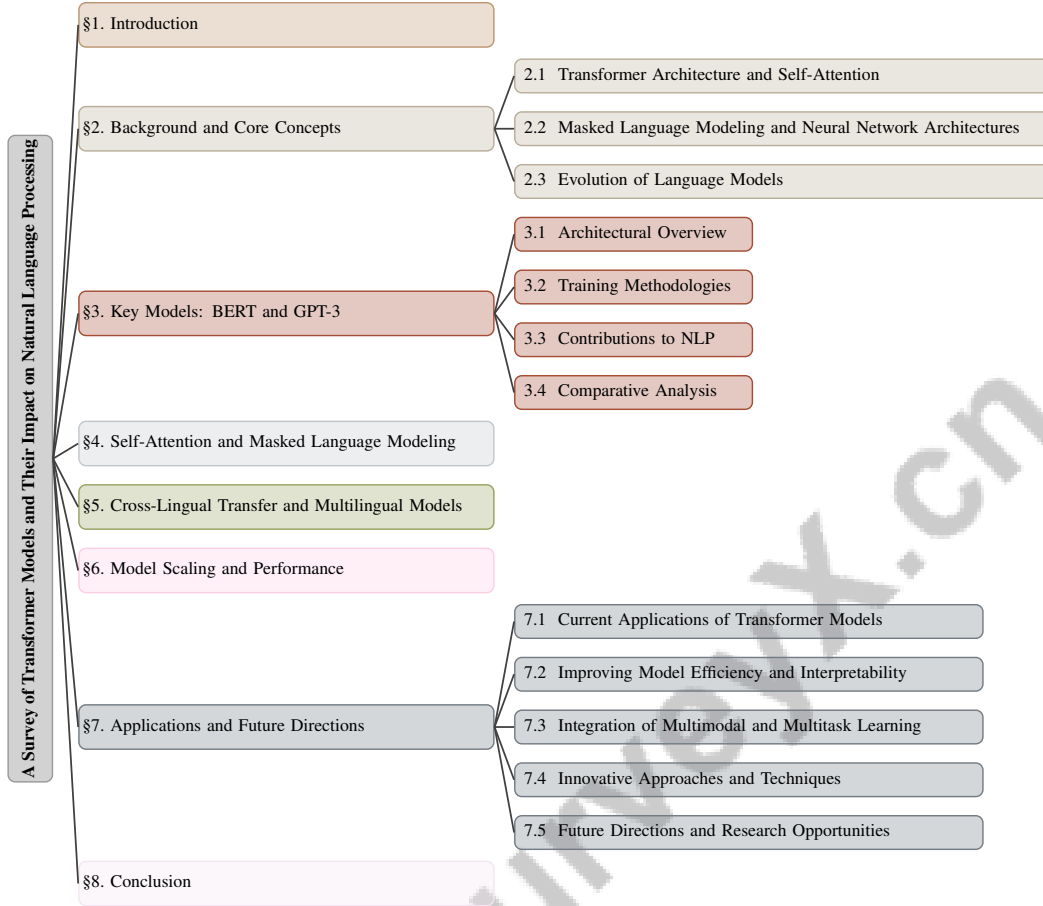


Figure 1: chapter structure

mechanisms are essential for enhancing the in-context learning capabilities of transformers. As NLP evolves, the role of transformers will continue to expand, driven by the growing significance of transfer learning and the development of larger models [7].

In speech recognition, transformers have significantly advanced Automatic Speech Recognition (ASR) systems by integrating linguistic paradigms that improve recognition accuracy across various languages [8]. Ongoing research and innovations will amplify the influence of transformers in NLP, addressing existing limitations and expanding their applicability across diverse domains.

## 1.2 Scope of the Survey

This survey provides a thorough examination of transformer-based models, focusing on large language models (LLMs) such as BERT and GPT, which have profoundly influenced NLP tasks [9]. The survey emphasizes architectural components, operational procedures, training methodologies, and practical applications, intentionally excluding non-transformer models to maintain a focused discourse [5].

Key aspects of LLMs are analyzed in detail, including pre-training, adaptation tuning, utilization, and capacity evaluation, particularly for models with parameter sizes ranging from tens to hundreds of billions [7]. The survey also addresses challenges related to unlearning harmful responses and copyrighted content, which are critical for ethical and legal considerations in their deployment [3].

Moreover, it explores the limitations of fixed attention patterns in pre-trained language models (PLMs) and investigates mechanisms to guide attention, thereby enhancing information processing capabilities [8]. Innovative approaches to decompose the Transformer architecture into smaller, more efficient components are also discussed, essential for optimizing model performance while maintaining quality.

---

The survey includes a comprehensive analysis of GPT-4's features and applications, highlighting advancements over GPT-3, such as enhanced multilingual capabilities, improved contextual understanding, and larger model size, while identifying potential applications like chatbots and language translation. However, it does not delve into technical implementations of these models or broader advancements in artificial intelligence beyond NLP [10, 11, 12, 13]. It emphasizes the need for new approaches to fully utilize the capabilities of deep hidden representations in pre-trained models.

By encompassing a wide range of critical elements, this survey aims to deliver an in-depth exploration of the evolution, implementation, and broader implications of transformer-based models in NLP. It seeks to advance both theoretical research and practical applications by examining foundational principles like the self-attention mechanism, analyzing prominent models such as BERT, GPT, and T5, and discussing their architectural innovations and performance benchmarks across various NLP tasks. Additionally, it addresses efficiency considerations and trade-offs between accuracy and resource requirements, ultimately contributing to sustainable NLP technology development [10, 14, 15]. The survey underscores the importance of evaluating multilingual capabilities and addressing data scarcity challenges, exemplified by benchmarks designed to assess question-answering performance on French datasets and the multilingual skills of models like GPT-3 in languages such as Catalan.

### 1.3 Structure of the Survey

This survey is systematically organized to offer a comprehensive analysis of transformer-based models and their significant impact on NLP. The paper begins with an **Introduction** that establishes the significance of transformer models in NLP, delineates the survey's scope, and underscores the importance of selected keywords. Following the introduction, the **Background and Core Concepts** section delves into the foundational architecture of transformers, the self-attention mechanism, and the evolution of language models, providing critical insights into the core technologies underpinning these models.

The survey then transitions to an examination of **Key Models: BERT and GPT-3**, exploring their architectural designs, training methodologies, and substantial contributions to various NLP tasks. This section includes a comparative analysis of BERT and GPT-3, highlighting their respective impacts on the field.

In the **Self-Attention and Masked Language Modeling** section, the survey investigates the role of self-attention in enhancing contextual understanding, recent innovations in self-attention mechanisms, and the challenges and solutions associated with self-attention and masked language modeling.

The discussion on **Cross-Lingual Transfer and Multilingual Models** follows, focusing on methods and benefits of cross-lingual transfer, the development of multilingual encoders, and advancements in multilingual and cross-lingual models. This section emphasizes the importance of language understanding across different languages.

Next, the survey addresses **Model Scaling and Performance**, analyzing the significance of scaling models, the associated trade-offs, and empirical evidence demonstrating the benefits of scaling transformer models. This section is crucial for understanding the balance between computational resources and performance gains.

The penultimate section, **Applications and Future Directions**, highlights current applications of transformer models, efforts to enhance model efficiency and interpretability, and the integration of multimodal and multitask learning. It also identifies innovative approaches and techniques, suggesting potential research areas and future directions for transformer model development.

The survey concludes with a **Conclusion** that synthesizes key points discussed, reinforcing the transformative impact of transformer models on NLP and reflecting on potential future innovations. This structured approach facilitates a comprehensive examination of transformer-based models in NLP, elucidating their architectural innovations, training methodologies, and practical applications. It also highlights the significant impact of these models on various NLP tasks, such as machine translation, text summarization, and sentiment analysis, while addressing inherent challenges and limitations in their deployment. By analyzing the development and evaluation of models like BERT, GPT, and T5, this exploration provides critical insights into the evolution of transformer technologies and their implications for future advancements in the field [10, 14, 15, 13]. The following sections are organized as shown in Figure 1.

---

## 2 Background and Core Concepts

### 2.1 Transformer Architecture and Self-Attention

The transformer architecture has reshaped natural language processing (NLP), revolutionizing tasks like machine translation, speech recognition, and language understanding through its self-attention mechanism, which captures intricate dependencies within sequences to produce contextually enriched representations [1]. This mechanism excels at handling long sequences, surpassing the limitations of recurrent neural networks (RNNs) and convolutional neural networks (CNNs) that struggle with sequential data processing [2].

Self-attention enhances a model's ability to grasp long-range dependencies by assigning variable importance to tokens within a sequence, albeit with quadratic computational and memory costs, which pose challenges in applications like automatic speech recognition (ASR) [16]. To mitigate these issues, simplified models such as SGPT have been developed, retaining in-context learning (ICL) capabilities without complex trainable attention weights [17].

Recent advancements focus on improving transformer models' efficiency and adaptability. For instance, BERTuit, pre-trained on a large corpus of Spanish tweets, excels in misinformation detection, outperforming existing multilingual models [3]. Moreover, transfer learning techniques significantly boost translation performance by adapting pretrained models from high-resource to low-resource language pairs [18].

The versatility of transformers extends to cross-lingual and multilingual applications, where structured taxonomies classify systems based on linguistic capabilities, ranging from monolingual to cross-lingual frameworks [8]. As transformer architectures evolve, their ability to integrate novel elements and attention mechanisms will expand their applicability across diverse domains, cementing their role in advancing NLP technologies.

### 2.2 Masked Language Modeling and Neural Network Architectures

Masked language modeling (MLM) is pivotal in pre-training transformer models, enhancing their capacity to capture nuanced linguistic features and contextual dependencies [19]. By masking specific tokens and predicting them based on surrounding context, MLM facilitates the learning of rich semantic representations without extensive manual feature engineering [20]. This approach has driven advancements in language models, as seen in frameworks like AntLM, which integrate MLM with Causal Language Modeling to enhance performance [21].

The combination of MLM with various neural network architectures has led to significant efficiency and adaptability improvements. For example, the Segatron model replaces traditional position encodings with indices for paragraphs, sentences, and tokens, improving the model's ability to capture hierarchical text structures [6]. Self-augmentation strategies, such as the SAS method, enable contextualized data augmentation during pre-training, optimizing the learning process [22].

Despite its advantages, MLM is computationally intensive, prompting the exploration of efficient training methodologies. Techniques like iterative mask filling, which involves progressive masking and prediction of words, aim to streamline training [23]. Addressing the quadratic time complexity of the self-attention mechanism remains a critical research area [20].

MLM's adaptability is further demonstrated in cross-lingual transfer scenarios, where language adapters enable zero-shot transfer across various languages, enhancing natural language understanding [24]. Moreover, MLM's potential extends beyond text-based tasks, contributing to synthetic data generation and improving language generation capabilities through task-specific fine-tuning [25].

Challenges in optimizing token masking strategies persist, as random masking may not adequately reflect the varying importance of words within a sentence. Addressing these challenges requires refined data selection and preprocessing methodologies to enhance performance across diverse NLP tasks and domains [23]. Continued innovation in MLM techniques and their integration with neural network architectures is crucial for realizing the full potential of transformer models, driving advancements in NLP technologies.

---

## 2.3 Evolution of Language Models

The evolution of language models has seen significant advancements, culminating in transformer models that enhance NLP capabilities. Initially, recurrent neural networks (RNNs) were favored for sequential data processing but faced limitations such as vanishing gradients and challenges in capturing long-term dependencies, hindering their effectiveness in complex language tasks [26]. These limitations led to the exploration of sequence-to-sequence models and the introduction of attention mechanisms, paving the way for transformers [26].

Transformers, characterized by self-attention mechanisms, overcome RNN limitations by enabling parallel processing and improved management of long-range dependencies. This innovation has significantly enhanced model efficacy and scalability, making transformers foundational in modern NLP [27]. Their ability to model dependencies across sequences results in lower perplexity and improved language generation, benefiting various downstream tasks [28].

Hybrid architectures, such as the integration of BERT and GPT-4, leverage the strengths of different models to enhance text generation [29]. This adaptability addresses diverse NLP challenges. However, the high computational costs associated with training attention-based models remain a critical challenge, necessitating ongoing research into more efficient methodologies [30]. Benchmark studies highlight the benefits of combining multiple architectures in machine translation, achieving significant improvements in translation quality across various linguistic phenomena [2].

The rise of large-scale models like ERNIE 3.0 Titan, with 260 billion parameters, reflects the trend toward scaling models to incorporate extensive knowledge bases, enhancing capabilities in complex language tasks [31]. This trend is complemented by efforts to tackle zero-shot cross-lingual transfer challenges, particularly in scenarios lacking annotated data for target languages [4]. Existing multilingual models, while trained on diverse languages, may underperform on less-represented languages due to biased training sets [32]. The primary challenge remains the unclear relationship between model capacity, data size, and effective learning of cross-lingual representations [33].

Advancements in model architecture, including efficient sentence-level representations in transformers, are crucial for tasks requiring semantic understanding and reconstruction of input sentences [34]. Furthermore, exploring scaling laws for linear complexity models provides insights into their scalability compared to traditional transformers, particularly regarding performance metrics like training loss and validation perplexity [35].

The evolution of language models is driven by the need to balance complexity with performance gains and the pursuit of efficient training methodologies that leverage large datasets effectively. As research progresses, ongoing refinement of transformer architectures and innovative training techniques are set to enhance the performance and versatility of these models, establishing them as essential components in NLP and enabling excellence across applications such as machine translation, text summarization, sentiment analysis, and question answering [14, 15].

## 3 Key Models: BERT and GPT-3

The emergence of transformer-based models has significantly advanced natural language processing (NLP), with BERT and GPT-3 being particularly notable for their architectural innovations and operational principles. Figure 2 illustrates the hierarchical structure of these models, emphasizing their architectural features, training methodologies, and contributions to NLP. This figure provides a comparative analysis that highlights the distinct capabilities and innovations of BERT and GPT-3, showcasing their impact on tasks such as language understanding and generation. Models like T5 and GPT-3 are effective in tasks such as machine-paraphrased plagiarism detection, posing challenges for academic integrity [36, 37, 38, 39, 15].

### 3.1 Architectural Overview

BERT (Bidirectional Encoder Representations from Transformers) and GPT-3 (Generative Pre-trained Transformer 3) represent key advancements, each offering distinct architectural features that enhance NLP capabilities. BERT's bidirectional training mechanism captures context from both directions within a sentence, excelling in question answering and sentiment analysis [3]. It uses masked language modeling (MLM) to predict masked tokens based on context, learning deep representations.

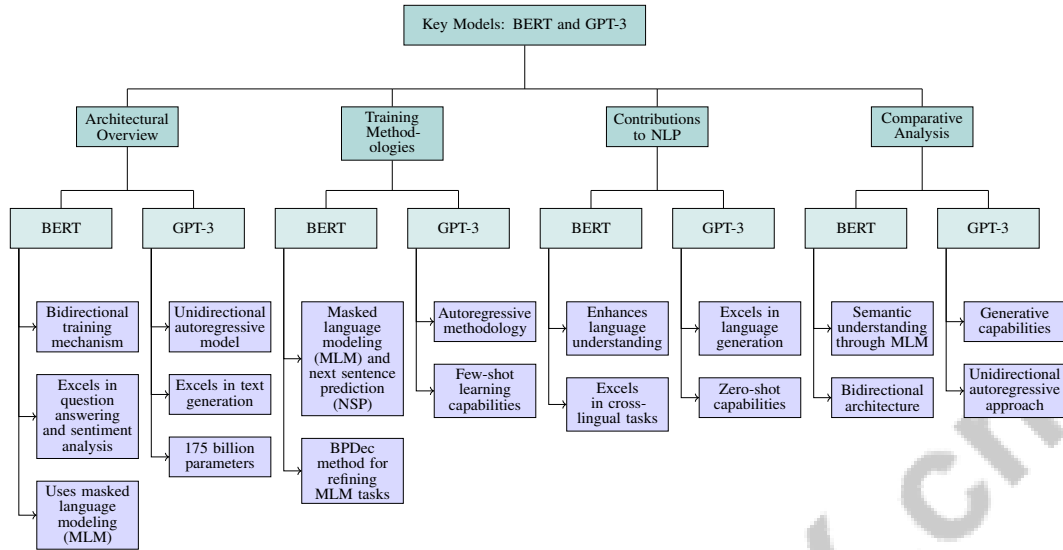


Figure 2: This figure illustrates the hierarchical structure of BERT and GPT-3, highlighting their architectural features, training methodologies, contributions to NLP, and comparative analysis. Each section emphasizes the distinct capabilities and innovations of these models, showcasing their impact on tasks like language understanding and generation.

GPT-3, using a unidirectional autoregressive model, predicts the next word based on preceding words, excelling in text generation despite potential limitations in coherence. Its 175 billion parameters enhance learning capacity from vast datasets, underscoring the transformative potential of transformers in processing large data volumes efficiently [5, 39, 21].

Both models employ self-attention mechanisms, capturing long-range dependencies crucial for performance across NLP tasks. Innovations like hierarchical attention mechanisms and integrating self-attention with convolutional layers enhance dependency management. Strategies like attention guiding and attention map reuse reduce computational complexity, optimizing models for real-time applications [16, 40, 41, 14, 42].

Alternative attention mechanisms, such as the Transformer with a Mixture of Gaussian Keys (Transformer-MGK), improve adaptability and performance by replacing redundant attention heads, reducing computational redundancy, and accelerating training [43, 44, 13, 45]. Hierarchical models like the Hourglass model exemplify efforts to optimize efficiency without compromising performance.

Recent research integrates BERT’s semantic capabilities with GPT-4’s generative strengths, creating a hybrid model that surpasses traditional architectures in metrics like Perplexity and BLEU, emphasizing the balance between accuracy and efficiency in future NLP technologies [46, 10, 29].

### 3.2 Training Methodologies

BERT and GPT-3’s training methodologies reflect diverse strategies enhancing model adaptability. BERT uses masked language modeling (MLM) and next sentence prediction (NSP) to develop deep bidirectional representations, crucial for capturing complex linguistic patterns [47]. Enhancements like the BPDdec method refine MLM tasks, improving contextual dependency learning [19].

GPT-3’s autoregressive methodology focuses on predicting the next token, facilitating coherent text generation. Its few-shot learning capabilities demonstrate effectiveness across numerous NLP tasks without extensive fine-tuning [48]. The model’s scalability, with its extensive parameter count, underscores the significance of its training strategy.

Innovative techniques, such as Self-Evolution learning (SE), enhance efficiency by focusing on informative tokens [49]. Methods like knowledge distillation and layer truncation, as seen in Compressed GPT-2 (C-GPT-2), optimize performance for downstream tasks [50].

Figure 3 illustrates the hierarchical structure of training methodologies in NLP, focusing on BERT and GPT-3, along with innovative techniques like self-evolution learning and the Fed-Grow framework. This visual representation complements the discussion by providing a clear overview of the relationships and advancements within these methodologies.

Integration of self-supervised adversarial loss and controllable language modeling loss in models like ERNIE 3.0 Titan illustrates advanced methodologies aimed at enhancing text generation [31]. These innovations, along with specialized encodings in models like Segatron, highlight ongoing advancements in training practices driving transformer evolution [28].

BERT and GPT-3’s adaptability and innovation in training methodologies optimize performance and efficiency, ensuring their significant influence in NLP. Their self-attention architecture captures long-range dependencies, leading to breakthroughs in tasks like machine translation and text summarization. Notable models like ChatGPT and BLOOM underscore the importance of continued research in this rapidly evolving landscape [14, 15].

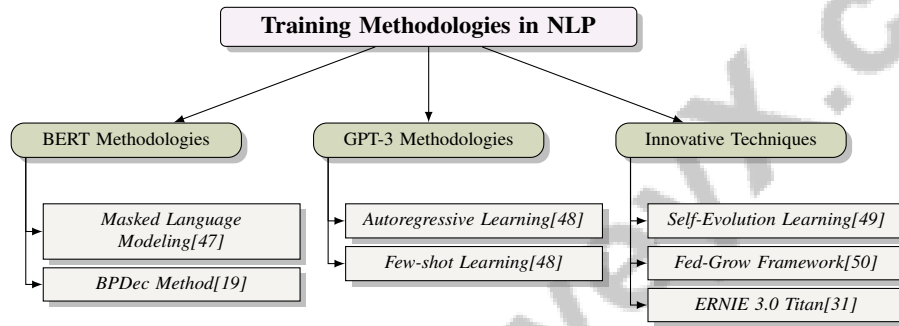


Figure 3: This figure illustrates the hierarchical structure of training methodologies in NLP, focusing on BERT and GPT-3, along with innovative techniques like self-evolution learning and the Fed-Grow framework.

### 3.3 Contributions to NLP

BERT and GPT-3 have substantially contributed to NLP, excelling in language understanding and generation. BERT’s bidirectional architecture enhances performance in question answering and sentiment analysis, outperforming traditional models like GloVe and Word2Vec, especially in cross-lingual tasks [51, 52, 53]. Its adaptability in low-resource settings is evident through language adapters enhancing cross-lingual transfer with minimal fine-tuning.

GPT-3 excels in language generation, leveraging its vast scale to produce coherent text across applications. Its zero-shot capabilities are notable in languages with limited representation in training data, like Catalan [54]. However, its few-shot performance may lag behind fine-tuned models [55].

Both models advance multimodal applications, adapting architectures for various data types. BERT’s versatility is evident in models like NoRBERT, which improve accuracy in missing data imputation tasks [30]. Their contributions extend to specialized domains, often surpassing state-of-the-art models in tasks and languages, as shown by fine-tuned XLM-R models improving NER performance in low-resource languages [27].

Techniques like self-conditioning and iterative mask filling improve performance in generating text with desired properties, including gender parity [23]. BERT’s high detection accuracy in identifying machine-generated texts underscores its robustness in content generation [56].

BERT’s distinct processing phases enhance its applicability in complex NLP tasks [46]. GPT-3’s architecture and few-shot learning capabilities set new benchmarks in text generation, solidifying its role in advancing NLP technologies [57]. These models’ integration into multi-task learning frameworks enhances text classification and summary generation performance, demonstrating their transformative impact on NLP [58].

---

### 3.4 Comparative Analysis

BERT and GPT-3 are influential transformer-based models, each showcasing distinct strengths. BERT excels in semantic understanding through masked language modeling, adept at tasks requiring deep comprehension, while GPT-3 demonstrates exceptional proficiency in generating coherent text due to its autoregressive capabilities and vast parameter size. Hybrid approaches integrating both models highlight enhanced performance potential by combining BERT’s semantic strengths with GPT-3’s generative capabilities [9, 13, 29, 20]. Both models use self-attention mechanisms, crucial for capturing contextual dependencies, yet their architectural differences yield distinct strengths.

BERT’s bidirectional architecture allows it to consider context from both directions, excelling in tasks requiring nuanced comprehension, particularly in cross-lingual scenarios where it outperforms models like mBERT [59].

GPT-3’s unidirectional autoregressive approach generates text based on preceding tokens, limiting future context incorporation. However, it excels in generating fluent text, with human evaluations rating its paraphrases comparably to originals [13, 39, 12, 21, 60]. This methodology excels in tasks requiring coherent text generation over extended passages, supported by its massive scale and zero-shot learning capabilities, though its few-shot performance may not always surpass fine-tuned models.

A notable distinction is their alignment with human moral judgments, revealing varying degrees of ethical considerations in NLP applications [61]. This highlights the importance of ethical dimensions in deploying language models, particularly in sensitive applications.

The choice between BERT and GPT-3 often depends on task requirements and resource availability. BERT’s efficiency suits tasks with limited data, while GPT-3’s generative capabilities excel in extensive language generation contexts. The decision to utilize mBERT or a language-specific model is similarly influenced by these factors [59].

## 4 Self-Attention and Masked Language Modeling

### 4.1 The Role of Self-Attention in Contextual Understanding

Self-attention mechanisms are crucial in transformer models, enhancing contextual understanding by dynamically assessing token relevance within sequences. This capability captures complex dependencies, crucial for various NLP tasks, and allows parallel processing, a significant improvement over traditional RNNs [8]. In ASR, self-attention improves accuracy by focusing on relevant input segments, though it may overemphasize current tokens, potentially neglecting historical context. This limitation necessitates architectures that balance attention across sequences for optimal predictions.

Innovations like the Segatron model address these issues by incorporating segment information, boosting performance across tasks. Self-attention also effectively decouples word positions and embeddings, capturing diverse linguistic information beneficial for cross-lingual tasks. Aligning representations of translated words significantly enhances cross-lingual transfer, with techniques like realignment improving performance, as seen in a 15.8-point accuracy increase in POS-tagging between English and Arabic [62, 14, 63].

Analyzing hidden states in transformer models, especially in QA tasks, reveals insights into reasoning processes often overlooked when focusing solely on attention weights. Our findings indicate that hidden states provide significant information on token representation evolution through processing layers, aligning with traditional pipeline tasks and integrating task-specific knowledge. Visualizations of hidden states enhance understanding of semantic capabilities, showing early layers can reflect prediction errors. This deeper understanding boosts interpretability and performance, highlighting self-attention’s importance in contextual understanding. Its refinement and integration with innovative designs promise further advancements, enabling models to emulate human reasoning and intelligence more effectively.

### 4.2 Innovations in Self-Attention Mechanisms

Recent advancements in self-attention mechanisms significantly enhance transformer models, improving performance across NLP tasks. Reweighting self-attention to amplify historical tokens’ influence



provides comprehensive context for predictions [64], mitigating traditional self-attention’s limitation of focusing excessively on immediate context. Adaptive sparse monotonic attention mechanisms, using adaptive  $\text{-entmax}$ , offer flexible attention distributions that capture salient features, while regularized monotonic multi-head attention structures improve sequence-to-sequence task performance [65].

The Progressive Transformer (PRO TST) exemplifies the progressive teacher-student paradigm, guiding models from foundational to complex tasks, enhancing learning, and reducing reliance on external tools [66]. This structured approach facilitates efficient training and generalization across diverse NLP applications. Additionally, integrating knowledge graph-based learning objectives within the self-attention framework, as seen in BERT models trained on free-text and structured data, refines performance and prediction accuracy in knowledge-intensive tasks [67].

These advancements illustrate the dynamic evolution of transformer models in NLP, enhancing adaptability to complex tasks like machine translation, text summarization, and sentiment analysis, while laying the groundwork for future breakthroughs. As researchers explore diverse transformer architectures and efficiencies, the potential for improved performance and broader applications becomes evident, indicating a promising trajectory for NLP technologies [13, 68, 10, 14, 15].

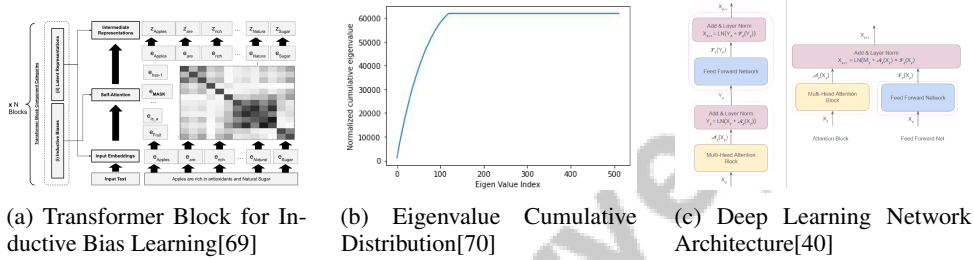


Figure 4: Examples of Innovations in Self-Attention Mechanisms

As shown in Figure 4, innovations in self-attention mechanisms and masked language modeling have significantly advanced deep learning. These examples highlight key advancements in understanding and implementing self-attention mechanisms, crucial for enhancing model performance and efficiency. The "Transformer Block for Inductive Bias Learning" demonstrates structuring transformer blocks to incorporate inductive biases, guiding models toward more generalizable patterns. The "Eigenvalue Cumulative Distribution" plot provides insights into dataset eigenvalues, influencing effective self-attention layer design. The "Deep Learning Network Architecture" image details components and layers in deep learning architectures, emphasizing multi-head attention and feed-forward networks' roles in data processing and transformation. Collectively, these examples underscore ongoing innovations in self-attention mechanisms driving progress in NLP and beyond [69, 70, 40].

### 4.3 Challenges and Solutions in Self-Attention and MLM

Implementing self-attention and masked language modeling (MLM) in transformers presents challenges related to computational complexity, scalability, and interpretability. The quadratic scaling of self-attention with input sequence length can be computationally prohibitive, limiting transformers' applicability in long sequences or resource-constrained environments. Reusing attention maps across layers offers a promising solution by reducing attention computations, enhancing efficiency without sacrificing performance [16].

Integrating causal and masked language modeling techniques presents another challenge. Existing models often struggle to balance both effectively without substantial architectural modifications, hindering coherent and contextually accurate language output generation [19]. Innovations like Parallel Decoding by Iterative Refinement (PaDIR) aim to refine target sequences in parallel, enhancing autoregressive model performance [21].

Interpretability is a critical concern, particularly in understanding feedforward network operations within transformers. Fine-grained interpretability can improve model prediction control, reducing harmful outputs and enhancing reliability [71]. Additionally, reliance on specific detection models

---

can limit generalization across different text types or domains, highlighting the need for robust interpretability frameworks [56].

In cross-lingual transfer contexts, challenges like catastrophic forgetting in multilingual models necessitate continual learning scenarios to maintain performance across languages [32]. The quality of pretrained transformers significantly impacts sentence representation tasks, requiring methods ensuring robust generalization across diverse applications [34].

Developing efficient augmentation strategies, like the Iterative Mask Filling (IMF) method, can enhance classification task performance by providing more effective augmentation than simpler methods [23]. Despite advancements, diffusion language models face challenges in complex reasoning tasks and generation order limitations, affecting output quality [25]. Addressing these limitations requires ongoing research into refining generation methods and improving task-specific adaptations.

Addressing challenges in self-attention and MLM necessitates a multifaceted approach, encompassing advancements in model architecture, training methodologies, and interpretability. These efforts are essential for enhancing transformer models' performance and versatility, significantly advancing NLP across tasks like machine translation, text summarization, sentiment analysis, and question answering. By leveraging architectural innovations and techniques like self-attention mechanisms, these initiatives aim to optimize model efficiency, making them more applicable in resource-constrained environments while maintaining high accuracy and effectiveness in diverse NLP applications [10, 72, 14, 15].

## **5 Cross-Lingual Transfer and Multilingual Models**

### **5.1 Cross-Lingual Transfer and Multilingual Representation Learning**

Cross-lingual transfer and multilingual representation learning are pivotal in advancing NLP by enabling linguistic knowledge transfer across languages, addressing linguistic diversity and data scarcity challenges. Models like multilingual BERT (mBERT) and XLM-Roberta excel in generating contextualized representations that support zero-shot cross-lingual transfer, leveraging high-resource languages to enhance low-resource contexts [33, 32]. This capability is crucial for inclusive language technologies across diverse linguistic landscapes.

Despite these advancements, multilingual models often face inefficiencies in specific languages, leading to resource wastage and deployment complexities [7]. Direct Transfer Learning (DTL) addresses these issues by allowing pre-trained models to adapt to new language pairs without altering architecture or hyperparameters, enhancing adaptability and efficiency [18]. Language adapters further facilitate parameter-efficient cross-lingual transfer, enabling multiple tasks with minimal additional parameters compared to full fine-tuning [8]. This broadens multilingual models' applicability across diverse linguistic contexts.

Examining cross-linguistic syntactic differences in models like mBERT highlights the significance of syntactic knowledge transfer among typologically diverse languages, improving translation accuracy and syntactic relationship handling in low-resource scenarios [3]. Benchmarks evaluating zero-shot and few-shot cross-lingual transfer provide insights into the performance of massively multilingual transformers across various languages and tasks [4, 54]. These evaluations emphasize the need for models capable of transferring linguistic features without explicit training data, facilitating more robust and accessible NLP technologies.

Future research should optimize fine-tuning methods for large language models (LLMs), particularly through intermediate task training to enhance performance across diverse NLP tasks. Improving model interpretability is crucial for understanding successful cross-lingual transfer mechanisms. Integrating external memory networks could advance handling low-resource language pairs and address challenges in neural machine translation, enhancing cross-lingual transfer capabilities [18, 73]. Leveraging cross-lingual transfer and multilingual representation learning can make NLP technologies more robust and accessible across different languages and cultural contexts.

### **5.2 Multilingual Encoders and Cross-Lingual Representations**

Multilingual encoders have been instrumental in advancing cross-lingual representation learning within transformer models, enabling effective processing and understanding of multiple languages.

---

Multilingual BERT (mBERT) exemplifies this capability by using a shared encoder to generate contextualized embeddings that facilitate cross-lingual transfer across languages and tasks [33, 74]. This shared architecture aligns multilingual representations effectively, particularly in lower layers, while using language-agnostic predictors in upper layers to enhance cross-lingual capabilities.

Benchmarks like XQuAD provide resources for evaluating cross-lingual question answering through translated datasets [75]. These benchmarks are essential for assessing models like mBERT and XLM-R in cross-lingual tasks, as demonstrated by alignment strategy tests [63]. The integration of morphosyntactic properties in selecting source languages for cross-lingual transfer emphasizes the importance of linguistic features in enhancing model performance [53].

Innovative approaches like Neuron Specialization modularize feed-forward layers in multilingual models, enhancing task specificity and reducing interference [76]. This method identifies specialized neurons within these layers, facilitating more efficient and targeted processing of linguistic information. The METAXLM method, combining meta-pretraining with cross-lingual pretraining, illustrates advanced training techniques' potential to improve cross-lingual representation learning [74].

Exploring transfer learning benefits, particularly adapting pretrained models from high-resource to low-resource language pairs, further highlights multilingual encoders' adaptability and efficiency [18]. These advancements are crucial for developing robust NLP systems capable of operating effectively across languages and applications, driving the continued evolution of multilingual encoders and cross-lingual representations in transformer models.

### 5.3 Advancements in Multilingual and Cross-Lingual Models

Recent advancements in multilingual and cross-lingual models have significantly enhanced NLP systems, improving language understanding and translation across diverse linguistic contexts. Cross-lingual language model meta-pretraining (METAXLM) introduces a meta-pretraining phase focused on generalization prior to cross-lingual pretraining, enhancing the model's ability to generalize across languages and improving cross-lingual transfer performance [74].

Studies consistently show that multilingual models outperform monolingual counterparts in cross-lingual transfer tasks, particularly in low-resource languages like Bantu languages, where models like AfriBERT excel [32]. This underscores multilingual models' potential to bridge the technological gap for underrepresented languages.

Despite these advancements, challenges remain in optimizing multilingual models for specific languages and tasks. Investigating factors influencing cross-lingual transfer, such as linguistic similarity and model architecture, reveals positive impacts on performance, while lexical overlap presents mixed results [53]. This emphasizes the importance of realignment methods for smaller multilingual models, evidenced by the strong correlation between multilingual alignment and cross-lingual transfer capabilities [63].

Transfer learning strategies enhance translation quality in low-resource scenarios through effective knowledge transfer from high-resource models [18]. This approach improves translation accuracy and efficiency without relying on extensive language-specific data.

Critical questions remain regarding optimal conditions for training multilingual models and various model architectures' specific contributions to cross-lingual performance [33]. Addressing these questions is essential for refining multilingual models and maximizing their utility across diverse linguistic landscapes.

## 6 Model Scaling and Performance

### 6.1 Importance of Model Scaling

Model scaling is pivotal for enhancing NLP systems by capturing complex linguistic patterns and long-distance dependencies essential for advanced language tasks. Larger models, with increased parameter counts, excel in understanding nuanced linguistic features and advancing multilingual comprehension [6]. They provide richer positional information, enhancing context comprehension and language modeling [6]. However, training these models requires substantial computational resources, posing challenges due to high costs and dependency on extensive labeled datasets [77].

---

Innovations like lightweight transformers, such as Segatron, demonstrate competitive performance while optimizing computational efficiency, underscoring the need for a balance between model size and resource utilization for diverse tasks [6].

Efficient inference advancements, such as attention map reuse, are significant for resource-limited applications [16]. Self-augmentation strategies like SAS integrate generation and discrimination tasks, enhancing computational efficiency and overall performance [22]. Optimizing transformer model depth further improves performance, highlighting architecture’s role in achieving state-of-the-art results [17]. Scaling models with both context and tasks enhances learning capabilities beyond traditional methods [17]. Continued research is essential for overcoming scaling challenges, ensuring large-scale models’ effective deployment across diverse tasks and domains. Balancing model complexity and computational efficiency is crucial for maintaining transformer models’ relevance in NLP, given their resource-intensive architectures and the demand for efficient solutions [10, 14, 36].

## 6.2 Trade-offs in Model Scaling

Scaling models in NLP involves navigating trade-offs between computational resources and performance enhancements. Larger models capture intricate linguistic patterns but incur significant computational demands, challenging resource-constrained environments [31]. The resources required for models like ERNIE 3.0 Titan exemplify these trade-offs, necessitating a careful evaluation of benefits versus costs [31]. Pruning techniques, such as A\* Pruning, manage these trade-offs by optimizing the removal of less significant attention heads, maintaining performance while reducing resource consumption [72]. However, multitask learning frameworks can introduce task conflicts and complexities in dynamically adjusting task weights, affecting scalability and adaptability [78].

The ‘curse of multilinguality’ complicates scaling, as multilingual models’ capacity can be stretched across languages, leading to negative interference and degraded performance for specific languages [79]. Specialized approaches like Neuron Specialization improve performance but still lag behind bilingual models for high-resource languages [76]. Fed-Grow enhances model accuracy and stability across clients while minimizing resource consumption during training [50]. Innovations like automated elastic pipelining optimize resource allocation and enhance training efficiency, addressing static training methods’ inefficiencies [80]. Challenges remain in managing out-of-vocabulary words and the computational complexity of training large models [26]. Fixed datasets in current benchmarks further complicate optimizing scaling by not fully capturing data distribution’s influence on scaling laws [35].

## 6.3 Empirical Evidence of Scaling Benefits

Empirical studies reveal substantial benefits from scaling transformer models, showing significant performance enhancements across NLP tasks. Scaling model sizes and data correlates with improved performance, though understanding emergent abilities remains challenging [81]. Larger models like GPT-3 exhibit advancements in few-shot learning, showcasing generalization across tasks with minimal fine-tuning [48]. The Hourglass model illustrates efficiency and performance improvements through optimized architectures, achieving state-of-the-art results in autoregressive image generation and competitive language modeling [82]. Strategic architectural scaling can yield high performance without extensive resource demands.

Research emphasizes optimizing communication performance across model sizes and parallelism strategies, highlighting efficient resource allocation in distributed training [83]. PipeTransformer exemplifies this by achieving significant training speed improvements, with up to 2.83-fold speedup without sacrificing accuracy, showcasing automated elastic pipelining benefits in scaling efforts [80]. Integrating symbolic knowledge distillation techniques enhances scaled models’ transparency and efficiency, addressing research gaps and opening exploration avenues [84]. GroupBERT offers a more efficient alternative to traditional transformer architectures, achieving up to 2.1× efficiency gains in FLOPs and training time while improving task performance across scales [30].

In domain-specific tasks, UMLS-KGI-BERT demonstrates structured knowledge’s significant enhancement of language models’ effectiveness in the biomedical domain, particularly improving named entity recognition (NER) tasks [67]. This underscores structured knowledge integration’s potential to optimize model performance in specialized domains. Empirical evidence supports scaling transformer models’ substantial performance improvements across NLP tasks, highlighting the need

---

for strategic scaling, efficient resource management, and continual architecture refinement to fully leverage advanced language models' benefits. Future research should explore alternative implementations and optimize model capacity allocation to enhance transformer models' scalability and applicability [79].

## 7 Applications and Future Directions

Recent advancements in natural language processing (NLP) have been significantly driven by transformer models, which have redefined performance standards across various tasks and introduced innovative methodologies that enhance their capabilities. This section delves into the current applications of transformer models, emphasizing their pivotal role in advancing NLP technologies and laying the groundwork for future research.

### 7.1 Current Applications of Transformer Models

Transformer models are central to diverse NLP applications, significantly enhancing tasks such as neural machine translation (NMT), where they outperform traditional word-based systems by effectively leveraging monolingual data and sentence representations [21]. Generative Masked Language Models (GMLMs) further improve translation tasks compared to autoregressive models [21].

In automatic speech recognition (ASR) and spoken language understanding (SLU), transformers have achieved notable improvements, particularly through speaker embeddings that enhance ASR accuracy [85]. Blockwise streaming Transformers show competitive performance in SLU and speech translation by employing ASR-based intermediate loss and joint CTC/attention translation methods [86, 16].

For sentiment analysis and named entity recognition (NER), transformer models like MaskAugment demonstrate versatility [87]. Datasets from the National Centre for Human Language Technology and the South African Centre for Digital Language Resources highlight transformers' role in advancing language technology [27].

Transformers also excel in automated writing and adaptive conversational agents, with large-scale generative architectures enhancing content generation and dialogue systems [29]. Experiments with GPT-3 emphasize its potential in creative writing, though its reliability varies across tasks [88, 5]. Innovations like Delta-KWT reduce inference time, making transformers more accessible for low-power devices [89].

Techniques such as HIRE enhance hidden representations, improving performance in natural language understanding tasks [90]. Larger models show significant improvements in few-shot learning, underscoring the transformative potential of scaled transformers [48]. Future research will focus on extending models like BERTuit to multilingual contexts and exploring novel applications in misinformation analysis [3, 31].

### 7.2 Improving Model Efficiency and Interpretability

Enhancing transformer models' efficiency and interpretability is crucial for diverse NLP tasks. Research focuses on optimizing efficiency to reduce computational demands, exploring training paradigms that maintain performance with limited data [37, 91]. Self-conditioning addresses model biases, aiding understanding of generative mechanisms in transformer language models (TLMs) [92].

Improving interpretability involves understanding decision-making processes within models. Tools visualize attention patterns and token representations, providing insights into models like BERT. Layer-wise analyses reveal how models transform token vectors, enhancing understanding of tasks like Question Answering [28, 37, 93]. Segment-Aware and Hypernym-Instructed Language Modeling elucidate language features, improving performance [28].

Feature engineering and learning approaches enhance interpretability by clarifying feature interactions with model decisions [94]. Addressing ethical implications and ensuring model safety in real-world applications are paramount. Benchmarks evaluating pre-trained language models' susceptibility to

---

harms are essential for model safety [95]. Future research should explore multi-modal integration and ethical deployment of large language models, ensuring responsible use and efficiency.

### 7.3 Integration of Multimodal and Multitask Learning

Integrating multimodal and multitask learning within transformers enhances their versatility, enabling simultaneous processing of diverse data types and execution of multiple tasks. Multimodal learning leverages transformers' ability to integrate information from text, image, and audio, improving understanding and generation capabilities [31].

Recent multimodal transformers excel in tasks like visual question answering and image captioning, using cross-attention mechanisms to align features from different modalities, enhancing performance [96]. Hierarchical attention structures further improve processing of complex dependencies.

Multitask learning allows concurrent training on related tasks, promoting knowledge sharing and improving generalization [78]. Dynamic task weight adjustments and task-specific adapters mitigate task interference, optimizing performance.

Combining multimodal and multitask learning enhances transformers' capability for complex applications, addressing data scarcity and task-specific training challenges. This integration paves the way for advancements in autonomous systems, human-computer interaction, and personalized content generation [96]. Further exploration of optimizing these strategies will fully harness transformers' potential in diverse applications, including efficient training methodologies, refined attention mechanisms, and novel architectures [91].

### 7.4 Innovative Approaches and Techniques

Innovative approaches in transformer development continue to expand NLP capabilities. Symbolic knowledge distillation enhances interpretability and efficiency by leveraging structured knowledge during training, addressing limitations in capturing complex semantic relationships [84].

Self-augmentation strategies, like SAS, optimize learning by combining data augmentation with pre-training, enhancing generalization in low-resource settings [22]. The Hierarchical Transformer with Mixture of Gaussian Keys (Transformer-MGK) introduces novel attention mechanisms, reducing redundancy and improving adaptability, facilitating efficient input sequence processing [82].

Automated elastic pipelining, as demonstrated by PipeTransformer, optimizes resource allocation, significantly improving training efficiency [80]. Parallel decoding techniques, such as PaDIR, enhance language generation by refining target sequences in parallel, improving coherence and fluency [21].

Surveys of transformer models highlight their dynamic evolution, showcasing their impact on tasks like machine translation, sentiment analysis, and text summarization. These advancements emphasize architectural innovations in models like BERT, GPT, and T5, paving the way for enhanced human-machine interactions and sophisticated applications across industries [97, 15, 13]. As research explores novel methodologies and architectures, transformers are poised for greater capabilities, driving the next wave of innovation in language understanding and generation.

### 7.5 Future Directions and Research Opportunities

Future development of transformer models in NLP will explore optimizing scalability, efficiency, and applicability across domains. Key areas include optimizing pretraining strategies and model architectures for improved few-shot learning, refining attention mechanisms, and integrating model components into visualization tools to enhance interactivity and interpretability [77].

Exploring hierarchical configurations and optimizing shortening mechanisms can enhance efficiency and performance [22]. Investigating optimal sublayer arrangements for tasks like translation and classification could leverage architecture modifications for performance gains. Expanding benchmarks to include diverse datasets and multimodal summarization techniques is crucial for effective deployment in real-world scenarios. Developing standardized benchmarks for language-specific models and improving data collection for low-resource languages are essential [7].

In token representation, applying insights about token evolution to other deep models and exploring additional learning objectives could streamline training. Further exploration of intricate translation

---

tasks and specialization of concepts to specific languages will enhance cross-lingual capabilities. Evaluating decomposition pipelines with larger digit operations and addressing multiplication task challenges could refine computational efficiency [16]. Expanding benchmarks to include additional languages and refining test suites for comprehensive syntactic coverage are promising research directions [17].

Optimizing learning rate strategies and exploring memory slots in complex tasks with longer input sequences could enhance performance [8]. Future research should investigate additional training strategies and expand benchmarks to encompass diverse domains and tasks, enhancing model versatility. Exploring pre-training techniques or datasets that promote general knowledge representations could lead to universal models excelling across multiple domains. Optimizing MTNs for hierarchical tasks, reducing memory footprints, and integrating them into mobile applications and conversational AI remain promising research avenues. Refining masking strategies and exploring complex visual tasks and datasets could enhance capabilities in visual domains [2].

Applying A\* Pruning to other architectures and tasks, enhancing heuristics for pruning, and investigating synthetic data's long-term effects on model performance, particularly in healthcare, are essential for reliability and accuracy [38, 98, 39]. Exploring alternative measures for latent space geometry and their implications for training and performance optimization could refine transformers. Understanding LLMs' emergent abilities, improving training efficiency, and aligning LLMs with human values and ethical considerations are crucial for advancing transformers.

Focusing on these research areas will drive future developments, ensuring transformers' relevance and impact across artificial intelligence and machine learning. Exploring additional tasks and datasets, enhancing robustness and understanding, and investigating interactions among factors, developing multilingual datasets, and exploring generative models' implications are important areas for further exploration [8]. Future work could aim at developing LLMs with human-like reasoning through enhanced training methodologies. Understanding emergent abilities, improving alignment with human values, and exploring efficient scaling strategies in data-constrained environments will be crucial for transformers' continued evolution. Future research should prioritize optimizing methods for specific tasks and exploring applicability across domains, including image and audio processing. Refining architecture search processes, incorporating sophisticated profiling techniques, and investigating applicability to other neural network architectures beyond transformers are promising avenues. Further investigation into sequence mixing layers in subquadratic models and their performance impact, along with optimizing the distillation process, are ripe for research. Advanced techniques for modifying model weights and diversifying evaluation metrics to include interpretability and fairness should be pursued. Finally, refining interface designs to promote active user engagement and exploring features that enhance narrative coherence are important considerations for future research. Scaling pre-training and enhancing MLP architectures to improve performance on vision-and-language tasks are also promising directions for exploration.

## 8 Conclusion

The survey highlights the transformative role of transformer models in natural language processing (NLP), establishing them as essential architectures that have significantly enhanced language understanding and generation. Utilizing self-attention mechanisms and scalable designs, transformers have achieved state-of-the-art results in various NLP tasks, such as machine translation, sentiment analysis, and automatic speech recognition [9]. Innovative strategies like MaskAugment have further improved model efficacy by enriching training data, thus validating the effectiveness of masked language models for data augmentation [87].

However, the survey also identifies the computational challenges posed by transformer models, necessitating enhanced efficiency. Recent studies suggest that lightweight transformer architectures can achieve competitive performance while improving computational efficiency [99]. Future research should explore the application of innovative architectures, such as GroupBERT, across various languages and domains, focusing on alternative sparsity techniques to boost model adaptability [30].

Additionally, the survey points to the potential advancements in generative masked language models (GMLMs), stressing the need to balance speed and quality for optimal performance [21]. The emergence of models like BERT-GPT-4, which exhibit enhanced fluency and semantic consistency, indicates promising pathways for generating coherent and contextually relevant language [29].

---

As transformer models continue to progress, ongoing research is crucial for optimizing their robustness and applicability across diverse languages and speech inputs [86]. Addressing these challenges while fostering innovation will ensure that transformer models remain at the forefront of NLP, propelling advancements in artificial intelligence and machine learning.

www.SurveyX.cn



---

## References

- [1] Dhanasekar Sundararaman, Vivek Subramanian, Guoyin Wang, Shijing Si, Dinghan Shen, Dong Wang, and Lawrence Carin. Syntax-infused transformer and bert models for machine translation and natural language understanding, 2019.
- [2] Felix Stahlberg, Adria de Gispert, and Bill Byrne. The university of cambridge’s machine translation systems for wmt18, 2018.
- [3] Javier Huertas-Tato, Alejandro Martin, and David Camacho. Bertuit: Understanding spanish language in twitter through a native transformer, 2022.
- [4] Anne Lauscher, Vinit Ravishankar, Ivan Vulić, and Goran Glavaš. From zero to hero: On the limitations of zero-shot cross-lingual transfer with multilingual transformers, 2020.
- [5] Manuel de Buenaga and Francisco Javier Bueno. Application of gpt language models for innovation in activities in university teaching, 2024.
- [6] He Bai, Peng Shi, Jimmy Lin, Yuqing Xie, Luchen Tan, Kun Xiong, Wen Gao, and Ming Li. Segatron: Segment-aware transformer for language modeling and understanding, 2020.
- [7] Zaid Alyafeai, Maged Saeed AlShaibani, and Irfan Ahmad. A survey on transfer learning in natural language processing. *arXiv preprint arXiv:2007.04239*, 2020.
- [8] Shruti Singh, Muskaan Singh, and Virender Kadyan. Speech recognition transformers: Topological-lingualism perspective, 2024.
- [9] Aman Pathak et al. Comparative analysis of transformer based language models. In *CS & IT Conference Proceedings*, volume 11. CS & IT Conference Proceedings, 2021.
- [10] Wazib Ansar, Saptarsi Goswami, and Amlan Chakrabarti. A survey on transformers in nlp with focus on efficiency. *arXiv preprint arXiv:2406.16893*, 2024.
- [11] Walid Hariri. Unlocking the potential of chatgpt: A comprehensive exploration of its applications, advantages, limitations, and future directions in natural language processing, 2025.
- [12] Jawid Ahmad Baktash and Mursal Dawodi. Gpt-4: A review on advancements and opportunities in natural language processing, 2023.
- [13] Chandra Sekhar Kolli, Pavan Kumar Vadrevu, and S Srinivasu. Comprehensive exploration of generative pre-trained transformer.
- [14] Pawan Sasanka Ammanamanchi. *Evaluation of Transformer Models on Summarization and Beyond*. PhD thesis, International Institute of Information Technology Hyderabad, 2024.
- [15] Ahmed El-Sayed, Laila El-Haddad, and Mohamed Ali. Natural language processing advancements: A survey of transformer models and beyond. *Artificial Intelligence and Machine Learning Review*, 2(2):1–9, 2021.
- [16] Kyuhong Shim, Jungwook Choi, and Wonyong Sung. Exploring attention map reuse for efficient transformer neural networks, 2023.
- [17] Amirhesam Abedsoltan, Adityanarayanan Radhakrishnan, Jingfeng Wu, and Mikhail Belkin. Context-scaling versus task-scaling in in-context learning, 2024.
- [18] Tom Kocmi. Exploring benefits of transfer learning in neural machine translation, 2020.
- [19] Wen Liang and Youzhi Liang. Bpdec: Unveiling the potential of masked language modeling decoder in bert pretraining, 2024.
- [20] Lucas Georges Gabriel Charpentier and David Samuel. Gpt or bert: why not both?, 2024.
- [21] Yuchen Li, Alexandre Kirchmeyer, Aashay Mehta, Yilong Qin, Boris Dadachev, Kishore Papineni, Sanjiv Kumar, and Andrej Risteski. Promises and pitfalls of generative masked language modeling: Theoretical framework and practical guidelines, 2024.

- 
- [22] Yifei Xu, Jingqiao Zhang, Ru He, Liangzhu Ge, Chao Yang, Cheng Yang, and Ying Nian Wu. Sas: Self-augmentation strategy for language model pre-training, 2023.
- [23] Himmet Toprak Kesgin and Mehmet Fatih Amasyali. Iterative mask filling: An effective text augmentation method using masked language modeling, 2024.
- [24] Jenny Kunz and Oskar Holmström. The impact of language adapters in cross-lingual transfer for nlu, 2024.
- [25] Jiasheng Ye, Zaixiang Zheng, Yu Bao, Lihua Qian, and Quanquan Gu. Diffusion language models can perform many tasks with scaling and instruction-finetuning, 2023.
- [26] Rohan Jagtap and Sudhir N. Dhage. An in-depth walkthrough on evolution of neural machine translation, 2020.
- [27] Ridewaan Hanslo. Deep learning transformer architecture for named entity recognition on low resourced languages: State of the art results, 2022.
- [28] He Bai. Novel methods for natural language modeling and pretraining. 2023.
- [29] Jiajing Chen, Shuo Wang, Zhen Qi, Zhenhong Zhang, Chihang Wang, and Hongye Zheng. A combined encoder and transformer approach for coherent and high-quality text generation, 2024.
- [30] Ivan Chelombiev, Daniel Justus, Douglas Orr, Anastasia Dietrich, Frithjof Gressmann, Alexandros Kolioussis, and Carlo Luschi. Groupbert: Enhanced transformer architecture with efficient grouped structures, 2021.
- [31] Shuohuan Wang, Yu Sun, Yang Xiang, Zhihua Wu, Siyu Ding, Weibao Gong, Shikun Feng, Junyuan Shang, Yanbin Zhao, Chao Pang, Jiaxiang Liu, Xuyi Chen, Yuxiang Lu, Weixin Liu, Xi Wang, Yangfan Bai, Qiuliang Chen, Li Zhao, Shiyong Li, Peng Sun, Dianhai Yu, Yanjun Ma, Hao Tian, Hua Wu, Tian Wu, Wei Zeng, Ge Li, Wen Gao, and Haifeng Wang. Ernie 3.0 titan: Exploring larger-scale knowledge enhanced pre-training for language understanding and generation, 2021.
- [32] Harish Thangaraj, Ananya Chenat, Jaskaran Singh Walia, and Vukosi Marivate. Cross-lingual transfer of multilingual models on low resource african languages, 2024.
- [33] Chi-Liang Liu, Tsung-Yuan Hsu, Yung-Sung Chuang, and Hung yi Lee. What makes multilingual bert multilingual?, 2020.
- [34] Ivan Montero, Nikolaos Pappas, and Noah A. Smith. Sentence bottleneck autoencoders from transformer language models, 2021.
- [35] Xuyang Shen, Dong Li, Ruitao Leng, Zhen Qin, Weigao Sun, and Yiran Zhong. Scaling laws for linear complexity language models, 2024.
- [36] Candida M Greco and Andrea Tagarelli. Bringing order into the realm of transformer-based language models for artificial intelligence and law. *Artificial Intelligence and Law*, 32(4):863–1010, 2024.
- [37] Sandra Wankmüller. Introduction to neural transfer learning with transformers for social science text analysis, 2022.
- [38] Zhi Hong, Aswathy Ajith, Gregory Pauloski, Eamon Duede, Kyle Chard, and Ian Foster. The diminishing returns of masked language models to science, 2023.
- [39] Jan Philip Wahle, Terry Ruas, Frederic Kirstein, and Bela Gipp. How large language models are transforming machine-paraphrased plagiarism, 2022.
- [40] Shashank Sonkar and Richard G. Baraniuk. Investigating the role of feed-forward networks in transformers using parallel attention and feed-forward net design, 2023.

- 
- [41] Shanshan Wang, Zhumin Chen, Zhaochun Ren, Huasheng Liang, Qiang Yan, and Pengjie Ren. Paying more attention to self-attention: Improving pre-trained language models via attention guiding, 2022.
  - [42] Yueyao Yu and Yin Zhang. Why "classic" transformers are shallow and how to make them go deep, 2024.
  - [43] Tam Nguyen, Tan M. Nguyen, Dung D. Le, Duy Khuong Nguyen, Viet-Anh Tran, Richard G. Baraniuk, Nhat Ho, and Stanley J. Osher. Improving transformers with probabilistic attention keys, 2022.
  - [44] Benjamin Hoover, Yuchen Liang, Bao Pham, Rameswar Panda, Hendrik Strobelt, Duen Horng Chau, Mohammed J. Zaki, and Dmitry Krotov. Energy transformer, 2023.
  - [45] Jesse Vig. A multiscale visualization of attention in the transformer model, 2019.
  - [46] Immanuel Trummer. From bert to gpt-3 codex: Harnessing the potential of very large language models for data management, 2023.
  - [47] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp, 2019.
  - [48] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.
  - [49] Qihuang Zhong, Liang Ding, Juhua Liu, Bo Du, and Dacheng Tao. Self-evolution learning for discriminative language model pretraining, 2023.
  - [50] Shikun Shen, Yifei Zou, Yuan Yuan, Yanwei Zheng, Peng Li, Xiuzhen Cheng, and Dongxiao Yu. Federating to grow transformers with constrained resources without model sharing, 2024.
  - [51] Beiduo Chen, Wu Guo, Quan Liu, and Kun Tao. Feature aggregation in zero-shot cross-lingual transfer using multilingual bert, 2022.
  - [52] Philipp Dufter and Hinrich Schütze. Identifying necessary elements for bert’s multilinguality, 2021.
  - [53] Ningyu Xu, Tao Gui, Ruotian Ma, Qi Zhang, Jingting Ye, Menghan Zhang, and Xuanjing Huang. Cross-linguistic syntactic difference in multilingual bert: How good is it and how does it affect transfer?, 2022.
  - [54] Jordi Armengol-Estapé, Ona de Gibert Bonet, and Maite Melero. On the multilingual capabilities of very large-scale english language models, 2021.
  - [55] Bonan Min, Hayley Ross, Elinor Sulem, Amir Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz, Eneko Agirre, Ilana Heintz, and Dan Roth. Recent advances in natural language processing via large pre-trained language models: A survey. *ACM Computing Surveys*, 56(2):1–40, 2023.
  - [56] Sinclair Schneider, Florian Steuber, Joao A. G. Schneider, and Gabi Dreo Rodosek. How well can machine-generated texts be identified and can language models be trained to avoid identification?, 2023.
  - [57] Shen Nie, Fengqi Zhu, Chao Du, Tianyu Pang, Qian Liu, Guangtao Zeng, Min Lin, and Chongxuan Li. Scaling up masked diffusion models on text, 2024.
  - [58] Narsimha Chilukuri, Eric Hunsberger, Aaron Voelker, Gurshaant Malik, and Chris Eliasmith. Language modeling using lmus: 10x better data efficiency or improved scaling compared to transformers, 2021.

- 
- [59] Debora Nozza, Federico Bianchi, and Dirk Hovy. What the [mask]? making sense of language-specific bert models, 2020.
- [60] Chapter 3 improving pre-trained.
- [61] Patrick Schramowski, Cigdem Turan, Nico Andersen, Constantin A. Rothkopf, and Kristian Kersting. Large pre-trained language models contain human-like biases of what is right and wrong to do, 2022.
- [62] Ofir Press, Noah A. Smith, and Omer Levy. Improving transformer models by reordering their sublayers, 2020.
- [63] Félix Gaschi, Patricio Cerda, Parisa Rastin, and Yannick Toussaint. Exploring the relationship between alignment and cross-lingual transfer in multilingual transformers, 2023.
- [64] Lei Sha, Yuhang Song, Yordan Yordanov, Tommaso Salvatori, and Thomas Lukasiewicz. Bird-eye transformers for text generation models, 2022.
- [65] Chendong Zhao, Jianzong Wang, Wen qi Wei, Xiaoyang Qu, Haoqian Wang, and Jing Xiao. Adaptive sparse and monotonic attention for transformer-based automatic speech recognition, 2022.
- [66] Hanxiao Lu, Hongyu Cai, Yiming Liang, Antonio Bianchi, and Z. Berkay Celik. A progressive transformer for unifying binary code embedding and knowledge transfer, 2024.
- [67] Aidan Mannion, Thierry Chevalier, Didier Schwab, and Lorraine Geouriot. Umls-kgi-bert: Data-centric knowledge integration in transformers for biomedical entity recognition, 2023.
- [68] Sainbayar Sukhbaatar, Edouard Grave, Piotr Bojanowski, and Armand Joulin. Adaptive attention span in transformers, 2019.
- [69] Kaushik Roy, Yuxin Zi, Vignesh Narayanan, Manas Gaur, and Amit Sheth. Knowledge-infused self attention transformers, 2023.
- [70] Madhusudan Verma. Revisiting linformer with a modified self-attention with linear complexity, 2020.
- [71] Mor Geva, Avi Caciularu, Kevin Ro Wang, and Yoav Goldberg. Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space, 2022.
- [72] Archit Parnami, Rahul Singh, and Tarun Joshi. Pruning attention heads of transformer models using a\* search: A novel approach to compress big nlp architectures, 2021.
- [73] Sovesh Mohapatra and Somesh Mohapatra. The (in)effectiveness of intermediate task training for domain adaptation and cross-lingual transfer learning, 2022.
- [74] Zewen Chi, Heyan Huang, Luyang Liu, Yu Bai, and Xian-Ling Mao. Cross-lingual language model meta-pretraining, 2021.
- [75] Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. On the cross-lingual transferability of monolingual representations, 2020.
- [76] Shaomu Tan, Di Wu, and Christof Monz. Neuron specialization: Leveraging intrinsic task modularity for multilingual machine translation, 2024.
- [77] Johannes von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent, 2023.
- [78] Zhen Qi, Jiajing Chen, Shuo Wang, Bingying Liu, Hongye Zheng, and Chihang Wang. Optimizing multi-task learning for enhanced performance in large language models, 2024.
- [79] Alan Ansell, Edoardo Maria Ponti, Anna Korhonen, and Ivan Vulić. Distilling efficient language-specific models for cross-lingual transfer, 2023.

- 
- [80] Chaoyang He, Shen Li, Mahdi Soltanolkotabi, and Salman Avestimehr. Pipetransformer: Automated elastic pipelining for distributed training of transformers, 2021.
  - [81] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.
  - [82] Piotr Nawrot, Szymon Tworkowski, Michał Tyrolski, Łukasz Kaiser, Yuhuai Wu, Christian Szegedy, and Henryk Michalewski. Hierarchical transformers are more efficient language models, 2022.
  - [83] Quentin Anthony, Benjamin Michalowicz, Jacob Hatef, Lang Xu, Mustafa Abduljabbar, Aamir Shafi, Hari Subramoni, and Dhabaleswar Panda. Demystifying the communication characteristics for distributed transformer models, 2024.
  - [84] Kamal Acharya, Alvaro Velasquez, and Houbing Herbert Song. A survey on symbolic knowledge distillation of large language models, 2024.
  - [85] Vishwas M. Shetty, Metilda Sagaya Mary N J, and S. Umesh. Investigation of speaker-adaptation methods in transformer based asr, 2021.
  - [86] Keqi Deng, Shinji Watanabe, Jiatong Shi, and Siddhant Arora. Blockwise streaming transformer for spoken language understanding and simultaneous speech translation, 2022.
  - [87] Ed S. Ma. Investigating masking-based data generation in language models, 2023.
  - [88] Laria Reynolds and Kyle McDonell. Multiversal views on language models, 2021.
  - [89] Zuzana Jelčicová and Marian Verhelst. Delta keyword transformer: Bringing transformers to the edge through dynamically pruned multi-head self-attention, 2022.
  - [90] Junjie Yang and Hai Zhao. Deepening hidden representations from pre-trained language models, 2020.
  - [91] Xavier Amatriain, Ananth Sankar, Jie Bing, Praveen Kumar Bodigutla, Timothy J Hazen, and Michael Kazi. Transformer models: an introduction and catalog. *arXiv preprint arXiv:2302.07730*, 2023.
  - [92] Xavier Suau, Luca Zappella, and Nicholas Apostoloff. Self-conditioning pre-trained language models, 2023.
  - [93] Betty van Aken, Benjamin Winter, Alexander Löser, and Felix A. Gers. How does bert answer questions? a layer-wise analysis of transformer representations, 2019.
  - [94] Daria Pylypenko, Kwabena Amponsah-Kaakyire, Koel Dutta Chowdhury, Josef van Genabith, and Cristina España-Bonet. Comparing feature-engineering and feature-learning approaches for multilingual translationese classification, 2021.
  - [95] Hezekiah J. Branch, Jonathan Rodriguez Cefalu, Jeremy McHugh, Leyla Hujer, Aditya Bahl, Daniel del Castillo Iglesias, Ron Heichman, and Ramesh Darwishi. Evaluating the susceptibility of pre-trained language models via handcrafted adversarial examples, 2022.
  - [96] Xiaoguang Tu, Zhi He, Yi Huang, Zhi-Hao Zhang, Ming Yang, and Jian Zhao. An overview of large ai models and their applications. *Visual Intelligence*, 2(1):1–22, 2024.
  - [97] Nour Eddine Zekaoui, Siham Yousfi, Maryem Rhanoui, and Mounia Mikram. Analysis of the evolution of advanced transformer-based language models: Experiments on opinion mining. *arXiv preprint arXiv:2308.03235*, 2023.
  - [98] Chancellor R. Woolsey, Prakash Bisht, Joshua Rothman, and Gondy Leroy. Utilizing large language models to generate synthetic data to increase the performance of bert-based neural networks, 2024.
  - [99] Chenguang Wang, Zihao Ye, Aston Zhang, Zheng Zhang, and Alexander J. Smola. Transformer on a diet, 2020.

---

**Disclaimer:**

SurveyX is an AI-powered system designed to automate the generation of surveys. While it aims to produce high-quality, coherent, and comprehensive surveys with accurate citations, the final output is derived from the AI's synthesis of pre-processed materials, which may contain limitations or inaccuracies. As such, the generated content should not be used for academic publication or formal submissions and must be independently reviewed and verified. The developers of SurveyX do not assume responsibility for any errors or consequences arising from the use of the generated surveys.

www.SurveyX.cn