

## A Comprehensive Survey of Text Encoders for Text-to-Image Diffusion Models

Shun Fang<sup>1,2,\*</sup>

<sup>1</sup> Peking University, Haidian, Beijing, 100091, China

<sup>2</sup> Lumverse Inc., Shijingshan, Beijing, 100043, China

### Abstract

In this comprehensive survey, we delve into the realm of text encoders for text-to-image diffusion models, focusing on the principles, challenges, and opportunities associated with these encoders. We explore the state-of-the-art models, including BERT, T5-XXL, and CLIP, that have revolutionized the way we approach language understanding and cross-modal interactions. These models, with their unique architectures and training techniques, enable remarkable capabilities in generating images from textual descriptions. However, they also face limitations and challenges, such as computational complexity and data scarcity. We discuss these issues and highlight potential opportunities for further research. By providing a comprehensive overview, this survey aims to contribute to the ongoing development of text-to-image diffusion models, enabling more accurate and efficient image generation from textual inputs.

**Keywords:** NLP, CLIP, T5-XXL, BERT, Text Encoder.

Received on 11 02 2024, accepted on 11 06 2024, published on 18 07 2024

Copyright © 2024 Fang *et al.*, licensed to EAI. This is an open access article distributed under the terms of the [CC BY-NC-SA 4.0](https://creativecommons.org/licenses/by-nc-sa/4.0/), which permits copying, redistributing, remixing, transformation, and building upon the material in any medium so long as the original work is properly cited.

doi: 10.4108/airo.5566

### 1. Introduction

In recent years, pre-training techniques that directly learn from raw text have transformed the landscape of the natural language processing (NLP) [1-7]. Language model pre-training has emerged as a highly effective approach for enhancing numerous natural language processing tasks [8-10]. Initially, early results on transfer learning for NLP made use of recurrent neural networks [2, 8]. However, in recent times, the utilization of models based on the “Transformer” architecture [11] has become increasingly prevalent, particularly BERT [4], T5-XXL [5], and CLIP [12]. The Transformer, initially demonstrated to be effective in machine translation, has subsequently found applications in a diverse range of NLP scenarios, further solidifying its position as a pivotal component in modern language processing [13, 14].

BERT, standing for Bidirectional Encoder Representations from Transformers, has established itself as a pivotal model in NLP. Its unique architecture allows it to capture rich contextual information from both directions of a text sequence, enabling it to achieve state-of-the-art performance on a wide range of NLP tasks. The success of BERT has spawned a flurry of research and innovations, further pushing the boundaries of language understanding.

T5-XXL, an extension of the Transformer architecture, builds on the strengths of BERT while scaling up to even larger model sizes. This scalability allows T5-XXL to tackle more complex tasks and achieve even better performance. Its ability to process vast amounts of data and capture intricate patterns in language makes it a powerful tool for addressing challenging NLP problems.

CLIP, on the other hand, represents a significant leap in cross-modal learning. By jointly training on text and image data, CLIP learns to align visual and textual representations, enabling it to perform tasks that require understanding and interaction between different modalities. This capability

\*Corresponding author. Email: [fangshun@pku.org.cn](mailto:fangshun@pku.org.cn)

opens up new avenues for applications in areas such as image captioning, visual question answering, and even augmented reality.

In this review paper, we delve into the principles, challenges, and opportunities associated with these three ground-breaking models. We explore their underlying architectures, training techniques, and the mechanisms that enable them to achieve such remarkable performance. We also discuss the limitations and challenges faced by these models, as well as the potential opportunities for further research and development.

Moreover, we highlight the significant impact these models have had on various application domains. Whether it's improving the accuracy of sentiment analysis in social media monitoring, enhancing the efficiency of document classification in legal systems, or enabling more intuitive and engaging user interactions in augmented reality applications, these models are poised to revolutionize the way we interact with and understand language and visual data.

By providing a comprehensive overview of BERT, T5-XXL, and CLIP, this review paper aims to contribute to the ongoing dialogue in the field of NLP and cross-modal learning. We hope that it will serve as a valuable resource for researchers, practitioners, and enthusiasts alike, inspiring further innovations and advancements in these exciting areas of artificial intelligence.

In the realm of NLP and text-to-image diffusion models, this paper introduces several pivotal contributions that advance the field's understanding and capabilities.

- a) We offer an in-depth analysis of state-of-the-art text encoders, particularly BERT, T5-XXL, and CLIP, which have revolutionized language comprehension and cross-modal interactions. This thorough exploration aims to provide readers with a comprehensive understanding of these models' capabilities and implications.
- b) We delve into the principles, architectures, and training techniques that underlie these text encoders. By examining their inner workings, we aim to reveal how they transform textual descriptions into vivid visual representations.
- c) We discuss the challenges and limitations that currently confront these models, such as computational complexity and data scarcity. Recognizing these obstacles is essential for guiding future research and development, as it highlights the areas that require further attention and innovation.
- d) We identify potential research opportunities by highlighting the challenges faced by the current models. We encourage the pursuit of more accurate and efficient methods for generating images from textual inputs, paving the way for novel advances in the field.
- e) This paper serves as a comprehensive overview of the current state of text encoders in the realm of NLP and cross-modal learning. By consolidating the knowledge scattered across various studies, we aim to contribute to the ongoing dialogue in this exciting area of research.
- f) We discuss the impact that these models have had on diverse application domains, ranging from social media

monitoring to legal systems and augmented reality applications. This examination underscores the practical significance and broad applicability of the research presented in this paper.

In this paper, we embark on a thorough exploration of the foundational principles and empirical results surrounding three influential text encoders: BERT, T5-XXL, and CLIP. Section 2 delves into the underlying mechanisms of these encoders, evaluates their performance, and offers insights into their respective strengths, as well as prospects for future advancements. In Section 3, we discuss the challenges confronting these models, including computational complexity and data scarcity, while highlighting avenues for further research and development. Finally, Section 4 provides a synthesis of the advancements made by BERT, T5-XXL, and CLIP in the fields of natural language processing (NLP) and cross-modal learning. We reflect on their common principles, unique characteristics, and the challenges they pose, while also peering into potential future research directions and the promise of more robust and versatile models.

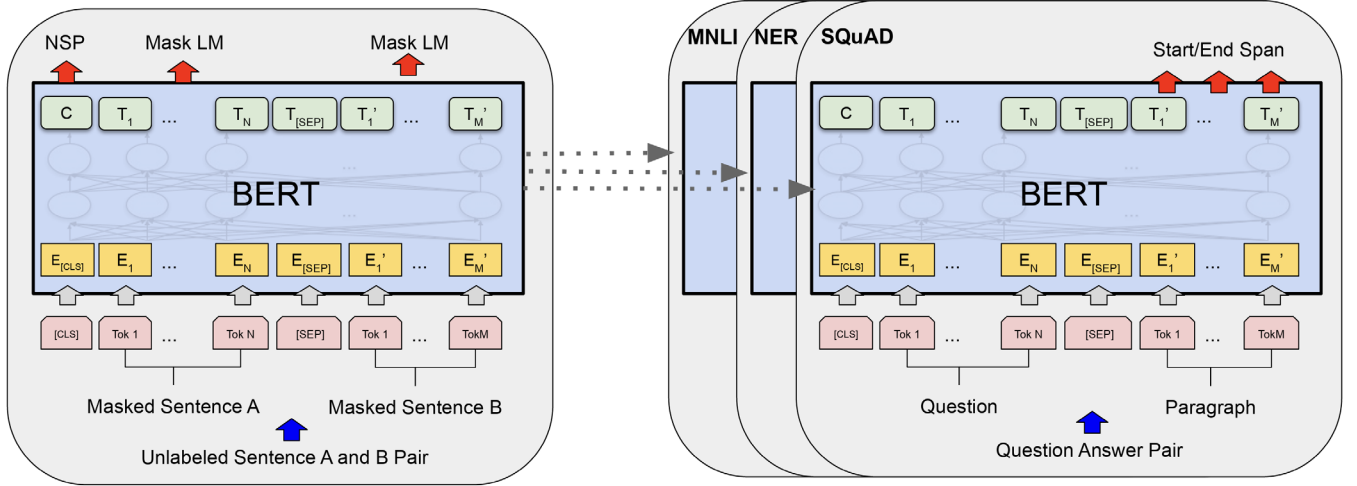
## 2. Methods

The current chapter delves into the core principles and experimental outcomes of three pivotal text encoders: BERT, T5-XXL, and CLIP. These models have revolutionized the landscape of text-to-image diffusion, enabling unprecedented capabilities in generating visually rich representations from textual inputs. We examine their underlying mechanisms and assess their performance, aiming to provide a comprehensive understanding of their strengths and potential for future advancements.

### 2.1 Adapting Deep Bidirectional Transformers through Pre-training for Enhanced Language Comprehension

#### Overview

BERT [4], or Bidirectional Encoder Representations from Transformers, is a groundbreaking language representation model tailored for deep bidirectional pretraining of unlabeled text. It uniquely conditions on both left and right context across all layers, enabling robust representation learning. Remarkably, BERT can be easily fine-tuned with an additional output layer, achieving cutting-edge performance across various natural language processing tasks, including question answering and language inference. Its architecture comprises a multi-layer bidirectional Transformer encoder, maintaining a consistent structure across diverse tasks (Figure 1). This flexibility extends to its input representation, which effortlessly handles both individual sentences and sentence pairs, facilitating its application to a broad spectrum of downstream tasks. The extensive pre-training data utilized by BERT, encompassing the BooksCorpus [15] and English Wikipedia, has resulted in significant advancements in tasks



**Figure 1.** The general pre-training and fine-tuning methodologies employed for BERT are outlined as follows. Notably, aside from the output layers, identical architectural frameworks are adopted for both pre-training and fine-tuning stages. Consequently, the pre-trained model parameters serve as the foundation for initializing models tailored to various downstream tasks. During the fine-tuning phase, all parameters undergo meticulous adjustment. Furthermore, a unique symbol designated as [CLS] is introduced at the commencement of each input instance, while [SEP] serves as a distinct separator token [4].

like General Language Understanding Evaluation (GLUE) score [16], Multi-Genre Natural Language Inference (MultiNLI) [17] accuracy, and the Stanford Question Answering Dataset (SQuAD) [18] question answering. In summary, BERT's bidirectional pretraining approach and its unified architecture underscore its simplicity and empirical prowess in advancing language understanding tasks.

### Neural Network

The BERT framework encompasses two fundamental stages: pre-training and fine-tuning. In the pre-training phase, the model undergoes rigorous training on unlabeled data [19], leveraging diverse pre-training tasks. Subsequently, for fine-tuning, the pre-trained BERT model initializes its parameters and undergoes further customization using labeled data tailored to specific downstream tasks. Notably, BERT boasts a uniform architecture that remains consistent across diverse tasks, exhibiting minimal architectural differences between its pre-trained and downstream configurations. The model's architecture is built upon a multi-layer bidirectional Transformer encoder, with key parameters such as the number of layers (L), hidden size (H), and the count of self-attention heads (A) tailored for different model sizes, namely  $BERT_{BASE}$  and  $BERT_{LARGE}$ .

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (1)$$

Eq.1 computes an attention score for each word in the input sequence with respect to a query word. The softmax function ensures that the model focuses on the most relevant words, and the scaling factor  $\sqrt{d_k}$  helps to stabilize the training process. The result is a contextually enriched representation

of the input sequence that the model can use for further processing.

- $Q, K, V$ : These represent the query, key, and value matrices respectively. In the context of the Transformer model, the query comes from the previous decoder layer, while the key and value come from the output of the encoder layer. These matrices are formed by concatenating the vectors for each word in the input sequence.
- $QK^T$ : This denotes the matrix multiplication of the query matrix  $Q$  with the transpose of the key matrix  $K$ . Each element  $i, j$  in the resulting matrix represents the dot product between the  $i$ -th query vector and the  $j$ -th key vector.
- $\sqrt{d_k}$ : This is the square root of the dimension of the key vectors  $d_k$ . Dividing the dot product by  $\sqrt{d_k}$  serves to scale the attention scores, which helps prevent the softmax function from having very small gradients during training when the key vectors have large dot products.
- softmax function**: The softmax function is applied to each row of the scaled dot product matrix. It converts the scaled dot products into probabilities that sum to 1. This step is crucial as it allows the model to focus more on certain parts of the input sequence that are more relevant to the current word being processed.
- $V$ : Finally, the softmax-normalized attention weights are used to take a weighted sum of the values  $V$ . This results in a weighted sum of the input values based on their relevance to the current query, which is the output of the attention mechanism for a given position.

BERT's input representations are meticulously designed to accommodate both individual sentences and sentence pairs. It employs WordPiece embeddings [20], supported by a comprehensive 30,000-token vocabulary. The input

sequence incorporates specialized tokens, serving as markers for classification tasks and sentence separators. Additionally, learned embeddings are employed to distinguish between sentence A and sentence B. Fine-tuning BERT is a straightforward process, attributed to the Transformer's self-attention mechanism. This mechanism enables the model to seamlessly adapt to a range of downstream tasks by adjusting its inputs and outputs accordingly. Fine-tuning is relatively cost-effective compared to pre-training, with replicable results achieved within a brief timeframe using the same pre-trained model. Overall, the intricate explanation of BERT's methodology underscores its adaptability and effectiveness in tackling a diverse array of natural language processing tasks.

The method incorporates a loss function composed of two principal elements: the Masked Language Model (MLM) [37-58] objective and the Next Sentence Prediction task.

With regard to the MLM objective, a random proportion of input tokens are masked, and the model is trained to reconstruct the original vocabulary id of the masked tokens by drawing upon contextual cues. The hidden vectors corresponding to the masked tokens are then fed into a softmax function across the vocabulary to generate predictions. This approach allows the model to capture bidirectional representations by merging the left and right context, thus overcoming the limitations of traditional unidirectional language models.

For the Next Sentence Prediction task, text-pair representations are jointly pre-trained by assessing whether a given pair of sentences follows one another in the original text. This task enhances the model's comprehension of the relationship between two sentences and bolsters its proficiency in handling tasks that involve text pairs, such as question answering.

During fine-tuning, the self-attention mechanism within the Transformer architecture enables the model to adapt to specific downstream tasks by modifying the inputs and outputs accordingly. The fine-tuning process involves substituting the appropriate inputs and outputs for each task and updating all parameters end-to-end. Compared to pre-training, this process is relatively cost-effective and can be efficiently executed on diverse hardware configurations.

Overall, the loss function employed in this method integrates the objectives of predicting masked tokens and determining the relationship between text pairs to pre-train a deep bidirectional Transformer model, which can be subsequently fine-tuned for a range of natural language processing tasks.

## Experiments

The BERT model has demonstrated its efficacy in numerous natural language processing tasks through experimental results (Table 1). Utilizing the GLUE benchmark [16], the performance of BERT was evaluated across 11 NLP tasks, encompassing natural language inference, question answering, and sentiment analysis, among others.

During fine-tuning, BERT adapts to specific downstream tasks by leveraging its self-attention mechanism and

adjusting inputs and outputs accordingly. This process is relatively cost-effective compared to pre-training and can be efficiently conducted across diverse hardware configurations. By harnessing bidirectional cross attention and integrating text pair encoding with self-attention, the model is capable of handling tasks involving both single text and text pairs.

Table 1. SQuAD2.0 results. We exclude entries that employ BERT as a constituent element in their composition [4].

System	Dev		Test	
	EM	F1	EM	F1
Top Leaderboard Systems(Dec 10 <sup>th</sup> 2018)				
Human	86.3	89.0	86.9	89.5
#1 Single-MIR-MRC(F-Net)	-	-	74.8	78.0
#2 Single-nlnet	-	-	74.2	77.1
Published				
unet(Ensemble)	-	-	71.4	74.9
SLQA+(Single)	-	-	71.4	74.4
<b>BERT<sub>LARGE</sub>(Single)</b>	<b>78.7</b>	<b>81.9</b>	<b>80.0</b>	<b>83.1</b>

The experiments further examined the influence of model size on fine-tuning task accuracy. Models with larger dimensions, including more layers, hidden units, and attention heads, exhibited improved accuracy across various tasks, even on datasets with limited labeled training examples. These results indicate that scaling to larger model sizes can significantly enhance performance on both large-scale and small-scale tasks, given sufficient pre-training.

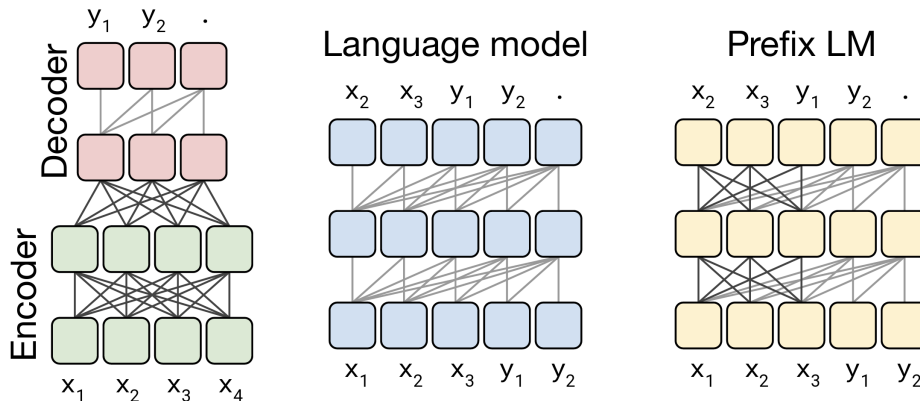
Additionally, the study investigated the impact of various masking strategies during pre-training on the performance of the BERT model. While fine-tuning appears robust to different masking strategies, exclusive reliance on the MASK or RND strategy can result in suboptimal performance in certain tasks, such as named entity recognition. This underscores the significance of pre-training tasks and the bidirectional nature of BERT in achieving superior performance across a wide range of NLP tasks.

## 2.2 Delving into the Boundaries of Transfer Learning Utilizing a Unified Transformer for Text-to-Text Conversion

### Overview

The T5-XXL [5] model delves into the realm of transfer learning techniques in NLP through the introduction of a unified Text-to-Text Transfer Transformer (T5) framework. This model compares a range of pre-training objectives, architectures, unlabeled datasets [21-23], and transfer strategies across various language understanding tasks. By capitalizing on the insights gained from its exploration, combined with its scale and a novel dataset named the Colossal Clean Crawled Corpus (C4), T5-XXL attains





**Figure 2.** Illustrations of the Transformer architecture variants we explore are presented here. In this visualization, blocks symbolize elements within a sequence, while lines depict the visibility of attention. Distinct color groupings of blocks represent varying stacks of Transformer layers. Dark grey lines signify full visibility masking, whereas light grey lines indicate causal masking. We employ the symbol “.” to represent a unique end-of-sequence token, marking the termination of a prediction. The input and output sequences are designated as  $x$  and  $y$ , respectively. On the left, a standard encoder-decoder architecture incorporates full visibility masking in both the encoder and encoder-decoder attention, along with causal masking in the decoder. In the middle, a language model comprises a single stack of Transformer layers and receives the concatenation of input and target, utilizing a causal mask throughout. On the right, appending a prefix to a language model allows for full visibility masking over the input [5].

cutting-edge performance on diverse benchmarks, encompassing tasks such as summarization, question answering, and text classification. The objective is to present a comprehensive overview of the current state of transfer learning in NLP and to foster future research in this domain by disseminating its dataset, pre-trained models, and associated code. T5-XXL strives to enhance the comprehension and utilization of transfer learning in NLP, exemplifying the prowess and promise of this technique in addressing a diverse array of language understanding challenges.

### Neural Network

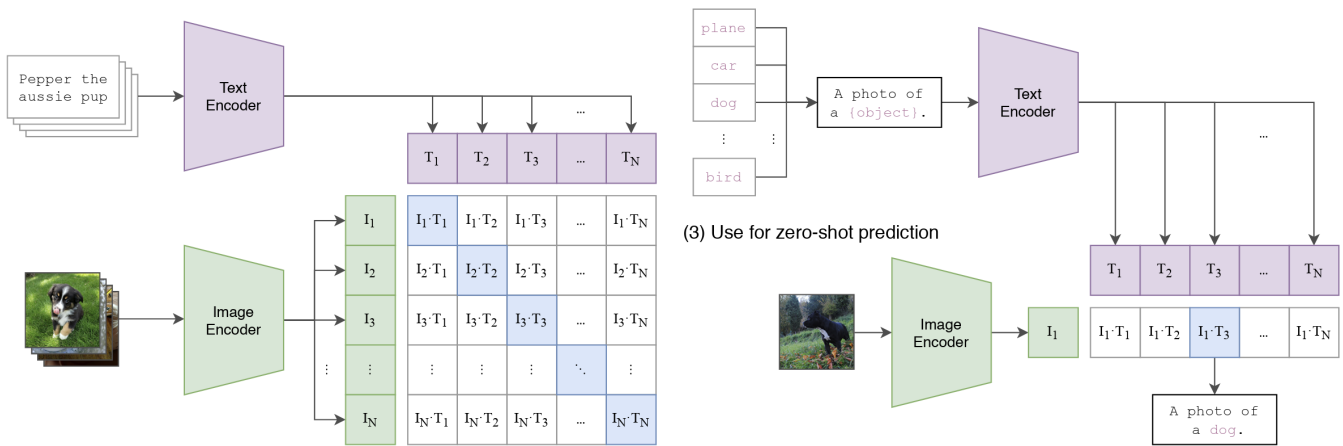
The Method delves exhaustively into the principles, processes, and steps inherent in their transfer learning approach within the realm of NLP, utilizing the T5 framework. This approach revolutionizes the landscape of NLP by introducing a unified text-to-text format, which transforms diverse text-based language challenges into a standardized structure (Figure 2). This standardization enables a methodical exploration of transfer learning techniques in NLP. The methodology involves a rigorous comparison of diverse pre-training objectives, architectures, unlabeled datasets, and transfer methods across a broad array of language understanding tasks. By leveraging profound insights derived from this extensive analysis, coupled with the immense scale of the newly introduced C4 dataset, the T5-XXL achieves state-of-the-art performance across a range of NLP benchmarks, encompassing summarization, question answering, and text classification. Additionally, the process involves experimenting with model scaling through the augmentation of parameters and computational resources, exploring the impact of prolonged training durations, and optimizing training objectives to enhance the performance of

NLP tasks. Furthermore, the T5-XXL underscores the crucial role of model sizes and training objectives in achieving superior results, emphasizing the importance of scaling up models and engaging in longer training sessions to push the boundaries of transfer learning in NLP.

The neural network architecture is grounded in a standard encoder-decoder Transformer model. This architecture comprises two layer stacks: the encoder and the decoder. The encoder is responsible for processing the input sequence, whereas the decoder is tasked with generating the output sequence. Central to the model is a fully-visible attention mask, enabling the self-attention mechanism to focus on any input entry during output generation. This masking pattern is well-suited for attending to a "prefix" context provided to the model for predictive purposes.

The model is designed with encoder and decoder structures that are comparable in size and configuration to a "**BERT<sub>BASE</sub>**" stack. Each structure consists of 12 blocks, encompassing self-attention, optional encoder-decoder attention, and a feed-forward network. The feed-forward networks within each block possess distinct dimensions and are followed by ReLU nonlinearity layers. The attention mechanisms are equipped with 12 heads, resulting in a model with approximately 220 million parameters.

During the training phase, the model undergoes pre-training on the C4 dataset for a specified number of steps, utilizing a text-to-text format and maximum likelihood training with AdaFactor [24] optimization. Subsequently, fine-tuning is conducted on downstream tasks for a designated number of steps, maintaining a constant learning rate. The architecture and training process are optimized to achieve cutting-edge results in natural language processing tasks.



**Figure 3.** An overview of the CLIP methodology is as follows: Unlike traditional image models that concurrently train an image feature extractor and a linear classifier to forecast specific labels, CLIP concurrently trains both an image encoder and a text encoder. The objective is to accurately predict the matching pairs within a batch of (image, text) training examples. During testing, the trained text encoder creates a zero-shot linear classifier by embedding the names or descriptions of the target dataset's classes. This allows CLIP to associate images with their corresponding textual descriptions without the need for explicit label-based training [12].

### Experiments

The experimental outcomes offer insights into the baseline model's performance across diverse downstream tasks. Multiple training iterations from scratch were conducted to assess the performance variance across different runs. The findings reveal that pre-training significantly enhances performance across most benchmarks, apart from tasks involving vast amounts of data, such as WMT [25] English to French translation. The study further underscores the significance of pre-training in boosting performance on tasks with limited data, emphasizing the benefits of transfer learning in optimizing model efficiency. Inter-run variance was generally low across most tasks, with exceptions being low-resource tasks like CoLA [26], CB [27], and COPA [28], which exhibited higher variance. The results tables comprehensively detail the scores achieved on each task, reflecting the model's performance across various experimental settings. In summary, the experimental findings underscore the effectiveness of pre-training and highlight the influence of model size and training steps on task performance.

### 2.3 Acquiring Transferable Visual Models through Natural Language-Guided Supervision

#### Overview

The CLIP [12] framework introduces an innovative method for training computer vision systems through natural language supervision. This approach involves a pre-training task where the model aims to predict the caption that corresponds to a given image, leveraging a vast dataset of 400 million (image, text) pairs sourced from the internet. By

extracting visual concepts directly from raw text, the model is able to achieve cutting-edge image representations and enable seamless zero-shot transfer to diverse downstream tasks without the need for task-specific training. The study evaluates the model's performance across over 30 computer vision datasets, exhibiting competitive results even compared to fully supervised baselines. In particular, it demonstrates comparable accuracy to ResNet-50 [29] on ImageNet [30, 31] zero-shot. This approach underscores the scalability, efficiency, and robustness of learning from natural language supervision, pointing to its promising potential for broader applications and ethical considerations within the realm of computer vision research.

#### Neural Network

The CLIP framework introduces a novel technique for training computer vision models using natural language supervision (Figure 3). Central to this approach is the pre-training phase, which involves utilizing a vast dataset of 400 million pairs sourced from the internet. During this phase, the model learns to associate captions with corresponding images, enabling it to capture state-of-the-art image representations directly from raw text, without relying on predefined object categories. Once pre-trained, the model can harness the power of natural language to reference and describe visual concepts, facilitating zero-shot transfer to a wide range of downstream tasks without the need for task-specific training.

The CLIP approach is evaluated across over 30 computer vision datasets, encompassing diverse tasks such as optical character recognition, action recognition in videos, and fine-grained object classification. Its performance is competitive with fully supervised baselines, and it even matches the accuracy of ResNet-50 on ImageNet for zero-shot learning.

Furthermore, the framework underscores the significance of prompt engineering and ensembling techniques in enhancing zero-shot performance, particularly in scenarios where class labels lack contextual information or exhibit polysemy issues.

By effectively leveraging natural language supervision and addressing these challenges, the CLIP framework demonstrates scalability, efficiency, and robustness in learning transferable visual models. This approach opens new avenues for broader applications and ethical considerations in the field of computer vision research.

The neural network architecture involves the training of a composite setup consisting of five ResNet models and three Vision Transformers. Among the ResNet models, diverse versions are utilized, including ResNet-50, ResNet-101, as well as scaled versions designated as RN50x4, RN50x16, and RN50x64. The Vision Transformers employ ViT-B/32, ViT-B/16, and ViT-L/14 [32]. The ResNet image encoders adopt a strategy that distributes additional computational resources across width, depth, and resolution to optimize performance.

The text encoder is based on a modified Transformer architecture and operates on a lowercased byte pair encoding representation of text. The training process involves 32 epochs, leveraging the Adam optimizer with decoupled weight decay regularization and a cosine learning rate decay schedule. To expedite training and conserve memory, various techniques are implemented, such as mixed-precision training, gradient checkpointing, and half-precision stochastic rounding for the text encoder weights.

## Experiments

Table 2. The accuracy rate achieved for classifying images of the 'White' category in FairFace, pertaining to Race, Gender, and Age, is expressed as a percentage [12].

Model	Age	Race	Gender
Linear Probe Instagram	54.2	90.8	93.2
Zero-Shot CLIP	57.1	58.3	95.9
Linear Probe CLIP	63.8	93.4	96.5
FairFace Model	59.7	93.7	94.2

The experimental outcomes reveal the CLIP model's proficiency across diverse tasks and datasets (Table 2). This study encompassed training models on tailored datasets, such as YFCC100M [33] and WIT, and comparing their performances in both zero-shot and linear probe settings. While the average performance was comparable between YFCC and WIT, there were noteworthy disparities in performance on specific fine-grained classification datasets, reflecting differences in data density and relevance within the pre-training datasets. Additionally, the study included an exploration of the model's sensitivity to classification terms carrying a high risk of representational harm, specifically focusing on denigration. When tested on classifying images from the FairFace [34] dataset, disparities in misclassification rates were observed across various

demographic subgroups, including race, age, and gender. Furthermore, experiments were conducted to assess how biases manifest in the model's output labels across various label probability thresholds, highlighting gendered associations in label distributions. The study also explored the model's performance in tasks such as coarse classification, fine-grained detection, and identity recognition using datasets like CelebA [35], demonstrating the model's strengths and limitations in diverse scenarios. Overall, these experimental findings offer valuable insights into the CLIP model's performance, biases, and potential societal implications across a range of tasks and datasets.

## 2.4 Comparison

This review paper has delved into the principles, challenges, and opportunities presented by BERT, T5-XXL, and CLIP. To provide a comparative analysis of these three models, it is important to consider their similarities and differences across various dimensions.

Firstly, in terms of principles, all three models are based on the transformer architecture and utilize pre-training on large-scale datasets to learn general-purpose representations. BERT and T5-XXL are primarily focused on language understanding tasks, while CLIP incorporates visual information, enabling it to perform tasks that require cross-modal understanding. This difference in scope underscores the versatility of CLIP and its potential to handle a wider range of tasks.

Moving on to challenges, all three models face computational resource constraints. Training and fine-tuning large-scale models such as BERT, T5-XXL, and CLIP require significant compute power and time. Additionally, dataset overlap and generalization capabilities are concerns shared by these models. Ensuring that the pre-training dataset does not overlap with evaluation datasets is crucial to avoid biased results, while the ability to generalize to unseen data is a key indicator of model robustness.

In terms of opportunities, there is significant potential for further research and development. As technology and compute power advance, more efficient training methods can be explored to scale these models even further. Additionally, addressing the challenges of dataset overlap and generalization can lead to more reliable and robust models. Moreover, there is a growing interest in multi-modal learning, where models like CLIP hold promise for combining different modalities such as text and images to solve complex tasks.

Finally, it is worth noting that each model has its unique strengths and weaknesses. BERT and T5-XXL excel in language-related tasks, while CLIP demonstrates impressive cross-modal capabilities. The choice of which model to use would depend on the specific task and domain of application.

In conclusion, BERT, T5-XXL, and CLIP are powerful models that have transformed the landscape of NLP and cross-modal learning. By comparing and contrasting their principles, challenges, and opportunities, we can gain a

deeper understanding of their strengths and weaknesses, and identify areas for future research and development.

### 3. Challenges and Opportunities

Optimization algorithms occupy a paramount position in the refinement and fine-tuning of text encoders, positioning them as a crucial and promising area of research. For instance, one approach [59] advocates the enhancement of computational efficiency in an incomplete data analysis model by utilizing an advanced Particle Swarm Optimization algorithm for hyper-parameter adjustment. Additionally, a novel Hybrid-of-Evolutionary-Schemes algorithm [60] is introduced for precision calibration of industrial robots, aiming for high accuracy. Furthermore, the utilization of diversified regularization techniques based on the Levenberg–Marquardt algorithm [61] is proposed to bolster the training process for effective manipulator calibration. These advancements demonstrate the significance and potential of optimization algorithms in the field of text encoder training and fine-tuning.

#### 3.1 The Challenges and Opportunities of the BERT

BERT demonstrates its proficiency in pre-training profound bidirectional representations across a range of natural language processing tasks. Nevertheless, the technique is not without its drawbacks and complexities. A significant constraint lies in the substantial computational resources necessary for training and refining large-scale BERT models, particularly for researchers constrained by limited access to high-performance computing facilities. Furthermore, the pre-training data sources, such as BooksCorpus [36] and English Wikipedia, might introduce biases or constraints in the model's comprehension of nuanced language expressions from diverse origins.

Looking ahead, there exist promising avenues for further exploration and enhancement of the BERT model. One promising direction is to delve into more diverse and exhaustive pre-training datasets, aiming to broaden the model's generalization abilities across various domains and languages. Another potential lies in refining the fine-tuning process to enhance its efficiency and accessibility to a more inclusive audience. Additionally, exploring methods to reduce the computational burden of training large models while preserving performance could pave new paths for deploying BERT in resource-limited settings. In summary, addressing these challenges and seizing these opportunities could propel advancements in leveraging BERT for a wider array of NLP applications.

#### 3.2 The Challenges and Opportunities of the T5-XXL

There are numerous constraints and difficulties encountered in the current methodological approach. A significant constraint lies in the ineffectiveness of the denoising objective employed for pre-training, potentially not being the most proficient approach to imparting generalized knowledge to the model. This underscores the necessity for more proficient methods of knowledge extraction, capable of enhancing fine-tuning performance without the prerequisite of extensive pre-training on vast datasets. Furthermore, the paper underscores the importance of formulating a more rigorous comprehension of task similarity between pre-training and downstream tasks, enabling more informed decisions when selecting unlabeled data for pre-training.

Another challenge identified is the absence of state-of-the-art results in translation tasks, particularly with English-only pre-training, indicating the need for deeper exploration of language-agnostic models capable of delivering robust performance across diverse languages.

Looking ahead, there exist avenues for future research to tackle these limitations and challenges. One potential direction involves developing more efficient and effective methods for imparting generalized knowledge to models during pre-training, possibly by exploring alternative techniques, such as distinguishing between authentic and machine-generated text. Additionally, there is an opportunity to push the boundaries of the field by formalizing a more robust concept of task similarity, drawing inspiration from related work in computer vision, to guide the selection of unlabeled data for pre-training. Moreover, the exploration of language-agnostic models offers an exciting prospect to enhance model performance across various languages and overcome the logistical complexities associated with language-specific pre-training. By addressing these areas, future research stands to improve the efficacy and applicability of transfer learning techniques in natural language processing.

#### 3.3 The Challenges and Opportunities of the CLIP

The approach outlined herein encounters numerous constraints and obstacles. A prime obstacle is the challenge of efficiently training extensive models for natural language supervision, given the immense computational resources required by contemporary computer vision systems. Another constraint lies in the potential overlap of data between the pre-training and evaluation datasets, which can undermine the model's generalization capabilities. Additionally, the reliance on prompt engineering and ensembling to boost zero-shot classification performance underscores the need for more robust and automated strategies for dataset labeling and task description handling.

Looking ahead, there are avenues for further exploration and enhancement in several domains. Firstly, developing more efficient training methodologies that can scale natural language supervision without resorting to extensive computational resources could broaden the accessibility and utility of such models. Secondly, addressing the challenges pertaining to dataset overlap and ensuring models' resilience



to out-of-distribution data can bolster the reliability and generalization abilities of the models. Lastly, exploring innovative methods for task description and label handling, beyond prompt engineering, can streamline the model training process and enhance performance across a broader range of tasks. By tackling these challenges and seizing these opportunities, future research can propel the capabilities and practical applications of models like CLIP in diverse fields.

## 4. Conclusion

The BERT, T5-XXL, and CLIP represent significant advancements in the field of natural language processing and cross-modal learning, leveraging the transformer architecture and large-scale pre-training to achieve state-of-the-art performance on a wide range of tasks. The analysis highlights the shared principles underlying these models, including their reliance on transformer-based architectures and the utilization of pre-training on vast datasets. However, each model also exhibits unique characteristics and strengths, reflecting the diversity and complexity of language and vision tasks. The paper also explores the challenges faced by these models, including computational resource constraints, dataset overlap, and generalization capabilities. Despite these challenges, the opportunities presented by these models are immense, offering new avenues for research and innovation. Looking ahead, future research is likely to focus on addressing the existing challenges while exploring new applications and extensions of these models. With advances in computing technology and the availability of larger datasets, we can expect even more powerful and versatile models to emerge in the future.

## References

- [1] Dai, A.M. and Le, Q.V., 2015. Semi-supervised sequence learning. *Advances in neural information processing systems*, 28.
- [2] Howard, J. and Ruder, S., 2018. Universal language model fine-tuning for text classification. *arXiv preprint arXiv:1801.06146*.
- [3] Radford, A., Narasimhan, K., Salimans, T. and Sutskever, I., 2018. Improving language understanding by generative pre-training.
- [4] Devlin, J., Chang, M.W., Lee, K. and Toutanova, K., 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- [5] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W. and Liu, P.J., 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, Vol.21, No.140, pp.1-67.
- [6] Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A. and Agarwal, S., 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33, pp.1877-1901.
- [7] Lyu, Q., Apidianaki, M. and Callison-Burch, C., 2024. Towards faithful model explanation in nlp: A survey. *Computational Linguistics*, pp.1-70.
- [8] Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018a. Deep contextualized word representations. In NAACL.
- [9] Radford, A., Narasimhan, K., Salimans, T. and Sutskever, I., 2018. Improving language understanding with unsupervised learning.
- [10] Howard, J. and Ruder, S., 2018. Universal language model fine-tuning for text classification. *arXiv preprint arXiv:1801.06146*.
- [11] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł. and Polosukhin, I., 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- [12] Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J. and Krueger, G., 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748-8763.
- [13] McCann, B., Keskar, N.S., Xiong, C. and Socher, R., 2018. The natural language decathlon: Multitask learning as question answering. *arXiv preprint arXiv:1806.08730*.
- [14] Yu, A.W., Dohan, D., Luong, M.T., Zhao, R., Chen, K., Norouzi, M. and Le, Q.V., 2018. Qanet: Combining local convolution with global self-attention for reading comprehension. *arXiv preprint arXiv:1804.09541*.
- [15] Zhu, Y., Kiros, R., Zemel, R., Salakhutdinov, R., Urtasun, R., Torralba, A., & Fidler, S., 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*, pp. 19-27.
- [16] Wang, A., Singh, A., Michael, J., Hill, F., Levy, O., & Bowman, S. R., 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*.
- [17] Williams, A., Nangia, N., & Bowman, S. R., 2017. A broad-coverage challenge corpus for sentence understanding through inference. *arXiv preprint arXiv:1704.05426*.
- [18] Rajpurkar, P., Zhang, J., Lopyrev, K., & Liang, P., 2016. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*.
- [19] Collobert, R., & Weston, J., 2008. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*, pp. 160-167.
- [20] Wu, Y., Schuster, M., Chen, Z., Le, Q. V., Norouzi, M., Macherey, W. & Dean, J., 2016. Google's neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.
- [21] Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R. R., & Le, Q. V., 2019. Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems*, 32.
- [22] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., & Stoyanov, V., 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- [23] Zellers, R., Holtzman, A., Rashkin, H., Bisk, Y., Farhadi, A., Roesner, F., & Choi, Y., 2019. Defending against neural fake news. *Advances in neural information processing systems*, 32.
- [24] Shazeer, N., & Stern, M., 2018. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, pp. 4596-4604.
- [25] Edunov, S., Ott, M., Auli, M., & Grangier, D., 2018. Understanding back-translation at scale. *arXiv preprint arXiv:1808.09381*.

- [26] Warstadt, A., Singh, A., & Bowman, S. R., 2019. Neural network acceptability judgments. *Transactions of the Association for Computational Linguistics*, 7, 625-641.
- [27] De Marneffe, M. C., Simons, M., & Tonhauser, J., 2019. The commitmentbank: Investigating projection in naturally occurring discourse. In *proceedings of Sinn und Bedeutung*, Vol. 23, No. 2, pp. 107-124.
- [28] Roemmele, M., Bejan, C. A., & Gordon, A. S., 2011. Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In *2011 AAAI Spring Symposium Series*.
- [29] Koonce, B., & Koonce, B., 2021. ResNet 50. Convolutional neural networks with swift for tensorflow: image recognition and dataset categorization, 63-72.
- [30] Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., & Fei-Fei, L., 2015. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115, 211-252.
- [31] Deng, J., Dong, W., Socher, R., Li, L. J., Li, K., & Fei-Fei, L., 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248-255.
- [32] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., & Houlsby, N., 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- [33] Thomee, B., Shamma, D. A., Friedland, G., Elizalde, B., Ni, K., Poland, D., & Li, L. J., 2016. Yfcc100m: The new data in multimedia research. *Communications of the ACM*, Vol.59, No.2, 64-73.
- [34] Karkkainen, K., & Joo, J., 2021. Fairface: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pp. 1548-1558.
- [35] Zhang, Y., Yin, Z., Li, Y., Yin, G., Yan, J., Shao, J., & Liu, Z., 2020. Celeba-spoof: Large-scale face anti-spoofing dataset with rich annotations. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XII 16*, pp. 70-85.
- [36] Zhu, Y., Kiros, R., Zemel, R., Salakhutdinov, R., Urtasun, R., Torralba, A., & Fidler, S., 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*, pp. 19-27.
- [37] Xia, Y., Sedova, A., de Araujo, P.H.L., Kougia, V., Nußbaumer, L. and Roth, B., 2024. Exploring prompts to elicit memorization in masked language model-based named entity recognition. *arXiv preprint arXiv:2405.03004*.
- [38] Hu, J. and Frank, M.C., 2024. Auxiliary task demands mask the capabilities of smaller language models. *arXiv preprint arXiv:2404.02418*.
- [39] Wiland, J., Ploner, M. and Akbik, A., 2024. BEAR: A Unified Framework for Evaluating Relational Knowledge in Causal and Masked Language Models. *arXiv preprint arXiv:2404.04113*.
- [40] Aguiar, M., Zweigenbaum, P. and Naderi, N., 2024. SEME at SemEval-2024 Task 2: Comparing Masked and Generative Language Models on Natural Language Inference for Clinical Trials. *arXiv preprint arXiv:2404.03977*.
- [41] Yu, J., Kim, S.U., Choi, J. and Choi, J.D., 2024. What is Your Favorite Gender, MLM? Gender Bias Evaluation in Multilingual Masked Language Models. *arXiv preprint arXiv:2404.06621*.
- [42] Thennal, D.K., Nathan, G. and Suchithra, M.S., 2024, May. Fisher Mask Nodes for Language Model Merging. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation*, pp. 7349-7355.
- [43] Labrak, Y., Bazoge, A., Daille, B., Rouvier, M. and Dufour, R., 2024. How Important Is Tokenization in French Medical Masked Language Models?. *arXiv preprint arXiv:2402.15010*.
- [44] Naguib, M., Tannier, X. and Névéol, A., 2024. Few shot clinical entity recognition in three languages: Masked language models outperform LLM prompting. *arXiv preprint arXiv:2402.12801*.
- [45] Zalkikar, R. and Chandra, K., 2024. Measuring Social Biases in Masked Language Models by Proxy of Prediction Quality. *arXiv preprint arXiv:2402.13954*.
- [46] Amirahmadi, A., Ohlsson, M., Etminani, K., Melander, O. and Björk, J., 2024. A Masked language model for multi-source EHR trajectories contextual representation learning. *arXiv preprint arXiv:2402.06675*.
- [47] Czinczoll, T., Hönes, C., Schall, M. and de Melo, G., 2024. NextLevelBERT: Investigating Masked Language Modeling with Higher-Level Representations for Long Documents. *arXiv preprint arXiv:2402.17682*.
- [48] Parra, I., 2024. UnMASKed: Quantifying Gender Biases in Masked Language Models through Linguistically Informed Job Market Prompts. *arXiv preprint arXiv:2401.15798*.
- [49] Toprak Kesgin, H. and Amasyali, M.F., 2024. Iterative Mask Filling: An Effective Text Augmentation Method Using Masked Language Modeling. *arXiv e-prints*, pp.arXiv-2401.
- [50] Liang, W. and Liang, Y., 2024. DrBERT: Unveiling the Potential of Masked Language Modeling Decoder in BERT pretraining. *arXiv preprint arXiv:2401.15861*.
- [51] Liu, Y., 2024. Robust Evaluation Measures for Evaluating Social Biases in Masked Language Models. *arXiv preprint arXiv:2401.11601*.
- [52] Velasco, A., Palacio, D.N., Rodriguez-Cardenas, D. and Poshyvanyk, D., 2024. Which Syntactic Capabilities Are Statistically Learned by Masked Language Models for Code?. *arXiv preprint arXiv:2401.01512*.
- [53] Velasco, A., Palacio, D.N., Rodriguez-Cardenas, D. and Poshyvanyk, D., 2024. Which Syntactic Capabilities Are Statistically Learned by Masked Language Models for Code?. *arXiv preprint arXiv:2401.01512*.
- [54] Jeong, M., Kim, M., Lee, J.Y. and Kim, N.S., 2024. Efficient parallel audio generation using group masked language modeling. *arXiv preprint arXiv:2401.01099*.
- [55] Bellamy, D.R., Kumar, B., Wang, C. and Beam, A., 2023. Labrador: Exploring the Limits of Masked Language Modeling for Laboratory Data. *arXiv preprint arXiv:2312.11502*.
- [56] Shi, B., Zhang, X., Kong, D., Wu, Y., Liu, Z., Lyu, H. and Huang, L., 2024, April. General phrase debiaser: Debiasing masked language models at a multi-token level. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 6345-6349.
- [57] Chen, T., Pertsemlidis, S., Watson, R., Kavirayuni, V.S., Hsu, A., Vure, P., Pulugurta, R., Vincoff, S., Hong, L., Wang, T. and Yudistyra, V., 2023. PepMLM: Target Sequence-Conditioned Generation of Peptide Binders via Masked Language Modeling. *arXiv preprint arXiv:2403.04187*.
- [58] Zhou, Y., Camacho-Collados, J. and Bollegala, D., 2023. A Predictive Factor Analysis of Social Biases and Task-Performance in Pretrained Masked Language Models. *arXiv preprint arXiv:2310.12936*.
- [59] Luo, X., Yuan, Y., Chen, S., Zeng, N. and Wang, Z., 2020. Position-transitional particle swarm optimization-incorporated latent factor analysis. *IEEE Transactions on*

Knowledge and Data Engineering, Vol.34, No.8, pp.3958-3970.

- [60] Chen, T., Li, S., Qiao, Y. and Luo, X., 2024. A Robust and Efficient Ensemble of Diversified Evolutionary Computing Algorithms for Accurate Robot Calibration. IEEE Transactions on Instrumentation and Measurement.
- [61] Li, Z., Li, S., Bamasag, O.O., Alhothali, A. and Luo, X., 2022. Diversified regularization enhanced training for effective manipulator calibration. IEEE Transactions on Neural Networks and Learning Systems.