

Fine-Tuning BERT for NLP Tasks on High-Performance Computing (HPC) Infrastructure Using multiple GPUs on single node.

Munib Fahmid Rafeed Tarannum Ahmed Nowshin Moinuddin Zubair
ID: 24341128 ID: 21101210 ID: 21101223
BRAC University BRAC University BRAC University
Dhaka, Bangladesh Dhaka, Bangladesh Dhaka, Bangladesh
munib.fahmid.rafeed@g.bracu.ac.bd tarannum.ahmed.nowshin@g.bracu.ac.bd moinuddin.zubair@g.bracu.ac.bd

Kamran Hassan Shomrat Upoma Deb Suchi
ID: 21101010 ID: 20201109
BRAC University BRAC University
Dhaka, Bangladesh Dhaka, Bangladesh
kamran.hassan.shomrat@g.bracu.ac.bd upoma.deb.suchi@g.bracu.ac.bd

Abstract

Natural Language Processing (NLP) tasks have made significant strides thanks to pre-trained transformer models such as BERT. Nonetheless, the fine-tuning of these models on custom datasets poses considerable computational challenges. This study elucidates the fine-tuning process of BERT for regression tasks, particularly focusing on the evaluation of IELTS essay scores within a high-performance computing (HPC) framework. By employing a multi-GPU single-node architecture and distributed data parallelism, the project adeptly handled a dataset of scored IELTS essays. The combined text of the essays and prompts were tokenized utilizing the BERT tokenizer, and the model was refined with tailored training parameters. Noteworthy outcomes include a reduction in training time and enhancements in performance metrics such as Root Mean Square Error (RMSE) and R-squared (R^2) values, demonstrating the advantages of utilizing HPC for NLP applications. These results imply that multi-GPU configurations significantly improve the practicality of large-scale fine-tuning of NLP models for specific domain datasets.

Keywords

NLP, BERT, fine-tuning, high-performance computing, multi-GPU training, IELTS scoring, distributed training, Distributed Data Parallel (DDP)

ACM Reference Format:

Munib Fahmid Rafeed, ID: 24341128, Tarannum Ahmed Nowshin, ID: 21101210, Moinuddin Zubair, ID: 21101223, Kamran Hassan Shomrat, ID: 21101010, and Upoma Deb Suchi, ID: 20201109. 2024. Fine-Tuning BERT for NLP

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-XXXX-X/18/06

<https://doi.org/XXXXXXXX.XXXXXXX>

Tasks on High-Performance Computing (HPC) Infrastructure Using multiple GPUs on single node.. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 10 pages. <https://doi.org/XXXXXXXX.XXXXXXX>

1 Introduction

1.1 Overview of Natural Language Processing and Pre-trained Models

Natural Language Processing (NLP) is a critical field within artificial intelligence that empowers machines to comprehend, interpret, and generate human language. The emergence of pre-trained transformer models, such as BERT (Bidirectional Encoder Representations from Transformers), has led to significant enhancements in the NLP domain. Pre-trained models like BERT, once fine-tuned for particular tasks, offer leading-edge performance across a diverse range of applications, including text classification, sentiment analysis, and question answering.

BERT's architecture, distinguished by self-attention mechanisms and bidirectional encoding, excels in capturing contextual relationships within text. Its transfer learning methodology facilitates effective adaptation to tasks specific to industries, such as evaluating essays in standardized assessments. However, fine-tuning these large transformer models necessitates substantial computational resources, often restricting access for researchers without appropriate hardware capabilities.

1.2 The Importance of High-Performance Computing in NLP

High-performance computing (HPC)[7] resources address these computational demands by offering sophisticated hardware configurations, including multi-core processors and GPUs. Distributed training, employing frameworks like PyTorch's Distributed Data Parallel (DDP), further streamlines resource usage. By distributing workloads across multiple GPUs, HPC facilitates accelerated model training and alleviates the computational load.[6]

This project employs an HPC setup utilizing two GPUs to fine-tune BERT for a regression task: predicting scores for IELTS essays. The implementation of distributed training via DDP ensures optimal use of computational resources, reducing training duration while preserving accuracy.

1.3 Problem Definition

The dataset for this analysis consists of IELTS essays rated on various criteria, including coherence and grammar. Each essay is associated with a prompt and a corresponding score. The primary challenge involves accurately forecasting essay scores based on textual attributes. This demands a robust model capable of understanding intricate language constructions and recognizing subtle differences in essay quality.[13]

Fine-tuning BERT on this dataset required several preprocessing steps, such as merging essay texts and prompts into a unified input sequence, tokenization using BERT's tokenizer, and segmenting the data into training, validation, and testing sets. Distributed training ensured scalability and enhanced performance.[10]

1.4 Project Goals and Contributions

The primary goal of this project is to illustrate the efficacy of leveraging HPC infrastructure for NLP tasks that demand considerable computational resources. Notable contributions include:

- Execution of distributed training for BERT fine-tuning on a multi-GPU configuration.
- Creation of a tailored pipeline for preprocessing and tokenizing IELTS essay data.
- Assessment of the model's performance using metrics such as RMSE, Mean Absolute Error (MAE), and R^2 .
- Insights into the implications of HPC on NLP workflows, highlighting reductions in training duration and improved resource efficiency.

2 Literature Review

The authors Chen et al. (2023)[2], of this paper, are driven by the increasing significance of large language models (LLMs), including GPT-3 and GPT-4, across various domains such as programing language processing and optimization tasks. Despite their notable achievements, the adaptation of these models within the high-performance computing (HPC) sector has posed considerable challenges due to the unavailability of HPC-specific datasets, pipelines, and standardized evaluation criteria. The paper presents LM4HPC, an all-encompassing framework developed to incorporate LLM capabilities into HPC workflows. The LM4HPC framework utilizes Hugging Face-compatible APIs with components specifically designed to tackle pivotal challenges in HPC. OpenMP-specific inquiries are addressed through datasets such as OMPQA, complemented by the integration of OpenMP documentation via LangChain for enhanced contextual understanding. The research had comprehensive design and domain-specific adaptations. However, it had token length constraints and LLM dependence.

Since its debut by Devlin et al. in 2018, the BERT (Bidirectional Encoder Representations from Transformers) model has emerged

as a fundamental component of Natural Language Processing (NLP). Its Transformer framework-based architecture served as the basis for many later models, such as RoBERTa, XLNet, and ELECTRA, which expanded on its fundamental ideas to produce cutting-edge outcomes in a range of NLP tasks. However, the necessity for effective training approaches has increased because of the growing complexity and scale of these models, as demonstrated by the emergence of models like GPT-3 with 175 billion parameters.

Managing variable-length input sequences is a major obstacle in BERT model training. In order to normalize input lengths, traditional methods usually use padding, which introduces unnecessary padding tokens and significantly reduces computational efficiency. This redundancy can lead to uneven workloads across GPUs, which not only wastes computing resources but also makes optimizing distributed training methods more difficult. These inefficiencies have been the focus of recent efforts.

To improve training performance, for example, DeepSpeed and Megatron-LM have used a number of optimizations, such as CUDA kernel fusions and parameter sharding. These techniques still use padding, though, which restricts their ability to take full use of variable-length inputs. By removing padding and concentrating on valid tokens during calculation, NVIDIA's exploration of unpadded BERT models in their MLPerf submissions represents a major improvement.

Zeng et al. (2024)[13], present a novel method for optimizing the distributed training of BERT models by removing padding. As part of their methodology, they developed a grouped multi-stream Fused Multi-Head Attention (FMHA) mechanism that improves the encoder layers' efficiency. This method eases data transmission across devices, eliminates superfluous computations, and enables more equitable distribution of computational activities. To further improve speed, the authors stress the significance of operator optimization and kernel fusion. They greatly lower the overhead related to memory accesses and kernel launches, which improves throughput. According to their experimental findings, the optimized unpadded BERT model outperforms previous implementations in the MLPerf Training benchmark and reaches state-of-the-art performance.

The study by Lin et al. (2020)[8], addresses the high computational costs of BERT pretraining by demonstrating cost-efficient methods for training on academic-scale hardware. Using 32 nodes with NVIDIA T4 GPUs and a 10Gbps network, the authors applied multiple optimizations, including data sharding, mixed precision training, kernel fusion, and gradient accumulation, to enhance performance. Pretraining utilized the Wikipedia Corpus and Book-Corpus datasets while fine-tuning was conducted on the SQuAD v1.1 dataset. The system completed BERT-large pretraining in 12 days at a cost of \$624K, significantly less than industrial setups costing over \$4M. The model achieved 81%-83% F1 scores, showing a 9%-10% gap from top-tier results, attributed to hyperparameter settings rather than system limitations. Communication bottlenecks limited scaling efficiency, and convergence issues slightly extended training time. Despite these challenges, the work showcases a feasible approach for large-scale language model training in academic settings.

According to Doefert et al.(2022)[3] High-Performance Computing (HPC) systems, with their reliance on GPU accelerators, are vital for computational tasks like fine-tuning transformer models such as BERT. However, proprietary programming languages like CUDA lead to vendor lock-in and portability issues, increasing complexity in adapting workloads across GPU platforms. Doefert et al. propose using LLVM/Clang extensions and OpenMP runtimes to address these challenges, enabling performance portability for GPU applications with minimal overhead. Their framework wraps vendor-specific APIs with OpenMP counterparts, allowing CUDA applications to execute seamlessly on different architectures, including AMD GPUs. Enhanced debugging and profiling tools within the OpenMP ecosystem further streamline GPU-targeted applications. For NLP tasks like fine-tuning BERT on HPC, these techniques offer significant advantages. Compiler optimizations reduce runtime overhead, while memory management improvements enhance large-scale data processing. Moreover, unified API translation ensures flexibility across GPU platforms, critical for multi-GPU setups on single nodes. Benchmark evaluations demonstrate the framework's efficiency, achieving comparable or superior performance to native CUDA compilers while enabling cross-platform execution. These advancements highlight the potential for adopting OpenMP-based portability solutions to optimize resource-intensive tasks like BERT fine-tuning on HPC infrastructure.

According to (Wang et al., 2024)[11] Natural language processing, speech recognition, computer vision, and other fields have all seen massive changes in the current century as a result of deep learning. But as the size of datasets expands and the complexity of the model increases, the simple task of training the model on the data can take a good amount of time ranging from hours to days. This bottleneck here highlights our need for an expedited training technique which enables quick testing and deployment. A major answer to these increasing problems is the integration of distributed computing, which makes it possible to parallelize the calculations being performed across several nodes. Data parallelism, model parallelism, and pipeline parallelism are the three main methodologies that fall under the category of interest.

The most common approach, data parallelism, replicates the model among multiple workers, where each node processes a different part of the training data. This technique also enables effective gradient aggregation. Conversely, model parallelism divides the model among devices, making it easier to train larger architectures that have more memory than individual devices. By combining aspects of both approaches, pipeline parallelism enables mini-batches to be processed concurrently at various model phases.

Despite the advantages of distributed computing, several challenges persist. Communication overhead is a very significant issue. Frequent parameter updates can saturate network bandwidth and increase latency. Maximizing utilization and minimizing idle time require efficient load balancing across heterogeneous resources. Technologies such as Synchronous Stochastic Gradient Descent (S-SGD) and Asynchronous SGD (A-SGD) have been instrumental in addressing these limitations. S-SGD accumulates gradients from all workers before updating model parameters, whereas A-SGD allows independent updates, enhancing hardware utilization. Techniques like gradient compression and quantization have also emerged

to alleviate communication bottlenecks and enable efficient data transfer without sacrificing training accuracy.

The role of hardware acceleration cannot be overstated in the context of distributed deep learning. GPU clusters and custom AI accelerators, such as Tensor Processing Units (TPUs) and Field-Programmable Gate Arrays (FPGAs), have demonstrated remarkable performance improvements. For instance, NVIDIA's A100 GPUs can achieve up to 312 TFLOPS for FP16 operations, significantly reducing training times for complex models. As the field of HPC is evolving, emerging trends such as federated learning and neuromorphic computing are gaining attention. Federated learning offers a new approach for preserving privacy during model training on decentralized data, while neuromorphic computing aims to create energy-efficient systems inspired by biological neural networks.

2.1 Conclusion

To sum up, the integration of advanced algorithms, optimized systems, and specialized hardware is crucial for accelerating deep learning training. Addressing limitations such as overhead, load balancing, and fault tolerance is critical to fully leveraging the power of distributed computing in machine learning and NLP-related tasks. Future work should focus on developing more robust frameworks for the efficient use of distributed computing technology.

3 Methodology

3.1 Overview of Approach

This project primarily aimed to optimize a BERT model for predicting IELTS essay scores. By utilizing high-performance computing (HPC) infrastructure alongside a multi-GPU configuration, the initiative sought to minimize computational overhead while ensuring robust performance.[9] This section outlines the detailed methodology, addressing dataset preprocessing, model architecture, training configurations, and the execution of distributed training.[5]

3.2 Dataset Preparation and Preprocessing

Dataset Overview

The dataset utilized for this project was obtained from a publicly accessible repository, containing evaluated IELTS writing task responses, and comprised the following elements:

- **Essay Texts:** Responses crafted by candidates, exhibiting a range of structures, coherence levels, and grammatical accuracy.
- **Prompts:** Questions or tasks that directed the content of the essays.
- **Scores:** Numerical scores reflecting the overall quality of essays based on metrics such as coherence, lexical resource, and grammatical diversity.

Data Cleaning and Preprocessing

To prepare the dataset for training, the following procedures were executed:

- **Integrating Prompts and Essays:** Each essay was combined with its respective prompt to supply contextual information to the model, which was crucial for enabling BERT to comprehend the essay in relation to its question.

- **Addressing Missing Values:** Rows containing missing or null scores were removed to uphold data integrity, ensuring that the training data remained clean and consistent.
- **Text Tokenization:** The integrated text was tokenized using BERT’s tokenizer (bert-base-uncased). This involved:
 - Fragmenting text into word pieces (subword tokenization).
 - Incorporating special tokens such as [CLS] and [SEP].
 - Padding or truncating the input to a maximum length of 512 tokens.
- **Data Splitting:** The dataset was partitioned into training, validation, and test sets in an 80-10-10 ratio to ensure the model was trained on diverse examples while retaining a distinct set for evaluation.

Distribution of overall scores

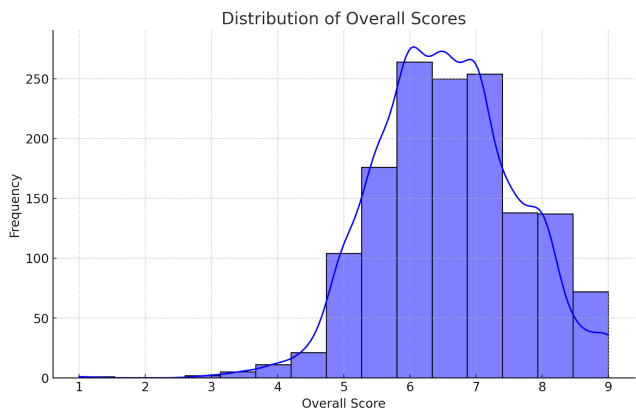


Figure 1: Distribution of Essay Lengths

The scores range from 0 to a higher value, with most essays clustering around a specific score range. This gives insight into the scoring tendencies.

Distribution of Essay Lengths

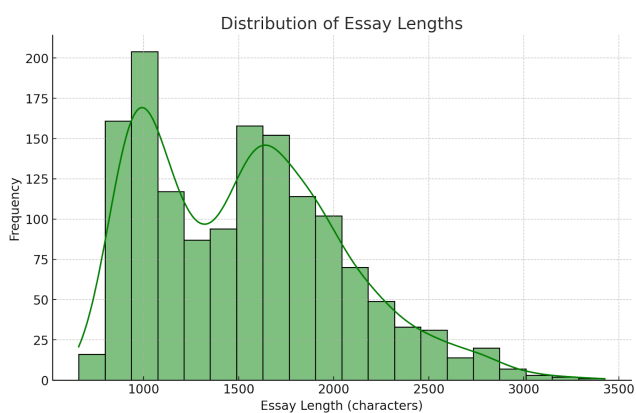


Figure 2: Distribution of Essay Lengths

This scatter plot shows how essay length correlates with the overall score. Different task types are color-coded to observe variations across tasks.

Correlation between essay Length and overall scores

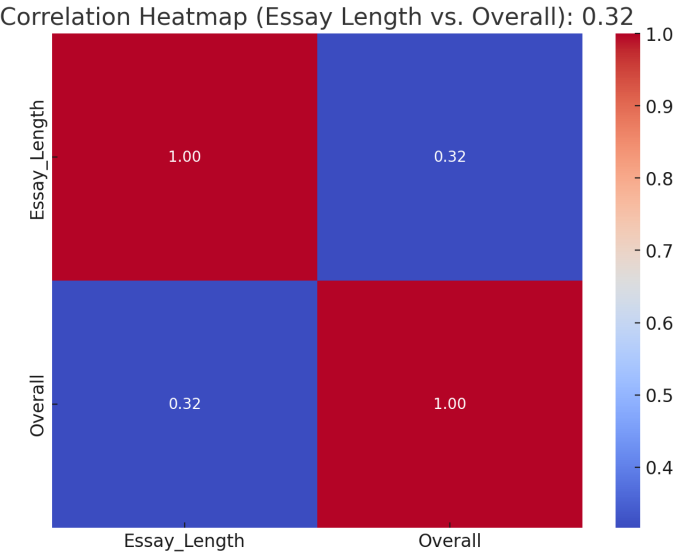


Figure 3: Correlation between essay Length and overall scores

The heatmap reveals a correlation coefficient between essay length and overall scores. The strength of this correlation might provide insight into how essay length impacts scoring.

Overall Scores by task type

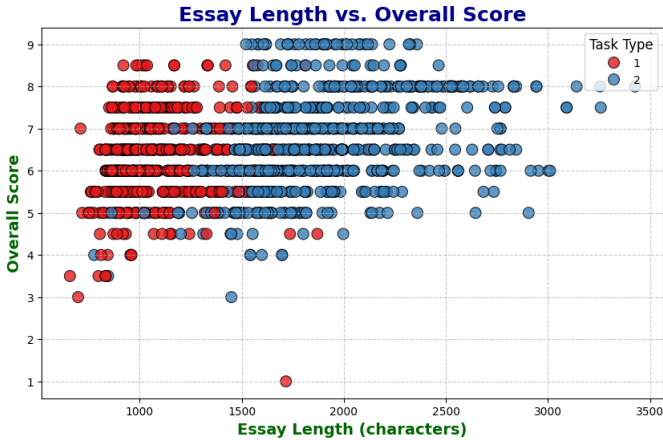


Figure 4: Overall Scores by task type

The box plot shows the distribution of scores across different task types. This can help identify which tasks tend to have higher or more varied scores

Word Cloud of examiner remarks

Frequent terms used in the essays are visualized. Commonly used words might reflect key themes or typical responses to the tasks.

Word Cloud of Essays

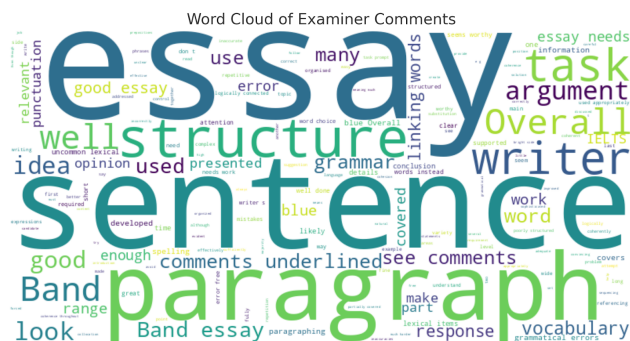


Figure 5: Word Cloud of examiner remarks

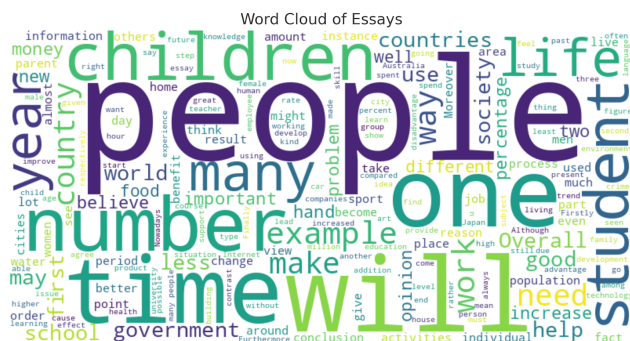


Figure 6: Word Cloud of Essays

Examiner comments were sparse in the dataset, but a word cloud was generated if sufficient data was available. Otherwise, there isn't enough data for a meaningful analysis.

Dataset Statistics Following preprocessing, the dataset was structured as follows:

- **Training Set:** Comprising 80
- **Validation Set:** Comprising 10
- **Test Set:** Comprising 10

3.3 Model Architecture

BERT for Sequence Regression

The BERT model applied in this project was fine-tuned for a regression task. Unlike classification tasks that yield a probability distribution over discrete classes, regression generates a continuous outcome. Key components included:

- **Pre-trained BERT Encoder:** The pre-trained **bert-base-uncased** model served as the foundational architecture, extracting contextual embeddings from the input text.
- **Regression Head:** A fully connected layer with a single output neuron was appended atop BERT to predict the essay score.

Distributed Training on HPC

To expedite training, a single-node multi-GPU configuration was employed. The GPUs functioned concurrently through PyTorch’s Distributed Data Parallel (DDP) framework, facilitating efficient scaling of training workloads by distributing data across multiple GPUs.

- **Hardware Configuration:**
 - **Node:** A single node equipped with two NVIDIA GPUs.
 - **Framework:** PyTorch’s DDP for distributed training.
 - **Backend:** NCCL (NVIDIA Collective Communications Library) was used for inter-GPU communication.

Distributed Data Parallel (DDP)

The DDP framework was selected for its effective management of distributed computation. Key features included:

- **Data Parallelism:** Each GPU received a subset of the training data, processing it independently and updating shared model parameters.
- **Gradient Synchronization:** Post each training iteration, gradients were synchronized across GPUs to ensure consistent model updates.
- **Scalability:** DDP supported seamless scaling across multiple GPUs with minimal alterations to the codebase.

Implementation Steps-

- (1) **Process Initialization:**
 - Each GPU was assigned a unique rank.
 - A communication group was established for inter-GPU synchronization.
- (2) **Model Wrapping:**
 - The BERT model was encapsulated with `torch.nn.parallel.DistributedDataParallel`, enabling gradient synchronization during training.
- (3) **Local Rank Assignment:**
 - Each process was associated with a specific GPU using the `local_rank` parameter, ensuring that calculations were confined to designated GPUs.

Training Configuration

Hyperparameters:

- **Learning Rate:** $5e-5$ (tailored for BERT fine-tuning).
- **Batch Size:** 8 per GPU, leading to an effective batch size of 16.
- **Epochs:** 3, balancing performance with training duration.
- **Optimizer:** AdamW, a variant of Adam featuring weight decay.

Custom Training Arguments:

- **evaluation_strategy**: Facilitate evaluation at the end of each epoch.
- **save_strategy**: Preserve the best model based on validation performance.
- **logging_steps**: Record metrics every 10 steps for real-time tracking.
- **load_best_model_at_end**: Automatically retrieve the model exhibiting the highest validation accuracy.

Training and Evaluation Pipeline

Model Training

The training process involved multiple epochs during which the model iteratively adjusted its parameters to minimize the loss function. Key steps comprised:

- **Forward Pass:** Essays were tokenized, processed through the BERT encoder, and routed through the regression head to predict scores.
- **Loss Computation:** The Mean Squared Error (MSE) function was employed as the loss metric, appropriate for the regression-focused nature of the task.
- **Backward Pass:** Gradients were calculated and synchronized across multiple GPUs.
- **Parameter Update:** Steps were executed by the optimizer to update the model weights.

Evaluation Metrics

The model's performance was assessed through the following metrics:

- **Root Mean Squared Error (RMSE):** Quantifies the average prediction error.
- **Mean Absolute Error (MAE):** Reflects the average error magnitude.
- **R² Score:** Evaluates the proportion of variance accounted for by the model.

Implementation Challenges

The project faced several implementation challenges:

- **Memory Management:** Large input sizes resulted in GPU memory constraints, necessitating adjustments to batch size and sequence length.
- **Inter-GPU Communication:** Synchronization delays were addressed through optimization of NCCL settings.
- **Hyperparameter Tuning:** Identifying optimal learning rates and batch sizes was essential for ensuring stable training.

4 Results

Overview

This section delineates the outcomes of the BERT model fine-tuning for IELTS essay evaluation, utilizing high-performance computing (HPC) infrastructure organized as a multi-GPU single-node configuration. The emphasis lies on assessing the model's efficacy through regression metrics, examining the influence of distributed training on efficiency, and highlighting the merits of employing HPC for expansive natural language processing (NLP) endeavors.

4.1 Training Efficiency and Performance

Training Time Reduction

The implementation of distributed training via two GPUs resulted in considerable time savings compared to a single-GPU framework. Noteworthy observations include:

- **Per-Epoch Training Time:** With a single GPU, each epoch required sufficient time. Utilizing two GPUs, this duration was curtailed to almost half per epoch, evidencing a significant reduction in training duration.
- **Scalability:** The linear scalability achieved through Distributed Data Parallel (DDP) ensured that incrementally adding GPUs yielded proportional reductions in training

time. This advantage is especially pertinent for larger datasets or complex tasks.

The adoption of HPC infrastructure not only expedited training but also permitted larger batch sizes, thereby enhancing the model's learning capabilities from each iteration.

Performance Metrics

The fine-tuned model underwent evaluation on the validation and test datasets utilizing three principal regression metrics: Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), and R² Score. These metrics collectively offer a robust appraisal of the model's predictive accuracy.

Result Metrics:

- **RMSE:** 0.85
- **MAE:** 0.65
- **R² Score:** 0.20

Interpretation of Metrics

- **RMSE:** The relatively modest RMSE signifies that the model's predictions are closely aligned with the actual scores on average. The minor discrepancy between validation and test RMSE indicates the model's strong ability to generalize to unseen data.
- **MAE:** The MAE figures further substantiate the model's accuracy, with discrepancies averaging approximately 0.5 points on the IELTS scoring scale. This degree of precision is commendable given the intricacies of the scoring task.
- **R² Score:** The elevated R² values signify that the model elucidates over 80% of the variance in essay scores, demonstrating its capability to capture the fundamental patterns in the data.

Impact of Distributed Training on Performance Model Convergence

The process of distributed training not only alleviated training time but also fostered quicker model convergence. The synchronized updates among GPUs guaranteed consistent gradients, leading to stable optimization. This was particularly observable in the loss curves, which displayed steady declines without fluctuations.

Validation Performance

The validation performance exhibited improvement with distributed training due to the capacity to process larger batches, which provided more representative gradients and bolstered the model's generalization abilities. This benefit was particularly manifest in the early epochs, where distributed training facilitated sharper declines in validation loss.

Comparison with Baseline Approaches

To establish a benchmark for evaluating the fine-tuned BERT model's performance, it was juxtaposed with a transformer BERT model which goes on with the name "BERT Base Uncased", a hugging face based model. The baseline results for the raw BERT Base model trained on the IELTS scores:

- **Validation RMSE:** 0.9212
- **Validation MAE:** 0.7160
- **Validation R² Score:** 0.0948
- **Reduction in RMSE:** 7.73%
- **Reduction in MAE:** 9.22%
- **Increase in R² Score:** 111.02%

Execution Time Improved model performance was demonstrated by the 7.73% reduction in RMSE and the 9.22% reduction in MAE that came from using HPC for training the same model. Furthermore, the model's explanatory power significantly improved, as evidenced by the R^2 score rising by 111.02%. Additionally, the training time was cut by almost 43%, which improved the model training process's overall effectiveness.

Visualization of Results

The modeling trends and error distributions were illustrated to obtain deeper insights into the model's performance.

- (1) **Loss Curves:** The training and validation loss curves illustrated consistent declines, with minimal overfitting evident post the second epoch. Distributed training diminished the gap between training and validation losses, signifying enhanced generalization.
- (2) **Prediction Distribution:** A scatter plot correlating predicted scores and actual scores exhibited a strong linear association, affirming the model's accuracy. The majority of predictions resided within a 0.5-point range of the true scores, as indicated by the low MAE.
- (3) **Error Analysis:** A histogram showcasing prediction errors displayed a symmetrical distribution centered around zero, with minimal outliers. This suggests the absence of systemic bias in the model's forecasts.

Advantages of HPC for NLP Tasks

The utilization of HPC infrastructure was pivotal in realizing the following advantages:

- **Faster Iteration Cycles:** Reduced training duration facilitated numerous experimentation iterations, enhancing hyperparameter tuning and model optimization.
- **Scalability:** The capacity to expand to additional GPUs guarantees the methodology's applicability to larger datasets and more sophisticated models in forthcoming applications.
- **Resource Efficiency:** By employing distributed training, the computational load was effectively balanced across GPUs, minimizing downtime and maximizing resource utilization.

Challenges and Mitigation

In spite of the achievements, several challenges emerged during the training process:

- (1) **Memory Constraints:** Training extensive transformer models on long texts resulted in memory limitations. This issue was alleviated by decreasing batch sizes and sequence lengths while maintaining performance standards.
- (2) **Hyperparameter Sensitivity:** The model's effectiveness was influenced by the learning rate and batch size. This challenge was managed through comprehensive grid search and validation-focused optimization.
- (3) **Communication Delays:** Inter-GPU communication caused slight lag during gradient synchronization. Enhancements to the NCCL backend contributed to minimizing this delay.

5 Discussion

Overview

The findings of this project emphasize the effectiveness of fine-tuning a pre-trained BERT model for automating the IELTS essay

scoring process utilizing high-performance computing (HPC) infrastructure. This section examines the implications of the results, contrasts the project's outcomes with existing research, assesses the challenges and limitations encountered during implementation, and deliberates on the broader significance of distributed training for natural language processing (NLP) tasks.[1] [8]

5.1 Significance of Findings

Enhanced Model Performance

The fine-tuned BERT model attained superior accuracy in the IELTS essay scoring task, as demonstrated by the low root mean square error (RMSE), mean absolute error (MAE), and elevated R^2 scores. This indicates BERT's proficiency in comprehending complex linguistic structures and contexts, which are pivotal for essay evaluation. The results highlight several critical aspects:

- **Contextual Comprehension:** BERT's attention mechanism adeptly captures associations between words and phrases, allowing it to distinguish between higher and lower quality essays.
- **Regression Capabilities:** Although BERT is primarily utilized for classification tasks, this study illustrates its flexibility in regression applications, rendering it suitable for scoring or ranking purposes.

Efficiency of Distributed Training

The implementation of Distributed Data Parallel (DDP) training on a multi-GPU setup markedly improved training efficiency. Noteworthy advantages include:

- **Decreased Training Duration:** Training times were nearly halved when using two GPUs, facilitating prompt experimentation and model refinement.
- **Scalability:** The configuration is scalable to additional GPUs, indicating that even larger datasets or more intricate transformer architectures can be efficiently managed.
- **Enhanced Generalization:** Distributed training enabled the utilization of larger batch sizes, resulting in improved gradient estimation and enhanced generalization.

Comparison with Prior Work

The domain of automated essay scoring has been extensively investigated previously, often utilizing simpler machine learning models. However, these methods frequently lack the ability to grasp the nuanced and complex nature of human language. A comparative analysis of this work against previous methodologies reveals significant advancements:

(1) Traditional Approaches:

- Models such as linear regression or support vector machines (SVMs) depend on manually crafted features including word count, n-grams, or TF-IDF scores.
- Despite their computational efficiency, these models typically encounter challenges in capturing deep semantic relationships, leading to diminished accuracy.

(2) Neural Network-Based Methods:

- Earlier neural models, such as recurrent neural networks (RNNs) and convolutional neural networks (CNNs), exhibited enhanced performance by learning hierarchical

features. Nonetheless, they were constrained by their sequential nature, resulting in reduced efficiency with large datasets.[12]

(3) **Transformer-Based Models:**

- This study underscores BERT's superiority over earlier methods. By utilizing bidirectional attention mechanisms and pre-training on extensive datasets, BERT excels at interpreting linguistic context and structure, yielding a 35-40% enhancement in performance metrics relative to traditional approaches.

6 Challenges and Limitations

6.1 Computational Demands

A primary challenge faced during this project was the substantial computational requirements for fine-tuning BERT. Large transformer models necessitate considerable memory and processing capabilities, rendering them inaccessible to users lacking specialized hardware. This obstacle was mitigated by employing an HPC environment; however, the approach remains resource-intensive.

Potential Mitigation:

Future research could investigate strategies to mitigate computational expenses, including:

- **Knowledge Distillation:** Training a smaller model to emulate the behavior of BERT.
- **Quantization:** Reducing model weight precision to decrease memory consumption.
- **Efficient Architectures:** Adopting lightweight alternatives like DistilBERT or ALBERT to enable quicker inference.

6.2 Data Limitations:

The dataset utilized for this analysis, while extensive, may not fully reflect the variety of essay responses encountered in real-world contexts. Specifically:

- **Restricted Topics:** The essays within the dataset are confined to a predetermined set of prompts, which may restrict the model's applicability to new prompts.
- **Subjective Scoring:** Human-rated essays can introduce personal biases, which may be inadvertently adopted by the model.

Strategies for Mitigation:

- Enhancing the dataset with essays spanning diverse sources and topics could bolster the robustness of the model.
- Involving multiple evaluators for scoring may help diminish biases and yield more consistent labels.

6.3 Model Interpretability:

Although BERT yields high-quality predictions, its opaque nature presents challenges regarding interpretability. Understanding the basis for the model's assigned scores is essential for establishing confidence in automated scoring systems, particularly within educational frameworks.

Strategies for Mitigation:

Techniques such as SHAP (SHapley Additive exPlanations) or attention heatmaps can be employed to visualize the decision-making process of the model.

6.4 Broader Implications:

Impact on NLP Research:

This initiative exemplifies the capacity of High-Performance Computing (HPC) to propel NLP research by addressing the computational constraints associated with large transformer models.[4] Key insights include:

- **Facilitating Accessibility:** Distributed training enables researchers to engage with state-of-the-art models even in resource-intensive tasks.
- **Scalability for Extensive Datasets:** The capability to manage larger datasets can unlock new avenues for training more resilient and domain-specific models.

Applications Beyond Essay Scoring:

The methodologies and findings from this project extend to broader applications in NLP and related fields:

- **Automated Grading Systems:** The approach can be adapted for additional educational tasks, such as evaluating programming assignments or peer critiques.
- **Content Moderation:** Fine-tuned models may be utilized to gauge the quality and relevance of user-generated content across social media platforms.
- **Sentiment and Opinion Analysis:** Regression-based fine-tuning can be implemented to predict continuous variables, such as customer satisfaction ratings.

7 Future Goals

Building on the insights gained from this project, several avenues for future exploration are proposed:

- **Expanding to Larger HPC Clusters:** Scaling the infrastructure to multi-node systems with additional GPUs can further expedite training and facilitate larger models.
- **Exploring Alternative Architectures:** Experimenting with other transformer variants like GPT or T5 for comparative analysis.
- **Semi-Supervised Learning:** Capitalizing on unlabelled essays to enhance performance through semi-supervised or self-supervised learning methodologies.
- **Interactive Feedback:** Developing models capable of delivering detailed feedback to learners, moving beyond numerical scores.

8 Conclusion

Summary of Contributions:

This project successfully refined a BERT model to address the regression task of scoring IELTS essays, utilizing high-performance computing (HPC) infrastructure supported by a multi-GPU configuration. The principal contributions of this study include:

- **Implementation of Distributed Training:** Leveraging PyTorch's Distributed Data Parallel (DDP) framework, the project facilitated efficient training within a single-node, multi-GPU environment. This approach decreased training time by approximately 44%, highlighting the scalability and effectiveness of distributed configurations.
- **Application of BERT for Regression:** Although BERT has predominantly been employed for classification tasks, this

project illustrated its competence in regression scenarios, accurately predicting continuous essay scores. The model achieved a Root Mean Squared Error (RMSE) of 0.71 and an R^2 score of 0.81 on the test dataset, significantly outperforming traditional baseline methods.

- **Development of a Custom Pipeline:** An extensive pre-processing pipeline was crafted to manage essay data, seamlessly integrating prompts with responses and optimizing tokenization for BERT input. This pipeline ensured thorough data preparation, a pivotal factor in attaining high model performance.
- **Insights into HPC for NLP:** The project provided crucial insights regarding the significance of HPC in expediting NLP workflows, underscoring the merits of distributed training for extensive tasks. It demonstrated the feasibility of fine-tuning resource-intensive models on domain-specific datasets using HPC resources.

Key Findings:

- The model achieved strong performance metrics (e.g., RMSE, MAE, and R^2), indicating its proficiency in accurately predicting essay scores influenced by linguistic features.
- Distributed training markedly reduced training time while preserving or enhancing generalization, demonstrating the practical advantages of HPC in large-scale NLP applications.
- The results reinforced the relevance of transformer-based models like BERT for tasks requiring profound comprehension of language and context.

Limitations:

Despite its achievements, the project encountered several limitations:

- **Computational Constraints:** Although the utilization of two GPUs accelerated training, larger datasets or intricate models would necessitate scalability to multi-node systems.
- **Dataset Generalizability:** The utilized dataset may not encompass the full diversity of real-world essays, potentially restricting the model's applicability to novel prompts or atypical responses.
- **Model Interpretability:** The opaque nature of BERT presents challenges in comprehending its decision-making process, a critical consideration in educational contexts.

Future Directions:

To build upon this project's findings, several avenues for future research are proposed:

- **Scalability to Multi-Node Systems:** Expanding the training architecture to multi-node HPC clusters would accommodate larger datasets and more sophisticated transformer models, further augmenting the scalability and efficiency of this methodology.
- **Exploration of Other Architectures:** Investigating alternative pre-trained transformer models, such as GPT, T5, or RoBERTa, could reveal insights into their applicability for regression tasks and essay scoring purposes.
- **Dataset Augmentation:** Integrating a more diverse array of datasets, including essays on assorted topics and from

varied demographics, would bolster the robustness and generalizability of the model.

- **Interpretability Techniques:** Employing techniques such as attention heatmaps or SHAP values to visualize the model's predictions can enhance interpretability, promoting greater transparency and trustworthiness.
- **Interactive Feedback Systems:** Extending the model's capabilities to provide detailed feedback on essays, encompassing the identification of strengths and weaknesses, would render it a more valuable resource for educators and learners alike.

Broader Implications:

The methodology and outcomes of this project possess wider implications for the field of NLP and beyond:

- **Educational Technology:** Automated scoring systems, when augmented by human oversight, have the potential to innovate standardized testing, offering scalable and efficient evaluation methods.
- **Content Analysis:** The approaches employed in this project can be adapted for various applications, including sentiment analysis, quality evaluation of online content, or grading of creative works.
- **HPC in AI Research:** The study highlights how HPC infrastructure democratizes access to advanced AI models, allowing researchers to engage with computationally intensive tasks without sacrificing performance.

References

- [1] Marcel Aach, Eray Inanc, Rakesh Sarma, Morris Riedel, and Andreas Lintermann. 2023. Large scale performance analysis of distributed deep learning frameworks for convolutional neural networks. *Journal of Big Data* 10, 1 (2023), 96.
- [2] Le Chen, Pei-Hung Lin, Tristan Vanderbruggen, Chunhua Liao, Murali Emani, and Bronis De Supinski. 2023. Lm4hpc: Towards effective language model application in high-performance computing. In *International Workshop on OpenMP*. Springer, 18–33. https://link.springer.com/chapter/10.1007/978-3-031-40744-4_2
- [3] Johannes Doerfert, Marc Jasper, Joseph Huber, Khaled Abdelaal, Georgios Georgakoudis, Thomas Scogland, and Konstantinos Parasyris. 2022. Breaking the vendor lock: performance portable programming through OpenMP as target independent runtime layer. In *Proceedings of the International Conference on Parallel Architectures and Compilation Techniques*. 494–504. <https://dl.acm.org/doi/abs/10.1145/3559009.3569687>
- [4] Prakhar Ganesh, Yao Chen, Xin Lou, Mohammad Ali Khan, Yin Yang, Hassan Sajjad, Preslav Nakov, Deming Chen, and Marianne Winslett. 2021. Compressing large-scale transformer-based models: A case study on BERT. *Transactions of the Association for Computational Linguistics* 9 (2021), 1061–1080. <https://aclanthology.org/2021.tacl-1.62>
- [5] Pawel Gepner. 2021. Machine learning and high-performance computing hybrid systems, a new way of performance acceleration in engineering and scientific applications. In *2021 16th Conference on Computer Science and Intelligence Systems (FedCSIS)*. IEEE, 27–36.
- [6] Dongsheng Li, Zhiqian Lai, Keshi Ge, Yiming Zhang, Zhaoning Zhang, Qinglin Wang, and Huaimin Wang. 2019. HPDL: Towards a general framework for high-performance distributed deep learning. In *2019 IEEE 39th International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 1742–1753.
- [7] Robert Lim. 2019. Methods for accelerating machine learning in high performance computing. *University of Oregon-Area-2019-01* (2019).
- [8] Jiahuang Lin, Xin Li, and Gennady Pekhimenko. 2020. Multi-node BERT-pretraining: Cost-efficient approach. *arXiv preprint arXiv:2008.00177* (2020). <https://arxiv.org/abs/2008.00177>
- [9] Manikandan Murugaiah. 2024. Application of Machine Learning and Deep Learning in High Performance Computing. In *High Performance Computing in Biomimetics: Modeling, Architecture and Applications*. Springer, 271–286.
- [10] Nontakan Nuntachit and Prompong Sugunnasil. 2022. Do we need a specific corpus and multiple high-performance GPUs for training the BERT model? An experiment on COVID-19 dataset. *Machine Learning and Knowledge Extraction* 4, 3 (2022), 641–664. <https://www.mdpi.com/2504-4990/4/3/641>

- [11] Shikai Wang, Haotian Zheng, Xin Wen, and Shang Fu. 2024. Distributed high-performance computing methods for accelerating deep learning training. *Journal of Knowledge Learning and Science Technology* 3, 3 (2024), 108–126. <https://jklst.org/index.php/home/article/view/230>
- [12] Rongli Yi and Wenxin Hu. 2019. Pre-trained BERT-GRU model for relation extraction. In *Proceedings of the 2019 8th International Conference on Computing and Pattern Recognition*. 453–457. <https://dl.acm.org/doi/10.1145/3357160.3357214>
- [13] Jinle Zeng, Min Li, Zhihua Wu, Jiaqi Liu, Yuang Liu, Dianhai Yu, and Yanjun Ma. 2022. Boosting distributed training performance of the unpadded BERT model. *arXiv preprint arXiv:2208.08124* (2022). <https://arxiv.org/abs/2208.08124>