



# Parallelizing Precision: Leveraging Mixed-Precision and Knowledge Distillation for Efficient BERT Training

Exploring Computational Strategies in NLP Model Optimization  
with Advanced GPU Technologies

Hadi Saghir

Computer and Information Science

Bachelor Thesis

15 hp

2024

Supervisor:

## Introduction

In the dynamic realm of deep learning, the advancement of modern GPUs has been pivotal in enabling techniques such as mixed-precision computing, marking a significant breakthrough in the field. This advancement offers a critical balance between computational efficiency and model accuracy. The technique of mixed precision, blending lower-precision formats with high precision where crucial, is reshaping the training of deep learning models. Crucially, it reduces memory demands and hastens the training process, becoming vital for machine learning experts. This study embarks on an exploration of mixed-precision computing intertwined with parallelization for training complex models, with a special focus on NLP. It particularly examines the role of Knowledge Distillation (KD) in this context, where a simpler "student" model learns from a more complex "teacher" model, aligning with mixed-precision's goals to achieve high accuracy with optimized resource use. The research delves into KD's application in parallelized environments and its impact on the efficiency and accuracy of NLP models, a territory less traversed in current studies.

## Background

The shift to mixed-precision computing signifies a pivotal change in the computational strategy for deep learning. Traditionally, the field relied on high-precision arithmetic for stability but at the cost of computational resources. Mixed-precision disrupts this by smartly leveraging lower-precision calculations where possible without compromising the essential high precision. This strategy offers substantial benefits like reduced memory use and faster processing, enabling the training of more sophisticated models. However, as models like BERT scale up, concerns about the precision-related degradation in accuracy arise. Here, the integration of KD into mixed-precision training becomes crucial. KD, focusing on a smaller 'student' model's training to replicate a larger 'teacher' model, may offer solutions to these challenges. By potentially reducing computational demands during key training phases, KD could complement mixed-precision computing, especially in large-scale NLP tasks, suggesting a promising new direction for research in this field.

## Problem

As deep learning models, particularly in NLP, grow more intricate, the importance of parallelized training has surged. This shift from traditional single-machine training to parallel processing aims to improve efficiency and performance. However, it introduces new challenges, especially in managing complex models like BERT. These challenges include efficient data distribution, consistent model initialization across nodes, and maintaining precision in gradient calculations during the backward pass in a parallelized setting.

Furthermore, while mixed-precision computing has proven effective in balancing computational efficiency and accuracy, concerns arise when scaling such models. Potential accuracy degradations due to reduced computational precision become a critical issue. Here, the incorporation of Knowledge Distillation (KD) presents a unique opportunity. This research investigates the synergistic potential of mixed-precision computing and KD in training complex NLP models like BERT. By training a less complex

student model through KD, the study aims to explore how this approach can mitigate the computational and precision challenges inherent in mixed-precision settings. The goal is to maintain, or even enhance, efficiency and scalability while balancing computational efficiency and accuracy, a vital aspect in advanced deep learning applications. This intersection presents a novel approach to advancing the field of deep learning, particularly in large-scale NLP tasks, making it a vital area of exploration.

The research will delve into how KD can complement mixed-precision training, potentially mitigating accuracy losses, thereby presenting a novel approach to model training in the realm of deep learning.

## Research Question

How does the integration of Knowledge Distillation in mixed-precision, parallelized training environments impact the efficiency, scalability, and accuracy of BERT models in large-scale NLP tasks?

## Research Objectives

The objectives include assessing the impact of mixed-precision computing on the performance of BERT models using KD, investigating optimization strategies for reduced-precision gradient calculations, analyzing training efficiency, scalability, and fault tolerance, and providing guidelines for implementing mixed-precision computing in parallelized environments.

## Scope

The research will concentrate on using mixed-precision computing techniques, particularly emphasizing Tensor Cores in NVIDIA GPUs, for training BERT models in parallelized environments. By focusing on Knowledge Distillation within this context, the study will provide insights into the practical implementation and challenges of mixed-precision computing in distributed systems for NLP tasks.

## Methodology

This section will delve into the methodology of the research.

## Overview

This study focuses on training BERT models for question-answering tasks using the WikiQA corpus, comparing three distinct training approaches to assess the impact of Knowledge Distillation (KD) and mixed-precision computing on model performance.

## Experiment Design

### Training Setup:

- **Model Selection:** Utilize BERT as the 'teacher' model in a Knowledge Distillation (KD) framework. Select an appropriate smaller model as the 'student'.
- **Dataset:** Employ the WikiQA corpus for training and evaluating the models. This dataset is specifically designed for question-answering tasks and is well-suited for evaluating the performance of NLP models in educational contexts.

### Training Procedures:

1. Baseline Model Training:
  - Traditional Training: Train a standard BERT model using single-precision computations, without KD and parallelization.
2. Full Precision Training with KD:
  - Model Selection: BERT as the 'teacher' and a smaller model as the 'student'.
  - Training Procedure: Train using full precision (32-bit) in a parallelized environment.
3. Mixed Precision Training with KD:
  - Model Selection: BERT as the 'teacher' and a smaller model as the 'student'.
  - Training Procedure: Implement training using mixed precision, utilizing the NVIDIA 3080 GPU's lower precision capabilities.

### Parallelization Strategy:

- Implement both training procedures in a parallelized environment, distributing the computational workload across the available processing units in the NVIDIA 3080 GPU.

## Metrics for Evaluation

### Speed:

- Measured as the total time taken from the start to the end of the training process.
- Use a timing tool to record the duration of training for both approaches.

### Efficiency:

- Evaluate GPU utilization and memory consumption during training.
- Measure power usage and computational resource allocation using system monitoring tools.

### Accuracy:

- Assess the model's performance on a held-out portion of the WikiQA dataset.
- Calculate metrics like precision, recall, and F1 score to evaluate the effectiveness of the model in answering questions correctly.

## Data Analysis

- **Comparative Analysis:** Perform statistical comparisons between the two training approaches across the defined metrics.
- **Interpretation:** Analyze results to understand the trade-offs between speed, efficiency, and accuracy in both training scenarios.

## Potential Challenges and Solutions

- **Optimal Performance Achievement:** A significant challenge lies in ensuring that each method – traditional training, full precision with KD, and mixed precision with KD – is implemented to achieve near-optimal performance. This involves fine-tuning the training parameters and carefully setting up each environment to accurately compare their effectiveness, leveraging best practices
- **Implementation Accuracy:** Ensuring accurate and effective implementation of each method is crucial. This includes appropriate configuration of the GPU, effective utilization of KD techniques, and precise adjustment of precision levels in mixed precision training.
- **Precision Loss:** In mixed precision training, there may be concerns about the loss of numerical precision. Implement strategies such as loss scaling to mitigate this.
- **Resource Limitations:** Ensure the parallelized training environment is optimally configured to handle the computational demands of both full and mixed precision training.