# Literature Review: Parallelization Strategies for Neural Network

## Introduction

Parallelization strategies in neural network training have become increasingly crucial due to the rising complexity and size of datasets. Efficient computational techniques are essential for timely model training and deployment. This review categorizes advancements into GPU-based parallelization, model parallelism, sparse neural networks, convolutional neural network (CNN) parallelization, and asynchronous optimization techniques.

## Recent Work on Parallelization Strategies in Neural Network Training

The field of neural network training has seen considerable recent advancements in parallelization strategies, reflecting growing efforts to tackle the challenges posed by increased data complexity and model size. This section outlines the contributions of various researchers and their impact on enhancing the efficiency and effectiveness of neural network training through innovative parallelization techniques.

### GPU-Based Parallelization

Recent advancements have underscored the effectiveness of GPUs in enhancing the computational efficiency of neural networks. Amin et al. (2019) demonstrated how GPU clusters can accelerate training tasks significantly, suggesting an improved allocation and utilization of GPU resources. Similarly, Bonvallet et al. (2014) discussed a CUDA-based implementation that optimizes real-time financial predictions, indicating the pivotal role of GPU-based parallelization in practical applications.

### Model Parallelism

Feng (2023) provided an insightful analysis into the application of model parallelism within deep residual networks, illustrating how distributing model parameters across multiple GPUs can accelerate training processes and manage larger models effectively. This approach is particularly advantageous in mitigating communication overhead and enhancing overall training speed.

Sparse Neural Networks

The focus on optimizing sparse computations in neural networks has led to notable research contributions, such as those by Gajurel et al. (2020), who explored the acceleration of sparse neural network training using GPUs. This research highlights the critical role of GPU architecture in enhancing the efficiency of sparse computations, which are essential for the scalability of neural network training.

CNN Parallelization

Krizhevsky (2014) introduced techniques for parallelizing convolutional neural networks, which are essential for tasks like image recognition. By distributing the workload of CNNs across multiple GPUs, significant reductions in training time were achieved, enhancing the practicality and efficiency of these networks in real-world applications.

GPU Asynchronous Optimization Techniques

Paine (2013) investigated the asynchronous stochastic gradient descent (SGD) methods on GPUs, a critical optimization technique that accelerates convergence and improves the efficiency of training neural networks. This method is particularly relevant for scaling neural network training to handle complex and large datasets efficiently.

Scalability Challenges and Solutions

Pourghassemi et al. (2020) and Singh and Bhatele (2021) explored the scalability challenges in parallelizing CNNs on GPUs and proposed Axonn, a framework designed for extreme-scale deep learning tasks. This research is pivotal in understanding the constraints and potentials of scaling deep learning models on advanced GPU architectures.

Advancements in Optimization Algorithms

Lavin and Gray (2016) focused on optimizing convolution operations for CNNs, significantly enhancing their execution efficiency on GPUs. Their work has been instrumental in advancing the widespread adoption of CNNs in high-performance applications requiring efficient image processing.

By categorizing the literature into these thematic areas, the review not only highlights the breadth of research in neural network parallelization but also underscores the specific advancements and challenges within each category. This structured approach facilitates a clearer understanding of where the field stands and where it is headed.

**Rationale**

The decision to pursue the parallelization of neural networks is driven by several key factors.

· Firstly, the exponential growth in dataset sizes and model complexities necessitates efficient computational strategies to ensure timely model training and deployment.

· Parallelization offers a scalable solution to this challenge, enabling the handling of larger datasets and more intricate models. Secondly, the advancements in GPU technology, such as the introduction of the Apple M1 chip, have significantly enhanced the computational power available for parallelization, making it a viable and efficient approach.

· Lastly, the implications of neural network parallelization extend beyond computational efficiency, driving innovation and enabling the development of more accurate and sophisticated models with applications in various domains.

Conclusion

In conclusion, the proposal to explore advancements in parallelization strategies for neural network training is supported by recent developments in the field and the clear benefits it offers in terms of computational efficiency and model sophistication. By leveraging GPU technology and optimizing parallelization frameworks, researchers and practitioners can unlock the full potential of machine learning technologies, paving the way for new advancements and applications in artificial intelligence.