

Efficient Training Techniques for Large-scale Neural Networks on Distributed Systems

Carole Antonio*

Department of Business Information Systems, Humboldt University of Berlin, Berlin, Germany

Introduction

The recent advancements in deep learning and the availability of vast amounts of data have led to the emergence of large-scale neural networks that require extensive computational resources for training. To address this challenge, researchers have explored various techniques to efficiently train these networks on distributed systems. This research article provides an overview of the state-of-the-art techniques employed for training large-scale neural networks on distributed systems, highlighting their advantages and limitations. The article discusses parallelization strategies, parameter server architectures, communication protocols, and synchronization mechanisms. Furthermore, it explores the trade-offs between computational efficiency and communication overhead. Finally, the article presents promising future directions for efficient training of large-scale neural networks on distributed systems. Large-scale neural networks, such as deep convolutional neural networks (CNNs) and recurrent neural networks (RNNs), have demonstrated remarkable success in various domains, including computer vision, natural language processing, and speech recognition. However, training these networks can be computationally expensive due to the increasing model sizes and the massive datasets involved. Distributed systems have become a popular solution for accelerating the training process by leveraging multiple computational resources in parallel [1-3]. This article investigates the efficient training techniques employed for large-scale neural networks on distributed systems.

Description

Parallelization strategies

Parallelization is a key technique for training large-scale neural networks on distributed systems. This section discusses data parallelism, model parallelism, and hybrid parallelism. Data parallelism involves dividing the training data across multiple workers and performing parallel computations on different subsets. Model parallelism divides the neural network model into subcomponents and assigns each component to different workers. Hybrid parallelism combines both data and model parallelism to achieve optimal performance.

Parameter server architectures

Parameter servers are critical components in distributed training systems that store and synchronize model parameters across workers. This section explores different parameter server architectures, including centralized parameter servers, decentralized parameter servers, and hybrid architectures. Each architecture has its advantages and considerations in terms of communication overhead, fault tolerance, and scalability [4,5].

*Address for Correspondence: Carole Antonio, Department of Business Information Systems, Humboldt University of Berlin, Berlin, Germany, E-mail: CaroleAntonio3@gmail.com

Copyright: © 2023 Antonio C. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution and reproduction in any medium, provided the original author and source are credited.

Received: 17 April, 2023, Manuscript No. jcsb-23-99542; Editor Assigned: 19 April, 2023, Pre QC No. P-99542; Reviewed: 03 May, 2023, QC No. Q-99542; Revised: 09 May, 2023, Manuscript No. R-99542; Published: 17 May, 2023, DOI:10.37421/0974-7230.2023.16.464

Communication protocols

Efficient communication protocols play a crucial role in reducing the training time of large-scale neural networks on distributed systems. This section discusses various communication protocols, such as parameter-based communication, gradient-based communication, and compressed communication. These protocols aim to minimize the communication overhead while ensuring accurate synchronization of model updates.

Synchronization mechanisms

Synchronization mechanisms ensure consistency across distributed workers during training. This section examines synchronous and asynchronous training approaches. Synchronous training achieves global synchronization at every iteration, while asynchronous training allows workers to update parameters independently, leading to potential staleness but improved scalability. Additionally, techniques like gradient compression and quantization are explored to mitigate the communication and synchronization costs.

Trade-offs and challenges

Efficient training of large-scale neural networks on distributed systems involves several trade-offs and challenges. This section discusses the trade-offs between computational efficiency and communication overhead, the impact of network topology on training performance, and the challenges of fault tolerance and straggler mitigation. It also highlights the importance of efficient memory utilization and load balancing in distributed training.

Future directions

To further enhance the efficiency of training large-scale neural networks on distributed systems, this section presents promising future directions. It explores the potential of hardware accelerators, such as GPUs and specialized AI chips, for distributed training. Additionally, research areas like automatic parallelization, dynamic load balancing, and adaptive communication strategies are identified as avenues for future exploration.

Conclusion

Efficiently training large-scale neural networks on distributed systems is crucial to meet the computational demands of modern deep learning applications. This research article provided an overview of efficient training techniques, including parallelization strategies, parameter server architectures, communication protocols, and synchronization mechanisms. It highlighted the trade-offs and challenges associated with these techniques.

Acknowledgement

None.

Conflict of Interest

Authors declare no conflict of interest.

References

1. Liang, Zhengfa, Yulan Guo, Yiliu Feng and Wei Chen, et al. "Stereo matching

- using multi-level cost volume and multi-scale feature constancy." *IEEE Trans Pattern Anal Mach Intell* 43 (2019): 300-315.
2. Guo, Yulan, Mohammed Bennamoun, Ferdous Sohel and Min Lu, et al. "3D object recognition in cluttered scenes with local surface features: A survey." *IEEE Trans Pattern Anal Mach Intell* 36 (2014): 2270-2287.
 3. Sayadnavard, Monireh H., Abolfazl Toroghi Haghighat and Amir Masoud Rahmani. "A multi-objective approach for energy-efficient and reliable dynamic VM consolidation in cloud data centers." *Eng Sci Technol Int J* 26 (2022): 100995.
 4. Wooller, Sarah K., Graeme Benstead-Hume, Xiangrong Chen and Yusuf Ali, et al. "Bioinformatics in translational drug discovery." *Biosci Rep* 37 (2017).
 5. Eltayieb, Nabeil, Rashad Elhabob, Alzubair Hassan and Fagen Li. "A blockchain-based attribute-based signcryption scheme to secure data sharing in the cloud." *J Syst Archit* 102 (2020): 101653.

How to cite this article: Antonio, Carole. "Efficient Training Techniques for Large-scale Neural Networks on Distributed Systems." *J Comput Sci Syst Biol* 16 (2023): 464.