**Programa de Mestrado e Doutorado em Engenharia de Telecomunicações**
**TP547 – Princípios de Simulação de Sistemas de Comunicação**
**Trabalho 2: Revisar um artigo e reproduzir um resultado**
**Alunos: Christopher Lima & Gabriel Pivoto**

Article Name: Quantized Compressive Sampling of Stochastic Gradients for Efficient Communication in Distributed Deep Learning [1].

# 1    Article summary

The article presents an exploration of enhancing communication efficiency in distributed deep learning systems. The authors introduce a novel technique called Quantized Compressive Sampling (QCS) to reduce the amount of data exchanged between nodes during the training process. Their approach addresses the challenge of high communication overhead in distributed settings, particularly when dealing with large-scale models.

# 2    Reading process

The authors meticulously detail the QCS method, with a well-written process, elucidating its underlying principles and demonstrating its effectiveness through mathematical equations. They provide clear explanations of the equations, ensuring accessibility for readers from various backgrounds. By leveraging compressive sensing and quantization, QCS enables a significant reduction in communication overhead while preserving the essential information required for model convergence.

In the results section, the authors showcase the practical implications of their approach through empirical evaluations. They compare QCS against existing communication compression methods, demonstrating superior performance in terms of communication efficiency without compromising model accuracy. Through extensive experiments on various datasets and deep learning architectures, the authors substantiate the efficacy of QCS in enabling efficient communication without sacrificing training quality.

While the article provides a comprehensive exploration of the proposed technique, there are a few potential issues or areas for improvement that the authors may have failed to address or clarify:

- **Scalability Concerns**: The scalability of the proposed approach to extremely large-scale distributed deep learning systems may not have been thoroughly discussed. While the authors demonstrate efficacy across various datasets and architectures, the performance of QCS in scenarios with a significantly higher number of nodes or more complex models could be a concern.

- **Sensitivity to Hyperparameters**: The sensitivity of QCS to specific hyperparameters or tuning settings is not extensively examined. Variations in compression ratio, quantization levels, or other parameters may have a significant impact on performance. A more in-depth analysis of the robustness of QCS to hyperparameter changes would enhance the reliability of the proposed technique.

- **Computational Overhead**: While the focus is primarily on reducing communication overhead, the potential increase in computational overhead associated with implementing QCS is not explicitly addressed. It would be beneficial to quantify any additional computational costs incurred by the compression and decompression processes, especially in resource-constrained environments.

- **Generalizability**: The generalizability of QCS across different types of deep learning tasks or optimization algorithms is not thoroughly discussed. Certain applications or optimization techniques may have specific characteristics that influence the effectiveness of communication compression methods differently. Providing insights into the generalizability of QCS would enhance its applicability across a wider range of scenarios.

All these limitations are mentioned by the authors but are not analyzed or explained.

Addressing these potential limitations or uncertainties would strengthen the proposed technique's credibility and applicability, ensuring its effectiveness across a broader range of distributed deep-learning scenarios.

Overall, the article makes a significant and relevant contribution to the field of distributed deep learning by introducing a novel technique for mitigating communication overhead. Its clear presentation, thorough explanation of equations, and compelling demonstration of results underscore its importance in advancing the efficiency of distributed deep learning systems.

# 3 Results Reproduction

Upon reviewing the article, Figure 1 and Figure 2 were deemed suitable for results reproduction.
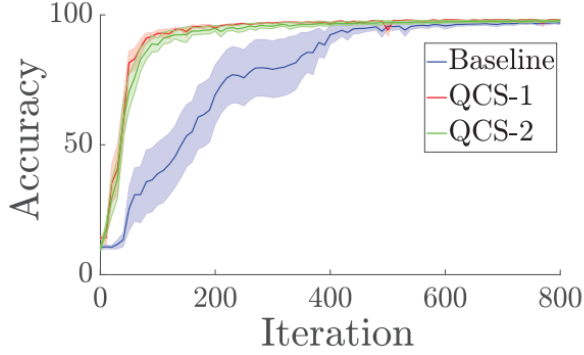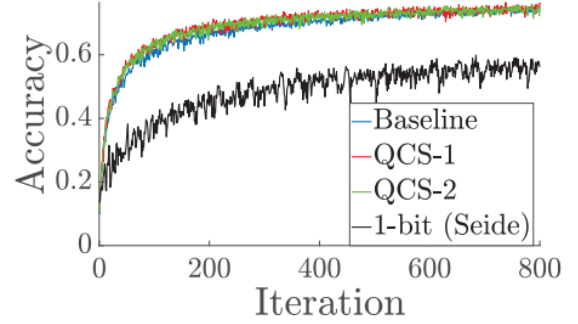


Figure 1: Curve 1



Figure 2: Curve 2

The curves were selected due to the complexity of the other graphs presented in this study, which feature intricate curves that cannot be represented by simple functions. This complexity significantly impeded the process of determining the corresponding functions. To extract the functions of the selected curves, the Web-PlotDigitizer [2] tool was employed. Upon loading each image, the x and y axes were carefully aligned, facilitating the acquisition of numerous points along the curves necessary for deriving their respective functions. Subsequently, these functions were utilized to develop a Python script capable of accurately reproducing the two curves. The results reproduction can be seen in Figure 3 and Figure 4. The implemented code is accessible on GitHub [1].
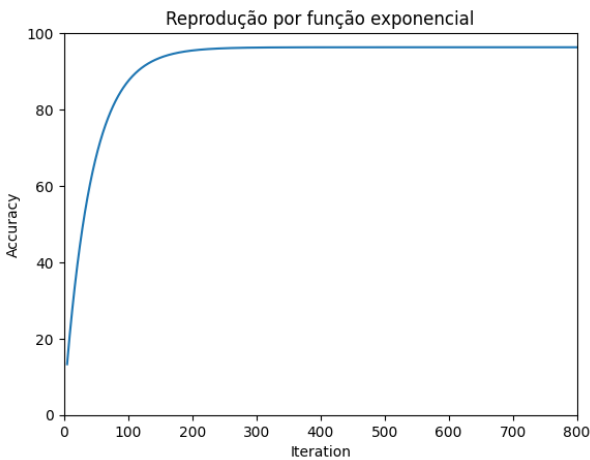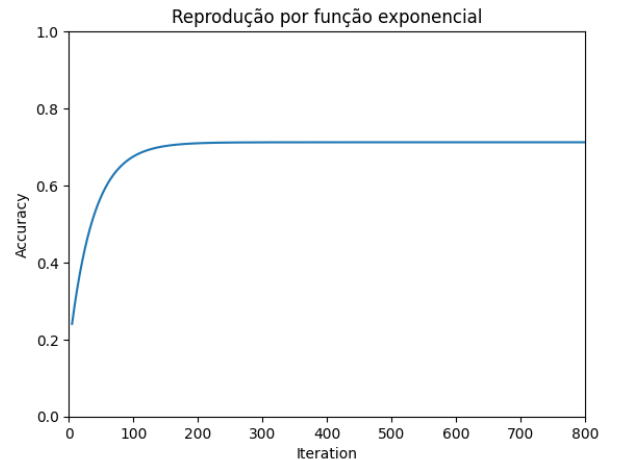


Figure 3: Result 1



Figure 4: Result 2

---

[1]https://github.com/GabrielPivoto/tp547/blob/master/TrabalhoFinalParte2/TrabalhoFinalParte2.ipynb

# References

[1] Shaohuai Shi, Xiaowen Chu, and Bo Li. "MG-WFBP: Merging Gradients Wisely for Efficient Communication in Distributed Deep Learning". In: *IEEE Transactions on Parallel and Distributed Systems* 32.8 (2021), pp. 1903–1917. DOI: 10.1109/TPDS.2021.3052862.

[2] Ankit Rohatgi. *WebPlotDigitizer*. URL: https://automeris.io/WebPlotDigitizer.html.