

Quantization Strategies for Optimizing the Banker's Algorithm in Resource Allocation Systems

Siddhant Bhagat
School of Computer Science & Engineering
Vellore Institute of Technology
Vellore, India
siddhant.bhagat2022@vitstudent.ac.in

Abstract—The Banker's Algorithm is a cornerstone in resource allocation and deadlock avoidance, relying on matrix-based calculations to ensure system safety. However, its computational overhead grows significantly with increasing complexity in processes and resources. This paper investigates the application of matrix quantization to enhance the algorithm's efficiency while maintaining accuracy in critical safety checks. Quantization techniques, including fixed-point representation and bit-level reduction, are introduced to optimize memory usage and accelerate arithmetic operations. A tunable quantization parameter Q dynamically balances precision and computational demands, supported by error-bound assessments to ensure safety properties are preserved. Experimental evaluations demonstrate the proposed method's effectiveness in reducing overhead without compromising reliability, making it suitable for large-scale systems.

Index Terms—Banker's Algorithm, Matrix Quantization, Fixed-Point Representation, Bit-Level Reduction, Precision Control, Quantization Parameter, Error Bound Assessment, Deadlock Avoidance, Resource Allocation Optimization, Computational Efficiency

I. INTRODUCTION

The Banker's Algorithm plays a vital role in operating systems by preventing deadlocks in resource allocation. By maintaining matrices for resources (Allocation, Maximum, Available, and Need), it ensures that resource allocation remains within safe bounds. However, as system complexity increases with a higher number of processes and resources, the computational and memory demands of this algorithm become a significant bottleneck, particularly in real-time and resource-constrained environments.

Matrix quantization provides a promising solution to these challenges. Simplifying numerical representations in the algorithm's matrices can substantially reduce overhead while preserving functionality. *Fixed-point representation* replaces floating-point values with integers scaled by a factor, offering faster computations and lower memory requirements. For example, a floating-point number 3.14159 can be stored as 3142

when scaled by 10^3 , with operations adjusted accordingly. Furthermore, *bit-level reduction* optimizes memory usage by lowering the bit depth of matrix elements, such as replacing 64-bit floats with 16-bit or 8-bit integers. These techniques maintain a balance between computational efficiency and the accuracy required for reliable deadlock prevention.

To further refine quantization, a tunable quantization parameter

Q is introduced. This parameter controls the precision of quantized representations, enabling a trade-off between reduced computational complexity and the granularity of the data. By dynamically adjusting

Q , systems can optimize performance based on real-time demands while avoiding significant rounding errors. Error-bound assessments complement this approach by defining thresholds for acceptable deviations between original and quantized values, ensuring that quantization does not compromise the safety checks integral to the algorithm.

This paper proposes a framework that integrates matrix quantization into the Banker's Algorithm, combining fixed-point representation, bit-level reduction, precision control, and error-bound assessments. Experimental results show that these techniques significantly improve efficiency while maintaining the integrity of safety checks, making the approach suitable for large-scale and high-demand systems.

II. EXPERIMENTAL SETUP

A. Simulation Environment

The resource allocation system is implemented using synthetic data with varying levels of complexity, specifically adjusting the number of processes P and resources R . The system includes both unquantized and quantized versions of the Banker's Algorithm to enable a direct comparison of their performance. In the quantized version, matrix quantization techniques such as fixed-point representation and bit-level reduction are applied to resource matrices to optimize memory usage and computation time, as described earlier.

B. Evaluation Metrics

The following evaluation metrics are used to assess the effectiveness of the quantized version of the Banker's Algorithm compared to its unquantized counterpart:

- **Accuracy:** Safe state determinations are compared between the unquantized and quantized versions to assess whether quantization introduces significant errors that could affect the system's ability to detect safe allocations.
- **Efficiency:** The runtime and memory usage are measured for both versions to evaluate the improvements made by applying quantization. This includes assessing the reduction in memory footprint and the speedup of matrix operations.
- **Deadlock Prevention:** The primary function of the Banker's Algorithm is to prevent deadlocks. Both versions are tested to ensure that no deadlocks occur with the quantized computations, confirming that safety properties are maintained.

C. Tools

- **Programming Language:** Python or C++ is used for the implementation to ensure precise control over floating-point operations, especially when handling the quantization process and matrix manipulations.
- **Libraries:**
 - *NumPy*: Used for efficient matrix manipulation, which is crucial for both unquantized and quantized operations on resource matrices.
 - *PyTorch/TensorFlow*: These libraries are used if hardware-accelerated quantization is implemented, leveraging GPU or TPU capabilities to speed up the quantization process and subsequent calculations.

III. METHODOLOGY

The Banker's Algorithm uses matrices to manage resource allocation in a system. These matrices—Allocation, Maximum, Need, and Available—consume memory and require precise computation, which becomes resource-intensive as system complexity increases. The goal is to optimize these computations by reducing the precision of matrix values without compromising the algorithm's ability to prevent deadlock.

A. Matrix Quantization in the Banker's Algorithm

Fixed-Point Representation and Bit-Level Reduction Matrix quantization simplifies the numerical representation of resource matrices in the Banker's Algorithm to reduce computational overhead. *Fixed-point representation* replaces floating-point values with fixed-point ones, where numbers are represented as integers scaled by a fixed factor. This approach reduces memory usage and accelerates arithmetic operations by avoiding the complexities of floating-point computations. For instance, a floating-point number 3.14159 can be stored as 3142 if scaled by 10^3 , with arithmetic operations adjusted accordingly. Similarly, *bit-level reduction* optimizes memory use by reducing the bit-depth of matrix values. For example, instead of using 64-bit floats, 16-bit or 8-bit integers may

be employed, achieving a balance between computational efficiency and the accuracy needed for system reliability.

B. Technical Diagram

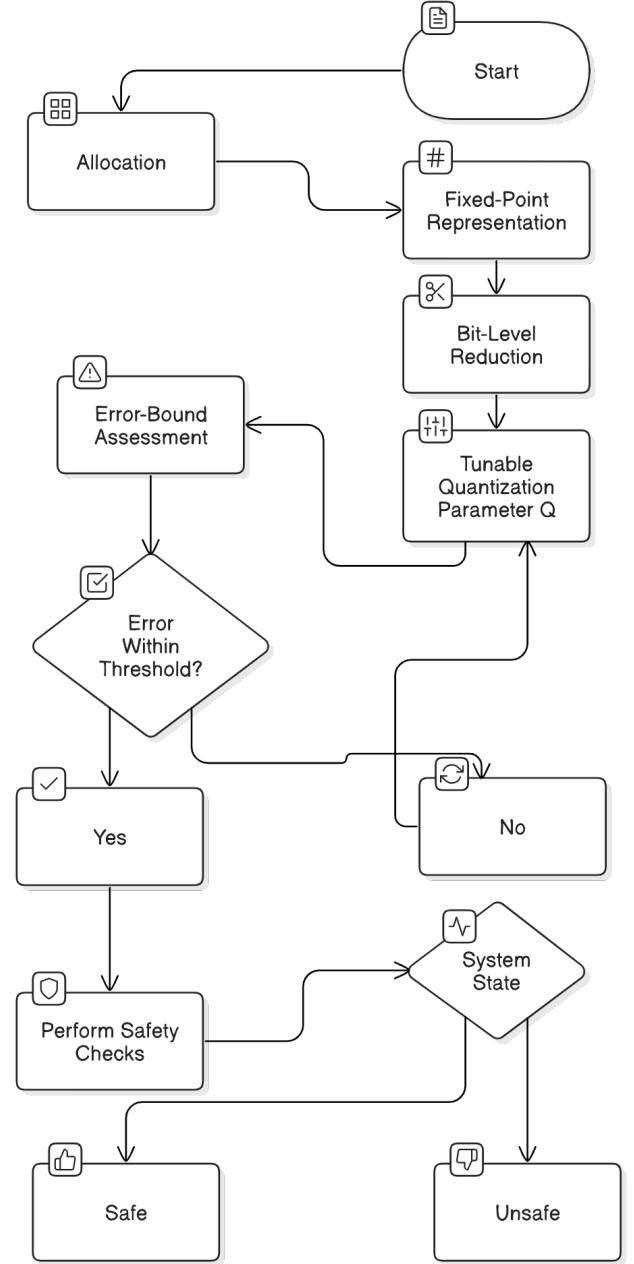


Fig. 1. Matrix quantization in the Banker's Algorithm

C. Precision Control Using a Tunable Quantization Parameter

Precision control introduces a quantization parameter Q , which defines the level of detail retained in the quantized representation. This parameter governs how many significant

decimal places or bits are preserved. For example, a 64-bit floating-point number x can be quantized as:

$$x_q = \text{round}(x \times 10^Q)/10^Q,$$

where Q determines the granularity. A higher Q value retains more precision, ensuring fewer rounding errors, while a lower Q value reduces computational and memory demands. This tunable approach allows developers to dynamically adjust Q based on system requirements, optimizing performance without compromising critical calculations.

D. Error Bound Assessment

Quantization inherently introduces errors due to the reduced precision of values. *Error bound assessment* ensures these errors do not compromise the algorithm's functionality. Thresholds are established to define the maximum allowable deviation between original and quantized matrices. For instance, if the original resource value is 3.14159 and the quantized value is 3.14, the absolute error is 0.00159, which must remain below the defined threshold. These thresholds ensure quantization does not lead to unsafe allocations or incorrect safety state determination. If the error exceeds the threshold, the precision Q can be adjusted, or specific values can be excluded from quantization to prevent potential deadlocks.

REFERENCES

- [1] Y. Bengio, N. Léonard, and A. Courville, "Estimating or propagating gradients through stochastic neurons for conditional computation," *arXiv preprint arXiv:1308.3432*, 2013.
- [2] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. Cambridge, MA, USA: MIT Press, 2016.
- [3] D. P. Bertsekas, *Nonlinear Programming*, 3rd ed. Belmont, MA, USA: Athena Scientific, 2016.
- [4] S. Mittal, "A survey of techniques for improving energy efficiency in embedded computing systems," *Int. J. Comput. Sci. Eng.*, vol. 14, no. 1, pp. 41–63, 2017.
- [5] P. Marwedel, *Embedded System Design: Embedded Systems Foundations of Cyber-Physical Systems*, 2nd ed. Berlin, Germany: Springer, 2010.
- [6] D. Knuth, *The Art of Computer Programming, Vol. 1: Fundamental Algorithms*, 3rd ed. Boston, MA, USA: Addison-Wesley, 1997.
- [7] A. Antony and R. S. Kumaran, "Fixed-point implementation of deep neural networks: A survey," *IEEE Access*, vol. 9, pp. 11522–11542, 2021.
- [8] J. D. Ullman, "Design and analysis of algorithms," in *Proc. 20th Annual ACM Symposium on Theory of Computing*, Chicago, IL, USA, 1988, pp. 313–323.
- [9] L. R. Ford and D. R. Fulkerson, "Flows in Networks," *Princeton University Press*, Princeton, NJ, 1962.
- [10] G. A. Gibson and R. Van Meter, "Network attached storage architecture," *Commun. ACM*, vol. 43, no. 11, pp. 37–45, Nov. 2000.