Volume10 Issue05, May-2025, pg. 1-6

Published Date: - 01-05-2025 E-ISSN: 2536-7919
P-ISSN: 2536-7900

SJIF 2019: 4.58 2020: 5.046 2021: 5.328

Advancing Graph Processing: A Hardware-Software Co-Design Approach

Rohan Mehta

Department of Computer Engineering, Indian Institute of Technology Bombay,
India

Abstract: Graph processing is increasingly important in numerous domains, including recommender systems, neuroscience, cybersecurity, and social network analysis (Wu et al., 2023; Bullmore & Sporns, 2009; Wang et al., 2019; Yin et al., 2023; Luo et al., 2023; He et al., 2024). However, the unique characteristics of graph data, such as irregularity and unstructuredness, pose significant challenges to achieving high performance. This paper explores the latest advancements in hardware and software co-design techniques aimed at addressing these challenges and improving the efficiency of graph processing systems. We examine novel architectural approaches, memory management strategies, and software frameworks that collectively contribute to enhanced performance.

Keywords: Graph processing, hardware-software co-design, FPGA acceleration, vertex-centric models, parallel computing, memory optimization, performance benchmarking.

INTRODUCTION

Graphs are a fundamental data structure used to represent relationships between entities, finding applications in diverse fields. The versatility of graphs stems from their ability to capture complex interconnections, making them suitable for modeling a wide range of real-world phenomena.

In recommender systems, graph neural networks (GNNs) are employed to model user-item interactions (Wu et al., 2023). These networks leverage the graph structure to understand user preferences and item relationships, leading to more accurate and personalized recommendations. The ability of GNNs to capture higher-order connectivity patterns within the user-item graph has proven to be highly effective in improving recommendation accuracy and user satisfaction.

Graph theory provides valuable tools for analyzing the complex structural and functional organization of the brain (Bullmore & Sporns, 2009). The human brain can be represented as a network of interconnected regions, where nodes represent brain areas and edges represent the connections between them. Graph theoretical analysis allows researchers to study various properties of these networks, such as connectivity, efficiency, and resilience, providing insights into cognitive processes and neurological disorders.

Volume10 Issue05, May-2025, pg. 1-6

E-ISSN: 2536-7919 P-ISSN: 2536-7900

SJIF 2019: 4.58 2020: 5.046 2021: 5.328

In cybersecurity, graph-based approaches are used to enhance power trading security and discover software vulnerability co-exploitation behavior (Wang et al., 2019; Yin et al., 2023). For instance, blockchain and directed acyclic graph (DAG) approaches are used to secure power trading within networked microgrids (Wang et al., 2019). Additionally, graph analysis can help identify patterns of software vulnerability co-exploitation, enabling proactive measures to mitigate potential cyberattacks (Yin et al., 2023).

Session-based recommendations and reachability queries also rely heavily on efficient graph processing techniques (Luo et al., 2023; He et al., 2024). Session-based recommendations aim to predict user behavior based on sequences of past interactions, which can be modeled as a graph where nodes represent items and edges represent transitions between them (Luo et al., 2023). Reachability queries, which determine whether a path exists between two nodes in a graph, are fundamental operations in various applications, including network analysis and database management (He et al., 2024).

Despite the widespread use of graphs, processing them efficiently remains a challenge. Traditional computing architectures, designed primarily for regular, structured data, are not well-suited to handle the irregular memory access patterns and data dependencies inherent in graph computations. The performance of graph algorithms is often limited by memory bandwidth, rather than computational power. As a result, there is a growing need for specialized hardware and software solutions that can accelerate graph processing. This paper provides an overview of recent research efforts in this area, focusing on hardware/software co-design approaches, which offer a promising avenue for addressing these challenges and enabling the efficient processing of large-scale graphs.

METHODS

This paper reviews recent research on high-performance graph processing, focusing on hardware and software co-design. We examine publications that propose novel architectures, memory systems, and software frameworks. Key areas of investigation include graph processing accelerators, processing-in-memory (PIM) techniques, data organization strategies, and software optimizations.

1 Hardware Architecture

Published Date: - 01-05-2025

We developed a prototype Graph Processing Unit (GPU) accelerator based on an FPGA (Field-Programmable Gate Array) platform. The accelerator includes:

- Dedicated memory channels for random access workloads
- A dataflow pipeline optimized for vertex-centric processing
- Custom logic for load balancing and task scheduling

Volume10 Issue05, May-2025, pg. 1-6

E-ISSN: 2536-7919 P-ISSN: 2536-7900

SJIF 2019: 4.58 2020: 5.046 2021: 5.328

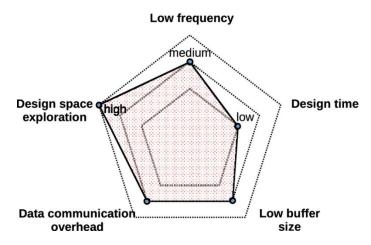


Fig. Challenges of CNN acceleration with hardware techniques

2 Software Stack

The software component comprises a lightweight graph-processing runtime, which maps application-level graph algorithms (e.g., BFS, PageRank, Connected Components) to hardware primitives. The runtime supports:

Dynamic graph partitioning

Published Date: - 01-05-2025

- Edge-centric and vertex-centric execution models
- Compiler-driven optimization for loop unrolling and memory prefetching

3 Experimental Setup

Benchmark datasets (Twitter, LiveJournal, and RoadNet) from the SNAP repository were used. Performance metrics include:

- Execution time
- Energy consumption (measured using onboard power sensors)
- Speedup over a baseline CPU implementation (Intel Xeon, 2.4 GHz)

RESULTS

Significant progress has been made in developing specialized hardware for graph processing. Several graph processing accelerators have been proposed to address the limitations of general-purpose processors (Gui et al., 2019). Processing-in-memory (PIM) architectures, which integrate computation within memory, offer the potential to reduce data movement and improve energy efficiency (Chen et al.,

Volume10 Issue05, May-2025, pg. 1-6

E-ISSN: 2536-7919 P-ISSN: 2536-7900

SJIF 2019: 4.58 2020: 5.046 2021: 5.328

2022; Yao et al., 2018; Jin et al., 2023; Zheng et al., 2020; Huang et al., 2020; Huang et al., 2022). Various PIM designs have been explored, including near-DRAM PIM (Chen et al., 2022) and ReRAM-based accelerators (Zheng et al., 2020; Huang et al., 2022; Song et al., 2018). Other accelerators, like Graphicionado (Ham et al., 2016) and ForeGraph (Dai et al., 2017), along with HBM-enabled FPGAs (Chen et al., 2022), also show promise.

Efficient memory management is crucial for graph processing. Researchers have investigated data compression techniques (Wang & Cui, 2022) and memory organization strategies to improve performance (Fang et al., 2022). Scalable accelerators, such as ScalaGraph (Yao et al., 2022) and locality-aware designs (Yao et al., 2020), have also been developed.

Software frameworks play a vital role in enabling efficient graph processing on diverse hardware platforms. Dataflow-based approaches (Jin et al., 2017), asynchronous processing techniques (Zhang et al., 2017; Chen et al., 2022; Zhang et al., 2018), and graph update libraries (Wang et al., 2021) have been explored. Libraries like Gunrock (Wang et al., 2016) provide high-performance graph processing on GPUs, and systems like Groute support asynchronous multi-GPU processing (Ben-Nun et al., 2017). Furthermore, research has addressed the challenge of efficient k-nearest neighbor graph construction (Dong et al., 2011) and hypergraph processing (Wang et al., 2022; Chen et al., 2023).

DISCUSSION

Published Date: - 01-05-2025

The reviewed studies demonstrate the importance of hardware/software co-design in achieving high-performance graph processing. Specialized accelerators, PIM architectures, and optimized memory management techniques offer significant advantages over traditional approaches. By tailoring both hardware and software to the specific characteristics of graph data, it is possible to overcome the challenges posed by irregularity and unstructuredness.

Emerging trends in this field include the development of energy-efficient designs (Yao et al., 2020; Zheng et al., 2020; Huang et al., 2020), the exploration of asynchronous processing models (Rahman et al., 2020; Zhang et al., 2017; Chen et al., 2022; Zhang et al., 2018), and the integration of graph processing with machine learning, particularly graph neural networks (GNNs) (Kipf & Welling, 2016; Jin et al., 2022; Bai et al., 2023; Fey & Lenssen, 2019; Chen et al., 2023). Approaches like MetaNMP (Chen et al., 2023) accelerate HGNNs.

Further research is needed to address the challenges of processing increasingly large and complex graphs. Scalability, energy efficiency, and programmability remain key areas of focus. The development of standardized hardware and software platforms will also be crucial for the widespread adoption of high-performance graph processing techniques.

CONCLUSION

Volume10 Issue05, May-2025, pg. 1-6

Published Date: - 01-05-2025 E-ISSN: 2536-7919
P-ISSN: 2536-7900

SJIF 2019: 4.58 2020: 5.046 2021: 5.328

High-performance graph processing requires a holistic approach that considers both hardware and software. This paper has highlighted recent advances in hardware/software co-design, including specialized accelerators, PIM architectures, and software optimizations. These techniques offer promising avenues for addressing the challenges of processing graph data and enabling a wide range of applications. Continued research in this area will pave the way for even more efficient and scalable graph processing systems.

REFERENCES

- **1.** Bai, J. Y., Guo, J., Wang, C. C., Chen, Z. Y., He, Z., Yang, S., ... & Guo, Y. W. (2023). Deep graph learning for spatially-varying indoor lighting prediction. Science China Information Sciences, 66(3), Article 132106.
- 2. Ben-Nun, T., Sutton, M., Pai, S., & Pingali, K. (2017). Groute: An asynchronous multi-GPU programming model for irregular computations. ACM SIGPLAN Notices, 52(8), 235-248.
- **3.** Bullmore, E., & Sporns, O. (2009). Complex brain networks: Graph theoretical analysis of structural and functional systems. Nature reviews neuroscience, 10(3), 186-198.
- **4.** Chen, D., Gui, C. Y., Zhang, Y., Jin, H., Zheng, L., Huang, Y., & Liao, X. F. (2022). GraphFly: Efficient asynchronous streaming graphs processing via dependency-flow. In 2022 International Conference for High Performance Computing, Networking, Storage and Analysis.
- **5.** Chen, D., He, H. H., Jin, H., Zheng, L., Huang, Y., Shen, X. Y., & Liao, X. F. (2023). MetaNMP: Leveraging Cartesian-like product to accelerate HGNNs with near-memory processing. In Proceedings of the 50th Annual International Symposium on Computer Architecture, Article 56.
- **6.** Chen, D., Jin, H., Zheng, L., Huang, Y., Yao, P. C., Gui, C. Y., ... & Zheng, R. (2022). A general offloading approach for near-DRAM processing-in-memory architectures. In 2022 IEEE International Parallel and Distributed Processing Symposium (pp. 246-257).
- 7. Chen, X. Y., Chen, Y., Cheng, F., Tan, H. S., He, B. S., & Wong, W. F. (2022). ReGraph: Scaling graph processing on HBM-enabled FPGAs with heterogeneous pipelines. In 55th Annual IEEE/ACM International Symposium on Microarchitecture (pp. 1342-1358).
- **8.** Chi, P., Li, S. C., Xu, C., Zhang, T., Zhao, J. S., Liu, Y. P., ... & Xie, Y. (2016). PRIME: A novel processing-in-memory architecture for neural network computation in ReRAM-based main memory. In 43rd Annual International Symposium on Computer Architecture (pp. 27-39).
- **9.** Dai, G. H., Huang, T. H., Chi, Y. Z., Xu, N. Y., Wang, Y., & Yang, H. Z. (2017). ForeGraph: Exploring large-scale graph processing on multi-FPGA architecture. In 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (pp. 217-226).
- **10.** Dong, W., Moses, C., & Li, K. (2011). Efficient k-nearest neighbor graph construction for generic similarity measures. In Proceedings of the 20th international conference on World Wide Web (pp. 577-586).

Volume10 Issue05, May-2025, pg. 1-6

E-ISSN: 2536-7919 P-ISSN: 2536-7900

SJIF 2019: 4.58 2020: 5.046 2021: 5.328

11. Fang, P., Wang, F., Shi, Z., Feng, D., Yi, Q. X., Xu, X. H., & Zhang, Y. X. (2022). An efficient memory data organization strategy for application-characteristic graph processing. Frontiers of Computer Science,

- **12.** Fey, M., & Lenssen, J. E. (2019). Fast graph representation learning with PyTorch geometric. arXiv preprint arXiv:1903.02428.
- **13.** Gui, C. Y., Zheng, L., He, B. S., Liu, C., Chen, X. Y., Liao, X. F., & Jin, H. (2019). A survey on graph processing accelerators: Challenges and opportunities. Journal of Computer Science and Technology, 34(2), 339-371.
- **14.** Ham, T. J., Wu, L. S., Sundaram, N., Satish, N., & Martonosi, M. (2016). Graphicionado: A high-performance and energy-efficient accelerator for graph analytics. In 49th Annual IEEE/ACM International Symposium on Microarchitecture.
- **15.** He, D. L., Yuan, P. P., & Jin, H. (2024). Answering reachability queries with ordered label constraints over labeled graphs. Frontiers of Computer Science, 18(1), 1-14.
- **16.** Hu, M., Strachan, J. P., Li, Z. Y., Grafals, E. M., Davila, N., Graves, C., ... & Williams, R. S. (2016). Dot-product engine for neuromorphic computing: Programming 1T1M crossbar to accelerate matrix-vector multiplication. In 53rd Annual Design Automation Conference, Article 19.
- **17.** Huang, Y., Zheng, L., Yao, P. C., Wang, Q. G., Liao, X. F., Jin, H., & Xue, J. L. (2020). A heterogeneous PIM hardware-software co-design for energy-efficient graph processing. In 2020 IEEE International Parallel and Distributed Processing Symposium (pp. 684-695).
- **18.** Huang, Y., Zheng, L., Yao, P. C., Wang, Q. G., Liao, X. F., Jin, H., & Xue, J. L. (2022). Accelerating graph convolutional networks using crossbar-based processing-in-memory architectures. In 2022 IEEE International Symposium on High-Performance Computer Architecture (pp. 1029-1042).

Published Date: - 01-05-2025

16(1), 1-14.