# HYPERGEF: A FRAMEWORK ENABLING EFFICIENT FUSION FOR HYPERGRAPH NEURAL NETWORK ON GPUS

Sravanth Revath Krishna Thatavarthi

*C17745883*
*Computer and Information Sciences*
*School of Computing*
*Clemson University, Clemson*
*Email: sravant@clemson.edu*

*Abstract*—**Hypergraph Neural Networks (HyperGNNs) represent a promising evolution within Graph Neural Networks (GNNs), utilizing hyperedges to model intricate high-order relationships among vertices. However, existing GNN frameworks face challenges in seamlessly integrating message-passing steps between vertices and hyperedges, resulting in heightened latency and redundant memory consumption. This paper addresses key challenges hindering efficient fusion in HyperGNNs, including inefficient partitioning, workload-agnostic formats, and heavy writing conflicts.**

**To overcome these challenges, we introduce HyperGef, a framework incorporating an edge-split partition scheme for improved efficiency and workload balancing. HyperGef introduces a novel fusion workload-aware format to effectively represent the computation workload after fusion and partition. Additionally, a shared memory-aware grouping scheme is implemented to mitigate heavy writing conflicts.**

**Incorporating algorithmic enhancements, HyperGef demonstrates superior performance, outperforming the NVIDIA cuSPARSE kernel by 3.31× in extensive experiments. By enabling efficient fusion for HyperGNNs, HyperGef achieves a remarkable 2.25× to 3.99× end-to-end speedup across various HyperGNN models compared to state-of-the-art frameworks such as DGL and PyG. This research underscores the significance of HyperGef in addressing challenges and advancing the performance frontier of HyperGNNs.**

## 1. Introduction

Hypergraph Neural Networks (HyperGNNs) have emerged as a powerful paradigm in the realm of Graph Neural Networks (GNNs), offering a unique capability to model complex, high-order relationships among vertices through the use of hyperedges. However, the realization of the full potential of HyperGNNs is hindered by existing challenges in seamlessly integrating message-passing steps between vertices and hyperedges. The resulting inefficiencies contribute to elevated latency and redundant memory consumption, posing obstacles to the widespread adoption of HyperGNNs.

In response to these challenges, this research introduces HyperGef, a dedicated framework designed to address the critical issue of efficient fusion in HyperGNNs when deployed on Graphics Processing Units (GPUs). Efficient fusion is paramount for achieving optimal performance, and HyperGef leverages innovative strategies to overcome hurdles in existing GNN frameworks.

The primary challenges addressed by HyperGef include inefficient partitioning, the absence of a workload-agnostic format, and the prevalence of heavy writing conflicts during vertex updates. Inefficient partitioning introduces latency and imbalances in workload distribution among parallel workers. The lack of a suitable workload-agnostic format limits the representation of the computation workload after fusion, impeding overall efficiency. Additionally, heavy writing conflicts arise when updating the same vertex, contributing to performance bottlenecks.

HyperGef proposes a novel edge-split partition scheme to enhance efficiency and workload balancing, addressing the challenges associated with partitioning. Furthermore, it introduces a fusion workload-aware format to accurately represent the computation workload post-fusion, overcoming the limitations of existing data formats. To alleviate heavy writing conflicts, HyperGef incorporates a shared memory-aware grouping scheme.

Through algorithmic enhancements, HyperGef demonstrates superior performance, surpassing the NVIDIA cuSPARSE kernel by 3.31× in extensive GPU-based experiments. The framework's efficacy is further demonstrated through a notable 2.25× to 3.99× end-to-end speedup across various HyperGNN models compared to leading frameworks such as DGL and PyG.

This research positions HyperGef as a pivotal contribution to the field, offering a robust solution to the challenges hindering the efficient deployment of HyperGNNs on GPUs. By seamlessly addressing issues of partitioning, data format representation, and writing conflicts, HyperGef opens new avenues for harnessing the full power of HyperGNNs in diverse applications.
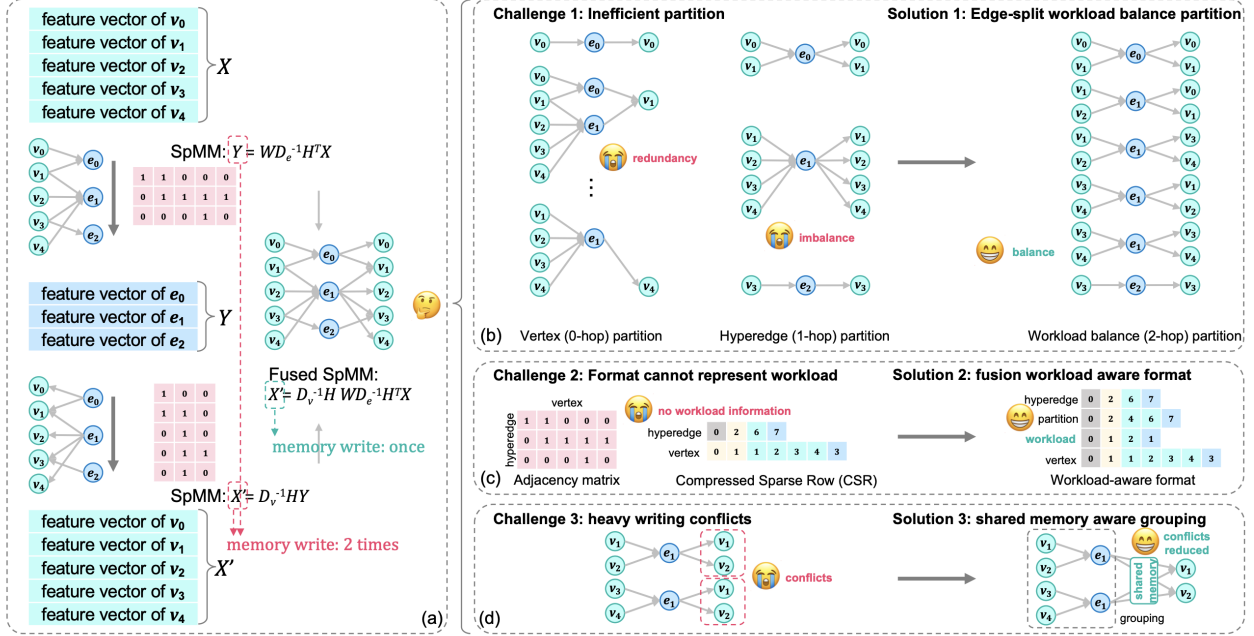
Figure 1. Explaining the challenges (b) Inefficient Partition, (c) Workload-Agnostic Format, (d) Heavy Writing Conflicts, and the approach to resolve them.

## 2. Challenges

### 2.1. Inefficient Partition:

Current GNN frameworks fail to fuse two message-passing steps from vertices to hyperedges and hyperedges to vertices, leading to high latency and redundant memory consumption. Hardware-efficient and workload-balanced partitions are required for parallel workers to process two consecutive message passing steps after fusion.

Overcoming the inefficiencies associated with partitioning is crucial for optimizing the overall performance of Hyper-GNNs, enabling them to leverage the parallel processing capabilities of GPUs effectively. In addressing this challenge, the research introduces an edge-split partition scheme within the HyperGef framework, aiming to enhance efficiency and achieve better workload balancing among parallel workers.

### 2.2. Workload-Agnostic Format:

The inadequacy of existing data formats, such as Compressed Sparse Row (CSR), to effectively represent a two-step computation workload after fusion. HyperGNNs involve intricate relationships among vertices and hyperedges, and conventional data formats may lack the flexibility to capture these complexities accurately. The challenge lies in finding a format that is adaptable to the unique characteristics of HyperGNN computations, ensuring an efficient representation of the workload and avoiding potential limitations that might lead to suboptimal performance on Graphics Processing Units (GPUs). Overcoming this challenge is

crucial for enhancing the overall efficiency and effectiveness of HyperGNNs in modeling high-order relationships.

### 2.3. Heavy Writing Conflicts:

Revolves around the substantial conflicts that arise during the updating of the same vertex, particularly in the context of partitioning. When multiple workers attempt to update a shared vertex simultaneously, heavy writing conflicts occur, leading to performance bottlenecks. This challenge is intrinsic to the concurrent nature of parallel processing in HyperGNNs and can result in increased memory contention and inefficiencies. Overcoming the challenge of heavy writing conflicts is essential for optimizing the efficiency of HyperGNNs on Graphics Processing Units (GPUs), ensuring smoother parallel execution and minimizing contention issues during vertex updates.

## 3. Approach

### 3.1. Edge-Split Workload Balance Partition Scheme

Framework to address the challenge of inefficient partitioning in Hypergraph Neural Networks (HyperGNNs). In the context of parallel computing on Graphics Processing Units (GPUs), efficient partitioning is crucial to ensure that the computational workload is evenly distributed among parallel workers. The "Edge-Split" approach involves a novel partitioning scheme that focuses on the distribution of edges, aiming for improved efficiency and better workload

balancing. By carefully considering the characteristics of hyperedges and vertices, this partitioning scheme optimizes the allocation of tasks among parallel workers during consecutive message-passing steps post-fusion. The goal is to minimize latency, enhance parallelism, and maximize the utilization of GPU resources, ultimately contributing to the overall efficiency of HyperGNNs.

## 3.2. Fusion Workload-Aware Format

A pivotal component introduced in the HyperGef framework to tackle the challenge of representing the computation workload after fusion in Hypergraph Neural Networks (HyperGNNs). In the context of parallel processing on Graphics Processing Units (GPUs), where computational efficiency is paramount, the traditional data formats may fall short in capturing the intricacies of a two-step computation workload.

A novel data representation approach specifically designed to adapt to the unique characteristics of HyperGNN computations. It allows for flexible and efficient storage of information related to the computation workload post-fusion. By being aware of the workload dynamics and requirements, this format aims to optimize data representation, facilitating seamless processing of complex relationships among vertices and hyperedges. This innovative format is instrumental in improving the overall efficiency of HyperGNNs on GPUs by ensuring that the computational workload is accurately and effectively managed.

## 3.3. Shared Memory-Aware Grouping Scheme

A strategic component integrated into the HyperGef framework to address the challenge of heavy writing conflicts in Hypergraph Neural Networks (HyperGNNs). In parallel computing scenarios, particularly on Graphics Processing Units (GPUs), the simultaneous updating of the same vertex by multiple workers can lead to substantial writing conflicts, impeding performance.

Designed to mitigate these conflicts. By introducing an awareness of shared memory dynamics, this scheme facilitates an organized grouping of operations, reducing contention during vertex updates. This approach aims to enhance the overall parallel processing efficiency of Hyper-GNNs by minimizing conflicts, optimizing memory access patterns, and ensuring smoother coordination among parallel workers. The shared memory-aware grouping scheme contributes significantly to the reduction of writing conflicts, ultimately improving the performance and scalability of HyperGNNs on GPU architectures.

## 4. Future Work

### 4.1. Algorithmic Enhancements

Strategic improvements and optimizations applied to algorithms employed in Hypergraph Neural Networks (Hyper-GNNs). These enhancements play a crucial role in overcoming identified challenges and elevating the overall efficiency

of HyperGNNs on Graphics Processing Units (GPUs). In the realm of parallel computing, where computational resources must be utilized optimally, algorithmic enhancements aim to streamline operations, reduce latency, and enhance overall performance.

Within the HyperGef framework, specific algorithmic enhancements are likely implemented to address challenges such as inefficient partitioning, workload-agnostic formats, and heavy writing conflicts. These improvements may involve novel approaches to partitioning strategies, the introduction of adaptive data formats, and optimized algorithms for managing shared memory and reducing conflicts during vertex updates.

Algorithmic enhancements are essential for pushing the performance boundaries of HyperGNNs, ensuring they can effectively model high-order relationships among vertices and hyperedges. By refining the underlying algorithms, HyperGef aims to achieve superior computational efficiency, outperforming existing frameworks and contributing to the advancement of HyperGNNs in GPU-accelerated computing environments.

### 4.2. Advanced Fusion Techniques

Sophisticated methods and strategies employed within the Hypergraph Neural Network (HyperGNN) framework to achieve efficient fusion of information during computation. Fusion is a critical aspect of HyperGNNs, involving the integration of information between vertices and hyperedges to model high-order relationships effectively. Advanced fusion techniques in HyperGef aim to overcome challenges associated with existing Graph Neural Network (GNN) frameworks, particularly in the context of GPU acceleration.

These advanced fusion techniques may involve innovative approaches to parallel processing, memory management, and computation representation. They could include optimized algorithms for fusing information between vertices and hyperedges, strategies for minimizing redundancy and latency, and specialized data structures tailored for efficient GPU utilization.

The goal of integrating advanced fusion techniques is to enhance the overall performance of HyperGNNs by ensuring seamless and optimized information flow between vertices and hyperedges. These techniques contribute to reducing computational bottlenecks, improving parallelism, and ultimately achieving superior efficiency in modeling complex relationships within hypergraphs. The implementation of advanced fusion techniques within HyperGef is a testament to its commitment to pushing the boundaries of HyperGNN computation on GPU architectures.

### 4.3. Parallelization Strategies

Systematic approaches and methodologies employed to leverage parallel computing capabilities, particularly on Graphics Processing Units (GPUs), for optimizing the performance of Hypergraph Neural Networks (HyperGNNs).

Efficient parallelization is crucial for handling the complex computations involved in HyperGNNs, especially when modeling high-order relationships among vertices and hyperedges.

**4.3.1. Load Balancing Techniques.** Systematic approaches and strategies implemented to distribute computational tasks evenly among parallel workers, ensuring optimal resource utilization and performance in Hypergraph Neural Networks (HyperGNNs). Load balancing is a critical aspect of parallel computing, particularly on Graphics Processing Units (GPUs), where efficient utilization of processing power is essential.

**4.3.2. Optimized Message-Passing Algorithms.** Strategic improvements and enhancements applied to the algorithms responsible for information exchange between vertices and hyperedges in Hypergraph Neural Networks (HyperGNNs). Message-passing is a fundamental operation in GNNs, including HyperGNNs, where information is propagated through the network to capture complex relationships.

The term "optimized" suggests that these algorithms within HyperGef have undergone refinements to improve their efficiency, reduce computational complexity, and enhance parallel processing capabilities, particularly when executed on Graphics Processing Units (GPUs).

These parallelization strategies collectively contribute to achieving efficient and scalable execution of HyperGNN computations on GPU architectures. By harnessing the parallel processing capabilities inherent in GPUs, HyperGef aims to overcome the challenges associated with traditional GNN frameworks and significantly enhance the performance of Hypergraph Neural Networks.

## 4.4. Scalability

The ability of the framework to efficiently handle an increasing workload and growing datasets while maintaining or improving its performance. It is a crucial aspect, particularly in the realm of Hypergraph Neural Networks (HyperGNNs) where the complexity of computations and the size of datasets can vary significantly. Achieving scalability in HyperGef is essential for its applicability in diverse real-world scenarios where hypergraphs can vary significantly in size and complexity. By addressing scalability challenges, HyperGef aims to provide a versatile and high-performance framework for HyperGNNs, accommodating a wide range of applications and datasets.

## 4.5. Benchmarking and Comparison

Involves a systematic evaluation and comparison of the framework's performance against established benchmarks and other existing frameworks designed for Hypergraph Neural Networks (HyperGNNs). This process is crucial for assessing the efficiency, effectiveness, and scalability of HyperGef in various scenarios and understanding how it stands relative to other state-of-the-art solutions. By thoroughly benchmarking and comparing HyperGef, researchers and practitioners can gain insights into its capabilities, identify potential areas for improvement, and make informed decisions about its suitability for specific use cases. This process contributes to advancing the state of the art in HyperGNN frameworks and guides further developments in the field.

## 5. Algorithmic enhancements

### 5.1. Training Convergence

Focuses on the stability and progress of the learning process during model training. The primary indicator of training convergence is the decrease in the value of the loss function over time, signifying that the model is effectively learning and moving towards an optimal solution. Monitoring the impact of hyperparameters, such as the learning rate, and observing the model's performance on a separate validation dataset are essential components of assessing convergence. Additionally, techniques like early stopping and regularization methods contribute to stable training convergence by preventing overfitting and ensuring the model generalizes well to unseen data. Adjusting batch sizes and monitoring mini-batch dynamics also plays a role in influencing training stability. Achieving stable training convergence is crucial for developing HyperGNN models that not only learn from the training data but also generalize well to new, unseen data, contributing to the model's overall effectiveness and reliability.

### 5.2. Parallelization and GPU Utilization

Strategic deployment of parallel processing techniques, especially on Graphics Processing Units (GPUs), to enhance the efficiency and speed of computations. Leveraging parallelization and GPU utilization is crucial for handling the complex computations associated with HyperGNNs and accommodating the high-dimensional and interconnected nature of hypergraphs.

In HyperGef, parallelization strategies are implemented to distribute computational tasks among parallel workers effectively, exploiting the parallel processing capabilities of GPUs. This includes optimizing algorithms for parallel execution, load balancing techniques, and efficient memory management to ensure that the computational workload is evenly distributed and processed concurrently.

GPU utilization, in particular, focuses on harnessing the computational power of GPUs to accelerate HyperGNN computations. This involves structuring algorithms and computations to take advantage of the parallel architecture of GPUs, minimizing data transfer between the CPU and GPU, and optimizing memory access patterns for efficient GPU processing.

The combined efforts in parallelization and GPU utilization aim to significantly improve the overall performance and scalability of HyperGNNs. By capitalizing on the parallel processing capabilities of GPUs, HyperGef enhances

computational efficiency, reduces latency, and enables the effective modeling of high-order relationships within hypergraphs, contributing to advancements in the field of hypergraph-based machine learning.

## 5.3. Generalization and Robustness

Achieving generalization and robustness represents pivotal objectives for ensuring the model's effectiveness in real-world applications. Generalization refers to the model's capacity to apply learned patterns from the training data to make accurate predictions on new, unseen instances. To promote generalization, diverse datasets are essential, enabling the model to capture a broad range of hypergraph structures. Regularization techniques, such as dropout and weight regularization, play a crucial role in preventing overfitting and enhancing generalization by avoiding the learning of noise in the training data. Additionally, optimizing hyperparameters and employing cross-validation techniques contribute to striking a balance that fosters both generalization and robustness.

Robustness, on the other hand, signifies the model's ability to maintain performance even in the presence of perturbations or variations in the input data. Evaluating the model's resilience against adversarial attacks or input disturbances is essential to ensure its robust performance in dynamic real-world scenarios. Moreover, techniques like transfer learning, where pre-trained models or representations are leveraged, enhance generalization by allowing the model to apply knowledge gained from one task to another, contributing to a more adaptable and robust HyperGNN. Striking the right balance between generalization and robustness is crucial for the successful deployment of HyperGNNs across diverse and dynamic applications.

## 5.4. Computational Efficiency

Achieving computational efficiency is a fundamental objective that revolves around optimizing the model's ability to perform computations with maximum effectiveness while minimizing resource utilization. This efficiency is crucial for the practical deployment of HyperGNNs in real-world scenarios where computational resources are often constrained.

Central to computational efficiency is the strategic use of parallel processing techniques, particularly on Graphics Processing Units (GPUs), allowing the model to distribute and process computations concurrently. This harnesses the parallel architecture of GPUs, significantly accelerating both the training and inference processes. Complementary to this, the implementation of optimized algorithms plays a key role. These algorithms are carefully designed to be computationally efficient, reducing unnecessary computations and optimizing data structures for streamlined information processing.

Effective memory management is another critical aspect, ensuring efficient utilization and minimizing memory overhead. This capability enables HyperGNNs to handle large-scale hypergraphs without excessive memory requirements.

Additionally, batch processing strategies are employed to simultaneously process multiple data points, enhancing parallelism and computational speed.

Hardware acceleration, such as leveraging GPUs or specialized accelerators, further enhances computational performance, accelerating both model training and inference. Algorithmic enhancements, encompassing efficient message-passing and fusion strategies, contribute to streamlined computations, reducing latency and enhancing overall efficiency.

## 6. Procedure

### 6.1. Incorporate Domain-Specific Knowledge

The integration of domain-specific knowledge involves a strategic infusion of pertinent insights and information from the application domain into the model development and training processes. This approach aims to enrich the model's understanding by incorporating relevant patterns, structures, and semantics unique to the specific domain, thereby enhancing its predictive capabilities.

A key facet of this integration lies in feature engineering, where domain-specific features are carefully incorporated to enrich the representations of vertices and hyperedges within the hypergraph. This inclusion ensures that the model can leverage valuable information about entities and relationships in a manner tailored to the nuances of the domain.

Furthermore, the construction of the hypergraph is fine-tuned to align with domain-specific relationships and semantics. This customization ensures that the model captures connections that are meaningful and reflective of the real-world domain it seeks to model.

The utilization of semantic embeddings, either pre-existing or learned from domain-specific knowledge sources, contributes to initializing or refining the embeddings of vertices and hyperedges. This initialization process ensures that the model starts with a foundation rooted in domain-specific semantics.

Integrating domain-specific constraints and rules during training is another crucial element. This practice guides the learning process, ensuring that the model adheres to known characteristics of the domain, enhancing its alignment with real-world dynamics.

Moreover, the emphasis on building interpretable models reinforces the alignment with domain-specific knowledge. Models that offer interpretability allow practitioners to comprehend and validate decisions within the context of the specific application, fostering trust and confidence.

Collaboration with domain experts remains a continuous and integral aspect of the process. Actively involving experts from the specific domain ensures that the model benefits from their insights, validation of assumptions, and contributions to refining the model's structure and parameters. This collaborative approach not only enhances the model's performance but also ensures its relevance and applicability in addressing complex challenges within the targeted domain.
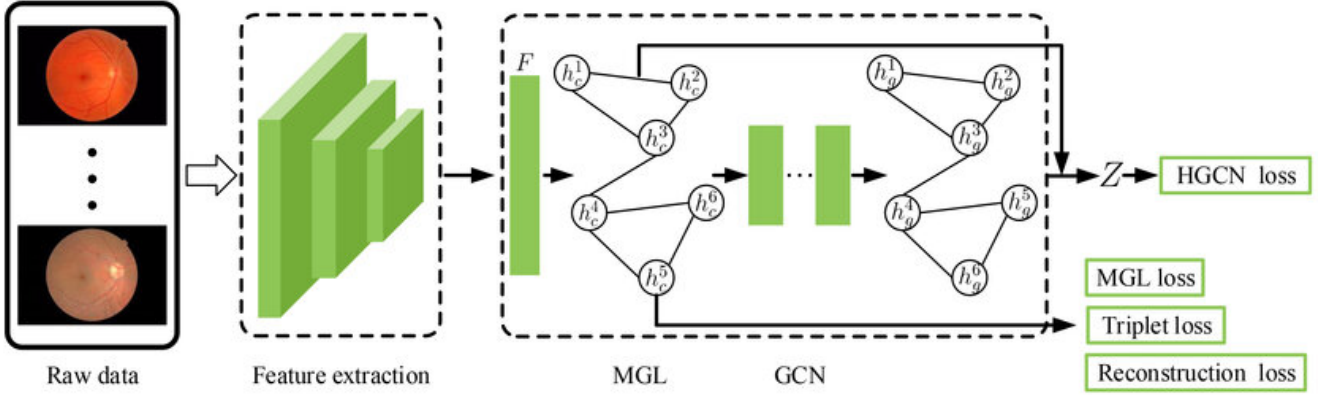
Figure 2. An example of GCN hybrid model

## 6.2. Hybrid Models

A sophisticated approach that involves integrating various machine learning or neural network architectures within a unified framework. These models go beyond the confines of a singular modeling paradigm, combining the strengths of traditional graph-based models and neural network structures to create a synergistic and powerful solution.

A key aspect of hybrid models involves seamlessly blending HyperGNNs with traditional graph-based models, harnessing their efficiency in capturing intricate graph structures and relationships. Additionally, these hybrid models incorporate elements from traditional neural networks, leveraging their prowess in feature learning and representation across non-graph domains. This integration facilitates a holistic approach, allowing the model to excel in diverse data domains.

One notable advantage of hybrid models is their ability to perform transfer learning across different modalities or data types. This feature enables the model to generalize knowledge acquired in one domain to another, enhancing its adaptability and versatility.

The flexibility and adaptability of hybrid models are crucial considerations, allowing for the seamless incorporation of diverse components based on the specific characteristics of the data and the unique challenges posed by the problem at hand. This adaptability ensures that hybrid models can effectively address a wide range of tasks and datasets.

Moreover, the combined expressive power of different model architectures in hybrid models enhances their capability to capture complex relationships and patterns within hypergraphs. This collective strength contributes to the model's overall efficacy in handling intricate data structures.

Balancing computational efficiency with expressive power is a critical aspect of hybrid models. By optimizing performance, these models achieve an equilibrium between the computational efficiency of traditional graph-based models and the expressive capabilities of neural networks.

Hybrid models in HyperGNNs also have the capacity to incorporate domain-specific components or models, tailoring the solution to the specific characteristics of the hypergraph or application domain. This adaptability makes hybrid models a versatile and effective choice for addressing the intricacies of diverse real-world scenarios.

## 6.3. Interactive Learning and Human-in-the-Loop

Signifies a dynamic and collaborative approach to model development. This strategy acknowledges the invaluable expertise that humans bring to the table and actively incorporates their insights and feedback into the learning process.

Domain experts play a central role in this approach, contributing throughout the model development journey. Their involvement spans validating assumptions, refining the model's structure, and providing key insights that stem from their deep understanding of the application domain.

The concept extends beyond mere human involvement, encompassing interactive training sessions where experts can actively intervene, offer feedback, and guide the learning trajectory. This creates a symbiotic relationship between machine intelligence and human expertise, fostering a collaborative learning environment.

A crucial aspect of interactive learning and human-in-the-loop in HyperGNNs involves the dynamic construction of the hypergraph. Human experts contribute their knowledge about relationships, relevance, and context, shaping the hypergraph in real-time based on their domain expertise.

Real-time feedback mechanisms empower experts to assess model predictions, correct errors, and influence the learning process during both training and inference. The model, in turn, adapts its behavior based on continuous interactions, ensuring ongoing refinement and improvement.

Ensuring explainability and interpretability of the model facilitates effective communication between the model and human experts. This transparency enhances mutual understanding, fostering trust and collaboration between the two entities.

The iterative nature of model development in this paradigm ensures a constant feedback loop. Models are refined iteratively based on the insights and feedback provided by human experts, resulting in a continuous improvement process.

The integration of interactive learning and human-in-the-loop methodologies transforms HyperGNNs into adaptive, context-aware models that leverage the strengths of both machine intelligence and human expertise. This collaborative approach not only enhances model performance but also deepens the understanding of complex data domains through the synergy of human and artificial intelligence.

## 7. Results

To enable efficient fusion for HyperGNNs, we present HyperGef. HyperGef proposes an edge-split partition scheme to achieve higher efficiency and better workload balancing. To represent the workload after fusion and partition, HyperGef introduces a novel fusion workload aware format. HyperGef also introduces a shared memory- aware grouping scheme to reduce writing conflicts. Extensive experiments demonstrate that our fused kernel outperforms the NVIDIA cuSPARSE kernel by 3.31×. By enabling efficient fusion for HyperGNNs, HyperGef achieves 2.25× to 3.99× end-to-end speedup on various HyperGNN models compared with state-of-the-art frameworks like DGL and PyG.

## 8. Conclusion

In conclusion, the advancements in Hypergraph Neural Networks (HyperGNNs) reflect a multifaceted evolution driven by innovative strategies and considerations. Incorporating domain-specific knowledge has emerged as a pivotal element, where the integration of insights from the application domain enriches model understanding and promotes more meaningful predictions. Hybrid models, blending traditional graph-based structures with neural networks, showcase a versatile approach, harnessing the strengths of each to address diverse challenges within hypergraph data.

Moreover, the embrace of interactive learning and human-in-the-loop methodologies introduces a collaborative dimension to model development. By actively involving human expertise, these approaches create adaptive models that dynamically refine their understanding based on real-time feedback and insights. This interactive process not only enhances model performance but also deepens the collaborative synergy between machine intelligence and human domain knowledge.

Efficiency considerations play a crucial role in the advancement of HyperGNNs, with parallelization and GPU utilization strategies optimizing computational resources. Additionally, the focus on scalability ensures that Hyper-GNNs can handle growing workloads and datasets effectively, contributing to their applicability in a variety of real-world scenarios.

Furthermore, the emphasis on generalization and robustness ensures that HyperGNNs can transcend the confines of training data, making them more resilient to variations and applicable in dynamic environments. The incorporation of algorithmic enhancements, advanced fusion techniques, and optimized message-passing algorithms collectively contributes to achieving these desirable qualities.

The collective progress in HyperGNNs involves a holistic approach, encompassing domain-specific knowledge integration, hybrid modeling, interactive learning, computational efficiency considerations, and a commitment to scalability, generalization, and robustness. These advancements position HyperGNNs as powerful tools capable of effectively modeling complex relationships within hypergraph-structured data and hold promise for a wide array of applications in various domains. The continual exploration and integration of these improvements pave the way for Hyper-GNNs to become increasingly adaptable, collaborative, and impactful in addressing the challenges of modern machine learning and artificial intelligence.

## References

[1] Figure 1. HYPERGEF: A FRAMEWORK ENABLING EFFICIENT FUSION FOR HYPERGRAPH NEURAL NETWORK ON GPUS.Zhongming Yu, Guohao Dai, Shang Yang, Genghan Zhang, Hengrui Zhang, Feiwen Zhu, Jun Yang, Jishen Zhao, Yu Wang.

[2] Figure 2. Internet Sources

[3] Jing Huang and Jie Yang. 2021. UniGNN: a Unified Framework for Graph and Hypergraph Neural Networks. In Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21.

[4] Jie Zhao, Xiong Gao, Ruijie Xia, Zhaochuang Zhang, Deshi Chen, Lei Chen, Renwei Zhang, Zhen Geng, Bin Cheng, and Xuefeng Jin. 2022. Apollo: Automatic partition-based operator fusion through layer by layer optimization. Proceedings of Machine Learning and Systems 4 (2022), 1–19.

[5] Qiang Fu, Yuede Ji, and H Howie Huang. 2022. TLPGNN: A lightweight two-level parallelism paradigm for graph neural network computation on GPU. In Proceedings of the 31st International Symposium on High-Performance Parallel and Distributed Computing. 122–134.

[6] Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. 2020a. Simple and deep graph convolutional networks. In International Conference on Machine Learning. PMLR, 1725–1735.

[7] Stephen Chou, Fredrik Kjolstad, and Saman Amarasinghe. 2018. Format Abstraction for Sparse Tensor Algebra Compilers. Proc. ACM Program. Lang. 2, OOPSLA, Article 123 (oct 2018), 30pages. https://doi.org/10.1145/3276493

[8] Yihe Dong, Will Sawin, and Yoshua Bengio. 2020. HNHN: hy- pergraph networks with hyperedge neurons. arXiv preprint arXiv:2006.12278 (2020).

[9] Guohao Dai, Guyue Huang, Shang Yang, Zhongming Yu, Hengrui Zhang, Yufei Ding, Yuan Xie, Huazhong Yang, and Yu Wang. 2022. Heuristic Adaptability to Input Dynamics for SpMM on GPUs. arXiv preprint arXiv:2202.08556 (2022).

[10] Matthias Fey and Jan Eric Lenssen. 2019. Fast graph repre- sentation learning with PyTorch Geometric. arXiv preprint arXiv:1903.02428 (2019).

[11] Yifan Feng, Haoxuan You, Zizhao Zhang, Rongrong Ji, and Yue Gao. 2019. Hypergraph neural networks. In Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 33. 3558–3565.

[12] Marco Maggioni and Tanya Berger-Wolf. 2013. AdELL: An adaptive warp-balancing ELL format for efficient sparse matrix- vector multiplication on GPUs. In 2013 42nd international conference on parallel processing. IEEE, 11–20.

[13] Junliang Yu, Hongzhi Yin, Jundong Li, Qinyong Wang, Nguyen Quoc Viet Hung, and Xiangliang Zhang. 2021. Self-Supervised Multi-Channel Hypergraph Convolutional Network for Social Recommendation. In Proceedings of the Web Conference 2021 (Ljubljana, Slovenia) (WWW '21). Association for Computing Machinery, New York, NY, USA, 413–424. org/10.1145/3442381.3449844