

λ Grapher: A Resource-Efficient Serverless System for GNN Serving through Graph Sharing

Haichuan Hu

Huazhong University of Science
and Technology
Wuhan, China
huhc@hust.edu.cn

Yongjie Yuan

Huazhong University of Science
and Technology
Wuhan, China
jayayuan@hust.edu.cn

Fangming Liu*

Peng Cheng Laboratory
Huazhong University of Science
and Technology
fangminghk@gmail.com

Zichen Xu

Nanchang University
Nanchang, China
xuz@ncu.edu.cn

Qiangyu Pei

Huazhong University of Science
and Technology
Wuhan, China
peiqiangyu@hust.edu.cn

Lin Wang

Paderborn University
Paderborn, Germany
lin.wang@uni-paderborn.de

ABSTRACT

Graph Neural Networks (GNNs) have been increasingly adopted for graph analysis in web applications such as social networks. Yet, efficient GNN serving remains a critical challenge due to high workload fluctuations and intricate GNN operations. Serverless computing, thanks to its flexibility and agility, offers on-demand serving of GNN inference requests. Alas, the request-centric serverless model is still too coarse-grained to avoid resource waste.

Observing the significant data locality in computation graphs of requests, we propose λ Grapher, a serverless system for GNN serving that achieves resource efficiency through graph sharing and fine-grained resource allocation. λ Grapher features the following designs: (1) adaptive timeout for request buffering to balance resource efficiency and inference latency, (2) graph-centric scheduling to minimize computation and memory redundancy, and (3) resource-centric function management with fine-grained resource allocation catered to the resource sensitivities of GNN operations and function orchestration optimized to hide communication latency. We implement a prototype of λ Grapher based on the representative open-source serverless platform Knative and evaluate it with real-world traces from various web applications. Our results show that λ Grapher can achieve an average savings of 61.5% in memory resource and 47.2% in computing resource compared with the state of the arts while ensuring GNN inference latency.

*Corresponding author: Fangming Liu (fangminghk@gmail.com). H. Hu, Q. Pei, and Y. Yuan are with the National Engineering Research Center for Big Data Technology and System, the Services Computing Technology and System Lab, Cluster and Grid Computing Lab in the School of Computer Science and Technology, Huazhong University of Science and Technology, 1037 Luoyu Road, Wuhan, China. F. Liu is with Peng Cheng Laboratory, and Huazhong University of Science and Technology, China. Z. Xu is with School of Mathematics and Computer Science, Nanchang University, China. L. Wang is with Paderborn University, Germany.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

WWW '24, May 13–17, 2024, Singapore, Singapore

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 979-8-4007-0171-9/24/05...\$15.00
<https://doi.org/10.1145/3589334.3645383>

CCS CONCEPTS

• Computer systems organization → Cloud computing.

KEYWORDS

serverless computing, graph neural networks, model serving

ACM Reference Format:

Haichuan Hu, Fangming Liu, Qiangyu Pei, Yongjie Yuan, Zichen Xu, and Lin Wang. 2024. λ Grapher: A Resource-Efficient Serverless System for GNN Serving through Graph Sharing. In *Proceedings of the ACM Web Conference 2024 (WWW '24), May 13–17, 2024, Singapore, Singapore*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3589334.3645383>

1 INTRODUCTION

Graphs, as a fundamental data structure, are prevalent in various domains including social networks [31, 52], financial networks [5, 41], and transportation networks [15, 34]. The rise of deep learning has empowered graph neural networks (GNNs) to be a powerful tool to extract features from graph structures [45, 46, 56]. Today, GNNs have been widely used in online web services, e.g., social network analysis [9, 23], short-video recommendation [26, 53], shopping recommendation [40, 50], and financial fraud detection [28, 41].

However, efficient serving of GNNs—running GNNs for time-sensitive inference tasks—remains a critical challenge, for the following reasons: (1) GNN inference is resource-hungry due to the large graph size and complex operations, while applications impose stringent service-level objectives (SLOs) on GNN inference latency [54]. (2) The arrival of GNN inference requests in web services is typically bursty and hard to predict [51]. (3) GNN execution intricately interleaves graph and tensor operations that show diverging resource sensitivities [38]. The resource inefficiency of GNN deployment leads to high operational costs for web services.

To deal with workload fluctuations, web services typically adopt autoscaling techniques to adjust the provisioned resources vertically and horizontally. Specifically, the system monitors a metric such as the CPU or memory utilization and applies a threshold-based scaling policy [4, 13]. Upon workload increases and the utilization exceeds the threshold, a more powerful service instance (e.g., with more CPU cores or memory) is launched to replace the current one in the case of vertical scaling, or more service instances are added

to serve requests in the case of horizontal scaling. The opposite will be applied when the workload decreases and the utilization drops below the threshold. While such autoscaling techniques can absorb workload variations at large time scales, the long delay in changing the provisioned resources (e.g., launching new virtual machines) limits their capability of handling short-term request spikes.

Serverless computing (and its popular implementation function as a service) offers new opportunities for efficient provisioning of web services thanks to its agile event-driven model [16]. However, a direct request-centric serverless deployment of GNN inference, i.e., invoking a separate function to process each arriving request, as done in financial fraud detection systems on AWS Lambda [6], may not provide us with the promised efficiency gain. There are two major reasons: (1) The fixed resource allocation for a function invocation per request ignores the diverging resource sensitivities of operations in different GNN execution stages, leading to low overall resource utilization. (2) Per-request function innovation leads to repeated computation and redundant memory usage across requests that potentially share parts of their computation graphs.

In this paper, we present λ Grapher, a scalable, resource-efficient serverless system for GNN inference. Our key observation is that GNN inference requests arriving in a given period show high spatial data locality, i.e., their computation graphs overlap significantly. Following this observation, λ Grapher features the following designs to achieve high resource efficiency. First, λ Grapher buffers requests and processes them in batches to exploit the data locality and reduce computation and memory redundancy. As request buffering introduces extra delay, to strike a good balance between resource efficiency and latency, λ Grapher incorporates *adaptive timeout configuration* to decide when the batch of requests in a buffer must be dispatched to avoid latency SLO violation. Second, λ Grapher adopts *graph-centric scheduling* to perform GNN inference computation. Specifically, we use multiple queues and distribute arriving requests to these queues, aiming to maximize the spatial data locality of requests in the same queue. To execute the aggregate computation of batched requests, we merge the computation graphs of all these requests and partition the merged graph accounting for locality so that resources allocated for a partition can be released immediately once the local computation is completed, leveraging the agility of serverless functions. Finally, λ Grapher employs *resource-centric function management* which allocates resources to functions catering to the resource sensitivities of the GNN operations performed by each function and orchestrates functions into a pipeline to reduce inter-function communication time overhead.

In short, this paper makes the following contributions. After conducting a thorough empirical analysis of GNN workload variations, data locality, and resource sensitivities of GNN operations (§2), we

- present a resource-efficient serverless system for GNN inference (§3) featuring an adaptive timeout mechanism for request buffering to balance resource efficiency and end-to-end latency.
- propose a graph-centric request scheduler that exploits data locality to minimize computation and memory redundancy and maximize resource elasticity.
- introduce a resource-centric function manager that caters the resource allocation to the specific resource sensitivities of GNN operations and orchestrates functions in pipelines to reduce inter-function communication latency.

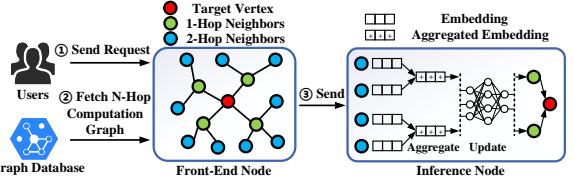
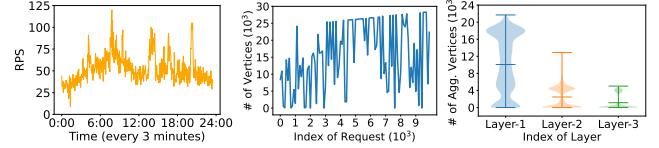


Figure 1: Typical GNN inference workflow.



(a) Request-Level (b) Graph-Level (c) Layer-Level
Figure 2: Workload fluctuations in GNN inference.

- implement λ Grapher on the serverless platform Knative. Our evaluation with real-world traces shows that λ Grapher achieves an average savings of 61.5% in memory resource and 47.2% in computing resource when compared to the state of the arts (§4).

§5 discusses related work and §6 concludes the paper.

2 BACKGROUND AND MOTIVATION

This section describes the fundamentals of GNN and the inference workflow in current systems, empirically studies the workload fluctuations of GNN inference, motivates a graph-centric serverless approach for GNN inference, and discusses the challenges in building an efficient graph-centric serverless system for GNN inference.

2.1 Fundamentals of GNN Inference

GNN basics. Denote the input graph as $G = (V, E)$, where V is the set of vertices representing specific entities and E is the set of edges representing relationships between entities. Each vertex $v \in V$ has a feature representation $h_v \in \mathbb{R}^d$, where d is the feature dimension. A GNN contains multiple layers, each comprising Aggregate and Update operations. In each layer, every vertex v aggregates information from its neighboring vertices with

$$h_v^{l+1} := \Phi^l \left(\{h_u^l : u \in N(v)\} \right), \quad (1)$$

where h_v^{l+1} is the representation of vertex v in layer $l + 1$, $N(v)$ is the set of neighbors of vertex v , h_u^l is the representation of neighboring vertex u in layer l , and Φ^l is the aggregation function. The representation of each vertex v is updated after each layer l with

$$h_v^{l+1} := Y^l(h_v^{l+1}, h_v^l), \quad (2)$$

where the update function Y^l typically includes neural network layers used to integrate information from the current layer and the previous layer, resulting in a new representation for the vertex.

GNN inference workflow. GNN inference has been employed by various time-sensitive online services, such as GraphLearn [1] and PlatoGL [26]. Figure 1 shows a typical GNN inference workflow. First, the request content is extracted as vertices and edges. Next, the platform sets the vertex to predict as the target vertex and extracts an n -hop computation graph. Then, the feature vectors are extracted following the vertices/edges in this graph. Finally, this graph and feature vectors are used as inputs for inference.

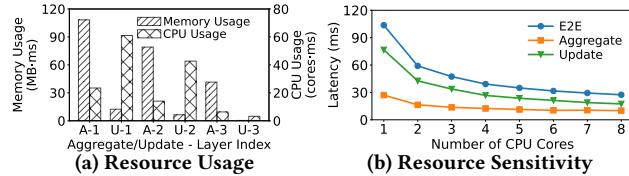


Figure 3: Varying resource sensitivity of GNN operations.
The experiment is conducted with a 3-layer GCN.

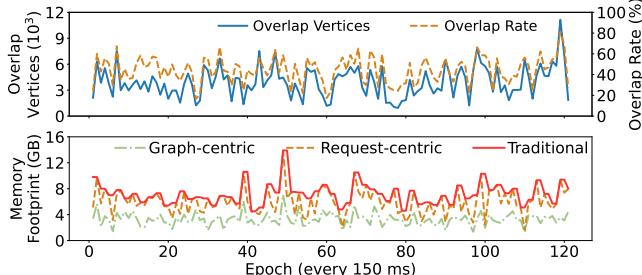


Figure 4: Computation graph overlap among requests over a period and comparative resource consumption analysis of traditional, request-centric, and graph-centric approaches.

2.2 Resource Inefficiency in Current Systems

Current GNN inference systems fall into two types: traditional elastic cloud systems and request-centric serverless systems. The former has pre-configured resources and applies autoscaling in coarse grains based on monitoring metrics, as explained before. Examples of this type include Alibaba’s GraphLearn [1] and Tencent’s PlatoGL [26]. Request-centric serverless systems handle each request by triggering a function invocation, allowing on-demand processing based on the specific computation graph of the request. AWS’s financial fraud detection system operates in this way [6]. Unfortunately, both types of systems suffer from resource inefficiency for one or both of the following two reasons.

Multi-scale workload fluctuations in GNN inference. Using widely recognized datasets of user request arrival traces from Twitter [2] and datasets of requests on social network graphs from Twitter [8], we show that the GNN inference workload fluctuates at three levels: request, graph, and layer. Request-level fluctuations are represented by burstiness in the user request intensity, measured by requests per second (RPS), as shown in Figure 2a. Graph-level fluctuations concern the size of the extracted computation graph of each request. We use a typical setup of a 3-hop computation graph from the target vertex for real-time inference and compare the graph size difference between any two consecutively arriving requests. Figure 2b shows the difference can be as large as 98.6 \times . Layer-level fluctuations are represented by the difference in the number of vertices at each GNN layer, demanding varying resources to perform computation. Figure 2c shows that this difference can reach 4 \times between Layers 1 and 2 and 9 \times between Layers 1 and 3.

Varying resource sensitivity of GNN operations. Each GNN layer is composed of two main operations alternatively executed: Aggregate and Update. Figure 15 (in Appendix A) shows the the structure of three classic GNN layers, namely GCN [18], GraphSAGE [11], and GIN [48]. Taking a 3-layer GCN model as an example, we investigate the demands and sensitivities of Aggregate and Update to different resource types. Figure 3a shows that Aggregate, a graph-based operation, is memory-bound, whereas Update,

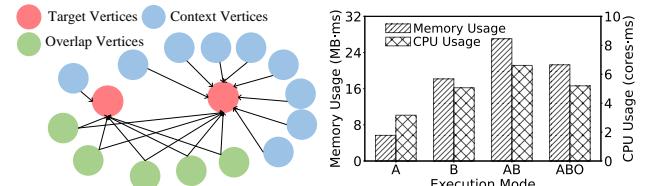


Figure 5: Resource utilization under different execution modes. “A” and “B” represent individual executions, “AB” represents batch processing without sharing, and “ABO” represents batch processing with sharing exploiting the overlap.

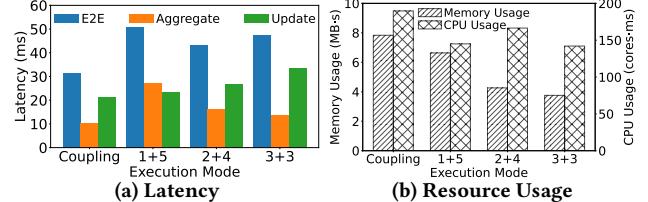


Figure 6: Latency and resource comparison between decoupled and coupled groups. Decoupled group “i+j” means allocating i CPU cores to Aggregate and j CPU cores to Update. a tensor-based operation, is CPU-bound. Figure 3b shows that with the increase of CPU cores, Aggregate shows a continuous latency reduction (up to 4.5 \times) while the latency for Update quickly plateaus with a maximum reduction of 2.2 \times . This implies that Aggregate is more sensitive to CPU resource than Update.

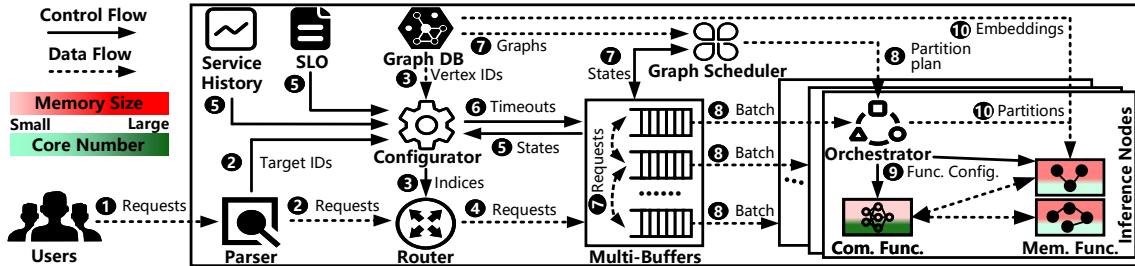
Existing systems consider request-level workload fluctuations at best, but none of them consider multi-scale workload fluctuations and varying resource sensitivities of GNN operations.

2.3 New Opportunities

The above analysis motivates us to switch from the request-centric serverless design to a *graph-centric* one. This design choice offers the following new opportunities.

Exploiting data locality for graph sharing. Based on the Twitter trace [2, 8], we observe a significant overlap between the computation graphs of requests arriving within a period. Figure 4 shows that the overlap rate can reach 44.2% for epochs of 150 ms, leading to considerable redundant computation and memory usage, which can be avoided by batching requests and sharing intermediate results across requests [46]. We show in Figure 4 that a graph-centric serverless approach could save, on average, 55.3% and 46.5% memory resource compared with the traditional and request-centric serverless approaches, respectively. Figure 5 shows the resource consumption of two consecutive requests under different execution modes. It shows that batching requests and eliminating redundancy reduces 21.3% of memory usage and 22.7% of CPU usage.

Decoupling GNN operations for fine-grained resource allocation. The sensitivity of Aggregate and Update to resources differs, suggesting a resource-centric approach to function management. Specifically, we can manage functions in resource groups, decoupling memory-sensitive Aggregate and compute-sensitive Update and customizing fine-grained resource allocation for each of them. Figure 6 shows that with this approach up to 52% memory reduction and 25% CPU reduction can be achieved (see the “3+3” mode). On the other hand, we pay the cost of slight latency increases, primarily caused by the inter-function communication overhead.

Figure 7: A system overview of λ Grapher.

2.4 Design Challenges

The graph-centric serverless approach with fine-grained resource allocation offers tremendous benefits, but also raises challenges.

C1: How to batch requests to exploit data locality? Request batching is a de-facto optimization in inference serving systems for improving resource efficiency. However, due to the heterogeneity of request graphs and irregular memory access in the Aggregate operation in GNN inference (see Figure 2b), batch processing can be inefficient if not treated carefully. As we have shown significant resource efficiency improvement can be achieved by reusing intermediate results among batched requests. The challenge is on quickly grouping requests to maximize the chance of reuse.

C2: How to efficiently execute batched requests? When batching requests, the computation graphs of these requests are merged into a big graph, e.g., with millions of vertices. The memory needed to host the merged graph can easily exceed the memory limit of serverless functions, leading to scalability concerns. A quick idea is to break down the merged graph into pieces and allocate a function for each piece. The challenge is on partitioning the merged graph at a suitable granularity to ensure scalability and take advantage of the agility of serverless functions to achieve resource efficiency.

C3: How to conceal inter-function communication overhead?

Decoupling GNN operations and enabling fine-grained resource allocation offers efficiency gains, but at the cost of extra inter-function communication overhead. One typical approach is to construct a pipeline to overlap function execution with communication. The challenge is to fine-tune this pipeline so that all the functions in the pipeline achieve load balancing to maximize overhead hiding.

3 SYSTEM DESIGN

We present λ Grapher and its design in this section.

3.1 System Overview

To address the shortcomings of the request-centric serverless service model discussed in Section 2.2, we develop a resource-efficient serverless GNN inference system with a graph scheduling and resource management engine. The main idea behind the engine lies in two aspects: (1) graph-centric scheduling which leverages the graph sharing of consecutively arriving requests to reduce computation and memory redundancy, and (2) resource-centric function management which involves fine-grained resource allocation for functions in the form of compute function groups and memory function groups, catering to the compute-sensitive and memory-sensitive operations, respectively, thereby maximizing resource efficiency. λ Grapher aims to optimize the resource efficiency during GNN serving while ensuring the latency SLOs of GNN requests.

Figure 7 illustrates the system overview of λ Grapher. At the beginning of the GNN serving, (1) a continuous stream of user requests arrives at the serverless platform. Then, (2) the *Parser* analyzes the content of the user requests, and the target IDs are dispatched to the *Configurator* and the *Router*. Next, (3) the *Configurator* queries the vertex IDs of the computation graph from the graph database and generates the data indices to send to the *Router*. According to the data indices, (4) the *Router* routes incoming requests to the buffer with the highest degree of graph sharing. While the requests are waiting, (5) the *Configurator* collects the buffer states and queries the built-in SLO and the service history to (6) periodically adjust the timeouts. As the requests are added to the *Multi-Buffers*, (7) the *Graph Scheduler*, with a global perspective, re-schedules the requests to enhance the data locality, and extracts the computation graphs, performing dynamic graph partitioning on each buffer. When a buffer times out or is full, (8) the batched requests are sent to the newly created *Orchestrator*. The *Orchestrator*, based on the graph partitions, (9) scales and maps the workloads to the compute functions and memory functions. (10) The compute functions load the neural network, while the memory functions load the graph partitions and features. Finally, the functions perform collaborative inference as per the orchestrated process.

3.2 Parser and Router

Target IDs. The *Parser* analyzes user requests to retrieve graph structure IDs for inference. Graph analysis tasks can be categorized into three types: vertex/edge/graph-level prediction. In this paper, vertex-level prediction task is taken as an example, where the ID of the vertex to be inferred is the target ID.

Routing strategy. The *Router* is responsible for routing each request to the buffer that has the highest graph-sharing degree for that request to enhance data locality. Each request corresponds to an n -hop computation graph $G(V, E)$ based on its target ID. The data index U_{r_i} for a request r_i is the vertex set V of its computation graph, denoted as $U_{r_i} = V_{r_i}$. The data index for a buffer b_i is the union of data indices for all requests it contains:

$$U_{b_i} = U_{r_0} \cup U_{r_1} \cup \dots \cup U_{r_j}, r_j \in b_i. \quad (3)$$

The routing strategy directs requests to the buffer with the highest graph sharing degree, determined by intersecting the buffers' data index with the request's data index:

$$S_{b_i}^{r_i} = |U_{b_i} \cap U_{r_i}| / |U_{r_i}|, |U_{r_i}| \neq 0, b_j = \arg \max_B S_B^{r_i}, \quad (4)$$

where $S_{b_i}^{r_i}$ is the graph sharing degree of request r_i with respect to buffer b_i in Multi-Buffers B , and b_j is the buffer with the highest graph sharing degree for the request r_i .

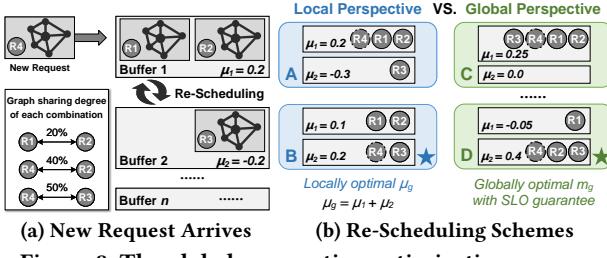


Figure 8: The global perspective optimization process.

3.3 Multi-Buffers and Configurator

Multi-Buffers. We observe a significant overlap among the computation graphs corresponding to user requests arriving continuously over a period as discussed in Section 2.2. Therefore, we design the *Multi-Buffers* which provides requests with an opportunity for graph sharing with other requests which have same subgraphs, by allowing requests to wait in the buffer for a certain period. The *Multi-Buffers*, denoted as B , consists of multiple individual buffers. The batched requests are sent to the functions for inference when the buffer times out or is full. Each buffer possesses a 4-tuple (R, S, Q, K) to characterize the state of the buffer at the current moment, where R denotes the requests per second for the buffer, $S \in [0, 1]$ represents the average graph sharing degree of all requests in the buffer, $Q \in [0, 1]$ indicates the ratio between the remaining time and the configured timeout of the buffer, and $K \in [0, 1]$ represents the ratio between the buffer’s configured timeout and the maximum allowable timeout setting. The buffer timeout is a crucial determinant of system performance. Optimal timeout configuration ensures efficient graph sharing to achieve high resource efficiency without compromising request violations. In the evolving inference service landscape, configuring buffer timeouts judiciously is essential. The batch size is also dynamically adjusted, determined by the resource limits of the functions and SLO slack.

Adaptive timeout configuration. The *Configurator* adapts buffer timeouts dynamically, using the decision tree algorithm [49] to swiftly balance the benefits of graph sharing and inference timeliness, ultimately achieving a comprehensive performance. We employ real-world traces mentioned in Section 4 and utilize the built-in SLO to conduct authentic service runs, thereby gaining service history. First, we determine the initial timeout T_0 and the maximum timeout T_{max} for the buffer based on the SLO:

$$T_0 = \gamma \times SLO, T_{max} = \delta \times SLO, 0 < \gamma < \delta < 1. \quad (5)$$

When the buffer reaches a threshold of new requests or a specific time interval elapses, its state shifts, prompting a decision from the decision set $X = \{0\} \cup \{1 \times \tau, \dots, i \times \tau\}$. Two types of decisions exist: $x_i = 0$ preserves the current timeout, and $x_i = i \times \tau$ extends it, with τ as the unit time interval. To quantify the impact of each decision, we propose a metric that measures the trade-off between graph-sharing benefits and inference timeliness:

$$\mu_{b_i} = \alpha \times S_{b_i} - \beta \times D_{b_i}, \alpha \in [0, 1], \beta \in [0, 1], \quad (6)$$

where μ_{b_i} represents the total performance gain of the buffer b_i , S_{b_i} signifies average graph sharing degree of buffer b_i , D_{b_i} represents the average delay ratio in buffer b_i due to waiting (i.e., average time each request is delayed relative to the timeout), and α and β are fixed coefficients set by the developers. We record the buffer’s state and

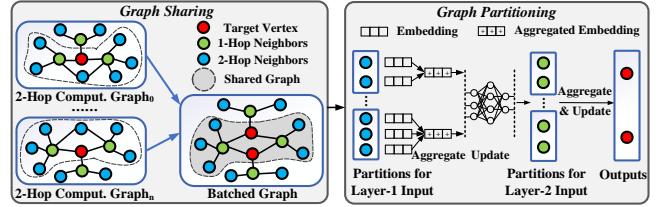


Figure 9: Demonstration of the dynamic graph scheduling.

the decision with the maximum performance gain when a decision is made, utilizing this as historical experience. We employ a decision tree regression model to fit the buffer state as independent variables and the corresponding decisions as the dependent variable:

$$Model = DecisionTreeRegressor.fit((R, S, Q, K), X). \quad (7)$$

The generated *Configurator* brings us some decision-making heuristics: 1) When benefits remain stable, the *Configurator* tends to favor $x_i = 0$, conservatively maintaining the current timeout. 2) If the buffer consistently receives requests that notably improve the average graph sharing degree, the *Configurator* typically chooses $x_i = i \times \tau$, greedily extending the timeout, the degree of which depends on the magnitude of the benefit increase.

3.4 Graph Scheduler

The *Graph Scheduler* is responsible for scheduling of the computation graphs corresponding to the requests, which involves three specific parts: 1) Globally schedule requests between the buffers to achieve the optimal graph sharing; 2) Preprocess graphs to reduce computation and memory redundancy through graph sharing; 3) Dynamically partition the computation graphs to improve resource efficiency and provide scalability for inference.

Global perspective optimization. New requests are routed to the buffer with the highest graph sharing degree according to the strategy. However, this can cause graph sharing results to converge towards local rather than global optima, as illustrated in Figure 8. Hence, we introduce a global perspective optimization algorithm to dynamically re-schedule remaining requests for graph sharing, aiming for the global optimum, as demonstrated in Algorithm 1 in Appendix B. Whenever a new request enters the *Multi-Buffers*, the *Graph Scheduler* places this request in the appropriate buffer based on the routing strategy and evaluates the current buffer’s performance gain (Line 1-Line 3). The *Graph Scheduler* identifies requests in other buffers eligible for graph sharing with the new arrival and computes their respective graph sharing degrees (Line 4-Line 7). Next, the *Graph Scheduler* calculates the performance gain if requests are moved in or out of the buffer (Line 8-Line 12). If the performance gain improves after the scheduling, the *Graph Scheduler* adopts this decision by transferring the requests into the buffer of the new request and removing them from their original buffer (Line 13-Line 16). After global perspective optimization, subsequent inference can fully benefit from graph sharing.

Graph sharing. The *Graph Scheduler* uses a hierarchically aggregated computation graph (HAG), building upon prior research [14], to merge redundant vertices and facilitate intermediate result sharing, reducing computational and memory redundancy in batch processing. The process of graph sharing primarily involves three steps: 1) Expand the computation graph of the target vertex into

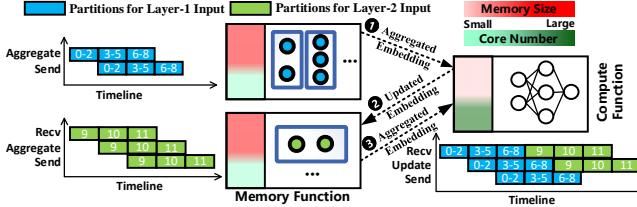


Figure 10: Demonstration of the collaborative inference.

a computation tree; 2) Traversal the computation tree to merge vertices at the same depth between different computation trees; 3) Conduct the aggregation operation on the merged vertices only once, and the intermediate results from the aggregation operation can be reused in subsequent steps, as illustrated in Figure 5a.

Dynamic graph scheduling. When requests move between buffers, their associated computation graph are transferred as well, with dynamic incremental graph partitioning optimizing inference for subsequent tasks. The dynamic graph scheduling serves two main objectives, as shown in Figure 9: 1) Not all graph vertices participate in the computation at every layer, as illustrated in Figure 2c. The *Graph Scheduler* partitions the graph for each GNN layer, utilizing serverless functions efficiently created and destroyed as needed. 2) Serverless functions have resource limitations and cannot load all graphs stored in buffers. Dynamic graph scheduling provides scalability for inference, addressing this constraint. Algorithm 2 in Appendix B describes the specific dynamic graph partitioning process. First, combine the graphs in the buffer with the arrived request graph to generate the HAG, which is the data structure resulting from shared graph scheduling (Line 1). Next, begin traversing from the task vertex to its predecessor vertices (note that even in the case of an undirected graph, it is represented as a directed graph), i.e., the vertices required for its aggregation, which are formed as a partition (Line 2-Line 9). The predecessor vertices visited in the previous iteration are treated as new task vertices for the subsequent traversal, and this process continues until the set of task vertices becomes empty, at which point the algorithm concludes (Line 10-Line 12). Finally, we obtain a two-dimensional list of graph partitions, where each row represents the input for each GNN layer, and the granularity of these graph partitions is fine, providing scalability for subsequent inference.

3.5 Orchestrator

The *Orchestrator* coordinates a set of serverless functions to perform GNN inference on batched requests, as shown in Figure 10, following resource-centric management that maximizes resource efficiency without violating SLOs, which comprises three stages: 1) The *Orchestrator* maps memory-sensitive graph workloads and compute-sensitive tensor workloads to memory functions and compute functions, respectively; 2) The *Orchestrator* employs a pipeline collaborative inference mechanism to distribute communication overhead among functions; 3) Based on the workloads and the remaining time, the *Orchestrator* scales memory functions and compute functions, customizing their resource allocation.

Workload mapping. The *Orchestrator* divides the GNN workload into graph workloads and tensor workloads and manages serverless functions with resource groups. The memory function group exclusively handles graph workloads, i.e., memory-sensitive Aggregate

operations, while the compute function group exclusively loads tensor workloads and handles computation-sensitive Update operations. To optimize resource usage, the *Orchestrator* maps GNN input graph partitions to memory functions, aiming to assign the same layer partitions to a single memory function. Exceeding memory limits for partitions from the same GNN layer prompts mapping excess partitions to new memory function instances. Tensor workloads, needing less memory, load neural networks for each GNN layer into a single compute function instance. Vertices finishing tasks early can exit batch processing and return results.

Collaboration between functions. The *Orchestrator* organizes collaborative inference between functions in a pipeline fashion, allowing the communication overhead between functions to be distributed within their respective computations, as illustrated in Figure 10. The entire pipeline process begins with the memory function inferring the first layer of GNN, and thus, the granularity of concurrent tasks in the pipeline is determined by the number of graph partitions and vertices processed in parallel at each step by the first layer memory function. The concurrent granularity needs to be considered when allocating resources for functions.

Function scaling. The *Orchestrator* customizes resources for memory functions and compute functions based on workload size and concurrent granularity, saving resources while ensuring SLO compliance. Specifically, the allocated memory resource amount for memory function F_i^m and compute function F_i^c are M_i^m and M_i^c :

$$M_i^m = M_r + M_{G_i} + M_{h_i}, M_i^c = M_r + M_{nn}. \quad (8)$$

where M_r represents the runtime memory size, M_{G_i} represents the memory size of loaded graphs, M_{h_i} represents the memory size of loaded embeddings, and M_{nn} represents the memory size of the neural network. the *Orchestrator* allocates the CPU cores to functions based the bayesian optimization [37]:

$$\text{BayesianOptimization}(\vec{F}, \vec{X}, \vec{\Gamma}, \vec{T}_l) \rightarrow \vec{C} \quad (9)$$

$$\text{Minimize : Cost} = \omega \sum F^m \times T_l^m + \eta \sum F^c \times T_l^c \quad (10)$$

$$\text{Constraints : } \sum T_l^m + \sum T_l^c \leq T_{slack} \quad (11)$$

where \vec{F} represents the function vector, \vec{X} represents the function workload size vector, $\vec{\Gamma}$ indicates the concurrent granularity vector, \vec{T}_l represents the inference time vector under different cores and task size, \vec{C} represents the core number vector, and ω and η indicate the cost per unit of memory usage and cpu usage respectively, which can be found in AWS Fargate Pricing [36]. T_l^m and T_l^c indicate the inference time of memory and compute functions, and T_{slack} indicate the SLO slack.

4 EVALUATION

In this section, we prototype λ Grapher and evaluate it with real-world traces from various web applications.

4.1 Experimental Setup

λ Grapher prototype. We prototype λ Grapher based on the open-source serverless platform Knative [19] with 3k LOC in Python and Go. Specifically, we implement the *Parser*, *Configurator*, *Router*, *Multi-Buffers*, *Graph Scheduler*, and *Orchestrator* in a VM instance as middleware between the request source and the Knative platform,

Table 1: Graph Datasets from Real-World Applications

Graph Datasets	$ V $	$ E $	Dim.	N-hop	SLO (s)
Bitcoin OTC [18]	5,881	35,592	1,024	3	0.3
KuaiRec [10]	4,738	4,676,570	58	2	0.6
Higgs Twitter [11]	456,626	14,855,842	1,024	3	1.0

and we deploy function instances through Knative Serving Service. We build an in-memory database service for fast graph queries.

Baselines. We compare λGrapher with the state-of-the-art GNN serving systems, including GraphLearn [1], representing the traditional cloud service architecture, and a financial fraud detection system based on AWS Lambda [6], denoted as AWSGNN, representing the request-centric serverless architecture. GraphLearn relies on monitoring memory occupancy threshold metrics to scale instances up or down, as most traditional elastic cloud services do [4]. AWSGNN dynamically allocates functions for each request based on its computation graph size.

Web application traces. We utilize real-world traces from Twitter [2] to generate the inter-arrival time of user requests, which is widely used for evaluating inference systems. We use three graph datasets from real-world applications to generate request contents, including KuaiRec [10] from the video-sharing mobile app Kuaishou [20], Bitcoin OTC [22] from Bitcoin transaction network, and Higgs Twitter [8] from Twitter network, which are widely applied in evaluating the GNN model designed for short-video recommendation [29], financial fraud detection [21] and social network analysis [35], respectively. The SLOs are set based on the requirements of the application scenario, as described in the previous work [54]. The details of graph datasets are shown in Table 1.

GNN workloads. We select three common GNN models with three layers using the Deep Graph Library (DGL) [42] API, including GCN [18], GraphSAGE [11], and GIN [48]. The structures of GNN layers are shown in Figure 15 in Appendix A. The three-layer model is used to evaluate both the BitCoin and Higgs Twitter traces, while the two-layer model is used to evaluate the KuaiRec trace.

Testbed. We implement λGrapher on a cluster with 10 kc1.8xlarge.2 machines, each of which includes 32 CPU cores at 2.3GHz and 64 GB RAM (Ubuntu 18.04). We collect real service data on physical machines, such as inference latency under various configurations. To expedite the experimental process, we transform the prototype implementation into a simulation mode as in [25].

4.2 Performance Comparison

We compare λGrapher with the state-of-the-art GNN serving systems, GraphLearn and AWSGNN, in terms of memory and computing resource efficiency, as well as end-to-end (E2E) latency and SLO violation rate. Figure 11 indicates that, compared to the state-of-the-art, λGrapher can achieve an average of 61.5% in memory resource and 47.2% in computing resource savings. Across three web application traces, including Bitcoin OTC, KuaiRec, and Higgs Twitter, the average graph sharing degrees of each buffer are 54.6%, 97.5%, and 58.6% respectively. The average graph sharing degree across KuaiRec is close to 1 because of its graph density of 99.6%, where each computation graph approximates the whole graph.

Memory resource efficiency. In different GNN workloads and across various traces, λGrapher reduces memory resource usage by

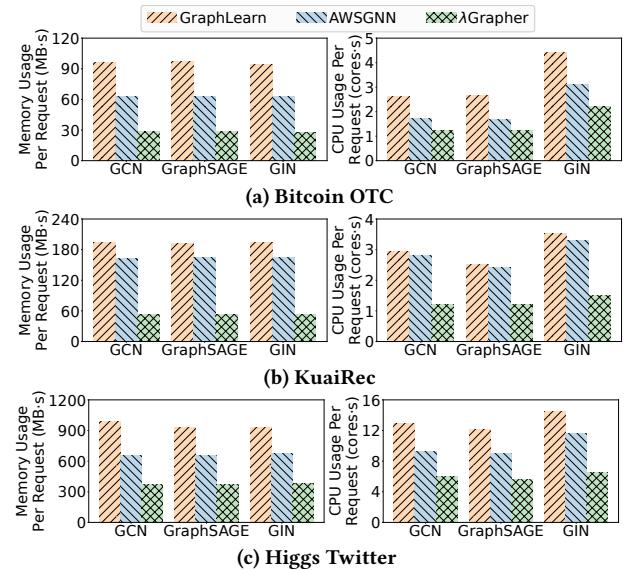


Figure 11: Resource efficiency between λGrapher and the state of the arts under different traces and GNN workloads.

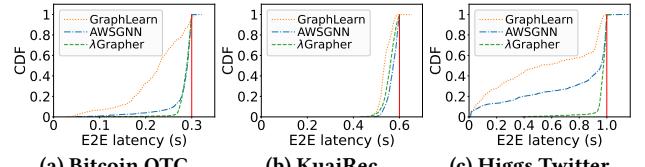
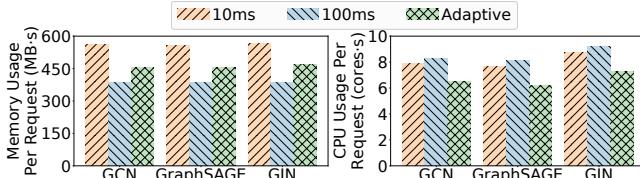


Figure 12: E2E latency between λGrapher and the state of the arts. The red solid lines represent the SLOs.

58.9% to 72.5% compared to GraphLearn and 42.9% to 67.9% compared to AWSGNN (on average 61.5%). GraphLearn, representing traditional cloud-based systems, over-allocates resources before service initiation and continuously monitors for request fluctuations and system availability. However, its coarse instance resource scaling granularity leads to significant memory wastage. On the other hand, AWSGNN, representing request-centric serverless systems, effectively manages request intensity fluctuations but lacks spatial data locality utilization due to serving individual requests with individual functions, resulting in memory redundancy during GNN inference. λGrapher adopts a graph-centric task scheduling approach, efficiently reducing memory redundancy by batching requests with common subgraphs. Besides, λGrapher employs a resource-centric approach, segregating graph workloads from tensor workloads to avoid memory overhead during tensor computations.

Computing resource efficiency. Under various GNN workloads and across different traces, λGrapher demonstrates a reduction in computing resource usage, achieving savings ranging from 49.3% to 57.2% compared to GraphLearn, and 27.8% to 56.9% compared to AWSGNN (on average 47.2%). As shown in Figure 3, GNN operations vary widely in resource sensitivities. Both GraphLearn and AWSGNN employ coarse-grained resource allocation for the entire GNN, resulting in suboptimal computing resource utilization. In contrast, λGrapher reduces computation redundancy through graph sharing and offers a resource-centric function management mechanism. By decoupling Aggregate and Update operations, λGrapher enables fine-grained resource allocation and orchestrates a refined

**Figure 13: Results of the adaptive timeout configuration.**

pipeline for customized functions to ensure load balancing, substantially enhancing computing resource efficiency.

E2E latency and SLO violation rate. Figure 12 depicts the cumulative distribution function (CDF) plot of the E2E latency for each system. Across the three traces, the average SLO violation rates of each system are 1.55%, 0.15%, and 0.09% respectively. GraphLearn’s strategy of over-provisioning resources ensures lower average E2E latency but its coarse-grained and homogeneous resource allocation policy struggles with GNN request fluctuations, resulting in a relatively high violation rate. AWSGNN sacrifices some inference performance to save resources, yet it cannot fully utilize SLO slack. λ Grapher maximizes resource efficiency by analyzing GNN workloads and fully leveraging SLO slack.

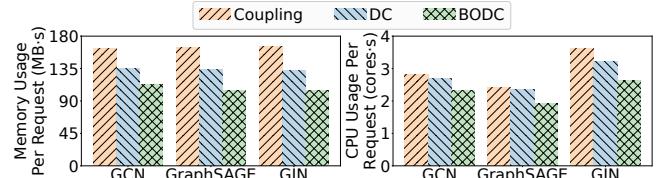
4.3 Module Analysis

Adaptive timeout module. To validate the adaptive timeout module’s performance, we set a fixed lower bound of 10ms and an upper bound of 100ms, allowing λ Grapher to dynamically adjust within this range. We conduct tests on the largest-scale graph datasets Higgs Twitter, as shown in Figure 13. The adaptive timeout module saves an average of 18.1% of memory and 17.4% of computing resources compared to the 10ms configuration, and achieves an average 21.6% reduction in computing resource usage compared to the 100ms configuration. The 10ms configuration overlooks data locality, hampering resource efficiency, while the 100ms timeout, despite optimizing memory resource efficiency through graph sharing, demands significant computing resources to meet SLO goals. The λ Grapher dynamically adjusts the timeout, balancing the benefits of graph sharing and the risk of violating SLOs.

Resource allocation module. We compare the resource allocation module’s performance, utilizing bayesian optimization (BODC), against coupling and decoupling (DC) solutions in KuaiRec. Figure 14 shows that BODC saves an average of 34.5% memory and 21.7% computing resources compared to coupling, and 19.9% memory and 16.6% computing resources compared to DC. Coupling results in wastage in computing resources during Aggregate operations and memory resources during Update operations. Under the DC approach, although separating graph and tensor workloads saves some resources, there is a mismatch in execution speed between memory and compute functions. By analyzing historical data, λ Grapher utilizes bayesian optimization to determine the optimal resource allocation ratio for memory and compute functions, maximizing resource efficiency.

5 RELATED WORK

GNN inference. In the traditional distributed environment, recent works focus on optimizing graph partitioning and resource mapping for acceleration [3, 17, 43, 55]. Wang et al. [43] propose an

**Figure 14: Results of the resource allocation module.**

adaptive and efficient system for GNN acceleration on GPUs, which preprocesses the model and input graph to achieve reasonable graph partitioning and resource mapping. In the cloud environment, in order to solve the problem of graph data distributed in different geographies, Zeng et al. [51] propose to conduct the GNN real-time inference by adopting the fog computing paradigm to reduce the communication overhead of the data collection before inference. The above works focus on inference of static GNN models, which pre-allocate computing node resources and provide services by continuous monitoring. This scheme is difficult to dynamically and adaptively allocate resources according to the fluctuation of user requests, resulting in waste of resources.

Serverless graph system. Due to the elastic scalability and flexibility of serverless computing [16, 27, 44], some web applications have been serverlessized, such as DNN inference or training [24, 33, 47], and IoT services [7, 32]. In particular, some scholars propose to migrate the graph processing system to the FaaS platform [12, 38, 39]. Toader et al. [39] implement the classic graph processing model Pregel [30] on the FaaS platform in a simple engineering manner. However, due to frequent data communication, the system performs poorly in performing large-scale graph algorithms. Thorpe et al. [38] make the GNN training process semi-serverless, introducing serverless threads to handle computation-sensitive tensor operations, while graph operations that are sensitive to memory resources are still executed on the CPU server. At present, there is a gap in the work of serverless-based GNN serving.

6 CONCLUSION

In this paper, we identify the resource inefficiency problem in current GNN serving systems. Through studying the web application traces, we observe the spatial data locality in computation graphs of requests. We propose a scalable, resource-efficient serverless system named λ Grapher for GNN serving. λ Grapher supports a graph-centric task scheduling strategy to reduce the computation and memory redundancy and facilitates a resource-centric function management mechanism which allocates resources to functions catering to the resource sensitivities of GNN fine-grained operations. Compared to the state of the arts, our λ Grapher prototype can save an average of 61.5% in memory resource and 47.2% in computing resource with real-world traces while meeting the SLOs.

ACKNOWLEDGMENTS

This work is supported in part by National Key Research & Development (R&D) Plan under grant 2022YFB4501703, in part by The Major Key Project of PCL (PCL2022A05), and in part by the NSFC under grant 61972158. Lin Wang has been funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – Project-ID 210487104 - SFB 1053.

REFERENCES

- [1] Alibaba. 2020. graph-learn: An Industrial Graph Neural Network. https://graph-learn.readthedocs.io/en/latest/index_en.html [Online Accessed, 12-Feb-2024].
- [2] ArchiveTeam. [n. d.]. Twitter streaming traces, 2017. <https://github.com/rickypinci/BATCH/tree/sc2020/traces> [Online Accessed, 12-Feb-2024].
- [3] Adam Auten, Matthew Tomei, and Rakesh Kumar. 2020. Hardware acceleration of graph neural networks. In *2020 57th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 1–6.
- [4] AWS. [n. d.]. AWS Auto Scaling. <https://aws.amazon.com/cn/autoscaling/> [Online Accessed, 12-Feb-2024].
- [5] Stefano Battiston, Guido Caldarelli, Robert M May, Tarik Roukny, and Joseph E Stiglitz. 2016. The price of complexity in financial networks. *Proceedings of the National Academy of Sciences* 113, 36 (2016), 10031–10036.
- [6] AWS Machine Learning Blog. [n. d.]. Build a GNN-based real-time fraud detection solution using Amazon SageMaker, Amazon Neptune, and the Deep Graph Library. <https://aws.amazon.com/cn/blogs/machine-learning/build-a-gnn-based-real-time-fraud-detection-solution-using-amazon-sagemaker-amazon-neptune-and-the-deep-graph-library/> [Online Accessed, 12-Feb-2024].
- [7] Gustavo André Setti Cassel, Vinicius Facco Rodrigues, Rodrigo da Rosa Righi, Marta Rosecler Bez, Andressa Cruz Nepomuceno, and Cristiano André da Costa. 2022. Serverless computing for Internet of Things: A systematic literature review. *Future Generation Computer Systems* 128 (2022), 299–316.
- [8] Manlio De Domenico, Antonio Lima, Paul Mougel, and Mirco Musolesi. 2013. The anatomy of a scientific rumor. *Scientific reports* 3, 1 (2013), 2980.
- [9] Fabrizio Frasca, Emanuele Rossi, Davide Eynard, Ben Chamberlain, Michael Bronstein, and Federico Monti. 2020. Sigma: Scalable inception graph neural networks. *arXiv preprint arXiv:2004.11198* (2020).
- [10] Chongming Gao, Shijun Li, Wenqiang Lei, Jiawei Chen, Biao Li, Peng Jiang, Xiangnan He, Jiaxin Mao, and Tat-Seng Chua. 2022. KuaiRec: A Fully-Observed Dataset and Insights for Evaluating Recommender Systems. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management* (Atlanta, GA, USA) (CIKM '22). 540–550. <https://doi.org/10.1145/3511808.3557220>
- [11] Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. *Advances in neural information processing systems* 30 (2017).
- [12] Chaoyang He, Emir Ceyani, Keshav Balasubramanian, Murali Annavaram, and Salman Avestimehr. 2021. Spreadgnn: Serverless multi-task federated learning for graph neural networks. *arXiv preprint arXiv:2106.02743* (2021).
- [13] Brendan Jennings and Rolf Stadler. 2015. Resource management in clouds: Survey and research challenges. *Journal of Network and Systems Management* 23 (2015), 567–619.
- [14] Zhihao Jia, Sina Lin, Rex Ying, Jiaxuan You, Jure Leskovec, and Alex Aiken. 2020. Redundancy-Free Computation for Graph Neural Networks. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (Virtual Event, CA, USA) (KDD '20). Association for Computing Machinery, New York, NY, USA, 997–1005.
- [15] Weiwei Jiang and Jiayun Luo. 2022. Graph neural network for traffic forecasting: A survey. *Expert Systems with Applications* (2022), 117921.
- [16] Eric Jonas, Johann Schleier-Smith, Vikram Sreekanth, Chia-Che Tsai, Anurag Khandelwal, Qifan Pu, Vaishala Shankar, Joao Carreira, Karl Krauth, Neeraja Yadwadkar, et al. 2019. Cloud programming simplified: A berkeley view on serverless computing. *arXiv preprint arXiv:1902.03383* (2019).
- [17] Kevin Kiningham, Philip Lewis, and Christopher Ré. 2022. GRIP: A graph neural network accelerator architecture. *IEEE Trans. Comput.* (2022).
- [18] Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* (2016).
- [19] Knative. [n. d.]. Knative is an Open-Source Enterprise-level solution to build Serverless and Event Driven Applications. <https://knative.dev/docs/> [OnlineAccessed, 12-Feb-2024].
- [20] KuaiShou. [n. d.]. KuaiShou is the video-sharing mobile app. <https://www.kuashou.com/en/> [OnlineAccessed, 12-Feb-2024].
- [21] Srijan Kumar, Bryan Hooi, Disha Makhija, Mohit Kumar, Christos Faloutsos, and VS Subrahmanian. 2018. Rev2: Fraudulent user prediction in rating platforms. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*. ACM, 333–341.
- [22] Srijan Kumar, Francesca Spezzano, VS Subrahmanian, and Christos Faloutsos. 2016. Edge weight prediction in weighted signed networks. In *Data Mining (ICDM), 2016 IEEE 16th International Conference on*. IEEE, 221–230.
- [23] Adam Lerer, Ledell Wu, Jiajun Shen, Timothee Lacroix, Luca Wehrstedt, Abhijit Bose, and Alex Peysakhovich. 2019. Pytorch-bigraph: A large scale graph embedding system. *Proceedings of Machine Learning and Systems* 1 (2019), 120–131.
- [24] Jie Li, Laiping Zhao, Yanan Yang, Kunlin Zhan, and Keqiu Li. 2022. Tetris: Memory-efficient Serverless Inference through Tensor Sharing. In *2022 USENIX Annual Technical Conference (USENIX ATC 22)*.
- [25] Zhuohan Li, Lianmin Zheng, Yinmin Zhong, Vincent Liu, Ying Sheng, Xin Jin, Yanping Huang, Zhifeng Chen, Hao Zhang, Joseph E Gonzalez, et al. 2023. AlpaServe: Statistical Multiplexing with Model Parallelism for Deep Learning Serving. *arXiv preprint arXiv:2302.11665* (2023).
- [26] Dandan Lin, Shijie Sun, Jingtao Ding, Xuehan Ke, Hao Gu, Xing Huang, Chonggang Song, Xuri Zhang, Lingling Yi, Jie Wen, et al. 2022. PlatoGL: Effective and Scalable Deep Graph Learning System for Graph-enhanced Real-Time Recommendation. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 3302–3311.
- [27] Fangming Liu and Yipei Niu. 2023. Demystifying the Cost of Serverless Computing: Towards a Win-Win Deal. *IEEE Transactions on Parallel and Distributed Systems* (2023).
- [28] Mingxuan Lu, Zhihao Han, Susie Xi Rao, Zitao Zhang, Yang Zhao, Yinan Shan, Ramesh Raghunathan, Ce Zhang, and Jiawei Jiang. 2022. BRIGHT-Graph Neural Networks in Real-time Fraud Detection. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 3342–3351.
- [29] Jingwei Ma, Kangkang Bian, Jiahui Wen, Yang Xu, Mingyang Zhong, and Lei Zhu. 2023. SRDPR: Social Relation-driven Dynamic network for Personalized micro-video Recommendation. *Expert Systems with Applications* 226 (2023), 120157.
- [30] Grzegorz Malewicz, Matthew H Austern, Aart JC Bik, James C Dehnert, Ilan Horn, Naty Leiser, and Grzegorz Czajkowski. 2010. Pregel: a system for large-scale graph processing. In *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data*. 135–146.
- [31] Seth A Myers, Aneesh Sharma, Pankaj Gupta, and Jimmy Lin. 2014. Information network or social network? The structure of the Twitter follow graph. In *Proceedings of the 23rd International Conference on World Wide Web*. 493–498.
- [32] Li Pan, Lin Wang, Shutong Chen, and Fangming Liu. 2022. Retention-aware container caching for serverless edge computing. In *IEEE INFOCOM 2022-IEEE Conference on Computer Communications*. IEEE, 1069–1078.
- [33] Qiangyu Pei, Yongjie Yuan, Haichuan Hu, Qiong Chen, and Fangming Liu. 2023. AsyFunc: A High-Performance and Resource-Efficient Serverless Inference System via Asymmetric Functions. In *Proceedings of the 2023 ACM Symposium on Cloud Computing*. 324–340.
- [34] Hao Peng, Hongfei Wang, Bowen Du, Md Zakirul Alam Bhuiyan, Hongyuan Ma, Jianwei Liu, Lihong Wang, Zeyu Yang, Linfeng Du, Senzhang Wang, et al. 2020. Spatial temporal incidence dynamic graph neural networks for traffic flow forecasting. *Information Sciences* 521 (2020), 277–290.
- [35] Huyen Trang Phan, Ngoc Thanh Nguyen, and Dosam Hwang. 2023. Fake news detection: A survey of graph neural network methods. *Applied Soft Computing* (2023), 110235.
- [36] AWS Fargate Pricing. [n. d.]. Serverless Compute Engine—AWS Fargate Pricing—Amazon Web Services. <https://aws.amazon.com/fargate/pricing/> [OnlineAccessed, 12-Feb-2024].
- [37] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P Adams, and Nando De Freitas. 2015. Taking the human out of the loop: A review of Bayesian optimization. *Proc. IEEE* 104, 1 (2015), 148–175.
- [38] John Thorpe, Yifan Qiao, Jonathan Eyolfson, Shen Teng, Guanzhou Hu, Zhihao Jia, Jinliang Wei, Keval Vora, Ravi Netravali, Miryung Kim, et al. 2021. Dorylus: Affordable, Scalable, and Accurate {GNN} Training with Distributed {CPU} Servers and Serverless Threads. In *15th USENIX Symposium on Operating Systems Design and Implementation (OSDI 21)*. 495–514.
- [39] Lucian Toader, Alexandra Uta, Ahmed Musaafir, and Alexandru Iosup. 2019. Graphless: Toward serverless graph processing. In *2019 18th International Symposium on Parallel and Distributed Computing (ISPDC)*. IEEE, 66–73.
- [40] Srinivas Virinchi, Anoop S V K K Saladi, and Abhirup Mondal. 2022. Recommending related products using graph neural networks in directed graphs. In *ECML-PKDD 2022*. <https://www.amazon.science/publications/recommending-related-products-using-graph-neural-networks-in-directed-graphs>
- [41] Daixin Wang, Jianbin Lin, Peng Cui, Quanhui Jia, Zhen Wang, Yanming Fang, Quan Yu, Jun Zhou, Shuang Yang, and Yuan Qi. 2019. A semi-supervised graph attentive network for financial fraud detection. In *2019 IEEE International Conference on Data Mining (ICDM)*. IEEE, 598–607.
- [42] Minjie Yu Wang. 2019. Deep graph library: Towards efficient and scalable deep learning on graphs. In *ICLR workshop on representation learning on graphs and manifolds*.
- [43] Yuke Wang, Boyuan Feng, Gushu Li, Shuangchen Li, Lei Deng, Yuan Xie, and Yufei Ding. 2021. GNNAdvisor: An adaptive and efficient runtime system for GNN acceleration on GPUs. In *15th USENIX symposium on operating systems design and implementation (OSDI 21)*.
- [44] Zhaojie Wen, Yishuo Wang, and Fangming Liu. 2022. StepConf: Slo-aware dynamic resource configuration for serverless function workflows. In *IEEE INFOCOM 2022-IEEE Conference on Computer Communications*. IEEE, 1868–1877.
- [45] Shiwen Wu, Fei Sun, Wentao Zhang, Xu Xie, and Bin Cui. 2022. Graph neural networks in recommender systems: a survey. *Comput. Surveys* 55, 5 (2022), 1–37.
- [46] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. 2020. A comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems* 32, 1 (2020), 4–24.
- [47] Fei Xu, Yiling Qin, Li Chen, Zhi Zhou, and Fangming Liu. 2021. λdnm: Achieving predictable distributed DNN training with serverless architectures. *IEEE Trans. Comput.* 71, 2 (2021), 450–463.

- [48] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. 2018. How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826* (2018).
- [49] Min Xu, Pakorn Watanachaturaporn, Pramod K Varshney, and Manoj K Arora. 2005. Decision tree regression for soft classification of remote sensing data. *Remote Sensing of Environment* 97, 3 (2005), 322–336.
- [50] Hongxia Yang. 2019. Aligraph: A comprehensive graph neural network platform. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*. 3165–3166.
- [51] Liekang Zeng, Peng Huang, Ke Luo, Xiaoxi Zhang, Zhi Zhou, and Xu Chen. 2022. Fograph: Enabling real-time deep graph inference with fog computing. In *Proceedings of the ACM Web Conference 2022*. 1774–1784.
- [52] Yanfu Zhang, Shangqian Gao, Jian Pei, and Heng Huang. 2022. Improving social network embedding via new second-order continuous graph neural networks. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 2515–2523.
- [53] Chengguang Zheng, Hongzhi Chen, Yuxuan Cheng, Zhezheng Song, Yifan Wu, Changji Li, James Cheng, Hao Yang, and Shuai Zhang. 2022. ByteGNN: efficient graph neural network training at large scale. *Proceedings of the VLDB Endowment* 15, 6 (2022), 1228–1242.
- [54] Hongkuan Zhou, Ajitesh Srivastava, Hanqing Zeng, Rajgopal Kannan, and Viktor Prasanna. 2021. Accelerating large scale real-time GNN inference using channel pruning. *arXiv preprint arXiv:2105.04528* (2021).
- [55] Hongkuan Zhou, Bingyi Zhang, Rajgopal Kannan, Viktor Prasanna, and Carl Busart. 2022. Model-Architecture Co-Design for High Performance Temporal GNN Inference on FPGA. In *2022 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*. 1108–1117.
- [56] Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. 2020. Graph neural networks: A review of methods and applications. *AI open* 1 (2020), 57–81.

A GNN LAYER STRUCTURES

Figure 15 shows the structure of three classic GNN layers, namely GCN [18], GraphSAGE [11], and GIN [48]. Each GNN layer is composed of two main operations alternatively executed: Aggregate and Update.

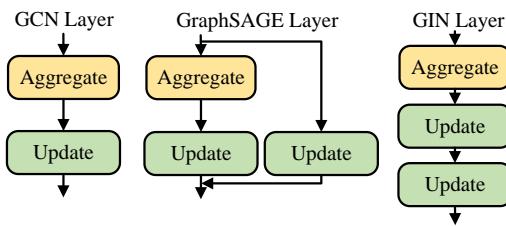


Figure 15: Classic GNN Layers.

B ALGORITHM DETAILS

Algorithm 1: Global Perspective Optimization Algorithm

Input	: Multi-Buffers $B = \{b_0, b_1, \dots, b_i\}$; Requests in the buffer $b_i = \{r_0^{b_i}, r_1^{b_i}, \dots, r_n^{b_i}\}$; Arrived requests $A = \{r_0, r_1, \dots, r_k\}$; A 's routing IDs $J = \{j_{r_0}, j_{r_1}, \dots, j_{r_k}\}$;
Output	: Modified Multi-Buffers B' ;
Parameters: Graph sharing degree of request with buffer S'_b ; Buffer to which the request is routed b_j ; Data index of request U_r ; Average sharing degree of buffer S_b , S'_b ; Average time ratio delayed in the buffer D_b, D'_b ; Performance gain of buffer μ_b, μ'_b ; Fixed coefficients α, β ;	

```

1 foreach  $r_k$  in  $A$  do
2   routeRequestToBuf( $r_k, j_{r_k}$ )  $\rightarrow b_j$ ;
3    $\alpha \times S_{b_j} - \beta \times D_{b_j} \rightarrow \mu_{b_j}$ ;
4   foreach  $b_i$  do
5     if  $b_i \neq b_j$  and  $S_{b_i}^{r_k} > 0$  then
6       foreach  $r_i$  in  $b_i$  do
7          $S_{r_i}^{r_k} = \frac{|U_{r_k} \cap U_{r_i}|}{|U_{r_k}|}$ ;
8         if  $S_{r_i}^{r_k} > 0$  then
9           modBufByDelRequest( $b_i, r_i$ )  $\rightarrow S'_{b_i}, D'_{b_i}$ ;
10           $\alpha \times S'_{b_i} - \beta \times D'_{b_i} \rightarrow \mu'_{b_i}$ ;
11          modBufByAddRequest( $b_j, r_i$ )  $\rightarrow S'_{b_j}, D'_{b_j}$ ;
12           $\alpha \times S'_{b_j} - \beta \times D'_{b_j} \rightarrow \mu'_{b_j}$ ;
13          if  $\mu'_{b_i} + \mu'_{b_j} > \mu_{b_i} + \mu_{b_j}$  then
14            transferRequestToBuf( $r_i$ )  $\rightarrow b_j$ ;
15            delRequestFromBuf( $r_i, b_i$ );
16 return  $B'$ ;

```

Algorithm 2: Dynamic Graph Scheduling Algorithm

Input	: Graph of the buffer $G_b(V_b, E_b)$; Graph of the arrived request $G_r(V_r, E_r)$; Target vertices IDs $W = [w_0, \dots, w_i]$;
Output	: Graph Partitions for each GNN layer $P = [[p_{00}, \dots, p_{0i}], \dots, [p_{n0}, \dots, p_{ni}]]$
Parameters: HAG $H(V_h, E_h)$; Traverse depth n ;	

```

1 GenerateHAG( $G_b(V_b, E_b), G_r(V_r, E_r)$ )  $\rightarrow H(V_h, E_h)$ ;
2 []  $\rightarrow P$  and  $0 \rightarrow n$ ;
3 while  $W \neq \emptyset$  do
4    $p_n = []$ ;
5   foreach  $w_i$  in  $W$  do
6     traversePredecessors( $w_i, H(V_h, E_h)$ )  $\rightarrow p_{ni}$ ;
7     append( $p_{ni}$ )  $\rightarrow p_n$ ;
8     append( $p_n$ )  $\rightarrow p$ ;
9      $p_{n0} \cup p_{n1} \cup \dots \cup p_{ni} \rightarrow W$ ;
10     $n + 1 \rightarrow n$ ;
11 return  $P$ ;

```
