# GraphSnapShot: Graph Machine Learning Acceleration with Fast Storage and Retrieval

> <sup>1</sup>University of Wisconsin-Madison <sup>2</sup>University of Notre Dame

## 1 Introduction

In our recent research, we have developed a framework called GraphSnapShot, which has been proven an useful tool for graph learning acceleration. The core idea of GraphSnapShot is to capture and update the state of local graph structures dynamically, just like taking snapshots of graphs.

GraphSnapShot is designed to efficiently capture, store and update the dynamic snapshots of graph data, enabling us to track patterns in the structure of graph networks. This technique is useful for most graph learning tasks that relies on topology analysis or networks are constantly evolving, such as social media analysis, biological networks, or any system where the relationships between entities change over time.

The key components of GraphSnapShot is the GraphSDSampler. GraphS-DSampler can efficiently capture, update, retrieve and store graph snapshots of topology while doing computation at the same time, which makes graph learning computation significantly faster.

In experiments, GraphSnapShot shows efficiency. It can promote computation speed significantly compared to traditional NeighborhodSamplers implemented in dgl, and it can reduce the GPU & memory usage during training, with little loss of accuracy. Experimental results show that the GraphSnapShot has potential to be a powerful tool for large graph training acceleration.

# 2 Background

Local structure analysis in graph learning applications are fundamental for understanding and analyzing complex networks in various domains. This concept pertains to the idea of capturing the localized patterns and relationships within a large graph structure. Efficient graph structure learning are crucial in var-

ious applications, from social network analysis to bioinformatics, where useful hidden pattern are discovered by graph structure learning.

The GraphSnapShot framework we proposed provide insight to study graph learning by using dynamically updated local graph snapshots. GraphSnapShot is particularly valuable in large-scale networks where global analysis can be computationally intensive and less informative for certain types of inquiries. By focusing on local structures, researchers can detect community patterns, identify influential nodes, and understand local clustering behaviors in a more efficient way.

Applications of local structure learning are diverse and impact several fields:

- Social Network Analysis: Understanding individual or entity interactions within a network to reveal social dynamics and community formation.
- **Bioinformatics:** In biological networks, local structures assist in identifying functional modules or predicting protein interactions.
- Recommendation Systems: Improving recommendation accuracy by focusing on immediate user-item interaction patterns.
- **Network Security:** Analyzing local patterns for detecting anomalies or potential security breaches.

## 3 Motivation

In machine learning, masks are powerful tools for reducing computation and enhancing model efficiency. In attention masks of long sequences, models selectively attend to relevant tokens within sequences, bypassing unnecessary computations and enabling more focused processing. Masks can also find utility in regularization techniques like dropout, where they aid in randomly deactivating units during training, thus preventing overfitting and promoting generalization. Additionally, in tasks such as instance segmentation in computer vision, masks efficiently delineate object boundaries, streamlining the computational load by restricting processing to specific regions of interest. Through these mechanisms, masks can reduce computational resources and contribute to the scalability of machine learning models. Notably, masked language modeling [STQ+19] [WGZC23], and masked vision learning [LFH+23] showcase the significance of masks.

Although masks are widely used in NLP and computer vision, their applications in graph learning are limited. In this study, we propose GraphSnapShot, a novel method for cache-masked graph structure learning. The aim of GraphSnapShot is to reduce computational resources and accelerate the speed of graph machine learning at large scale.

## 4 Model Construction

GraphSnapShot introduces caching masked graph structure for graph learning acceleration. The main process of GraphSnapShot involves two phases: preprocess sampling and dynamic resampling, both are methods for caching graph mask.

Traditional graph learning methods struggle with extensive memory usage and slow computation. Node-wise sampling methods, including well-known ones like GraphSAGE [HYL18], are not optimized for multi-hop scenarios due to the computational resources usage that grows exponentially with each additional hop— a phenomenon known as the "neighbor explosion." This makes them unsuitable for large-scale graphs where the capture of extended neighborhoods is crucial. Layer-wise sampling, such as in FastGCN [CMX18], although designed to mitigate this explosion by sampling nodes per layer, often loses valuable connectivity information, leading to a sparse and inaccurate representation of the graph's multi-hop structure. Furthermore, subgraph sampling methods such as Cluster-GCN [CLS+19] focus on localized computations and overlook vital inter-subgraph connections, failing to capture the broader topology necessary for a comprehensive multi-hop domain analysis.

GraphSnapShot's core idea is to achieve efficient cache management in graph learning, it is designed to maintain an up-to-date cache through the following strategies:

- Fully Cache Refresh (FCR): Periodically refresh the whole cache to get a new representative of the graph local structure.
- On-The-Fly (OTF): Retreive the new representation of graph local structure in each computation, with partially retrieve from the previous graph local structure and partially retrieve from the disk. In detail the OTF can be identified as follows based on their fetch and refresh methods:
  - Partial Cache Refresh and Full Fetch: Dynamically updates portions of the cache or retrieves all data from cache for computation as needed.
  - Partial Refresh with Partial Fetch: OTF method in partially updating and partially retrieving data from both disk and cache, balancing cache up-to-date cache with computational efficiency.
  - Only Fetch Mode: Fetch data for computation from cache by OTF mode with update existing cache periodically.
  - Only Refresh Mode: Updates the existing cache by OTF mode.
- Shared Cache Mode: Unlike traditional methods that cache individual nodes or edges, this mode involves caching significant portions of the graph that are frequently accessed by various processes. This strategy optimizes memory usage and ensures that the most frequently accessed data is readily available, enhancing the responsiveness and efficiency of graph machine learning tasks.

```
Algorithm 1 Static Sampling and Dynamic ReSampling
```

```
1: function PREPROCESS(G)
       sG \leftarrow \text{StaticSample}(G)
        return sG
 4: end function
 5: function DYNAMICSAMPLE(G, \alpha, N)
       num \leftarrow \alpha \times N
6:
        sNodes \leftarrow \text{SampleFromDisk}(G, num)
7:
       {f return}\ sNodes
9: end function
10: function RESAMPLE(sG, \alpha, N)
       num \leftarrow \alpha \times N
11:
        sNodes \leftarrow \text{SampleFromMemory}(sG, num)
12:
       return \ sNodes
13:
14: end function
```

## Algorithm 2 Local Snapshot Swap

```
function Add G, Nodes)

2: G \leftarrow G \cup Nodes
Nodes \leftarrow MoveToMemoryFromDisk(Nodes)

4: return G
end function

6: function RemoveG, Nodes)
G \leftarrow G - Nodes

8: Nodes \leftarrow MoveToDiskFromMemory(Nodes)
return G

10: end function
```

#### Algorithm 3 GraphSnapShot Process

```
function Process(G, \alpha, N)
 2:
        sG \leftarrow \text{Preprocess}(G)
        while (InComputation) do
            dNodes \leftarrow \text{DynamicSample}(G, \alpha, N)
 4:
            rNodes \leftarrow \text{ReSample}(sG, (1-\alpha), N)
            uG \leftarrow \text{Combine}(dNodes, rNodes)
 6:
            sG \leftarrow \text{Remove}(sG, rNodes)
            sG \leftarrow \text{Add}(sG, dNodes)
 8:
            result \leftarrow Compute(uG)
        end while
10:
        {\bf return}\ result
12: end function
```

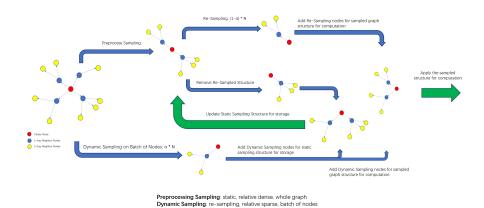


Figure 1: GraphSnapShot Model, code: https://github.com/NoakLiu/GraphSnapShot

GraphSnapShot is designed to solve the inefficacy for current graph masking methods on large graphs. Typical distributed graph processing System such as Marius [MWX<sup>+</sup>21], will re-sample all the local structure each time, and load corresponding node embeddings from disk, which is time-consuming. GraphSnap-Shot proposing a new method for quick local structure retrieval by integrating the benefits of both static and dynamic processing. In the preprocessing phase, a static graph snapshot is sampled and stored in memory from the disk. Then, in the computation phase, the algorithm re-samples the graph snapshot, using a mix of data partially from memory and partially from disk. Concurrently, a smaller, dynamic snapshot is sampled from the disk. This dynamic snapshot is crucial as it is used to swap with the portion of the graph stored in memory. This process ensures the computation can always learn "fresh" snapshot without overwhelming the memory capacity. GraphSnapShot allows the algorithm to continuously update and adapt to changes in the graph while maintaining a manageable memory footprint, making it an effective tool for processing large, dynamic graph data.

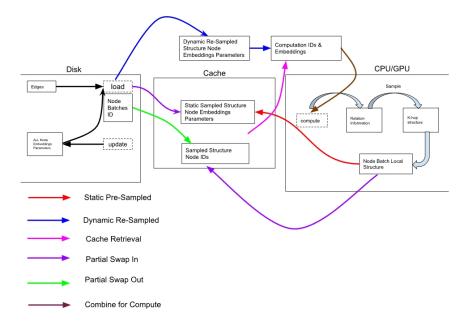


Figure 2: Disk-Cache-CPU/GPU Architecture for GraphSnapShot

In the GraphSnapShot system, disks are used to store edge and node embedding information, while computational resources (CPU/GPU) are utilized for neural network computations (such as backpropagation) and local structure masking. Additionally, the cache is employed to store the K-V pairs information of the sampled structure.

- Disk: Persistent storage of graph structure information, such as edges and node embedding parameters. Those embeddings are stored long-term and loaded into memory when needed for further processing and analysis.
- Cache: Store frequent K-V pairs of GraphSnapShot to accelerate graph learning. These K-V pairs include the identifiers of the sampled structure nodes and their embedding parameters, allowing for more efficient repeated access to the graph data.
- CPU/GPU: Computational Resources to responsible for performing local structure computations on the graph, including neural network forward and backward propagation as well as the dynamic resampling and masking of the graph's local structure.

# 5 Algorithmic Design

As for the cache retrieval and refreshing methods, we develop three methods to simulate different strategies.

#### 5.1 FBL

The first method is FBL, which is short for full batch load. FBL is the traditional method where most graph processing systems current are deployed. Those system does not caching any information and resample graph structure each time from disk to memory for computation.

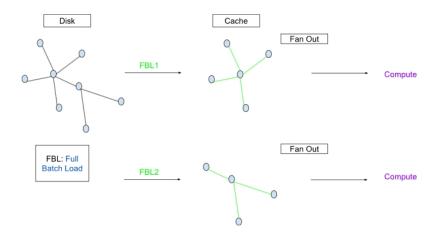


Figure 3: FBL Cache Snapshot Strategy

#### 5.2 OTF

The "On The Fly" (OTF) method, updates portions of the cache dynamically while leaving other segments of the cache information unchanged. In detailed implementation of OTF, we design 4 different mode of it for real world application. They are (Full Fetch, Partial Fetch)x(Full Refresh, Partial Refresh).

In \*OTF-refresh\* mode, we implement a partial cache refresh strategy, which dynamically updates only a fraction of the cache based on a defined refresh rate.

#### 5.2.1 Algorithmic Design for OTF-refresh

Let G = (V, E) denote a heterogeneous graph, where V and E represent the vertices and edges, respectively. The vertices are categorized by types, specified by hete\_label.

- fanouts =  $[f_1, f_2, ..., f_n]$ : An array indicating the desired number of neighbors to sample at each depth level.
- $\alpha$ : A scaling factor for fanouts to determine the initial cache size for each depth.
- T: The interval after which the cache is considered for partial refreshing.
- γ: Refresh rate defining the fraction of the cache to update during each refresh cycle.

#### Algorithm 4 Neighbor Sampler with OTF Cache Refresh

```
1: procedure INITIALIZE(Graph G, Array fanouts, \alpha, T, \gamma, hete_label)
       Calculate cache sizes: cache_size[i] \leftarrow f_i \times \alpha for each i
 3:
       Initialize each cache layer by pre-sampling using calculated sizes
 4: end procedure
 5: procedure RefreshCache(Layer ID, Cache Refresh Size)
       Sample a portion of the current cache to retain
 6:
 7:
       Sample new neighbors to replace the refreshed portion
       Merge retained and new samples to update the cache
 8:
 9: end procedure
   procedure SampleBlocks(Seed nodes S)
10:
       if current cycle modulo T=0 then
11:
          for each cache layer do
12:
              Refresh a portion of the cache based on \gamma
13:
          end for
14:
       end if
15:
       for each layer in reverse order do
16:
          Use the cache to sample blocks for seed nodes
17:
          Update seed nodes for the next layer
18:
       end for
19:
       return Sampled blocks
21: end procedure
```

The OTF refresh strategy enhances the algorithm's responsiveness to changes in the graph's structure by updating a subset of the cache, thus balancing between computational efficiency and sampling accuracy. This approach is particularly beneficial in dynamic environments where frequent full cache refreshes are impractical.

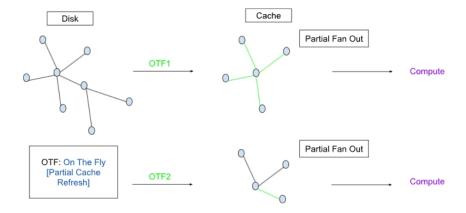


Figure 4: OTF Cache Snapshot Strategy

In \*OTF-fetch\* mode, we design an algorithm that can maintaining up-todate neighborhood information in dynamic graph environments. Traditional methods that rely solely on static pre-sampling can lead to stale data, affecting the model's performance. OTF-fetch mode introduces a fetching mechanism that balances between pre-sampled cache and real-time graph data, enhancing accuracy and efficiency.

#### 5.2.2 Algorithmic Design for OTF-fetch

Let G = (V, E) be a heterogeneous graph, where V and E represent the vertices and edges, respectively. Each vertex belongs to a specific type, categorized under hete\_label.

- fanouts =  $[f_1, f_2, \dots, f_n]$ : An array indicating the number of neighbors to sample at each layer.
- amp\_rate: Amplification rate that determines the size of the initial cache.
- $T_{\text{fetch}}$ : The period after which the cache is partially fetched to adjust the sample data dynamically.
- fetch\_rate: Determines the fraction of the cache to be dynamically fetched alongside real-time data from the graph.

The OTF fetching strategy introduces a significant enhancement to graph sampling by allowing dynamic data retrieval, which integrates cached and real-time data. This mechanism ensures that the graph samples are both current and computationally efficient to retrieve, catering particularly well to environments where the graph's structure may change frequently. Moreover, this approach reduces the overhead associated with frequent full cache refreshes, thereby optimizing both memory usage and processing time.

#### Algorithm 5 Neighbor Sampler with OTF Fetching Mechanism

```
1: procedure Initialize(Graph G, Array fanouts, amp_rate, T_fetch,
   fetch_rate, hete_label)
       Compute cache sizes: cache_size[i] \leftarrow f_i \times \text{amp\_rate} for each i
2:
      Initialize each cache layer by pre-sampling using these sizes
3:
 4: end procedure
5: procedure OTF_FETCH(Layer ID, Seed Nodes, Fetch Size)
       Compute the number of neighbors to fetch from the disk
6:
       Fetch required neighbors from the cache
7:
8:
      Fetch additional neighbors directly from the graph
      Merge both sets to form the updated sample
9:
10: end procedure
   procedure SampleBlocks(Seed nodes S)
       for each layer from last to first do
12:
          Determine the number of neighbors to fetch based on fetch_rate
13:
          if current cycle modulo T_{\text{fetch}} = 0 then
14:
              Call OTF_Fetch to update sampled data
15:
          else
16:
             Sample directly from the pre-filled cache
17:
          end if
18:
          Update seed nodes for the next layer
19:
       end for
20:
       return Sampled blocks
22: end procedure
```

# 5.3 Algorithmic Design for OTF-partial fetch & partial refresh

Let G = (V, E) represent a heterogeneous graph where V and E denote the sets of vertices and edges, respectively, with vertices categorized by types specified by hete\_label.

- fanouts =  $[f_1, f_2, ..., f_n]$ : Specifies the number of neighbors to sample at each depth.
- amp\_rate: Amplification rate determining the initial cache size for each layer.
- T: Period after which the cache is refreshed.
- refresh\_rate: Proportion of the cache to be refreshed dynamically during each cycle.

#### Algorithm 6 Neighbor Sampling with OTF Partial Cache Fetch and Refresh

- 1: **procedure** INITIALIZE(Graph G, fanouts, amp\_rate, refresh\_rate, hete\_label)
- 2: Compute cache sizes using amp\_rate and fanouts.
- 3: Initialize cache for each layer.
- 4: end procedure
- 5: **procedure** OTF\_RF\_CACHE(layer\_id, cached\_graph, seed\_nodes, refresh\_rate, fanout)
- 6: Calculate portions of the cache to refresh and fetch.
- 7: Fetch new data for specified portions while maintaining other data unchanged.
- 8: Update the cache dynamically to reflect recent changes.
- 9: end procedure
- 10: **procedure** SampleBlocks(Graph G, seed\_nodes)
- 11: **for** each layer in cache **do**
- 12: **if** refresh required **then**
- 13: Call OTF\_rf\_cache to refresh and fetch data.
- 14: end if
- 15: Use updated cache to sample blocks for seed nodes.
- 16: Update seed nodes for the next layer.
- 17: end for
- 18: **return** Sampled blocks and updated seed nodes.
- 19: end procedure

The OTF partial fetch and refresh strategy significantly enhances the responsiveness and accuracy of graph sampling in dynamic environments. By selectively updating the cache, the algorithm reduces computational overhead while maintaining high data freshness, crucial for downstream applications like machine learning where timely and accurate data is vital.

## 5.4 OTF partial refresh & full fetch

OTF partial refresh & full fetch is a mechanism for both partial cache refresh and full fetch. This methodology optimizes the sampling process by dynamically adjusting cache content to reflect the latest graph state, thus balancing computational efficiency with data freshness.

Efficient graph sampling is pivotal in applications involving dynamic graphs, especially in scenarios where graph structures evolve over time. Traditional static caching methods often lead to outdated data, adversely affecting the performance of graph-based algorithms. Our approach, through dynamic cache management, ensures timely updates of the sampled data, significantly enhancing the relevance of the data used in downstream processes.

Let G = (V, E) denote a heterogeneous graph where V and E represent vertices and edges, respectively, categorized by hete\_label.

- fanouts =  $[f_1, f_2, ..., f_n]$ : An array specifying the number of neighbors to sample at each depth.
- amp\_rate: Amplification rate that determines the initial size of the cache for each depth.
- T: Period after which the cache undergoes a refresh operation.
- refresh\_rate: Proportion of the cache that is refreshed in each cycle.

#### 5.5 FCR

The third method is FCR, which is short for "fully cache refreshing" in a batch. For FCR, in each batch, the structure information are sampled from cache for computation, and between each batch, the cache is fully refreshed.

The Fully Cache Refresh (FCR) strategy effectively scales the depth and breadth of pre-sampled neighborhoods stored in the cache. This mechanism reduces the need for frequent neighborhood sampling from the graph, thereby lowering computational overhead and enhancing performance, especially in dynamic and large-scale heterogeneous graph environments. The periodic refresh rate T ensures the cache remains up-to-date and reflects recent changes in the graph structure.

Let G = (V, E) be a heterogeneous graph where V and E represent vertices and edges respectively. The vertices may belong to different types, specified by hete label.

- fanouts =  $[f_1, f_2, ..., f_n]$ : An array indicating the number of neighbors to sample at each depth of the sampling process.
- $\alpha$ : Fractional caching rate, controlling the scaling of **fanouts** for presampling.

# **Algorithm 7** Neighbor Sampling with OTF Partial Cache Refresh and Full Fetch

- 1: **procedure** INITIALIZE(Graph G, fanouts, amp\_rate, refresh\_rate, hete\_label)
- 2: Compute cache sizes using amp\_rate and fanouts.
- 3: Initialize cache for each layer.
- 4: end procedure
- 5: **procedure** OTF\_REFRESH\_CACHE(layer\_id, cached\_graph, seed\_nodes, refresh\_rate)
- 6: Calculate the proportion of the cache to refresh.
- 7: Fetch new data for the refreshed portion while retaining other data unchanged.
- 8: Update the cache dynamically to incorporate recent graph changes.
- 9: end procedure
- 10: **procedure** SampleBlocks(Graph G, seed\_nodes)
- 11: **for** each layer in cache **do**
- 12: **if** refresh required **then**
- 13: Call OTF\_refresh\_cache to refresh the data.
- 14: **end if**
- 15: Fully fetch from the updated cache to sample blocks for seed nodes.
- 16: Update seed nodes for the next layer.
- 17: end for
- 18: **return** Sampled blocks and updated seed nodes.
- 19: end procedure

• T: The interval after which the cache is refreshed to incorporate recent graph changes.

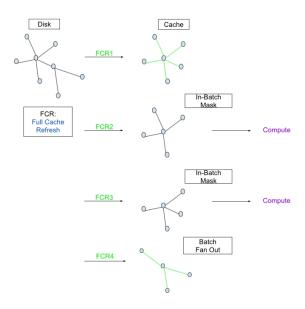


Figure 5: FCR Cache Snapshot Strategy

Traditional methods often suffer from high computational costs due to repetitive neighborhood sampling. Our approach mitigates these issues by integrating a Fully Cache Refresh, which dictates the extent of neighborhood pre-sampling and cache utilization.

The Fully Cache Refresh approach, while ensuring the freshness of the sampled data, can be computationally intensive compared to partial or incremental refresh strategies, as it involves re-sampling large portions of the graph frequently. Therefore, it's crucial to carefully configure the refresh period T to balance performance with computational overhead, particularly in very large or rapidly-changing graphs.

This approach is ideal for scenarios where the accuracy and timeliness of the graph data are critical, such as in real-time monitoring systems or dynamic social networks where relationships and attributes may frequently change. By refreshing the entire cache periodically, the FCR method ensures that the sampling process remains robust and relevant to the current graph structure, thus supporting high-quality graph-based analyses and applications.

## Algorithm 8 NeighborSampler with Fully Cache Refresh

```
1: procedure Initialize(Graph G, Array fanouts, \alpha, T, hete_label)
       amplified_fanouts \leftarrow [f \times \alpha \text{ for } f \text{ in fanouts}]
 3:
       Initialize cache to empty
       Pre-fill cache using amplified_fanouts
 4:
 5: end procedure
   procedure Cacheresh
       Clear cache
 7:
       Re-sample neighborhoods using amplified_fanouts
 9: end procedure
10: procedure SampleBlocks(Seed nodes S)
       if current cycle modulo T=0 then
11:
          Cacherefresh
12:
       end if
13:
       for each level k from n down to 1 do
14:
          Sample using cached neighborhoods
15:
          Update seed nodes for next level
16:
       end for
17.
       return Sampled blocks
18:
19: end procedure
```

## 5.6 Shared Cache Design

Graph sampling is often computationally intensive, especially when dealing with large-scale heterogeneous graphs. The shared cache design aims to mitigate these challenges by maintaining a cache of pre-sampled neighborhoods, thus facilitating rapid access and reducing the frequency of expensive graph traversal operations.

Let G = (V, E) denote a heterogeneous graph where V and E represent the sets of vertices and edges, respectively. fanouts defines the number of neighbors to sample at each depth of sampling, and  $\alpha$  represents the scaling factor for cache size determination.

- fanouts =  $[f_1, f_2, \dots, f_n]$ : Specifies the number of neighbors to sample at each layer.
- $\bullet$   $\alpha$ : Amplification factor applied to each fanout to determine the initial cache size.
- T: Period after which the cache is refreshed to ensure it remains representative of the graph's evolving structure.

The shared cache design effectively reduces the computational overhead associated with frequent and intensive graph sampling operations. By leveraging an amplified fanout strategy, the cache maintains a ready pool of neighborhood data that can be rapidly accessed during the sampling process, significantly enhancing performance and efficiency. The periodic refresh mechanism ensures the

## Algorithm 9 Initialization and Refresh of Shared Cache

```
1: procedure INITIALIZE(Graph G, fanouts, alpha)
2: sc_size \leftarrow \max(fanouts) \times \alpha \triangleright Compute cache size
3: Pre-sample and populate cache using G, fanouts, and sc_size
4: end procedure
5: procedure CACHEREFRESH
6: if current cycle modulo T = 0 then
7: Pre-sample neighborhoods to refresh the cache
8: end if
9: end procedure
```

cache's data remains current, accommodating changes in the graph's topology over time.

# 6 GraphSnapShot Overview

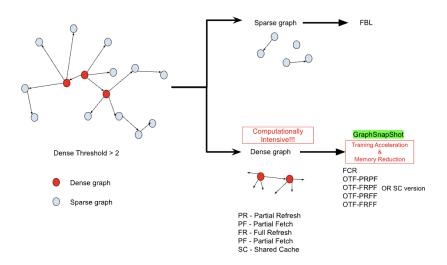


Figure 6: Graph Processing Overview of GraphSnapShot

Our approach involves splitting a graph into two subgraphs based on the node degree threshold. Nodes with a degree higher than a specified threshold are categorized into a dense subgraph, while those with a degree below the threshold form a sparse subgraph. This division allows for the application of different sampling techniques to each subgraph, optimizing the computational resources and processing time.

For dense subgraphs, we utilize advanced sampling methods in GraphSnap-Shot such as FCR, OTF, FCR-SC, OTF-SC, etc, which are designed to efficiently manage the memory cache while maintaining the quality of the sampled neighborhood. These samplers dynamically adjust their sampling rates and cache refresh rates based on the learning phase and node centrality, ensuring efficient processing of dense regions.

For sparse subgraphs, we employ a Full Batch Load (FBL) approach, akin to DGL's NeighborSampler. This method is suitable for areas of the graph with lower degrees, where the computational overhead is inherently less, allowing for direct processing without the need for intricate sampling techniques.

We use the Deep Graph Library (DGL [WZY<sup>+</sup>20]) framework to implement new samplers in GraphSnapShot, leveraging its efficient handling of graph data structures and operations. Our implementation involves loading the graph, applying the graph splitting based on node degree, and then processing each subgraph with the respective sampling technique. We evaluated the performance of our new models on datasets in ogbn-benchmark [HFZ<sup>+</sup>21]. The results show that GraphSnapShot can significantly reduce training time and memory usage without compromising the quality of GNN training.

Bridging the Gap Between Pure Dynamic Algorithms and Static Memory Storage with GraphSnapShot: Dynamic graph algorithms, such as GraphSAGE, which require resampling the entire graph for each computation, incur significant overhead. This is primarily due to the need for constant, full-graph resampling to accommodate the ever-changing nature of dynamic graphs. GraphSnapShot addresses this issue by establishing a static snapshot at preprocessing phase and use dynamic processing to update partial snapshot at computation. This approach significantly reduces the computational burden by providing a stable, memory-stored snapshot (sG) as a baseline. GraphSnapShot ensures efficient computation and maintains the flexibility needed to handle dynamic changes in the graph.

Tradeoff Between Dynamic Sampling and Resampling: In GraphSnap-Shot, the balance between dynamic sampling (DynamicSample) and resampling (ReSample) plays a pivotal role in maintaining an up-to-date representation of the graph. Dynamic sampling involves integrating new nodes from the disk into memory, capturing the latest updates and changes. Resampling, on the other hand, focuses on sampling and processing existing nodes from memory, which are then rotated back to disk storage. However, if DynamicSample is too large, it is close to traditional sampling, resulting in prolonged disk retrieval times. Conversely, while resampling processes memory-stored nodes, an excessive Re-Sample may overly rely on memory data, potentially leading to a decrease in accuracy. Striking this balance is essential: moderating dynamic sampling to prevent inefficiencies and ensuring resampling maintains graph accuracy without an undue dependence on memory. In short, this is a tradeoff between optimizing memory usage (quick accessing) and model accuracy.

**Local Snapshot Swap Strategy:** The *Add* and *Remove* functions are designed to simulate the swap of nodes between memory and disk, reflecting the dynamic shifts between active and less active parts of the graph. This local snapshot swapping strategy is crucial for handling large-scale dynamic graphs, as it reduces the demand on memory and computational resources.

In conclusion, the design of our GraphSnapShot Process algorithm considers the complexity and variability of dynamic graph analysis. By combining static sampling with dynamic resampling, the algorithm effectively manages and analyzes graphs whose structures evolve over time, demonstrating robust performance and flexibility in handling large-scale network analyses.

Dense & Sparse Graphs Processing: Our approach shows significant improvements in computational efficiency compared to traditional methods. The dual sampling strategy for dense and sparse graphs allows for targeted resource allocation, effectively handling dense and sparse regions within the graph. By using FBL, sparse graph processing can avoid resampled on same node set; by using GraphSnapShot, dense graph processing can have smaller cache size and quicker processing time. This results in faster training times without compromising the accuracy of the GNN models. Belows are figures to visulize the sample efficiency for Dense & Sparse Processing.

#### Algorithm 10 Calculate Total Degree

```
1: Input: Heterogeneous graph G, node type t
 2: Output: Array total_degrees
 3: total_degrees \leftarrow array of zeros, size equal to |G_t| (number of nodes of type
 4: for each (s, r, d) in G.edge_types do
       if d = t then
 5:
 6:
           total\_degrees \leftarrow total\_degrees + G.in\_degrees((s, r, d))
 7:
       end if
       if s = t then
 8:
           total\_degrees \leftarrow total\_degrees + G.out\_degrees((s, r, d))
 9:
10:
       end if
11: end for
12: return total_degrees
```

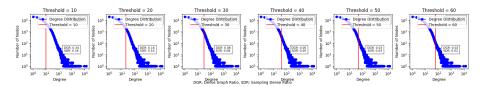


Figure 7: Dense and Sparse Processing for ogbn-arixv

## Algorithm 11 Split Graph by Degree

- 1: **Input:** Heterogeneous graph G, node type t, degree threshold  $\theta$
- 2: Output: Subgraphs  $G_{\text{high}}, G_{\text{low}}$
- 3: total\_degrees  $\leftarrow$  Calculate\_Total\_Degree(G, t)
- 4: high\_nodes  $\leftarrow$  indices where total\_degrees  $> \theta$
- 5: low\_nodes  $\leftarrow$  indices where total\_degrees  $\leq \theta$
- 6:  $G_{\text{high}} \leftarrow G.\text{subgraph(high\_nodes)}$
- 7:  $G_{low} \leftarrow G.subgraph(low\_nodes)$
- 8: **return**  $G_{\text{high}}, G_{\text{low}}$

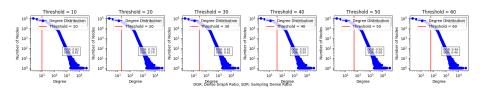


Figure 8: Dense and Sparse Processing for ogbn-products

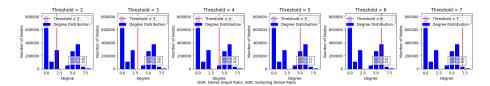


Figure 9: Dense and Sparse Processing for ogbn-mag

# 7 Experimental Result

# 7.1 Theoretial Analysis of Traditional Disk-Memory Strategy and Our Disk-Cache-Memory Strategy

In the traditional Graph System like Marius [MWX<sup>+</sup>21], a disk-cache model is deployed, where the graph structure is resampled each time from the disk into the memory for computation. In our current GraphSnapShot model, we deploy a disk-cache-memory model where K-V pair for reused graph structure are stored in the cache to promote the system performance.

Theoretically, if we denote the batch processing size as S(B), the cache size as S(C) and the cache refreshing rate as  $\alpha$ , and the processing speed of cache as  $v_c$ , and the processing speed of memory as  $v_m$ . Then we can get the batch average processing time for disk-memory model as  $\frac{S(G)}{v_m}$  and the batch average processing time for disk-cache-memory model as  $\frac{S(B)-S(C)}{v_m}+\frac{(1-\alpha)*S(C)}{v_c}$ .

# 7.2 Training Time Reduction Analysis

The provided table presents a detailed comparison of training time acceleration for various computational methods relative to the baseline method FBL. It quantitatively illustrates the percentage reduction in training times when employing FCR, FCR-shared cache, OTF, and OTF-shared cache methods across three distinct settings: [20, 20, 20], [10, 10, 10], and [5, 5, 5]. Each entry reflects the efficiency gains achieved by the respective method, highlighting how advancements in computational techniques can significantly decrease the time required for model training.

Table 1: Training Time Acceleration Percentage Relative to FBL

Method/Setting	[20,20,20]	[10,  10,  10]	[5, 5, 5]
FCR	7.05%	14.48%	13.76%
FCR-shared cache	7.69%	14.33%	14.76%
OTF	11.07%	23.96%	23.28%
OTF-shared cache	13.49%	25.23%	29.63%

The training time comparison table reveals that the OTF-shared cache method consistently provides the most substantial time accelerations, achieving up to 29.63% faster training times than the FBL method, particularly in the smallest setting. This performance underscores its effectiveness in scenarios demanding quick model training. The FCR methods, both with and without shared caching, offer steady improvements, enhancing training speeds by approximately 7% to 14.76%. Notably, the OTF method shines in smaller settings, suggesting increased efficiency in less complex computational tasks. The positive impact of shared caching across the FCR and OTF methods indicates that optimized resource management can lead to significant reductions in training times, making these methods highly beneficial in high-performance computing environments where reducing time is essential.

## 7.3 Runtime Memory reduction Analysis

From the visualizations and data provided in the table of runtime memory reduction analysis, it is evident that the computational methods FCR, FCR-shared cache, OTF, and OTF-shared cache demonstrate significant improvements in reducing runtime memory usage compared to the baseline FBL method up to 90%.

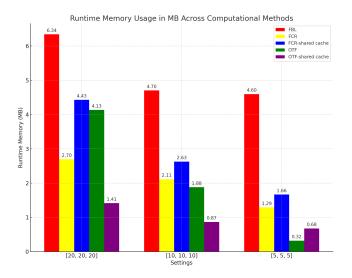


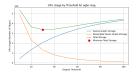
Figure 10: Runtime Memory Reduction

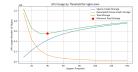
# 7.4 GPU Reduction Analysis

This significant reduction in GPU memory usage demonstrates the effectiveness of the GraphSnapShot methodology. By intelligently balancing the trade-off between full storage of sparse graph regions and sampling of dense graph regions, the proposed approach greatly enhances memory efficiency, which is particularly beneficial for large-scale graph processing tasks.

Table 2: GPU Storage Optimization Comparison

Dataset	Original	Optimized	Compression
ogbn-arxiv	1,166,243	552,228	52.65%
ogbn-products	123,718,280	20,449,813	83.47%
ogbn-mag	5,416,271	556,904	89.72%





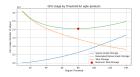


Figure 11: GPU Reduction Visualizations for Datasets (ogbn-arxiv, products and mag)

## 7.5 Fluent Generation Analysis

The partial fetch and partial refresh mode are designed to achieve fluent generation in GraphSnapShot. As in the GraphSnapShot, the most important part of the way is to develop a great opportunity to achieve ...

Method	Description	Time (seconds)	
Method	Description	Average	Worst
Full Refresh (FR) + OTF Partial Fetch (PF)	Most Common	-	-
Partial Fetch (PF) + Partial Refresh (PR)	Less Common	0.23	0.82
Partial Refresh (PR) + Partial Fetch (PF)	Less Common	0.54	0.54
Partial Refresh (PR) + Full Fetch (FF)	Less Common	0.31	0.31

Table 3: Performance Comparison of Refresh and Fetch Methods

Full Refresh (FR) + OTF Partial Fetch (PF) is the most commonly used method, leveraging full refresh alongside on-the-fly partial fetch, likely offering advantages in performance and efficiency, although no time data is provided. Partial Fetch (PF) + Partial Refresh (PR) exhibits a shorter average time of 0.23 seconds but a longer worst time of 0.82 seconds, indicating good performance in most cases but potential delays in certain scenarios due to its implementation of partial fetch and partial refresh. Partial Refresh (PR) + Partial Fetch (PF) has an average time of 0.54 seconds and a worst time of 0.54 seconds, showing more consistent performance across various scenarios, being closest to a graph structure query with each refresh. Partial Refresh (PR) + Full Fetch (FF) shows both an average and worst time of 0.31 seconds, similar to the second method in average time but excelling in worst-case scenarios, while also offering the best fluent generation due to its implementation of partial refresh and full fetch.

## 7.6 Experimental Analysis of SSDReS Performance

In our experiments, we explored the efficacy of the GraphSnapShot algorithm by testing the tradeoff between static and dynamic resampling. Static fetch (proportional to fr) focuses on quick access to graph topology information stored in memory. Dynamic resampling (proportional to rr), on the other hand, is a measurement for the memory swap, it addresses the real-time changes in the graph, frequently requiring access to disk-stored data and swap in the memory to update the in-memory topology.

Analyzing the outcomes of these experiments provides insights into the algorithm's performance:

• Accuracy and fetch rate: Our results showed a general increase in accuracy as fetch rate increased, which shows dynamic resampling has critical roles in enhancing model precision. A higher rate of dynamic resampling implies a greater reliance on the stable structure of the graph, aiding in capturing patterns in graph topology. A relative high dynamic resampling rate is beneficial for the model's training stability and generalizability.

- Loss and fetch rate: Loss decreased with an increase in fetch rate, further validating the importance of dynamic resampling in improving model performance. Lower loss indicates reduced error in data fitting, suggesting that static resampling helps the algorithm to learn and predict the graph's structure more accurately. This could also imply that dynamic resampling offers a stable learning environment, potentially reducing the risk of overfitting.
- Training Time and refresh rate: Our experiments indicated the existence of an optimal refresh rate value between 0 and 1, which minimizes training time, and with a little loss in acuracy. We call this point as "trade-off" point between static and dynamic resampling. Over-dependence on dynamic resampling (low refresh rate) could increase the computational burden due to frequent updates required for rapidly changing graph structures, while excessive reliance on static resampling (high refresh rate) might cause the model to overly learn from limited portion of the graph structure in memory, deduce its ability to effectively learn from the entire graph structure.

The effectiveness of the GraphSnapShot algorithm lies in the way that it deploys disk-cache-memory architecture to store and update intermediate graph local structure (GraphSnapShot). By caching the K-V pairs of the reused graph structure, the graph machine learning is significantly accelerated.

The GraphSnapShot excels in processing dynamic graph data and tasks that require exploration of graph structures, particularly for applications with complex and large-scale network structures.

# References

- [CLS+19] Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. Cluster-GCN. In Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery Data Mining. ACM, jul 2019.
- [CMX18] Jie Chen, Tengfei Ma, and Cao Xiao. Fastgcn: Fast learning with graph convolutional networks via importance sampling, 2018.
- [HFZ<sup>+</sup>21] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs, 2021.
- [HYL18] William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs, 2018.
- [LFH<sup>+</sup>23] Yanghao Li, Haoqi Fan, Ronghang Hu, Christoph Feichtenhofer, and Kaiming He. Scaling language-image pre-training via masking, 2023.

- [MWX<sup>+</sup>21] Jason Mohoney, Roger Waleffe, Yiheng Xu, Theodoros Rekatsinas, and Shivaram Venkataraman. Learning massive graph embeddings on a single machine. *CoRR*, abs/2101.08358, 2021.
- [STQ<sup>+</sup>19] Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. Mass: Masked sequence to sequence pre-training for language generation, 2019.
- [WGZC23] Alexander Wettig, Tianyu Gao, Zexuan Zhong, and Danqi Chen. Should you mask 15
- [WZY+20] Minjie Wang, Da Zheng, Zihao Ye, Quan Gan, Mufei Li, Xiang Song, Jinjing Zhou, Chao Ma, Lingfan Yu, Yu Gai, Tianjun Xiao, Tong He, George Karypis, Jinyang Li, and Zheng Zhang. Deep graph library: A graph-centric, highly-performant package for graph neural networks, 2020.

# A Appendix

# A.1 dgl Experiments

#### A.1.1 datasets

Table 4: Overview of OGBN Datasets

Feature	ARXIV	PRODUCTS	MAG
Type	Citation Net.	Product Net.	Acad. Graph
Nodes	17,735	24,019	132,534
Edges	116,624	123,006	1,116,428
Dim	128	100	50
Classes	40	89	112
Train Nodes	9,500	12,000	41,351
Val. Nodes	3,500	2,000	10,000
Test Nodes	4,735	10,019	80,183
Task	Node Class.	Node Class.	Node Class.

# A.1.2 Training Time Acceleration

This table provides a comparison of the training times ("training time") of different computational methods (FBL, FCR, FCR-shared cache, OTF, OTF-shared cache) across three settings ([20, 20, 20], [10, 10, 10], and [5, 5, 5]). It also calculates the time acceleration percentage relative to FBL for methods other than FBL, expressed as a percentage.

Table 5: Training Time Acceleration Through Each Method

Method	Setting	Time (s)	Acceleration (%)
FBL	[20, 20, 20]	0.2766	-
	[10, 10, 10]	0.0747	-
	[5, 5, 5]	0.0189	-
FCR	[20, 20, 20]	0.2571	7.05
	[10, 10, 10]	0.0639	14.48
	[5, 5, 5]	0.0163	13.76
FCR-shared cache	[20, 20, 20]	0.2554	7.69
	[10, 10, 10]	0.0640	14.33
	[5, 5, 5]	0.0161	14.76
OTF	[20, 20, 20]	0.2460	11.07
	[10, 10, 10]	0.0568	23.96
	[5, 5, 5]	0.0145	23.28
OTF-shared cache	[20, 20, 20]	0.2393	13.49
	[10, 10, 10]	0.0559	25.23
	[5, 5, 5]	0.0133	29.63

Table 6: Detailed Settings for Each Computational Method

Method	Detailed Settings
FCR	[20, 20, 20], alpha=2, T=50
	[10, 10, 10], alpha=2, T=50
	[5, 5, 5], alpha=2, T=50
FCR-shared cache	[20, 20, 20], alpha=2, T=50
	[10, 10, 10], alpha=2, T=50
	[5, 5, 5], alpha=2, T=50
OTF	[20, 20, 20], amp_rate=2, refresh_rate=0.15, T=50
	[10, 10, 10], amp_rate=2, refresh_rate=0.15, T=50
	[5, 5, 5], amp_rate=2, refresh_rate=0.15, T=50
OTF-shared cache	[20, 20, 20], amp_rate=2, refresh_rate=0.15, T=50
	[10, 10, 10], amp_rate=2, refresh_rate=0.15, T=50
	[5, 5, 5], amp_rate=2, refresh_rate=0.15, T=50

- OTF-shared cache shows the highest time acceleration at the smallest setting ([5, 5, 5]), with a 29.63% improvement, highlighting its efficiency for smaller datasets.
- Both **FCR** and **FCR-shared cache** methods consistently perform better than FBL, especially in medium and smaller settings, indicating significant improvements in training time.
- The time acceleration percentages are notably higher in the smallest setting for all methods, suggesting that these optimized approaches are particularly effective for smaller datasets.

#### A.1.3 Runtime Memory Reduction

Table 7: Comparison of Runtime Memory Reduction Across Computational Methods. RM stands for Runtime Memory, ar for amp\_rate/alpha, and rr for refresh\_rate.

Method	Setting	RM (MB)	Reduction (%)
FBL	[20]*3	6.33	0.00%
	[10]*3	4.70	0.00%
	[5]*3	4.59	0.00%
FCR	[20]*3, ar=2, T=50	2.69	57.46%
	[10]*3, ar=2, T=50	2.11	55.04%
	[5]*3, ar=2, T=50	1.29	71.89%
FCR-SC	[20]*3, ar=2, T=50	4.42	30.13%
	[10]*3, ar=2, T=50	2.62	44.15%
	[5]*3, ar=2, T=50	1.66	63.79%
OTF	[20]*3, ar=2, rr=0.15,	4.13	34.80%
	T=358		
	[20]*3, ar=2, rr=0.15,	4.19	33.86%
	T=50		
	[10]*3, ar=2, rr=0.15,	1.87	60.07%
	T=50		
	[5]*3, ar=2, rr=0.15,	0.32	93.02%
	T=50		
OTF-SC	[20]*3, ar=2, rr=0.15,	1.41	77.68%
	T=50		
	[10]*3, ar=2, rr=0.15,	0.86	81.58%
	T=50		
	[5]*3, ar=2, rr=0.15,	0.67	85.29%
	T=50		

The reduction percentage is calculated by comparing the runtime memory usage of a method under a specific setting to the baseline method's usage at the same setting. The formula used to compute the reduction is given by:

$$\label{eq:Reduction} \text{Reduction } (\%) = \left(\frac{\text{Base RM} - \text{Current RM}}{\text{Base RM}}\right) \times 100\%$$

where:

- Base RM represents the runtime memory (in MB) of the baseline method for a given setting, which is the FBL method in this case.
- Current RM represents the runtime memory of the current method under the same setting.

• Baseline Setting (FBL): All calculations are compared against the FBL method as the baseline, where the reduction is always 0% because it is compared against itself.

#### • FCR Method:

- For the setting [20, 20, 20], memory reduction is 57.46%, showing significant effectiveness of this method at larger settings.
- For the settings [10, 10, 10] and [5, 5, 5], memory reductions are 55.04% and 71.89% respectively, indicating effective memory reduction across various settings, especially at smaller scales.

#### • FCR-shared cache Method:

- This method typically uses more memory than FCR in a shared cache situation but still shows significant memory reductions, especially in the setting [5, 5, 5], with a reduction of 63.79%.

#### • OTF Method:

- For the setting [20, 20, 20], whether at a longer cycle (T=358) or a shorter cycle (T=50), memory reductions are around 34%, indicating that the change in response time does not significantly affect memory reduction.
- At smaller settings ([5, 5, 5]), memory is reduced by up to 93.02%, showing extremely high efficiency of the OTF method in handling small-scale tasks.

#### • OTF-shared cache Method:

This method shows very high memory reduction rates across all settings, particularly in [5, 5, 5] with a reduction of 85.29%, and 77.68% in [20, 20, 20], indicating that shared caching significantly enhances memory efficiency in the OTF method.

# A.1.4 GPU Usage Reduction

Table 8: GPU Usage for  ${\tt ogbn-arxiv}$  (edges)

Threshold	Sparse Graph	Resampled Graph	Total Storage
10	86,053	948,495	1,034,548
20	223,031	431,520	$654,\!551$
30	$333,\!599$	223,305	556,904
40	$423,\!513$	128,715	$552,\!228$
50	491,462	85,125	$576,\!587$
60	$543,\!682$	62,415	606,097
70	585,752	48,810	$634,\!562$
80	$619,\!468$	$39,\!855$	659,323
90	649,094	$33,\!255$	682,349
100	674,965	28,290	$703,\!255$
110	698,646	24,225	722,871
120	720,183	20,940	741,123
130	$738,\!536$	18,420	$756,\!956$
140	755,044	16,335	771,379

Table 9: GPU Usage Simulation for  $\operatorname{\mathsf{ogbn-products}}$  (edges)

Threshold	Sparse Graph	Resampled Graph	Total Storage
10	92,784	30,942,345	31,035,129
20	$491,\!202$	$26,\!278,\!605$	26,769,807
30	1,122,314	23,034,195	$24,\!156,\!509$
40	1,894,456	20,644,980	22,539,436
50	2,781,440	18,723,690	21,505,130
60	3,797,306	17,072,610	20,869,916
70	4,943,892	15,588,690	$20,\!532,\!582$
80	$6,\!230,\!698$	$14,\!219,\!115$	20,449,813
90	7,697,986	12,920,310	20,618,296
100	9,327,842	11,708,340	21,036,182
110	11,063,622	10,622,880	$21,\!686,\!502$
120	12,882,860	9,671,610	$22,\!554,\!470$
130	14,737,002	8,839,755	$23,\!576,\!757$
140	$16,\!667,\!254$	8,088,615	24,755,869

Table 10: GPU Usage Simulation for ogbn-mag (edges)

Threshold	Sparse Graph	Resampled Graph	Total Storage
10	432,912	4,934,910	5,367,822
20	1,335,788	2,321,400	3,657,188
30	2,059,531	1,307,445	3,366,976
40	2,610,841	813,090	3,423,931
50	3,055,313	$533,\!295$	3,588,608
60	3,422,953	363,120	3,786,073
70	3,717,725	$257,\!385$	3,975,110
80	3,962,156	187,365	4,149,521
90	4,159,659	140,670	4,300,329
100	4,322,855	107,835	4,430,690

# A.2 PyTorch Experiments

Dataset	Nodes	Edges	Features	Classes
PubMed	19,717	44,338	500	3
Cora	2,708	5,429	1,433	7
CiteSeer	3,312	4,732	3,703	6

Table 11: Comparison of Datasets

Operation	Duration (s)	Simulation Frequency
Simulated Disk Read	5.0011	0.05
Simulated Disk Write	1.0045	0.05
Simulated Cache Access	0.0146	0.05
In-Memory Computation	Real Computation	Real Computation

Table 12: Simulation Durations and Frequencies

Operation	k_h_sampling	k_hop_retrieval	k_h_resampling
Disk Read	✓		✓
Disk Write	✓		✓
Memory Access		✓	

Table 13: Function Access Patterns for Different Operations

Table 14: Experimental Settings - Setting 1

Dataset	Alpha	Presampled	Resampled	Sampled Depth
CiteSeer	0.1,, 0.9	100	40	1, 2, 3, 4
Cora	0.1,, 0.9	100	40	1, 2, 3, 4
PubMed	0.1,, 0.9	100	40	1, 2, 3, 4

Table 15: Experimental Settings - Setting 2  $\,$ 

Dataset	Alpha	Presampled	Resampled	Sampled Depth
CiteSeer	0.1,, 0.9	20	10	1, 2, 3, 4
Cora	0.1,, 0.9	20	10	1, 2, 3, 4
PubMed	0.1,, 0.9	20	10	1, 2, 3, 4

Table 16: IOCostOptimizer Functionality Overview

Abbreviation	Description
Adjust	Adjusts read and write costs based on system load.
Estimate	Estimates query cost based on read and write oper-
	ations.
Optimize	Optimizes query based on context ('high_load' or
	'low_cost').
Modify Load	Modifies query for high load optimization.
Modify Cost	Modifies query for cost efficiency optimization.
Log	Logs an I/O operation for analysis.
Get Log	Returns the log of I/O operations.

Table 17: BufferManager Class Methods

Method	Description
init	Initialize the buffer manager with capacity.
load	Load data into the buffer.
get	Retrieve data from the buffer.
store	Store data in the buffer.