# Marius: Learning Massive Graph Embeddings on a Single Machine

Paper # 143

## Abstract

We propose a new framework for computing the embeddings of large-scale graphs on a single machine. A graph embedding is a fixed length vector representation for each node (and/or edge-type) in a graph and has emerged as the de-facto approach to apply modern machine learning on graphs. We identify that current systems for learning the embeddings of large-scale graphs are bottlenecked by data movement, which results in poor resource utilization and inefficient training. These limitations require state-of-the-art systems to distribute training across multiple machines. We propose Marius, a system for efficient training of graph embeddings that leverages partition caching and buffer-aware data orderings to minimize disk access and interleaves data movement with computation to maximize utilization. We compare Marius against two state-of-the-art industrial systems on a diverse array of benchmarks. We demonstrate that Marius achieves the same level of accuracy but is up to one order-of magnitude faster. We also show that Marius can scale training to datasets an order of magnitude beyond a single machine's GPU and CPU memory capacity, enabling training of configurations with more than a billion edges and 550GB of total parameters on a single AWS P3.2xLarge instance.

## 1   Introduction

Graphs are used to represent the relationships between entities in a wide array of domains, ranging from social media and knowledge bases [36, 8] to protein interactions [3]. Moreover, complex graph analysis has been gaining attention in neural network-based machine learning with applications in clustering [28], link prediction [37, 30], and recommendation systems [35]. However, to apply modern machine learning on graphs one needs to convert discrete graph representations (e.g., traditional edge-list or adjacency matrix) to continuous vector representations [11]. To this end, learnable *graph embedding* methods [10, 5, 33] are used to assign each node (and/or edge) in a graph to a specific continuous vector representation such that the structural properties of the graph (e.g., the existence of an edge between two nodes or their proximity due to a short path) can be ap-
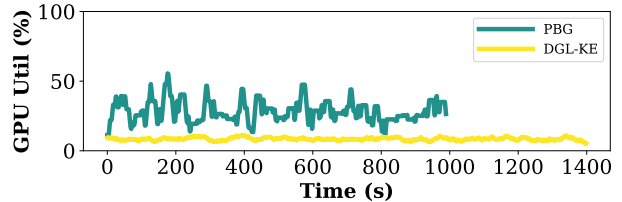


Figure 1: The GPU utilization of DGL-KE and PBG for one training epoch of ComplEx embeddings on the Freebase86m knowledge graph.

proximated using these vectors. In general, graph embedding models aim to capture the global structure of a graph and are complementary to graph neural networks (GNNs) [18]. Graph embedding models are primarily used in *link prediction* tasks and can also be used to obtain vector representations that form the input to GNNs.

However, learning a graph embedding model is a resource intensive process. First, training of graph embedding models can be compute intensive: many graph embedding models assign a high-dimensional continuous vector to each node in a graph [2, 34, 31], for example, it is common to assign a 200-dimensional continuous vector to each node [17, 38]. Consequently, the computational capabilities of GPUs and optimization methods such as mini-batch Stochastic Gradient Descent (SGD) are needed to accelerate training. Second, graph embedding models are memory intensive: the model from our previous example needs 800 bytes of storage per node and requires more than 32GB (the largest GPU memory) for a modest 50 million vertex graph. Thus, it is necessary to store the learnable parameters on off-GPU memory. Third, the training of graph embedding models requires optimizing over loss functions that consider the edges of the graph as training examples (e.g., the loss can enforce that the cosine similarity between the vector representations of two connected nodes is close to one, see Section 2.1) making training IO-bound for models that do not fit in GPU memory. This limitation arises due to irregular data accesses imposed by the graph structure. As a result, training of large graph embedding models is a non-trivial challenge.

Due to the aforementioned factors, scaling graph embedding training to instances that do not fit in GPU memory introduces costly data movement overheads that

1

can result in poor resource utilization and slow training. In fact, current state-of-the-art systems, including DGL-KE [38] from Amazon, and Pytorch BigGraph (PBG) [17] from Facebook, exhibit poor GPU utilization due to these overheads. Figure 1 shows the GPU utilization during a training epoch when using a single GPU for DGL-KE and PBG. As shown, DGL-KE only utilizes 10% of the GPU, and average utilization for PBG is less than 30%, dropping to zero during data movement.

GPU under-utilization can be attributed to how these systems handle data movement: To support out-of-GPU-memory training, DGL-KE stores parameters in CPU memory and uses synchronous GPU-based training over minibatches. However, the core computation during graph embedding training corresponds to dot-product operations between vectors (see Section 2), and thus, data transfers dominate the end-to-end run time. Moreover, DGL-KE is fundamentally limited by CPU memory capacity. To address this last limitation, PBG uses a different approach for scaling to large graphs. PBG partitions the embedding parameters into disjoint, node-based partitions (Figure 3) and stores them on external storage (e.g., SSDs) where they can be accessed sequentially. Partitions are then loaded from storage and sent to the GPU where training proceeds synchronously. Doing so avoids copying data from the CPU memory for every batch, but results in GPU underutilization when partitions are swapped. This problem is exacerbated if the storage device has low read throughput when compared to GPU processing. Thus, to scale to large instances both systems opt for distributed training over multiple compute nodes, making training resource hungry. However, the problems these systems face are not insurmountable and can be mitigated. *We show that one can train embeddings on billion-edge graphs using just a single machine.*

We introduce a new pipelined training architecture that can interleave data access, transfer, and computation to achieve high utilization. In contrast to prior systems, our architecture results in high GPU utilization throughout training: for the same workload shown in Figure 1, our approach can achieve an average $\sim 70\%$ GPU utilization while achieving the same accuracy (see Section 5).

To achieve this utilization, our architecture introduces asynchronous training of nodes with *bounded staleness*. We combine this with synchronous training for edge embeddings to handle graphs that may contain edges of different types, for example knowledge graphs where an edge may capture different relationship types. Specifically, we consider learning a separate vector representation for each edge-type. For clarity, we refer to edge-type embeddings as *relation embeddings*. This is because updates to the embedding vectors for nodes are sparse and therefore well suited for asynchronous training. However due to the small number of edge-types in

real-world graphs $(10, 000s)$, updates to relation embedding parameters are dense and require synchronous updates for convergence. We design the pipeline to maintain and update node embedding parameters in CPU memory asynchronously, allowing for staleness, while keeping and updating relation embeddings in GPU memory synchronously. Using this architecture, we can train graph embeddings for a billion-edge Twitter graph *one order of magnitude faster* than state-of-the-art industrial systems for the same level of accuracy: Using a single GPU, our system requires 3.5 hours to learn a graph embedding model over the Twitter graph. For the same setting, DGL-KE requires 35 hours.

To scale training beyond CPU memory, we propose a partition buffer to hide and reduce IO from swapping of partitions. Partitions of the graph and embedding parameters are kept on disk. They are swapped into a partition buffer in CPU memory and then used by the training pipeline. Our partition buffer supports pre-fetching and async writes of partitions to hide waiting for IO, resulting in a reduction of training time by up to $2\times$. Further, we observe that the order in which edge partitions are traversed can impact the number of IOs. Thus, we introduce a *buffer-aware* ordering that uses knowledge of the buffer size and what resides in it to minimize number of IOs. We show that this ordering achieves close to the lower bound and provides benefits when compared to locality-based orderings such as Hilbert space-filling curves [15].

Our design is implemented in Marius, a graph embedding engine that can train billion-edge graphs on a single machine. Using one AWS P3.2xLarge instance, we demonstrate that Marius improves utilization of computational resources and reduces training time by up to an order of magnitude in comparison to existing systems. Marius is 10x faster than DGL-KE on the Twitter graph with 1.46 billion edges, reducing training times from 35 hours to 3.5 hours. Marius is $1.5\times$ faster than PBG on the same dataset. On Freebase86m with 86 million nodes and 338 million edges, Marius trains embeddings 3.8x faster than PBG, reducing training times from 8.5 hours to 2.3 hours. We also show that Marius can scale to configurations where the parameter overhead exceeds CPU and GPU memory by an order of magnitude, training a configuration with 550GB of total parameters, $35\times$ and $9\times$ larger than CPU and GPU memory respectively.

## 2 Preliminaries

We first discuss necessary background on graph embeddings and related systems. Then, we review challenges related to optimizing data movement for training large scale graph embedding models. These are the challenges that this work addresses.
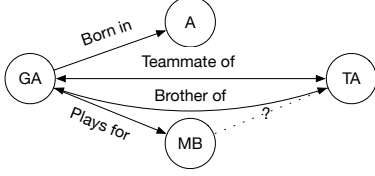
Figure 2: A sample knowledge graph.

## 2.1 Background and Related Work

**Graphs with Multiple Edge Types** We focus on graphs with multiple edge types defined as $G = (V, R, E)$ where $V$ is the set of nodes, $R$ is a set of edge-types or *relations*, and $E$ is the set of edges. Each edge $e = (s, r, d) \in E$ is defined as a triplet containing a source node, relation, and destination node. An example of such a graph is a *knowledge graph*, e.g., Freebase [9]. Here, the source node in a triplet defines a subject (an entity), the relation defines a predicate, and the destination node an object (an entity or literal) (see example in Figure 2. Knowledge graphs are commonly used both in industry and academia to represent real-world facts.

**Graph Embedding Models** A graph embedding is a fixed length vector representation for each node (and/or edge-type) in a graph. That is, each node and relation is represented by a corresponding $d$-dimensional vector $\theta$, also known as an *embedding* [11]. There are $d(|V|+|R|)$ total learnable parameters. To learn these vector representations, embedding models rely on *score functions* that capture structural properties of the graph. We denote the score function $f(\theta_s, \theta_r, \theta_d)$ where $\theta_s, \theta_r, \theta_d$ are the vector representations of the elements of a triplet $e = (s, r, d)$. For example, a score function can be the scaled dot product $f(\theta_s, \theta_r, \theta_d) = \theta_s^T \text{diag}(\theta_r)\theta_d$ with the requirement that the parameter vectors are such that $f(\theta_s, \theta_r, \theta_d) \approx 1.0$ if nodes $s$ and $d$ are connected via an edge of type $r$ and $f(\theta_s, \theta_r, \theta_d) \approx 0.0$ otherwise. There are several score functions proposed in the literature ranging from linear score functions [2, 21] to dot products [34, 31, 25] and complex models [12, 11].

Score functions are used to form loss function for training. The goal is to maximize $f(\theta_s, \theta_r, \theta_d)$ if $e \in E$ and minimize it if $e \notin E$. Triplets that are not present in $E$ are known as *negative edges*. A standard approach [38, 17] is to use the the score function $f(\theta_s, \theta_r, \theta_d)$ to form the cross entropy loss:

$$\mathcal{L} = - \sum_{s,r,d \in E} (f(\mathbf{e}_\theta) + \log(\sum_{s',r',d' \notin E} e^{f(\mathbf{e}_\theta')} + e^{f(\mathbf{e}_\theta)})) \quad (1)$$

where $\mathbf{e}_\theta = (\theta_s, \theta_r, \theta_d)$ and $\mathbf{e}_\theta' = (\theta_s', \theta_r', \theta_d')$.

The first summation term is over all true edges in the graph and the second summation is over all negative edges. There are a total of $|V|^2|R| - |E|$ negative edges in a knowledge graph; this makes it computationally infeasible to perform the full summation and thus is commonly approximated by *negative sampling*, in which a set of negatives edges is generated by taking a (typically uniform) sample of nodes from the graph for each positive edge. With negative sampling the term in the logarithm is approximated as $\sum_{s,r,d' \in N_e} e^{f(\mathbf{e}_\theta')} + e^{f(\mathbf{e}_\theta)})$. Where $N_e$ is the set of negative samples for $e$.

Graph embeddings are commonly used for *link prediction*, where the similarity of two node vector representations is used to infer the existence of a missing edge in a graph. For example, in the knowledge graph in Figure 2 we can use the vector representation of $TA$ and $MB$ and the relation embedding for *plays-for* to predict the existence of the edge $TA \xrightarrow{plays\text{-}for} MB$, marked with a questionmark in the figure.

**The Need for Scalable Training** The largest publicly available multi-relation graphs have hundreds of millions of nodes and tens of thousands of relations [32] (Table 1). Companies have internal datasets which are an order of magnitude larger than these, e.g., Facebook has over 3 billion users [4]. Learning a 400-dimensional embedding for each of the users will require the ability to store and access 5TB of embedding parameters efficiently, far exceeding the CPU memory capacity of the largest machines. Furthermore, using a larger embedding dimension has been shown to improve overall performance on downstream tasks [29]. For these two reasons it is important that a system for learning graph embeddings can scale beyond the limitations of GPU and CPU memory.

**Scaling Beyond GPU-Memory** We review approaches for scaling the training of graph embedding models out of GPU memory. Prior works follow in two categories: 1) Methods that leverage CPU Memory to store embedding parameters, and 2) Methods that leverage Block Storage and Partitioning of the model parameters. We discuss these two approaches in turn.

Following the first approach, systems such as DGL-KE [38] and GraphVite [39], store node embedding parameters in CPU memory and relation embedding parameters in GPU memory. Shown in Algorithm 1, training is performed synchronously and batches are formed and transferred on-demand. While synchronous training is beneficial for convergence, it is resource inefficient. The GPU will be idle while waiting for the batch to be formed and transferred; furthermore, gradient updates also need to be transferred from the GPU to CPU memory and applied to the embedding table, adding additional delays. The effect of this approach on utilization can be seen in Figure 1, where DGL-KE on average only utilizes about 10% of the GPU. This approach is also fundamentally limited by the size of the CPU memory, preventing the training of large graph embedding models.

**Algorithm 1:** Synchronous Embedding Training

**for** $i$ *in* $range(num\_batches)$ **do**

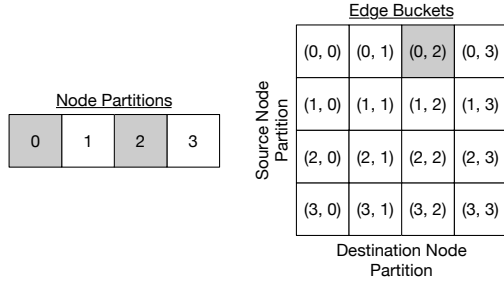| | |
|---|---|
| 1 | $\mathbf{B_i} = \text{getBatchEdges}(i)$; |
| 2 | $\mathbf{\Theta_n} = \text{getCpuParameters}(\mathbf{B_i})$; |
| 3 | $\text{transferBatchToDevice}(B_i, \mathbf{\Theta_n})$; |
| 4 | $\mathbf{\Theta_r} = \text{getGpuParameters}(\mathbf{B_i})$; |
| 5 | $\mathbf{B_\theta} = \text{formBatch}(\mathbf{B_i}, \mathbf{\Theta_n}, \mathbf{\Theta_r})$; |
| 6 | $\mathbf{G_n}, \mathbf{G_r} = \text{computeGradients}(\mathbf{B_\theta})$; |
| 7 | $\text{updateGpuParameters}(\mathbf{B_i}, \mathbf{G_r})$; |
| 8 | $\text{transferGradientsToHost}(\mathbf{G_n})$; |
| 9 | $\text{updateCpuParameters}(\mathbf{B_i}, \mathbf{G_n})$; |



Figure 3: Partitions and edge-buckets with $p = 4$. All edges in edge-bucket $(0, 2)$ have a source node in node-partition 0 and a destination node in node-partition 2.

The second approach is adopted by PyTorch BigGraph (PBG) [17]. PBG uses uniform partitioning to split up node embedding parameters into $p$ disjoint partitions and stores them on a block storage device (see example in Figure 3). Edges are then grouped according the partition of their source and destination nodes into $p^2$ edge buckets, where all edges in edge bucket $(i, j)$ have a source node which has an embedding in the $i$-th partition, and the destination node which has an embedding in the $j$-th partition. A single epoch of training requires iterating over all edge buckets while swapping corresponding pairs of node embedding partitions into memory for each edge bucket. This approach enables scaling to embedding models that exceed CPU memory capacity.

The major drawback of partitioning is that partition swaps are expensive and lead to the GPU being idle while a swap is happening. In fact, utilization goes towards zero during swaps as shown in Figure 1. We find that PBG yields an average GPU utilization of 28%. To best utilize resources, a system using partitioning to scale beyond the memory size of a machine, will need to mitigate overheads that arise from swapping partitions.

## 2.2 Data Movement Challenges

We discuss how to optimize data movement and related challenges that Marius' architecture addresses; we discuss the architecture in detail in Sections 3 and 4.

**Traditional Optimizations for Data Movement** Pipelining is a common approach used in a number of system designs to overlap computation with data movement, thereby improving utilization [14, 24, 26]. Using an image classifier as an example, a simple pipeline will consist of multiple worker threads that pre-process training images in parallel, forming batches and transferring them to the GPU. Once on the GPU, batches of training data are pushed onto a queue, with a training process constantly polling the queue for new batches. By keeping the queue populated with new batches, the GPU will be well utilized.

In IO-bound applications, buffer management can also be used to prevent unnecessary IO by caching data in memory. Buffer management is well studied in the area of databases and operating systems and has been applied to a myriad of applications and workloads [27, 13]. When using a buffer, the order in which data is accessed and swapped impacts end-to-end performance. When the data access pattern exhibits good locality, buffer managers typically yield good performance. Additionally, if the ordering is known ahead of time the buffer manager may prefetch data items and use Belady's optimal cache replacement algorithm to minimize reduce IO [1].

In graph processing, locality-aware data layouts of graph edges have been shown to improve locality of accesses and performance of common graph algorithms such as PageRank [23]. One such data layout, utilizes Hilbert space filling curves to define an ordering over the adjacency matrix of the graph. The ordering produced is a 1D index that preserves the locality of the 2D adjacency matrix. Storing and accessing edges according to this index improves OS cache hit rates [23, 22].

**Challenges for Graph Embeddings** For large graphs and embedding sizes, graph embedding models can be multiple orders of magnitude larger than the GPU's memory capacity, a key difference from deep neural network models that typically fit in a single GPU. To design a pipeline for graph embedding training, not only will training data (formed by considering edges) have to be piped to the GPU but also the corresponding model parameters (the node and relation embeddings of the endpoints and the type of each edge). Furthermore, model updates need to be piped back from the GPU and applied to the underlying storage. By pipelining model parameters and updates, we introduce the possibility of stale parameters, which must mitigated (see Section 3).

Buffer management techniques paired with data orderings can be used to buffer partitions in CPU memory to reduce IO. However, we find that prior locality-aware replacement policies fall short and still result in IO bound training ( Section 5.3). To address this challenge we propose a nearly-optimal data replacement policy, referred to as *elimination ordering*, in Section 4.

4

# 3 Pipelined Training Architecture

We review Marius' pipelined architecture for training graph embedding models. We first discuss the overall design, then the details of each stage, and finally discuss how staleness arises due to interleaving computation with data movement and how we can mitigate it.

**Pipeline Design** Our architecture follows Algorithm 1 and divides its steps into a five-stage pipeline with queues separating each stage (Figure 4). Four stages are responsible for data movement operations, and one for model computation and in-GPU parameter updates. The four data movement stages have a configurable number of worker threads, while the model computation stage only uses a single worker to ensure that relation embeddings stored on the GPU are updated synchronously.

We now describe the different stages of the pipeline and draw connections to the steps in Algorithm 1:

**Stage 1: Load** This stage is responsible for loading the edges (i.e., entries that correspond to a pair of node-ids and the type of edge that connects them) and the corresponding node embedding vectors that form a batch of inputs used for training. The edge payload constructed in this stage includes the true edges appearing in the graph and a uniform sample of negative edges (i.e., fake edges) necessary to form the loss function in Equation 1 (Line(s): 1-2 in Algorithm 1).

**Stage 2: Transfer** The input to this stage consists of the the edges (node-id and edge-type triples) and the node embeddings from the previous stage. Worker threads in this stage asynchronously transfer data from CPU to GPU using `cudaMemCpy` (Line 3 in Algorithm 1).

**Stage 3: Compute** The compute stage is the only stage that does not involve data movement. This stage takes place on GPU where the payload of edges and node embeddings created in Stage 1 is combined with relation embedding vectors (corresponding to the edge-type associated with each entry) to form a full batch. The worker thread then computes model updates and applies updates to relation embeddings stored in the GPU. The updates to node embeddings (i.e., the scaled gradients that need to be added to the previous version of the node embedding parameters) are placed on the output queue to be transferred from GPU memory (Lines: 4-7 in Algorithm 1)

**Stage 4: Transfer** The node embedding updates are transferred back to the CPU. We use similar mechanisms as in Stage 2 (Line(s): 8 in Algorithm 1)

**Stage 5: Update.** The final stage in our pipeline applies node embedding updates to stored parameters in CPU memory (Line(s): 9 in Algorithm 1).

This hybrid-memory architecture allows us to execute sparse parameter updates asynchronously (i.e., the node embedding parameter updates) and dense updates (i.e., the relation embedding parameter updates) synchronously, and optimize resource utilization as we show experimentally in Section 5.

**Bounded Staleness** The main challenge with using a pipelined design as described above, is that it introduces *staleness* due to asynchronous processing. To illustrate this, consider a batch entering the pipeline (Stage 1) with the embedding for node $A$. Once this batch reaches the GPU (Stage 3), the gradients for the embedding for $A$ will be computed. While the gradient is being computed, consider another batch that also contains the embedding for node $A$ entering the pipeline (Stage 1). Now, while the updates from the first batch are being transferred back to the CPU and applied to parameter storage, the second batch has already entered the pipeline, and thus it contains a stale version of the embedding for node $A$.

To limit this staleness, we bound the number of batches in the pipeline at any given time. For example, if the bound is $4$, embeddings in the pipeline will be at worst 4 updates behind. However, due to the sparsity of node embedding updates, it is unlikely a node embedding will even become stale. To give a realistic example, take the Freebase86m graph which has 86 million nodes. A typical batch size and staleness bound for this benchmark is 10,000 and 16 respectively. Each batch of 10,000 edges will have at most 20,000 node embeddings and given this staleness bound there can be at most 320,000 node embeddings in the pipeline at any given time, which is just about .4% of all node embeddings. Even with this worst case, only a very small fraction of node embeddings will be operated on at a given time. The same property does not hold for relation embeddings since there are very few of them (15K in Freebase86m), hence our design decision to keep relation embeddings in GPU memory and update them synchronously, bypasses the issue of stale relation embeddings. We study the effect of staleness and Marius' performance as we vary the aforementioned bound in Section 5.5.

# 4 Out-of-memory Training

As described in Section 2.1, to learn embedding models for graphs that do not fit in CPU memory, existing systems partition the graph into non-overlapping blocks. They correspondingly partition the parameters as well so that they can be loaded sequentially for processing. However as IO from disk can be slow (e.g, a partition can be around 100s of MB in size), it is desirable to hide the IO wait times and minimize the number of swaps from disk to memory. In this section, we describe how we can effectively hide IO wait time by integrating our training pipeline with a *partition buffer* that constitutes an in-memory cache of partitions. We also describe how we can minimize the number of swaps from disk to memory
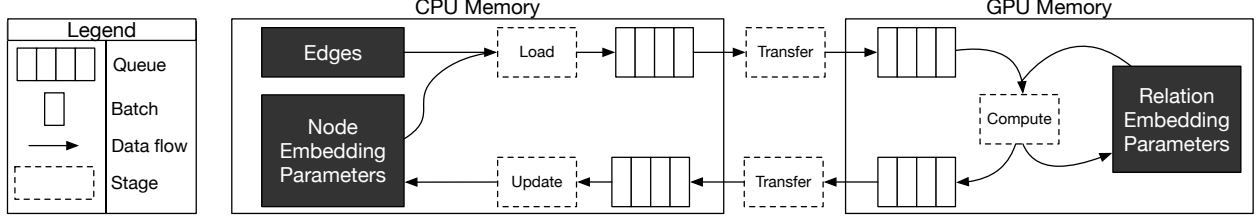
Figure 4: Marius training pipeline.

<div style="display:flex">

**Algorithm 2:** Training Using a Partition Buffer

1   Buffer = {};
2   **for** *k in range*($p^2$) **do**
3     $\mathbf{E_{ij}}, i, j$ = getEdgeBucket(Ordering[$k$]);
4     **if** *i not in Buffer* **then**
5       **if** *Buffer.size() == c* **then**
6         Buffer.evictFurthest(Ordering, $k$);
7       Buffer.admit($i$);
8     **if** *j not in Buffer* **then**
9       **if** *Buffer.size() == c* **then**
10         Buffer.evictFurthest(Ordering, $k$);
11       Buffer.admit($j$);
12     $\mathbf{\Theta_i}$ = Buffer.get($i$); // Source Node Partition
13     $\mathbf{\Theta_j}$ = Buffer.get($j$); // Destination Node Partition
14     trainEdgeBucket($\mathbf{E_{ij}}, \mathbf{\Theta_i}, \mathbf{\Theta_j}$);

</div>

by developing a new ordering for traversing graph data.

**Partition-based training** Consider a graph that is partitioned into $p^2$ edge buckets. Training one epoch requires iterating over all $p^2$ edge buckets, where each edge in a given bucket $(i, j)$, will have a source node in partition $i$ and destination node in partition $j$. When processing an edge bucket $(i, j)$, node partition $i$ and node partition $j$ must be present in the CPU partition buffer in order for learning to proceed using the pipelined training architecture (see Section 3). If either one is not present, it must be loaded from disk and *swapped* into the buffer, replacing an already present partition if the cache is full. Partition-based training is described in Algorithm 2.

Given a partitioned graph, there are a number of orderings that can be used for traversal. To minimize the number of times partitions need to be loaded from disk, *we seek an ordering over edge buckets which minimizes the number of required partition swaps*. We also observe that if we know the order in which edge buckets will be processed, we can further mitigate IO overhead by 1) prefetching to load node partitions as they are needed in the near future and 2) using a optimal cache eviction policy which removes partitions used farthest in the future.

We next discuss the problem of determining an optimal ordering over edge buckets and describe *elimination ordering*, a new ordering scheme that achieves near-optimal number of partition swaps.

## 4.1   Edge Bucket Orderings

We develop an edge bucket ordering scheme that minimizes the number of swaps. First, we first derive a lower bound on the number of swaps necessary to complete one training epoch for a cache of size $c$ and $p$ ($p >= c$) partitions. To derive the lower bound, we view an edge bucket ordering as a sequence of caches over time, where each item in the sequence describes what node partitions are in the buffer at that point. Each successive cache differs by one swapped partition as shown in Figure 5.

Given such a sequence, an edge bucket ordering can be constructed by processing edge bucket $(i, j)$ when partitions $i$ and $j$ are in the cache. For simplicity, we can do this the first time $i$ and $j$ appear together. Note that 1) $i$ and $j$ must appear together at least once otherwise no ordering over all edge buckets can be constructed, 2) self-edge buckets (i.e. $(i, i)$) can also be added to the ordering the first time $i$ appears in the cache as it does not contribute additional swaps, and 3) there are many edge bucket orderings with the same sequence of caches (depending on how the edge buckets in a particular cache are processed). Viewed in this light, *we seek the shortest (min. swaps) cache sequence where all node partitions pairs appear together in the cache at least once*.

**Lower bound** We assume that initializing the first full cache does not count as part of the total number of swaps. Thus there are $\frac{p(p-1)}{2}$ (the total number of pairs) minus $\frac{c(c-1)}{2}$ (the number of pairs we get in the first cache) remaining. On any given swap, the most new pairs we can cover is if the partition entering the cache has not been paired with anything already in the buffer (everything in the buffer has already been paired with everything else in the buffer). Thus, for each swap the best we can hope for is to get $c - 1$ pairs we have not already seen. With this in mind a lower bound on the minimum number of swaps required is $\lceil (p(p - 1)/2 - c(c - 1)/2)/(c - 1) \rceil$. We use this lower bound to evaluate the performance of different edge bucket orderings in the next section. We experimentally show that the new ordering strategy we propose is nearly optimal with respect to this bound.

**Elimination ordering** We describe elimination ordering, a edge bucket ordering that achieves close to optimal number of partition swaps and improves upon locality-
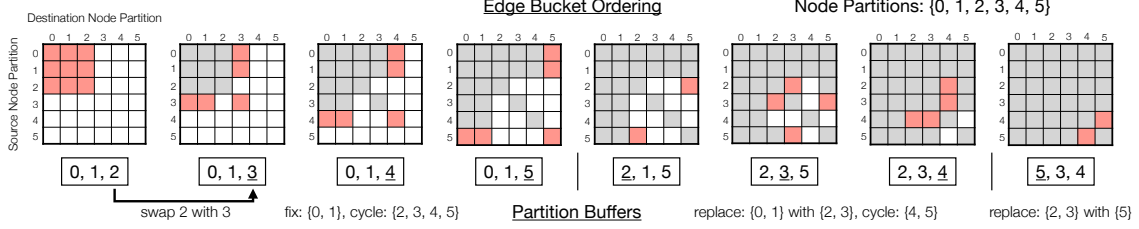
Figure 5: Example elimination ordering for $p = 6$ and $c = 3$. The sequence of partition buffers corresponds to first eliminating $\{0, 1\}$ then $\{2, 3\}$ and finally $\{5\}$. Each successive buffer differs by one swap. A corresponding edge bucket ordering is shown above the buffers. At each step, all previously unprocessed edge buckets which have their source and destination node partitions in the buffer are added to the ordering (red edge buckets). For each buffer, these edge buckets can be added in any order.



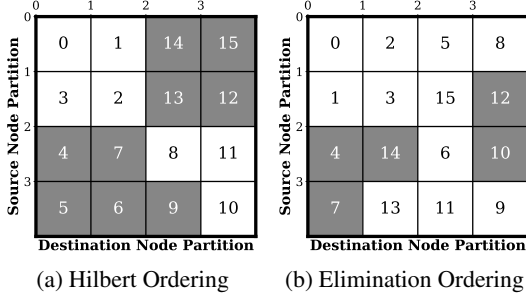(a) Hilbert Ordering  (b) Elimination Ordering

Figure 6: Hilbert and Elimination edge bucket orderings. Numbers indicate the order in which the bucket is processed. Gray cells indicate misses to the buffer.
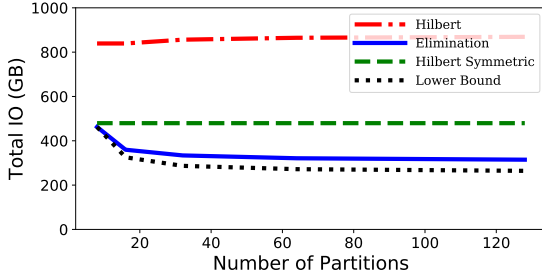


Figure 7: Simulated total IO performed during a single epoch of training Freebase86m with $d = 100$.

aware orderings such as Hilbert space-filling curves [15].

Consider a cache that was initialized with the first $c$ node-partitions of the graph. The key insight that guides our elimination ordering is to *swap in a partition that provides the highest number of new, previously unseen pairs*. We use this insight in the following fashion: once the first $\frac{c(c-1)}{2}$ edge buckets—created by considering pairs of node partitions in the cache—have been processed, we fix $c-1$ entries (node-partitions) in the cache, and bring in one new node partition $\hat{c}$ and process the pairs $(\hat{c}, c_i)$ and $(c_i, \hat{c}) \forall i \in [0, c-2]$. We note that all the $p - c$ partitions outside the cache are equally good candidates at this point, and we break ties randomly. Once

the pairs, formed by partition $\hat{c}$ and the $c - 1$ partitions in the cache, have been processed, we continue to fix the $c-1$ partitions and swap in the remaining $p - c - 1$ partitions, one at a time. With every swap we have $c - 1$ new pairs (edge buckets) to be processed. Once this is complete, all pairs for the fixed $c-1$ partitions have been processed and we need to refresh our cache with new node-partitions. We refer to this as one *round* of the algorithm. The previously-fixed $c - 1$ partitions can be retired and we can replace these $c - 1$ partitions by randomly sampling from the $p - c$ other partitions to form a new cache and then repeat the above steps. We show an example of this ordering in Figure 5. We also analyze the number of swaps generated by the above algorithm: given $p$ partitions and a cache of size $c$, with $x = \left\lfloor \frac{p-c}{c-1} \right\rfloor$ the number of swaps is $(p - c) + (x + 1) \left[(p - c) - x(c - 1)/2\right]$.

We observe that the elimination ordering has a number of useful properties that make it advantageous to implement in practice. Since all partitions are symmetrically processed we do not need to track any extra state or use any priority mechanisms other than knowing which partitions have been retired. Further, within a round for every disk IO, we have $c - 1$ new partition pairs to process meaning we can overlap most IO operations with compute. Finally, the only bottleneck arises in between rounds when $c - 1$ partitions are replaced but this only happens at most $\left\lfloor \frac{p-c}{c-1} \right\rfloor + 1$ number of times in one epoch.

**Comparison with Hilbert, lower bound** We illustrate how the elimination ordering compares to a Hilbert space-filling curve on a small $p = 4, c = 2$ case in Figure 6. We see that while the Hilbert ordering has nine cache misses the elimination ordering only has five misses. We also performed simulations that compare the number of IO operations incurred by the elimination algorithm with a random ordering, a space-filling curve based ordering, and the analytical lower bound. Figure 7 shows the number of IO accesses for each of them when varying $p$ and using a cache with $\frac{p}{4}$ partitions. Space

filling curve orderings like Hilbert attempt to define an ordering that preserves 2D locality. A key advantage of the elimination ordering when compared to these methods is that it is cache-aware, i.e., the algorithm knows the cache size and specifically aims to minimize partition swaps. In contrast, orderings like Hilbert are unaware of this information, aiming instead for locality. Figure 7 also shows that elimination ordering yields nearly optimal performance across partition configurations.

We leave an investigation of a provably-optimal ordering for future work. Our initial studies have shown that there exist cases of $p$, $c$ where no valid ordering can match the lower bound as well as cases where an ordering which requires slightly less swaps than the elimination ordering does exist. Thus,the optimal algorithm requires IO-swaps somewhere between the lower bound and the elimination ordering in Figure 7.

## 4.2 Partition Buffer

We next describe mechanisms that we use in the partition buffer to further minimize IO overhead. The partition buffer is a fixed sized memory region that has capacity to store $c$ embedding partitions in memory. We *co-design* the buffer replacement policies with the elimination ordering described above. Co-designing the edge traversal with buffer replacement policy means that we have knowledge about which partitions will be accessed in the future. This allows the buffer to use the optimal replacement policy: *evict the partition that will be used furthest in the future* [1]. Given this policy we also design a prefetching mechanism that can minimize the amount of time spent waiting for partitions to swap. Again, based on knowing the order in which partitions are used, we use a prefetching thread that reads the next partition in the background. Correspondingly when a partition needs to be evicted from memory, we perform asynchronous writes using a background writer thread.

## 5 Evaluation

We evaluate Marius on standard benchmarks using a single AWS P3.2xLarge instance and compare against SoTA graph embedding systems. We show that:

(1) Due to optimized resource utilization, Marius yields up to 10× faster training in comparison to SoTA systems.

(2) Elimination ordering reducesIO required by up to 2× when compared to other locality-based graph orderings, thus alleviating the IO bottleneck during training.

(3) Marius is able to scale to graph embedding models that rely on increased vector dimensions to achieve higher accuracy. Due to the increased vector dimensions these models exceed CPU memory size. For instance, we show that Marius can learn an embedding model using 800-d vector representations on a graph with 86M nodes on a single machine. In this configuration there are 550GB of total parameters and optimizer state, which is 35× GPU memory size and 9× CPU memory size.

## 5.1 Setup

**Implementation:** Marius is implemented in about 10,000 lines of C++. We use LibTorch [20], which is the C++ API of Pytorch, as the underlying tensor engine. LibTorch provides access to the wide-ranging functionality of PyTorch, making it easy to extend Marius to support more complex embedding models. We also implement an abstracted storage API, which allows for embedding parameters to be stored and accessed across a variety of backends under one unified API. This allows us to easily switch between storage backends, say from using a CPU memory-based backend to a disk-based backend.

**Hardware Setup:** All experiments are run on a single AWS P3.2xLarge instance which has: 1 Telsa V100 GPU with 16GB of memory, 8 vCPUs with 61 GB of memory, and an attached EBS volume with 400MBps of read and write bandwidth. The only exception is for running DGL-KE on the Twitter and Freebase86m datasets, where DGL-KE runs out of memory when training. We use a larger machine for these cases, which has: 1 Telsa V100 GPU with 32GB of memory, 200 CPUs with 500 GB of memory, block storage is not used for these cases.

**Datasets:** For our evaluation, we use standard benchmark datasets that include social networks(Twitter [16], Livejournal [19]) and knowledge graphs (FB15K and Freebase86m [17, 38] derived from Freebase [9]). A summary of the dataset properties is shown in Table 1. FB15k uses an 80/10/10 train, validation and test split. All others use a 90/5/5 split.

**Embedding Models:** For comparisons on FB15k, we use ComplEx [31] and DistMult [34]. On LiveJournal and Twitter we use *Dot* [18], which is a dot product between the node embeddings of an edge. On Freebase86m we use ComplEx embeddings. These models are widely used and produce the best quality embeddings [17, 38].

**Hyperparameters:** To ensure fair comparisons, we use the same hyperparameters across each system instead of tuning separately. Hyperparameter values for each configuration were chosen based on those used in the evaluation of the DGL-KE and PBG and is shown in Table 1. All systems use the Adagrad optimizer [7] for training, which empirically yields higher-quality emebeddings over SGD. One drawback of using this optimizer is that it effectively requires storing a learning rate per parameter, doubling the overall memory footprint of the embeddings during training.

**Evaluation Task and Metrics:** We evaluate the quality of the embeddings using the link prediction task. Link

| Name | Type | |E| | |V| | |R| | Size | Hyperparameters |
|------|------|-----|-----|-----|------|-----------------|
| FB15k | KG | 592k | 15k | 1.3k | 52MB | $d = 400, lr = .1, b = 10^4, n_t = 10^3, \alpha_{n_t} = .5, \text{FilteredMRR}$ |
| LiveJournal | Social | 68M | 4.8M | - | 1.9GB | $d = 100, lr = .1, b = 5 \times 10^4, n_t = 10^3, \alpha_{n_t} = .5, n_e = 10^4, \alpha_{n_e} = 0$ |
| Twitter | Social | 1.46B | 41.6M | - | 33.2GB | $d = 100, lr = .1, b = 5 \times 10^4, n_t = 10^3, \alpha_{n_t} = .5, n_e = 10^3, \alpha_{n_e} = .5$ |
| Freebase86m | KG | 338M | 86.1M | 14.8K | 68.8GB | $d = 100, lr = .1, b = 5 \times 10^4, n_t = 10^3, \alpha_{n_t} = .5, n_e = 10^3, \alpha_{n_e} = .5$ |

Table 1: Datasets used for evaluation. The size column indicates total size of embedding parameters with the embedding dimension $d$, including the Adagrad optimizer state. $lr$: learning rate, $b$: batch size, $n_t$: training negatives, $\alpha_{n_t}$: train degree-based negatives fraction, $n_e$: evaluation negatives, $\alpha_{n_e}$: eval degree-based negatives fraction.

prediction is a commonly used evaluation task in which embeddings are used to predict if a given edge is present in the graph. Link prediction metrics reported are Mean Reciprocal Rank (MRR) and Hits@$k$, which are derived from the rank of the score of each candidate edge, where the scores are produced from the embedding score function $f$. For a given candidate edge $i$, it has a rank $r_i$ which denotes the position of the score of candidate edge in descending sorted array $S_i$, where $S_i$ contains the score of the candidate edge and the scores of a set of negative samples. Given this, the MRR and Hits@$k$ can be computed from a set of candidate edges $C$ as follows: $\frac{1}{|C|} \sum_{i \in C} \frac{1}{r_i}$ and $\frac{1}{|C|} \sum_{i \in C} \mathbb{1}_{r_i <= k}$ respectively.

Metrics can be filtered or unfiltered. Filtered evaluation involves comparing candidate edges with $|N|$ negative samples, produced by using all of the nodes in the graph. Some of the produced negative samples will be false negatives, which will not be used in filtered evaluation. Because all nodes in the graph are used, filtered evaluation is expensive for large graphs. Unfiltered evaluation samples $n_e$ nodes from the graph, with a fraction $\alpha_{n_e} n_e$ by degree and $(1 - \alpha_{n_e})n_e$ uniformly. False negatives are not removed in unfiltered evaluation, but will not be common if $n_e << |V|$. Unfiltered evaluation is much less expensive and is well suited for large scale graphs. We use filtered metrics only on FB15k and unfiltered metrics elsewhere. The same evaluation approach is adopted by prior systems [17].

## 5.2 Comparison with Existing Systems

To demonstrate that Marius utilizes resources better than current SoTA systems leading to faster training, we compare Marius with PBG and DGL-KE on four benchmark datasets. We do not compare with GraphVite since it is significantly slower than DGL-KE as reported in Zheng et al. [38]. FB15k and LiveJournal fit in the machine's GPU memory and therefore do not have data movement overheads. Twitter exceeds GPU memory which introduces data movement overheads from storing parameters off-GPU. Freebase86m exceeds the CPU memory of the machine, which prevents DGL-KE from training these embeddings on a single P3.2xLarge instance, therefore we only compare against PBG.

| System | Model | Filtered MRR | Hits | | Time (s) |
|--------|-------|--------------|------|------|----------|
| | | | @1 | @10 | |
| DGL-KE | ComplEx | .795 | .766 | .848 | 35.6s ± .69 |
| PBG | ComplEx | .795 | .736 | .888 | 40.3s ± .1 |
| Marius | ComplEx | .795 | .736 | .888 | **27.7s** ± .12 |
| DGL-KE | DistMult | .792 | .766 | .848 | 32.8s ± .88 |
| PBG | DistMult | .790 | .728 | .888 | 46.2s ± .46 |
| Marius | DistMult | .790 | .727 | .889 | **28.7s** ± .15 |

Table 2: FB15k Results. All systems reach peak accuracy at about the same number of epochs with 30 and 35 epochs for ComplEx and DistMult respectively.

**FB15k** In this experiment, we compare Marius with PBG and DGL-KE on FB15k to show that Marius achieves similar embedding quality as the other systems on a common benchmark. We measure the FilteredMRR, Hits@$k$, and runtime of the systems when training ComplEx and DistMult embeddings with $d = 400$ to peak accuracy, averaged over five separate runs. DGL-KE requires slightly more epochs to reach peak accuracy with Results are shown in Table 2. It should be noted that all parameters and training data fit in GPU memory for this dataset. We find that Marius achieves near identical metrics as PBG when learning the same embeddings, this is expected as both systems have similar implementations for sampling edges and negative samples. DGL-KE on the other hand only achieves a similar FilteredMRR. DGL-KE has implementation differences for initialization and sampling which likely account for the difference in metrics. While Marius is not designed for small knowledge graphs, we can see that it performs comparably to SoTA systems, achieving similar embedding quality in lesser time.

**LiveJournal** To show that the systems are comparable on social graphs, we compare the quality of 100-dimensional embeddings learned by the three systems using a dot product score function. While Livejournal is two orders of magnitude larger than FB15k, all parameters still fit in GPU memory with a total of $2GB$. As before, we measure MRR, hits@k, and runtime, averaging over three runs; but we use unfiltered MRR instead of FilteredMRR. We do so because FilteredMRR is computationally expensive to evaluate on larger graphs, as discussed in 5.1. Instead of using all nodes in the graph to construct negative samples, we sample 10,000 nodes

| System | Model | MRR | Hits | | Time |
| --- | --- | --- | --- | --- | --- |
| | | | @1 | @10 | (min) |
| PBG | Dot | .751 | .671 | .873 | 23.6m +-.17 |
| DGL-KE | Dot | .750 | .672 | .874 | 17.5m +-.03 |
| Marius | Dot | .749 | .671 | .871 | **12.5m +-.01** |

Table 3: LiveJournal Results. Marius and PBG are trained to 25 epochs, and DGL-KE to 17.

| System | Model | MRR | Hits | | Time |
| --- | --- | --- | --- | --- | --- |
| | | | @1 | @10 | |
| PBG | Dot | .313 | .239 | .451 | 5h15m |
| DGL-KE | Dot | .220 | .153 | .385 | 35h3m |
| Marius | Dot | .383 | .316 | .508 | **3h28m** |

Table 4: Twitter Results. Trained to 10 epochs.

| System | MRR | Runtime |
| --- | --- | --- |
| PBG (N=16) | .730 | 514m |
| Marius (N=16, PBG Sim) | .685 | 631m |
| Marius (N=16, C=4) | .685 | 162m |
| Marius (N=16, C=8) | .687 | 136m |

Table 5: Freebase86m with embedding size 100. PBG is simulated on Marius(no partition buffer, no prefetching, and with synchronous training). Trained to 10 epochs.

uniformly for evaluation, as done in [17]. Results are shown in Table 3. We see that all three systems achieve near identical metrics when learning embeddings for this dataset. There are slight differences in runtimes that can be attributed to implementation differences. PBG checkpoints parameters after each epoch, while this is optional in Marius and DGL-KE. Without checkpointing, PBG would likely achieve similar runtimes to DGL-KE and Marius. Overall, we find that Marius performs as well or better than SoTA systems on this social graph benchmark.

**Twitter** We now move on to evaluating Marius on large-scale graphs for which embedding parameters do not fit in GPU memory. The Twitter follower network has approximately 1.4 billion edges and 41 million nodes. We train 100-dimensional embeddings on each system using a dot product score function. We report results for one run for each system since we observed that training times and MRR are stable between runs. In total, there are $16GB$ of embedding parameters with another $16GB$ of optimizer state, since all systems use the Adagrad optimizer, as discussed above. To construct negatives for evaluation we use the approach from Zheng et al. [38], where 1,000 nodes are sampled uniformly from the graph, and 1,000 nodes are sampled by degree.

Unlike the previous two datasets, each system uses a different methodology for training embeddings beyond GPU memory sizes. DGL-KE uses the approach described in Algorithm 1, storing parameters in CPU memory and processing batches synchronously while waiting for data movement. PBG does not utilize CPU memory and instead uses the partitioning approach with 16 partitions. Marius stores parameters in CPU memory, and utilizes its pipelined training architecture to overlap data movement with computation.

We compare the peak embedding quality learned by each of the systems after ten epochs of training in Table 4. We find that Marius is able to train higher quality embeddings faster than the other systems, 10× faster than

DGL-KE and 1.5× faster than PBG. DGL-KE's long training times can be attributed to data movement wait times inherent in synchronous processing. PBG on the other hand, only pays a data movement cost when swapping partitions. PBG achieves comparable runtimes because this dataset has a large amount of edges relative to the total number of parameters, meaning that computation times dominate partition swapping times.

Turning our attention to the embedding quality, we find that Marius learns embedding of comparable quality to the next-best system: Marius yields an MRR of 0.383 versus 0.313 PBG. On the other hand, DGL-KE only achieves an MRR of 0.220. We attribute this gap in quality to to implementation differences between the systems (since all three use the same hyperparameters for training). This gap between the embedding quality of each system merits further study by modifying PBG and DGL-KE to use the same initialization and sampling methods. However, we do not examine this further since it is not the focus of our contributions.

**Freebase86m** We now evaluate Marius on a large-scale knowledge graph for which embedding parameters do not fit in CPU or GPU memory. We train 100-dimensional ComplEx embeddings for each system. In total, there are about $32GB$ of embedding parameters with another $32GB$ of Adagrad optimizer state. We do not evaluate DGL-KE on this dataset since it is unable to process this configuration on a single P3.2xLarge instance. For evaluation, we sample 1000 nodes uniformly and 1000 nodes based on degree as negative samples.

We compare the peak embedding quality of Marius and PBG where both systems are trained to 10 epochs in Table 5. Both systems use 16 partitions for training and in Marius we vary the number of partitions we hold in the CPU memory buffer. We find that Marius is able to train to peak embedding quality 4× faster when the buffer has a capacity of 8 partitions. We note that PBG achieves a higher embedding quality in Marius, however we do not believe the cause is due to our contributions. To demonstrate this, we simulate the training scheme of PBG in Marius by performing synchronous training and limiting the capacity of the buffer to two partitions. This simulated configuration has a comparable runtime to PBG but does not achieve the same embedding quality as PBG.
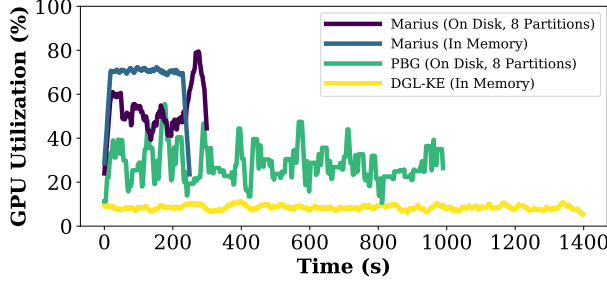
Figure 8: GPU utilization of Marius, DGL-KE and PBG during a single epoch of training $d = 50$ embeddings on Freebase86m. Utilization is smoothed over a 25-second window.
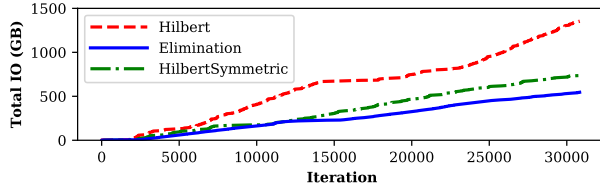


Figure 9: Total IO during a single epoch of training.

Because of this, we do not believe that our contributions impact the embedding quality and this gap between systems is likely due to an unidentified implementation difference. The runtime difference between the two systems can be attributed to the fewer number of partition swaps Marius performs and the ability to prefetch partitions.

**Utilization** We also include a comparison of GPU utilization during a single epoch of training $d = 50$ embeddings. Figure 8 shows the utilization of two configurations of Marius compared to DGL-KE and PBG. One configuration of Marius stores embeddings in CPU memory while the other uses eight partitions on disk with four partitions buffered in CPU memory. We see that Marius is able to utilize the GPU $8\times$ more than DGL-KE when training in memory and about $6\times$ more when using the partition buffer. Compared to PBG, our partition buffer design leads to nearly $2\times$ GPU utilization with fewer drops in utilization when waiting for partition swaps. In this configuration, Marius is utilizing $\sim 100\%$ of the CPUs.

## 5.3 Partition Orderings

We now evaluate our buffer-aware *Elimination* ordering and compare it to two Hilbert curve based orderings. The first, *Hilbert*, is the ordering generated directly from a Hilbert curve over the $n \times n$ matrix of edge buckets. The second, *HilbertSymmetric*, modifies the previous curve by processing edge buckets $(i, j)$ and $(j, i)$ together, which reduces the overall number of swaps that need to be performed by about $2\times$. All experiments use

| d | Size | Partitions | MRR | Runtime (Epoch) |
|---|---|---|---|---|
| 20 | 13.6GB | - | .665 | 5m |
| 50 | 34.4GB | - | .688 | 7.7m |
| 100 | 68.8GB | 32 | .691 | 12.9m |
| 400 | 275.2GB | 100 | .701 | 92.4m |
| 800 | 550.4GB | 250 | .701 | 396m |

Table 6: Freebase86m. $d = 400$ and $d = 800$ trained to 5 epochs, other cases are trained to 10. Buffer capacity is fixed to 16 partitions.

32 partitions and a buffer capacity of 8 partitions.

We compare the orderings on Freebase86m with $d = 50$ and $d = 100$ sized embeddings, where the latter configuration exceeds CPU memory size. For $d = 50$ we include an in-memory configuration which does not use partitioning as a baseline. Results are shown in Figure 10. We find that the *Elimination* ordering reduces training time to nearly in-memory speeds, while only keeping 1/4 of the partitions in memory at any given time. The runtime of the three orderings is directly correlated with the amount of IO required to train a single epoch. Since the *Hilbert* and *HilbertSymmetric* orderings require more IO, training stalls more often waiting for IO to complete. Results for $d = 100$, also in Figure 10, show that *Elimination* has the lowest training time, which is directly correlated with the amount IO performed. Overall, the *Elimination* ordering is well suited for training large-scale graph embeddings through reducing IO.

We also compare the orderings on Twitter with $d = 100$ and $d = 200$ sized embeddings. Results for $d = 100$ are shown in Figure 11. We find that the choice of ordering does not impact runtime for this configuration. Even though *Elimination* results in the smallest amount of total IO, the prefetching of partitions to the buffer always outpaces the speed of computation for the other orderings. We see this in Twitter and not Freebase86m, because Twitter has nearly $10\times$ the density of Freebase86m, i.e., more computation needs to be performed per partition. When we increase the embedding dimension to $d = 200$ (Figure 11) we see a difference in running time. By increasing the embedding dimension by $2\times$ we increase the total amount of IO by $2\times$, and now the prefetching of partitions is outpaced by the computation.

Overall, we see that certain configurations are *data bound* and others are *compute bound*. For data bound configurations like $d = 50$ and $d = 100$ on Freebase86m and $d = 200$ on Twitter, the choice of ordering will impact overall training time, with *Elimination* performing best. But for *compute bound* workloads such as $d = 100$ on Twitter, the choice of ordering makes little difference since the prefetching always outpaces computation.
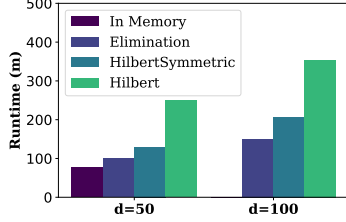
11

Figure 10: 10 epochs runtime per edge bucket ordering on Freebase86m.
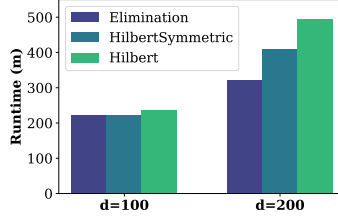

Figure 11: 10 epochs runtime per edge bucket ordering on Twitter.
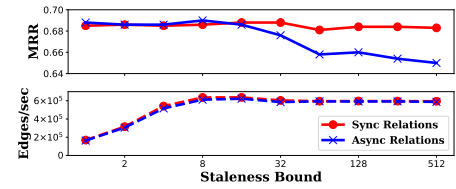

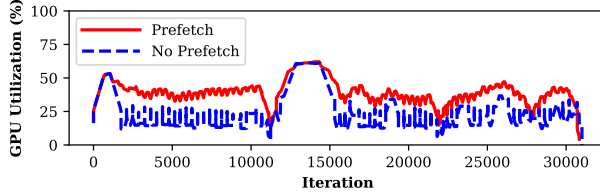Figure 12: Impact of staleness bound.


Figure 13: Effect of prefetching with Freebase86m.

## 5.4 Large Embeddings

We evaluate the ability of Marius to scale training beyond CPU memory sizes in this section. We vary the embedding dimension from a small dimension of $d = 20$, for which training fits in GPU memory, to a large embedding dimension $d = 800$, which is well beyond the memory capacity of a single P3.2xLarge instance. The results are shown in Table 6. We find that the embedding quality increases with increased embedding dimension. We also see that as the embedding size increases, the training time increases quadratically. We see this because the number of swaps and total IO scales quadratically with the number of partitions, if the buffer capacity is held fixed. And because training is bottlenecked by IO for large embedding sizes, we see quadratic runtime increases. It should be noted that with a faster disk we would observe improved runtimes with a $2\times$ faster disk leading to $2\times$ faster training for large embeddings. With NVMe-based SSDs becoming more common, the design of Marius will best be able to leverage future fast sequential storage mediums and scale training to embedding sizes beyond what we show here.

## 5.5 Microbenchmarks

**Bounded Staleness** We now show how our pipelined training architecture with bounded staleness effects the embedding quality and throughput of training. We train Marius on Freebase86m with $d = 50$, and vary the number of batches allowed into the pipeline at any given time. We evaluate how the performance and MRR vary as we vary the staleness bound. Results are shown in Figure 12. We see that embedding quality is only affected when the relations are not updated synchronously. The reason

for this is that updates to relation embeddings are dense, whereas node embedding updates are sparse. We also find that increasing the staleness bound improves overall throughput of training up to a bound of about eight. At this point, the system becomes CPU bound since there are only eight vCPUs on the instance we use.

**Prefetching Effects** We evaluate the effect of prefetching partitions to the buffer on GPU utilization. We train Marius on Freebase86m with $d = 100$, 32 partitions, and a buffer capacity of eight. We show the average utilization of the GPU during each iteration of a single epoch of training in Figure 13. We can see that prefetching results in a higher sustained utilization of the GPU since less time is spent waiting for partition swaps. Interestingly, both configurations see a utilization bump starting at about iteration 12,000. This is because the *Elimination* ordering does not require any swaps during this period. Overall, prefetching is able to mitigate wait times for partition swaps improving utilization and training times.

## 6 Conclusion

We introduced Marius, a new framework for computing large-scale graph embedding models on a single machine. We demonstrated that the key to scalable training of graph embeddings is optimized data movement; a task that is challenging as the learnable parameters in graph embedding models scale proportionally to the number of nodes and the dimensions of the embedding vectors. To optimize data movement and maximize GPU utilization, we proposed a pipelined architecture that leverages partition caching and *elimination ordering*, a novel buffer-aware data ordering scheme. We showed using standard benchmarks that Marius achieves the same accuracy but is up to one order-of magnitude faster than existing systems. We also showed that Marius can scale to graph instances with more than a billion edges and up to 550GB of model parameters on a single AWS P3.2xLarge instance. In the future, we plan to explore how the ideas behind Marius' design and our new data ordering can help speed up training of graph neural networks. We also aim to extend our system to embeddings that go beyond Euclidean geometry towards Hyperbolic Geometry [6].

# References

[1] Laszlo A. Belady. A study of replacement algorithms for a virtual-storage computer. *IBM Systems journal*, 5(2):78–101, 1966.

[2] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. *Advances in neural information processing systems*, 26:2787–2795, 2013.

[3] Sylvain Brohee and Jacques Van Helden. Evaluation of clustering algorithms for protein-protein interaction networks. *BMC bioinformatics*, 7(1):488, 2006.

[4] Nathan Bronson, Zach Amsden, George Cabrera, Prasad Chakka, Peter Dimov, Hui Ding, Jack Ferris, Anthony Giardullo, Sachin Kulkarni, Harry Li, et al. Tao: Facebook's distributed data store for the social graph. In *USENIX ATC 2013*, pages 49–60, 2013.

[5] Hongyun Cai, Vincent W Zheng, and Kevin Chen-Chuan Chang. A comprehensive survey of graph embedding: Problems, techniques, and applications. *IEEE Transactions on Knowledge and Data Engineering*, 30(9):1616–1637, 2018.

[6] Christopher De Sa, Albert Gu, Christopher Ré, and Frederic Sala. Representation tradeoffs for hyperbolic embeddings. *Proceedings of machine learning research*, 80:4460, 2018.

[7] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.

[8] Kimm Fairchild, Steven E Poltrock, and George W Furnas. Graphic representations of large knowledge bases. *Cognitive science and its applications for human-computer interaction*, page 201, 1988.

[9] Google. Freebase data dumps. https://developers.google.com/freebase, 2018.

[10] Palash Goyal and Emilio Ferrara. Graph embedding techniques, applications, and performance: A survey. *Knowledge-Based Systems*, 151:78–94, 2018.

[11] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. *CoRR*, abs/1607.00653, 2016.

[12] William L. Hamilton, Rex Ying, and Jure Leskovec. Representation learning on graphs: Methods and applications. *CoRR*, abs/1709.05584, 2017.

[13] Joseph M Hellerstein, Michael Stonebraker, and James Hamilton. *Architecture of a database system*. Now Publishers Inc, 2007.

[14] John L Hennessy and David A Patterson. *Computer architecture: a quantitative approach*. Elsevier, 2011.

[15] David Hilbert. Über die stetige abbildung einer line auf ein flächenstück. *Mathematische Annalen*, 38(3):459–460, 1891.

[16] Haewoon Kwak, Changhyun Lee, Hosung Park, and Sue Moon. What is twitter, a social network or a news media? In *Proceedings of the 19th International Conference on World Wide Web*, WWW '10, page 591–600, 2010.

[17] Adam Lerer, Ledell Wu, Jiajun Shen, Timothee Lacroix, Luca Wehrstedt, Abhijit Bose, and Alex Peysakhovich. Pytorch-biggraph: A large-scale graph embedding system. *arXiv preprint arXiv:1903.12287*, 2019.

[18] Jure Leskovec. WWW-18 Tutorial: Representation Learning on Networks. http://snap.stanford.edu/proj/embeddings-www/.

[19] Jure Leskovec and Andrej Krevl. Snap datasets: Stanford large network dataset collection. http://snap.stanford.edu/data, 2014.

[20] LibTorch: PyTorch C++ API. https://pytorch.org/cppdocs.

[21] Hailun Lin, Yong Liu, Weiping Wang, Yinliang Yue, and Zheng Lin. Learning entity and relation embeddings for knowledge resolution. *Procedia Computer Science*, 108:345–354, 2017.

[22] Steffen Maass, Changwoo Min, Sanidhya Kashyap, Woonhak Kang, Mohan Kumar, and Taesoo Kim. Mosaic: Processing a trillion-edge graph on a single machine. In *Proceedings of the Twelfth European Conference on Computer Systems*, EuroSys '17, pages 527–543, New York, NY, USA, 2017. ACM.

[23] Frank McSherry, Michael Isard, and Derek G Murray. Scalability! but at what {COST}? In *15th Workshop on Hot Topics in Operating Systems (HotOS {XV})*, 2015.

[24] Deepak Narayanan, Aaron Harlap, Amar Phanishayee, Vivek Seshadri, Nikhil R Devanur, Gregory R Ganger, Phillip B Gibbons, and Matei Zaharia. Pipedream: generalized pipeline parallelism for dnn training. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles*, pages 1–15, 2019.

[25] Maximilian Nickel, Volker Tresp, and Hans-Peter Kriegel. A three-way model for collective learning on multi-relational data. In *Proceedings of the 28th International Conference on International Conference on Machine Learning*, ICML'11, page 809–816, 2011.

[26] Shoumik Palkar and Matei Zaharia. Optimizing data-intensive computations in existing libraries with split annotations. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles*, pages 291–305, 2019.

[27] Raghu Ramakrishnan, Johannes Gehrke, and Johannes Gehrke. *Database management systems*, volume 3. McGraw-Hill New York, 2003.

[28] Satu Elisa Schaeffer. Graph clustering. *Computer science review*, 1(1):27–64, 2007.

[29] C. Seshadhri, Aneesh Sharma, Andrew Stolman, and Ashish Goel. The impossibility of low-rank representations for triangle-rich complex networks. *Proceedings of the National Academy of Sciences*, 117(11):5631–5637, 2020.

[30] Ben Taskar, Ming-Fai Wong, Pieter Abbeel, and Daphne Koller. Link prediction in relational data. *Advances in neural information processing systems*, 16:659–666, 2003.

[31] Théo Trouillon, Johannes Welbl, Sebastian Riedel, Eric Gaussier, and Guillaume Bouchard. Complex embeddings for simple link prediction. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48, pages 2071–2080, 20–22 Jun 2016.

[32] Denny Vrandečić and Markus Krötzsch. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*, 57(10):78–85, 2014.

[33] Quan Wang, Zhendong Mao, Bin Wang, and Li Guo. Knowledge graph embedding: A survey of approaches and applications. *IEEE Transactions on Knowledge and Data Engineering*, 29(12):2724–2743, 2017.

[34] Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. Embedding entities and relations for learning and inference in knowledge bases. *arXiv preprint arXiv:1412.6575*, 2014.

[35] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. Graph convolutional neural networks for web-scale recommender systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 974–983, 2018.

[36] Reza Zafarani, Mohammad Ali Abbasi, and Huan Liu. *Social media mining: an introduction*. Cambridge University Press, 2014.

[37] Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. *Advances in Neural Information Processing Systems*, 31:5165–5175, 2018.

[38] Da Zheng, Xiang Song, Chao Ma, Zeyuan Tan, Zihao Ye, Jin Dong, Hao Xiong, Zheng Zhang, and George Karypis. DGL-KE: Training knowledge graph embeddings at scale. *arXiv preprint arXiv:2004.08532*, 2020.

[39] Zhaocheng Zhu, Shizhen Xu, Meng Qu, and Jian Tang. Graphvite: A high-performance cpu-gpu hybrid system for node embedding. In *The World Wide Web Conference*, pages 2494–2504. ACM, 2019.