



# Accelerating Graph Sampling for Graph Machine Learning using GPUs

#### Abhinav Jangda

University of Massachusetts Amherst

Sandeep Polisetty

University of Massachusetts Amherst

Arjun Guha Northeastern University Marco Serafini University of Massachusetts Amherst

#### Graph Data is Everywhere



Social Networks



**Relational Data** 





Recommendation Systems



**Knowledge Bases** 

Fraud and Financial Data

#### Machine Learning on Graphs

Graph Machine Learning uses the network structure of the underlying data to improve predictive outcomes.

Personalization and Recommendation Systems



combine features from multiple users and products

#### Search, Q&A, semantic web



#### Machine Learning on Graphs



Graph Neural Network maps vertices of graphs to an embedding in a N-dimensional space.

#### Types of Graph Neural Networks (GNNs)



Input Graph

2-D Vertex Embeddings

Two types of Graph Neural Networks:

- Sampling based GNNs samples the input graph and train using these samples
- Whole Graph based GNNs train on the input graph directly

#### Workflow of a Sampling based GNNs



A Sampling based GNN first samples the graph using a Graph Sampling Algorithm and use these samples for data parallel training.

#### Graph Sampling is an Active Area of Research



DeepWalk and node2vec are random walks of fixed length starting from a vertex.



FastGCN and LADIES samples vertices for each layer in their Neural Network.

GraphSAGE samples 2-Hop Neighborhood of a vertex

2-hop

1-hop



ClusterGCN divides a graph into many clusters and samples one or more of these clusters ML Domain Experts Implement Graph Sampling on CPU



GNN implementations use CPU for Graph Sampling and GPU for Neural Network

#### Graph Sampling is a Major Overhead in GNNs



Fraction of time spent in Sampling and Training (per epoch)

Experiments Performed on 32-Core CPU with 1 NVIDIA Tesla V100

How can domain experts have best of both worlds: easy to implement and fast Graph Sampling?

### NextDoor: A System to Accelerate Graph Sampling on GPUs

- 1. A simple yet powerful API to express diverse graph sampling algorithms.
- 2. A new "Transit Parallel" approach to parallel graph sampling to achieve regularity.
- 3. Load balancing and caching to optimize GPU utilization.
- 4. Improves end-to-end training time of GNNs by upto 4x.



#### An Abstraction for Graph Sampling Applications

- A graph sampling application runs for *k* steps.
- In the beginning each sample has one or more root vertices.
- At each step *i* 
  - A *transit* vertex for *i* is a vertex whose neighbors may be added to the sample
  - Sample *m*<sub>i</sub> of those neighbors





✤ Transit Parallel Sampling

# Graph Sampling using NextDoor's API

\* Sample Parallel Sampling

- Simple yet powerful API based on the abstraction.
- Implementations need only few lines of code
- We implemented diverse algorithms
  - DeepWalk, PPR, and node2vec
  - K-hop Neighborhood Sampling
  - Importance Sampling

NextDoor's API

- ClusterGCN Sampling
- Minimal Variance Sampling
- Layer Sampling

```
Vertex next(s, trn, trnEdges) {
    int idx = randInt(0,
        trnEdges.size()-1);
    return trnEdges[idx];
}
int steps() {return 2;}
int sampleSize(int step)
{return (step == 0) ? 25:10;}
```

GraphSAGE's k-hop Neighborhood in NextDoor

The API provides information to perform effective load balancing and caching:

- It provides a distinction between a sample and a transit vertex
- Provides the number of steps and number of vertices sampled at each step.

How can we implement this API on a GPU?

### A GPU Can Execute Thousands of Threads Simultaneously



### Sample Parallel Graph Sampling on GPUs

Sample Parallel: Assign samples to consecutive threads.



6

#### Transit Parallel: Rethinking Parallel Graph Sampling

**Transit Parallel:** Assign samples with common transits to consecutive threads. Samples **S**<sub>2</sub> Samples Samples Transit 2 Transits step 1 1) Transits roup 4 6 1 à Ū 4 6 5 4 **GPU Sampling Kernel** Edge List in Global Memory 2-Hop Neighborhood of (1), (2) & (3)6 Transit Parallelism achieves regularity: 6 4 Helps to coalesce memory accesses S<sub>2</sub> S<sub>2</sub> S<sub>2</sub> S<sub>2</sub> S Can cache in shared memory and registers Consecutive threads access edges of same transit vertices

## Load Balanced Transit Parallelism in NextDoor

NextDoor performs load balancing for different transits and cache of neighbors of a transit.



#### Implementation

- 1. NextDoor utilizes efficient parallel radix sort and prefix scan for group by operation and load balancing.
- 2. NextDoor is implemented in C++14 and CUDA 11.
- 3. NextDoor is available at <u>https://plasma-umass.org/nextdoor-eurosys21/</u>.







#### Experiments

Benchmarks:

- DeepWalk
- Personalized Page Rank (PPR)
- node2vec
- Multi Dimensional Random Walks
- K-hop Neighborhood Sampling
- ClusterGCN Sampling
- FastGCN Sampling
- LADIES Sampling
- Minimal Variance Sampling (MVS)
- Layer Sampling

CPU only Baselines:

- KnightKing
- Samplers in existing GNN Systems

Graph Name	Abrv	Nodes	Edges
Protein-Protein Interactions	PPI	50K	1.4M
com-Orkut	Orkut	3M	117M
cit-Patents	Patents	3.77M	16.5M
soc-LiveJournal1	LiveJ	4.8M	68.9M

### Experiments

**Evaluation System:** 

- Dual Socket 16-Core Intel Xeon Silver CPUs
- 4x NVIDIA Tesla V100
- 128 GB of RAM
- Ubuntu 18.04
- CUDA 11.2

\* NextDoor's API \* Sample Parallel Sampling \* Transit Parallel Sampling \* NextDoor's Runtime \*

Experiments

### NextDoor: End-to-End Speedups with GNN Training



NextDoor's fast Graph Sampling implementations significantly improves training time of GNNs.

### NextDoor against existing Graph Sampling Implementations



Speedup over KnightKing for Random Walks

Speedup over GNN implementations

NextDoor achieves orders of magnitude speedup over CPU baselines

#### Summary

- Graph Sampling takes up to 62% of training time in Graph Neural Networks.
- Efficient Graph Sampling on GPUs is hard due to irregular nature of graphs.
- NextDoor: accelerate Graph Sampling using GPUs.
  - ✓ API to easily write efficient graph sampling applications.
  - Runtime that optimizes memory accesses in GPUs and efficiently balances load.
- NextDoor achieves orders of magnitude improvement over existing solutions.
- NextDoor improves end-to-end training time of GNNs by up to 4 times on large graphs.

#### Existing Implementations Have Limited Parallelism



Is it possible to increase parallelism in the standard approach to parallel Graph Sampling?

## Regular Code achieves High Performance on GPU







# Scheduling Transit Parallel in NextDoor

The NextDoor API exposes three degrees of parallelism that match the GPU architecture



#### Additional Overhead of Transit Parallel over Sample Parallel

- Transit Parallel uses a Group By operation.
- Random Walks spend up to 40% time in grouping operation.
- Despite overhead Transit Parallel still achieves up to 2x speedup over Sample Parallel in Random Walks.
- Other Application spend less than 10% time in grouping operation.
- Random Walks spends more time because they sample only one neighbor of the transit vertex at each step.



Percentage of time spent in group by operation over total time.

# Standard Approach to Parallel Graph Sampling

#### Sample Parallel Graph Sampling

- Samples can be expanded in parallel by assigning samples to a single thread.
- Approach adopted by existing systems.



Can we use this approach for a GPU based parallel graph sampling?

### Leverage GPUs for Graph Sampling is hard!

#### **Regular Computations** $\begin{vmatrix} a_{2} & a_{3} \\ a_{5} & a_{6} \end{vmatrix} \begin{vmatrix} b_{1} & b_{2} & b_{3} \\ b_{4} & b_{5} & b_{6} \end{vmatrix} = \begin{vmatrix} c_{1} & c_{2} & c_{3} \\ c_{4} & c_{5} & c_{6} \end{vmatrix}$ $a_4$ b a, a。 a11 **a**12 **a**13 b11 **b**12 **b**13 a11+b11 a12+b12 a13+b13 a21 + b21 b22 b23 = a21+b21 a22+b22 a23+b23 a22 a23 **a**31 **a**33 b32 **b**33 a31+b31 a32+b32 a33+b33 **a**<sub>32</sub>

Consecutive Memory Accesses
 Convergent Control Flow

Utilize faster shared memory





#### Irregular Computations

Different neighbors of vertices



Random Memory Access Divergent Control Flow Cannot Utilize faster memory

#### NextDoor speedup over Transit Parallel



#### Workflow of Graph Neural Network Training



Each Random Walk is a mini-batch

#### An Abstraction for Graph Sampling Applications

- A graph sampling application runs for k steps. step = 1
  Each execution of application produces one sample of the graph.
  In the beginning each sample has root vertice(s).
  At step i , the application samples m<sub>i</sub> vertices.
  Function next describes the sampling procedure.
  Step = 3
  - A *transit* vertex at a step *i* is a vertex whose neighbors may be sampled at step *i*.

