

Hello!

# CS 744: MARIUS

Shivaram Venkataraman

Spring 2024

# ADMINISTRIVIA

- Midterm grades ~ 95% done . TODAY
- Regrade requests
- Course Project: Check in by April 16<sup>th</sup>  
↳ next Tuesday

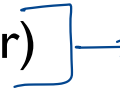
# PROJECT CHECK-INS

One page document that includes the following

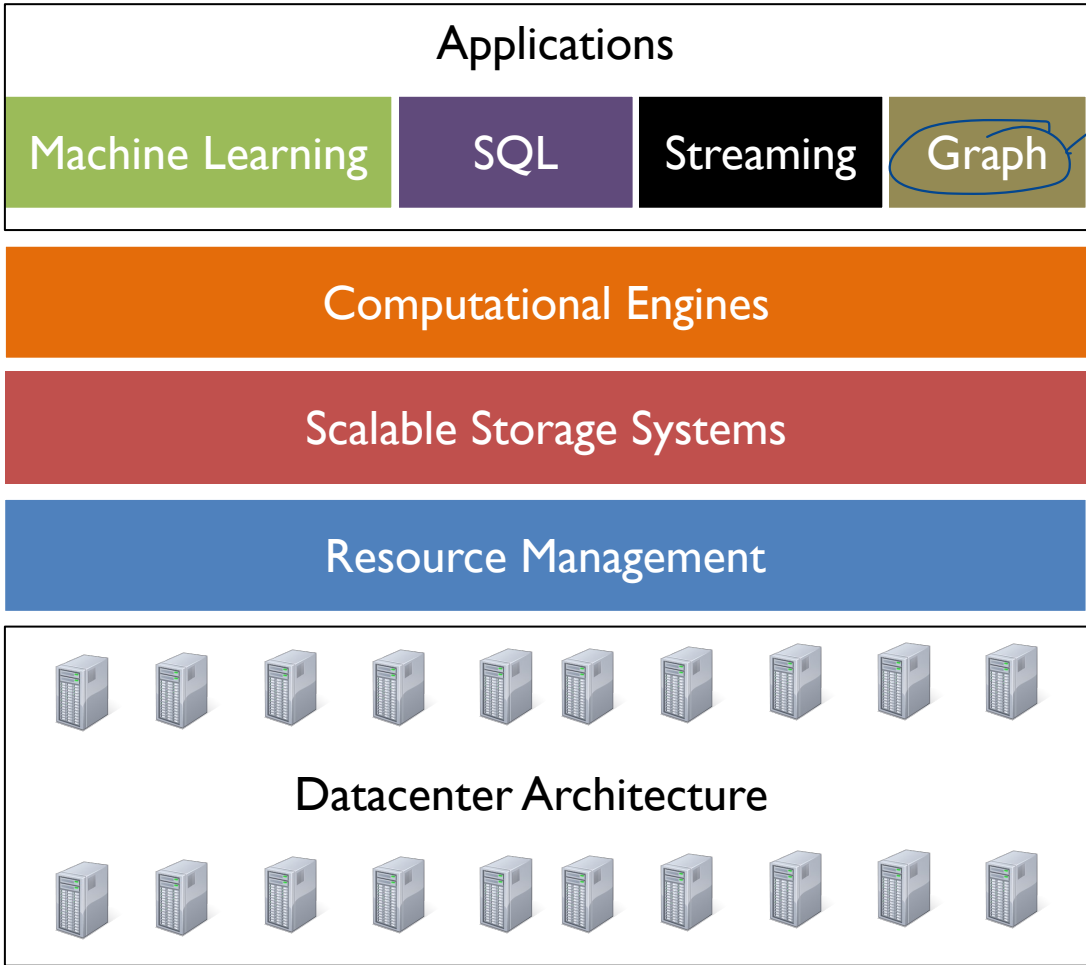
- What have you done so far
- Any challenges that you have faced so far
- Your timeline (from now till end of the semester)
- Things you need help from the course staff
- Any other comments/remarks

*cloud lab*

*hardware*



*changed a lot  
from project  
introduction?*



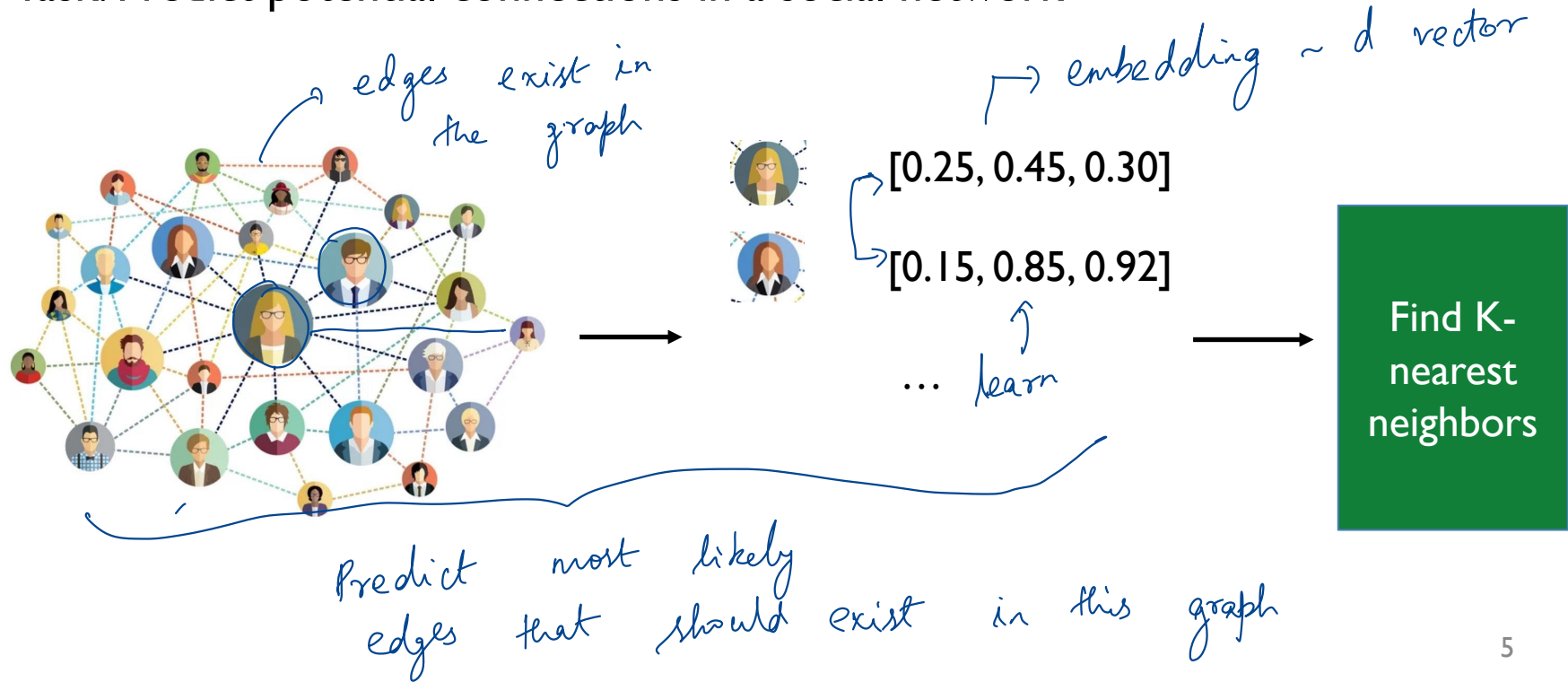
Graph Analytics

→ Page Rank  
Connected Components

Learning on graph structured data?

# EXAMPLE: LINK PREDICTION

Task: Predict potential connections in a social network



# BACKGROUND: GRAPH EMBEDDING MODELS

Score function  $\longrightarrow$  maximize score in order to learn embeddings

Capture structure of the graph given source, destination embedding

Loss function *Contrastive learning*

Maximize score for edges in graph  $\longrightarrow$  positive edges

Minimize for others (negative edges)  $\longrightarrow$  not present in graph

$$\mathcal{L} = \sum_{e \in G} \sum_{e' \in S'_e} \max(f(e) - f(e') + \lambda, 0)$$

$[0.1, 0.2]$   $[0.2, 0.3]$



*cosine similarity*  
*Distance function*

Initialize embeddings  
to random vector

# TRAINING ALGORITHM

SGD/AdaGrad optimizer

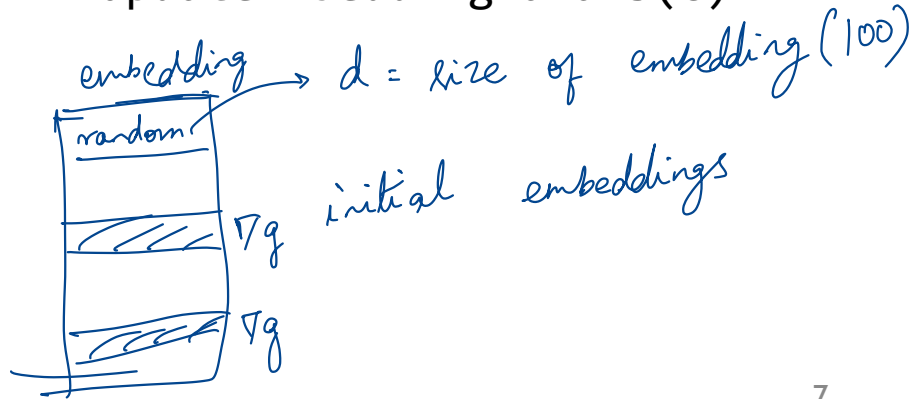
Sample positive, negative edges

Access source, dest embeddings for  
each edge in batch

Model is  
this!



vertex  
(n)  
(M)



sample "positive" edges from the graph

```
for i in range(num_batches)
    B = getBatchEdges(i)
```

```
E = getEmbeddingParams(B)
```

```
G = computeGrad(E, B)
```

```
updateEmbeddingParams(G)
```

# CHALLENGE: LARGE GRAPHS

Large graphs → Large model sizes

Example

3 Billion vertices, d = 400

Model size = 3 billion \* 400 \* 4 = 4.8 TB!

→ embedding table  
size

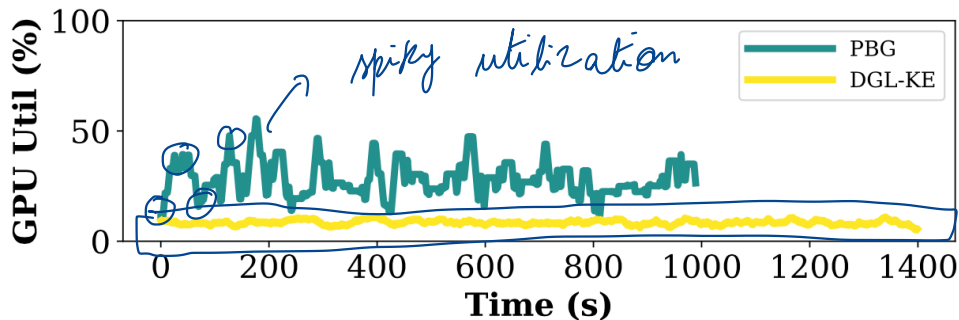
Need to scale beyond GPU memory, CPU memory!



# CHALLENGE: DATA MOVEMENT

DGL-KE: Sample edges, embeddings from CPU memory

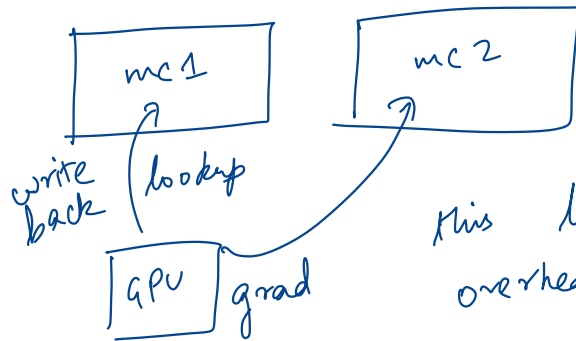
Pytorch-BigGraph: Partition embeddings so that one partition fits on GPU memory. Load sequentially



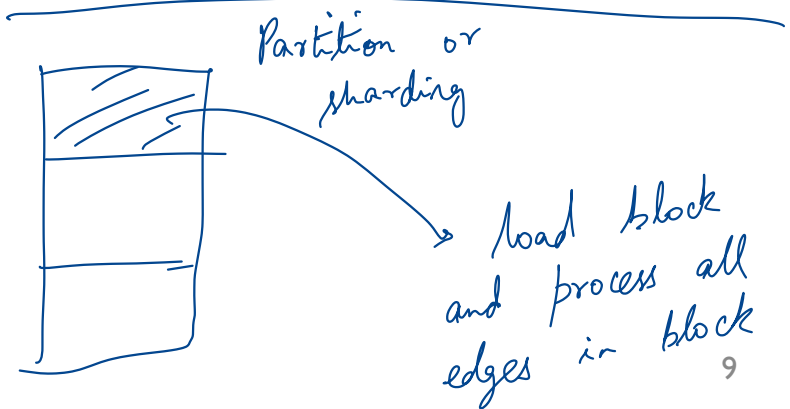
One epoch on the Freebase86m knowledge graph

shard emb. table

cpu memory



this lookup overhead is quite significant



# MARIUS

I/O efficient system for learning graph embeddings

## Marius Design

- Pipelined training → how to keep the GPU always busy
- Partition ordering  
↙ minimize I/O during 1 epoch

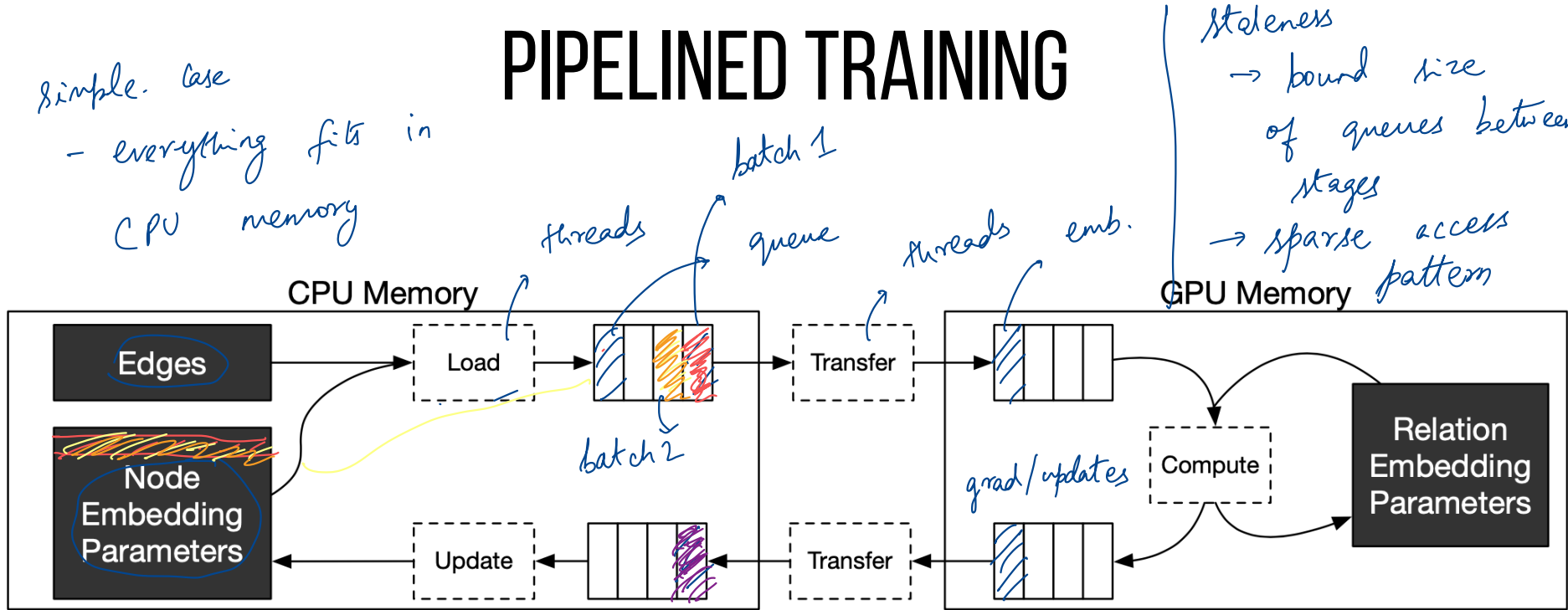


marius

# PIPELINED TRAINING

Simple case

- everything fits in CPU memory



Staleness  
 → bound size of queues between stages  
 → sparse access pattern

- ① sample positive, neg edges  
load \_emb parameters
- ② Transfer emb, edges to GPU

- ③ Compute grad GPU
- ④ Transfer back
- ⑤ Apply on emb. table

# OUT OF MEMORY TRAINING

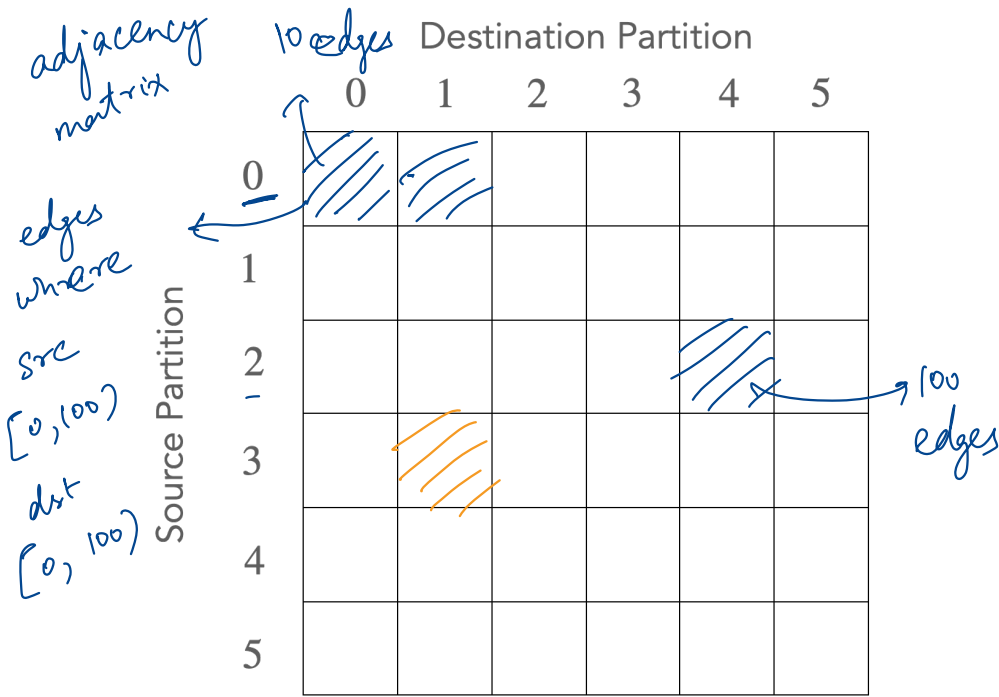
Key idea: Maintain a *cache* of partitions in CPU memory

Questions

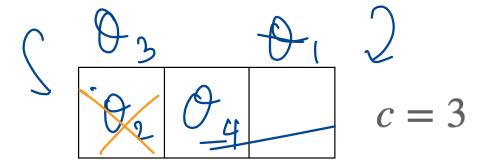
Order of partition traversal?

How to perform eviction?

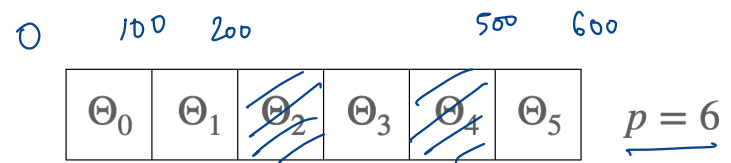
1 pass over all edges  
 ⇒ visit all buckets in edge partitions



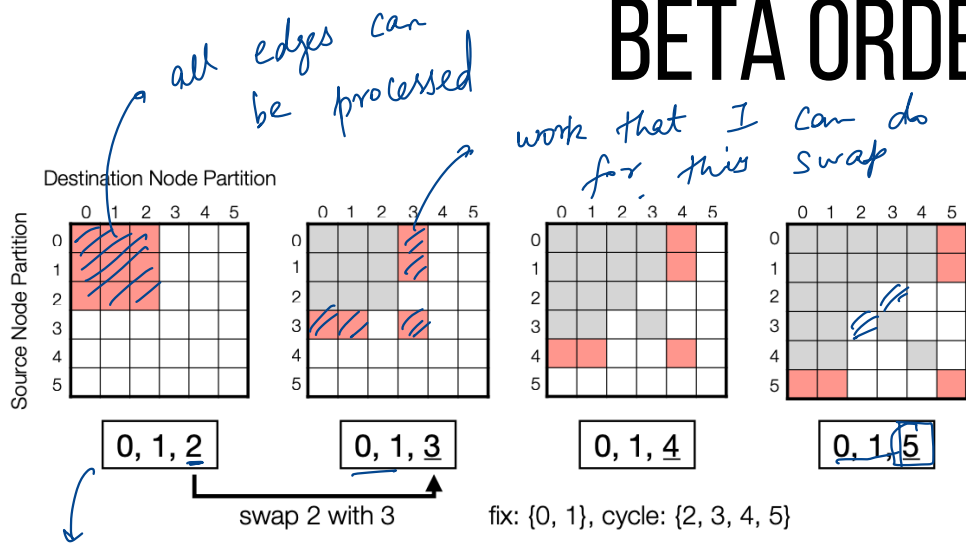
Partitions in Buffer



Partitions on disk



# BETA ORDERING



Initialize cache with  $c$  partitions

Swap in partition that leads to highest number of unseen pairs

Achieved by fixing  $c-1$  partitions and swap remaining in any order

in-memory =  $c$  partition ;  $c=3$

More work for every swap  $\Rightarrow$  minimize number of swaps

# SUMMARY

Graph Embeddings: Learn embeddings from graph data for ML

Marius: Efficient single-machine training

→ Pipelining to use CPU, GPU

→ Partition buffer, BETA ordering

↳ size of embedding table  
large  
→ access pattern sparse



# DISCUSSION

<https://forms.gle/9H6dhiiSUtjU29yd7>

How does the partitioning scheme used in this paper differ from partitioning schemes used in PowerGraph and why?

Powergraph → minimizing replicas / nbrs which are remote

Marius → ordering of computation

Similarity → each edge is in only 1 bucket !!  
↳ vertex may be visited many times



high utilization → cheaper

System	Deployment	Epoch Time (s)	Per Epoch Cost (\$)
Marius	1-GPU	727	.61
DGL-KE	2-GPUs	1068	1.81
DGL-KE	4-GPUs	542	1.84
DGL-KE	8-GPUs	277	1.88
DGL-KE	Distributed	1622	2.22
PBG	1-GPU	3060	2.6
PBG	2-GPUs	1400	2.38
PBG	4-GPUs	515	1.75
PBG	8-GPUs	419	2.84
PBG	Distributed	1474	2.02

more GPUs  
utilization  
goes down

cheaper?

faster than 1-GPU Marius

What are some shortcomings of Marius? What could the authors do to further improve the system?

# NEXT STEPS

Next class: Recommendation Models

Project check-ins next week