

Hi!

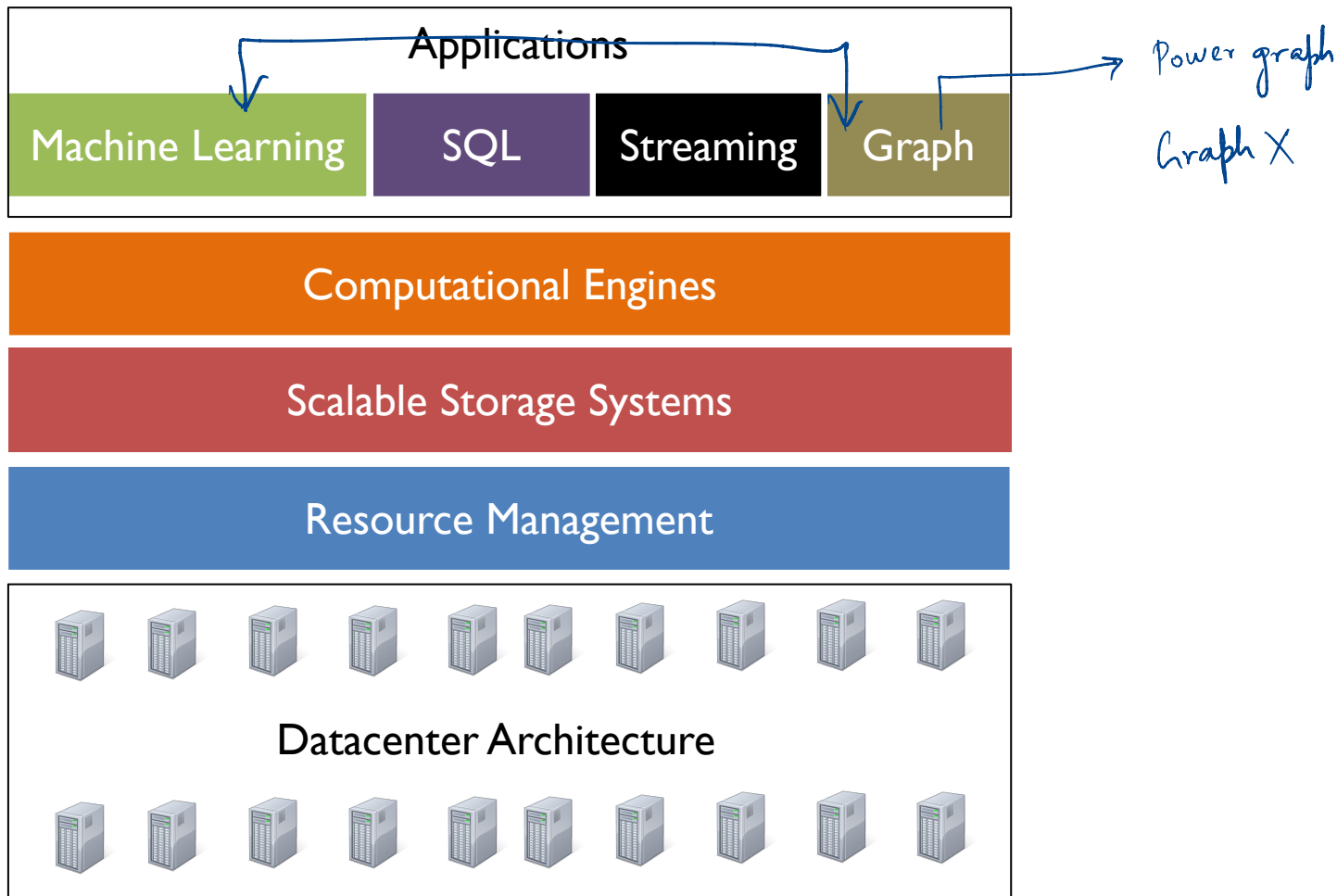
# CS 744: MARIUS

Shivaram Venkataraman

Fall 2021

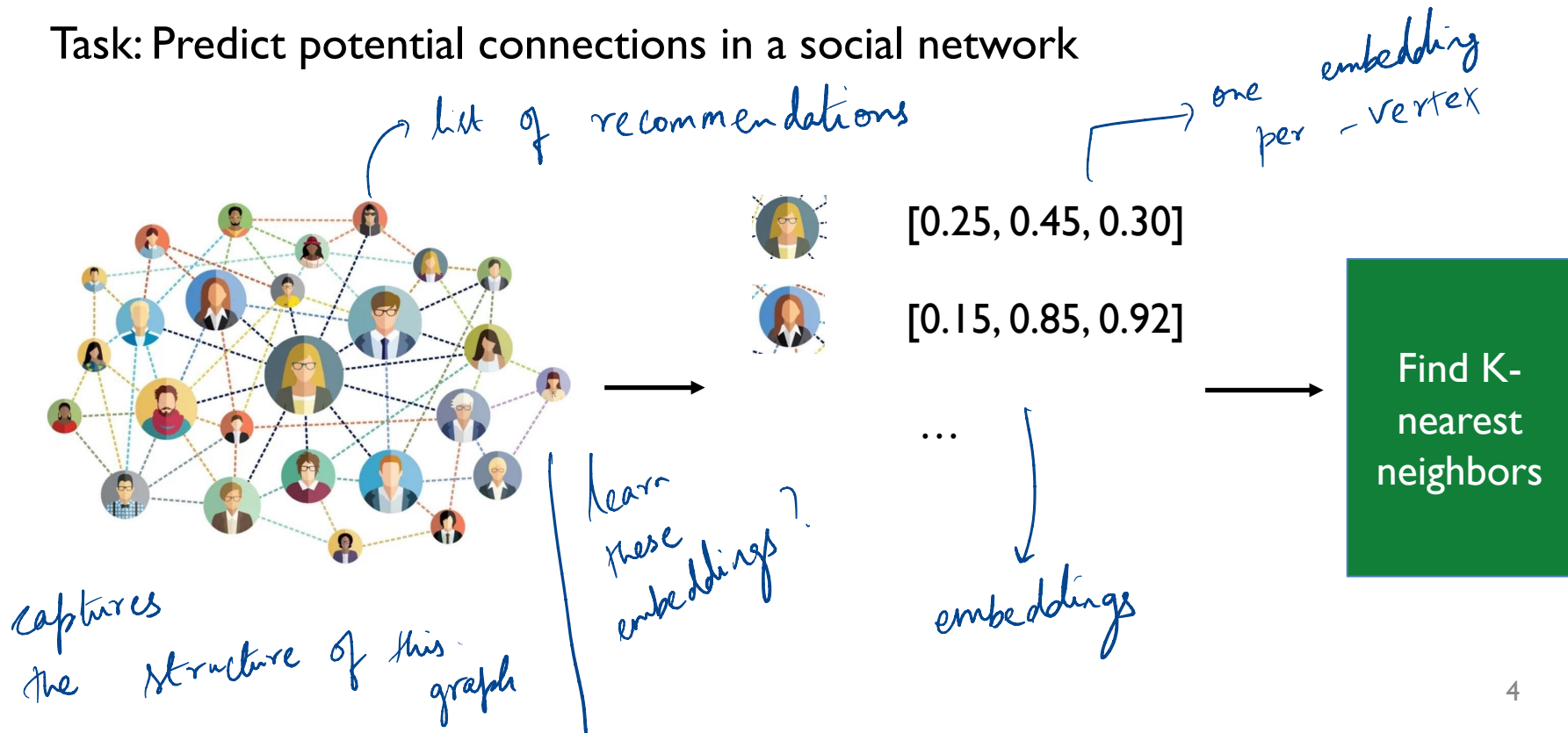
# ADMINISTRIVIA

- Midterm grades today! → Pick up papers in my office hours
- Course Project: Check in by Nov 30<sup>th</sup>  
Yien's O/H  
↳ Monday
- DON'T PANIC!
- On the class on Tuesday
  - skip the discussion
  - solution walk through



# EXAMPLE: LINK PREDICTION

Task: Predict potential connections in a social network



# BACKGROUND: GRAPH EMBEDDING MODELS

Score function  $\rightarrow$  for a particular edge

Capture structure of the graph given source, destination embedding

$\hookrightarrow$  cosine similarity

Loss function

Maximize score for edges in graph

Minimize for others (negative edges)

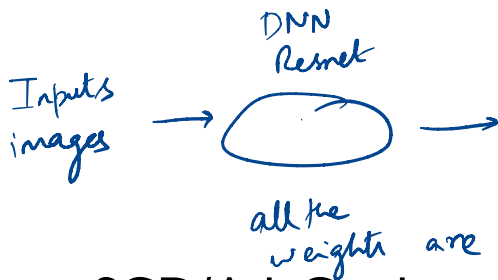
$$\mathcal{L} = \sum_{\underline{e \in G}} \sum_{\underline{e' \in S'_e}} \max(\underbrace{f(e)}_{\text{score for positive edge}} - \underbrace{f(e')}_{\text{score for negative edge}} + \lambda, 0)$$

score for positive edge

score for negative edge

regularization

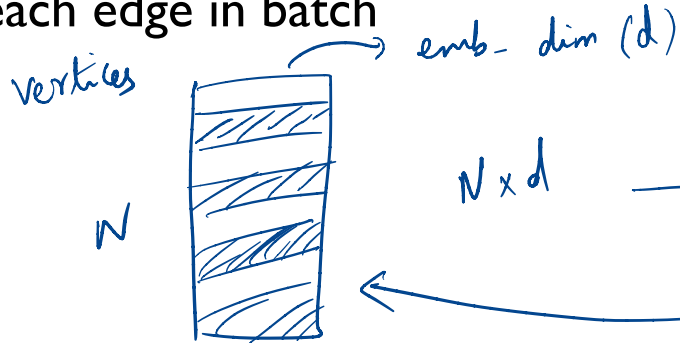
# TRAINING ALGORITHM



SGD/AdaGrad optimizer

Sample positive, negative edges

Access source, dest embeddings for each edge in batch



take a batch of edges

1 epoch = all positive edges

```
for i in range(num_batches)
    B = getBatchEdges(i)
    E = getEmbeddingParams(B)
    G = computeGrad(E, B)
    updateEmbeddingParams(G)
```

return the embeddings at the end.

update embedding params

# CHALLENGE: LARGE GRAPHS

Billions of vertices

Large graphs → Large model sizes

Example

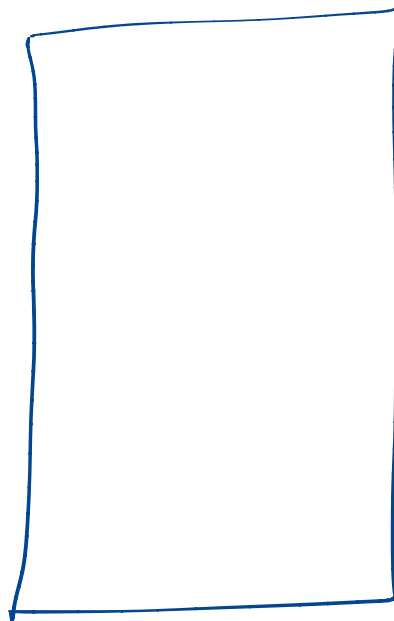
3 Billion vertices,  $d = 400$

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

Need to scale beyond GPU memory, CPU memory!

$d = 100$  to  $1000$

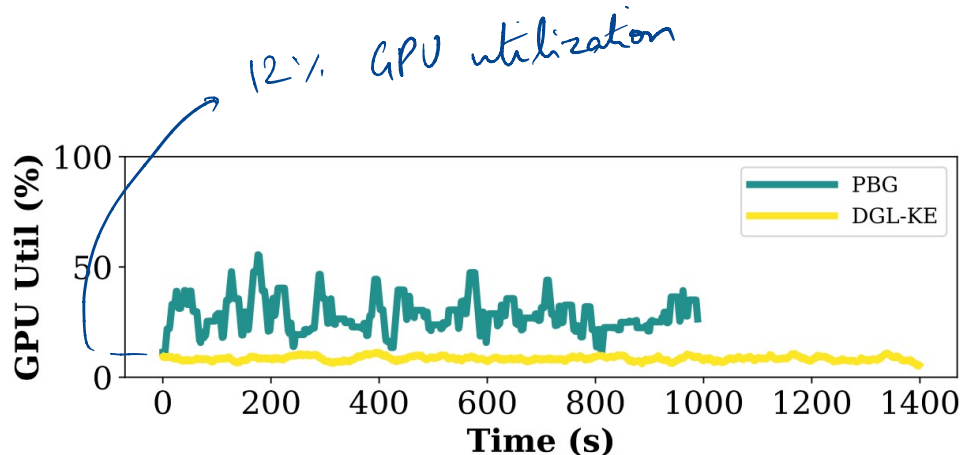
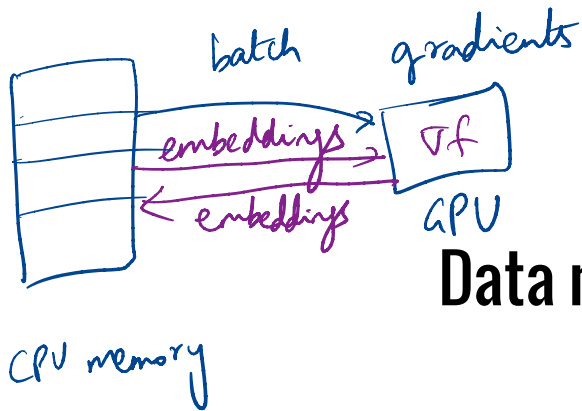
$N = 3^B$



# CHALLENGE: DATA MOVEMENT

(a) Sample edges, embeddings from CPU memory (DGL-KE)

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



One epoch on the Freebase86m knowledge graph

Data movement overheads → low GPU util



# MARIUS

I/O efficient system for learning graph embeddings

## Marius Design

- Pipelined training
- Partition ordering

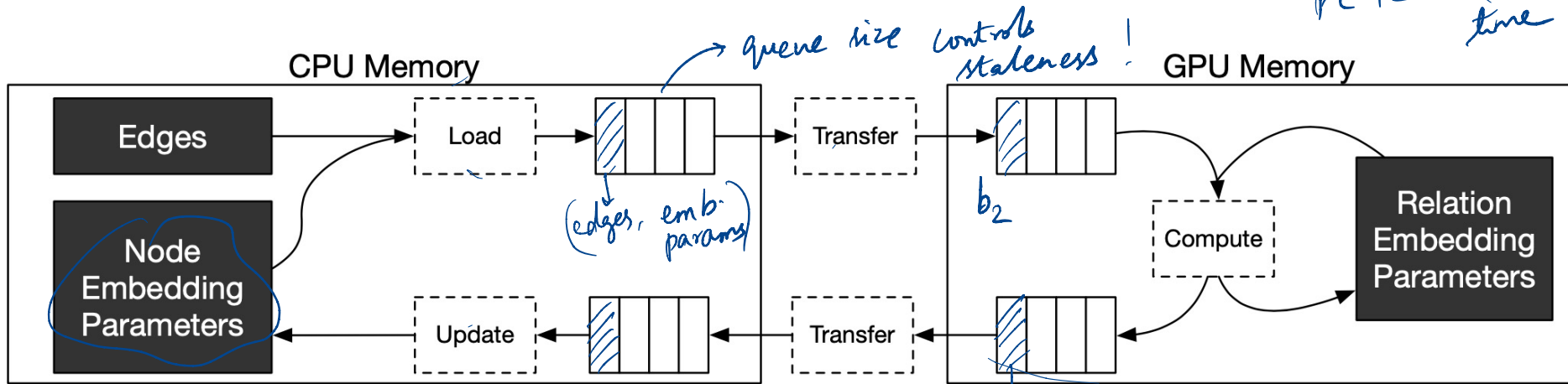


# PIPELINED TRAINING

splits the algorithm into 5 stages

and pipelines stages to

utilize CPU & GPU & PCIe at same time



Stale data!

- If  $b_2$  has shared embeddings with  $b_1$  then  $b_2$ 's embeddings are stale

updates to embeddings  
 Sparse access patterns  
 ⇒ not too many conflicts

# OUT OF MEMORY TRAINING

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

Questions

Order of partition traversal?

How to perform eviction?

adjacency  
matrix  
my graph

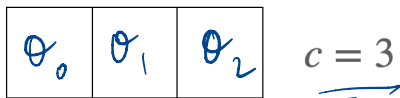
Source Partition

Destination Partition

	0	1	2	3	4	5
0						
1						
2						
3						
4						
5						

1 epoch  
is all  
the squares

Partitions in Buffer



→ Cache in CPU memory

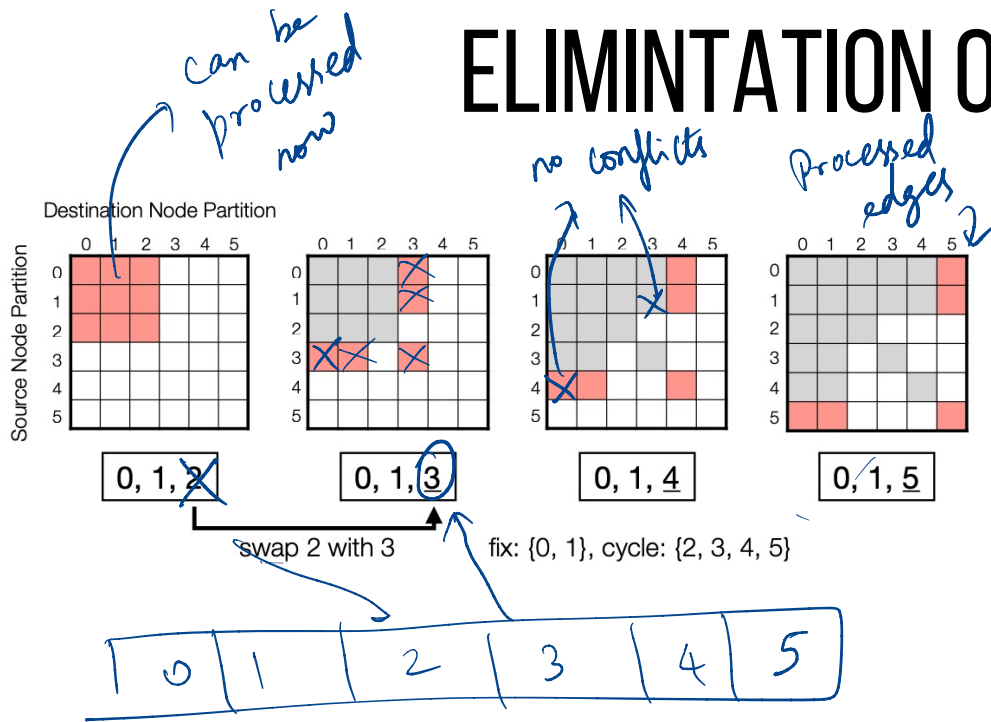
Partitions on disk



$p=6$

node embedding  
parameters

# ELIMINATION ORDERING



1. Initialize cache with  $c$  partitions
2. Swap in partition that leads to highest number of unseen pairs
3. Achieved by fixing  $c-1$  partitions and swap remaining in any order

number of swaps close to the lower bound

# SUMMARY

Graph Embeddings: Learn embeddings from graph data for ML

Marius: Efficient single-machine training

Pipelining to use CPU, GPU →

Partition buffer, BETA ordering → *minimize number of swaps*

# DISCUSSION

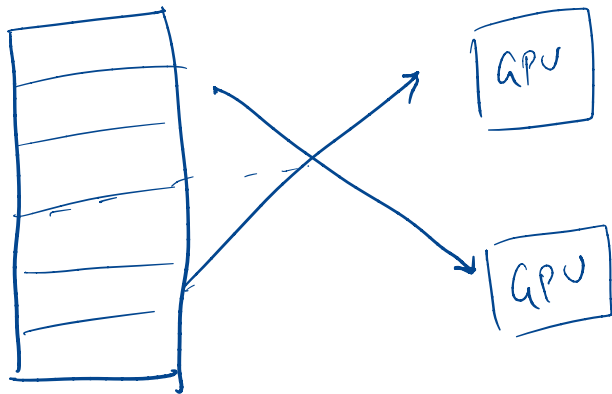
<https://forms.gle/LtoT8nEmw3oLvXuo9>

If you were going to repeat the COST analysis for knowledge graph embedding training, what would you expect to find and why?

→ Distributed systems need ~~4~~ - 8 GPUs to match Single GPU

→ Marinus 1 GPU ~ 3 GPUs of DGL

→ Communication cost if we go to many GPUs



KG Training

→ Compute  $\gg$  more compute for PageRank

→ I/O the bigger bottleneck ??

↳ data dependent

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

Workload was focused on

→ vertex and its neighbors → PageRank job

→ edge drives the computation

↳ source, dest for this edge



# NEXT STEPS

Next class: New module!

Project check-ins by Nov 30th