

*Good morning!*

# CS 744: BAGPIPE

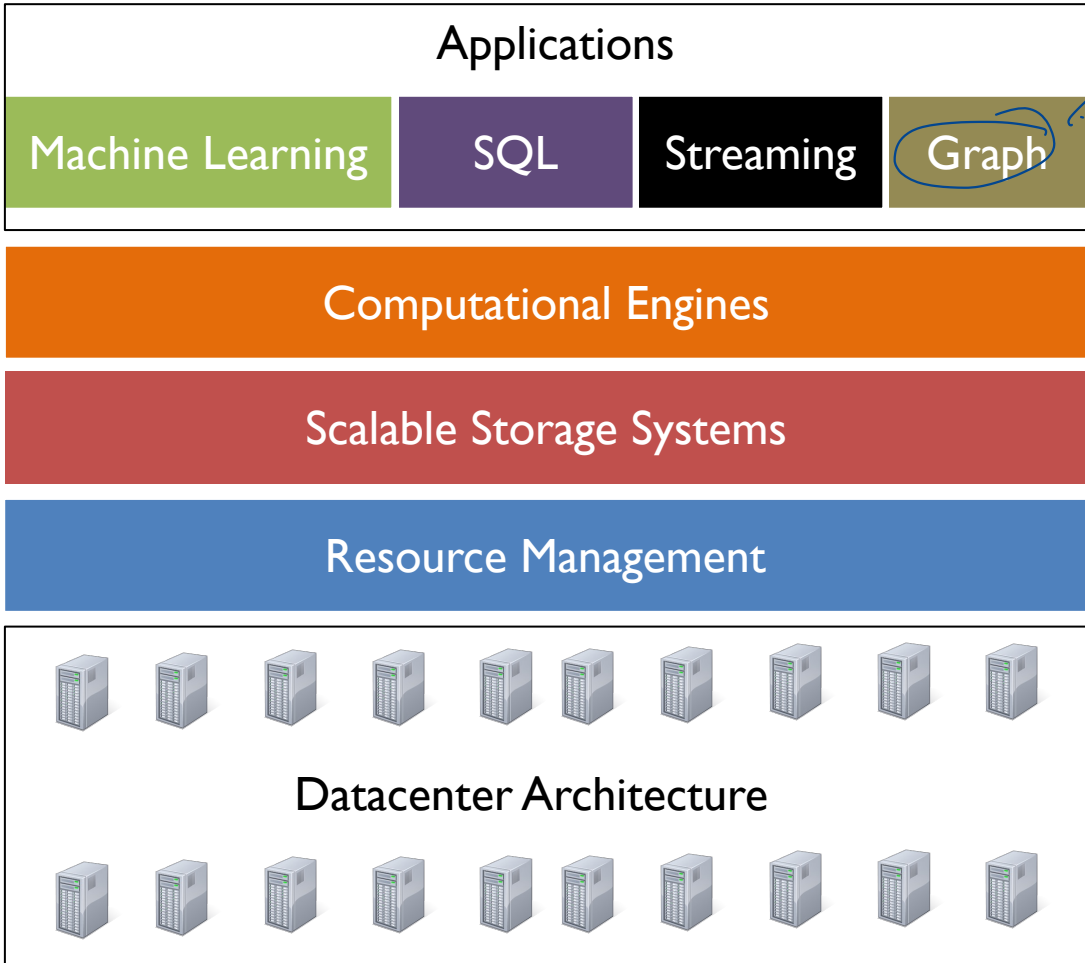
Shivaram Venkataraman

Spring 2024

# ADMINISTRIVIA

- Midterm grades on Gradescope!
  - Submit regrade requests through Gradescope
- Course Project: Check in by April 16<sup>th</sup>

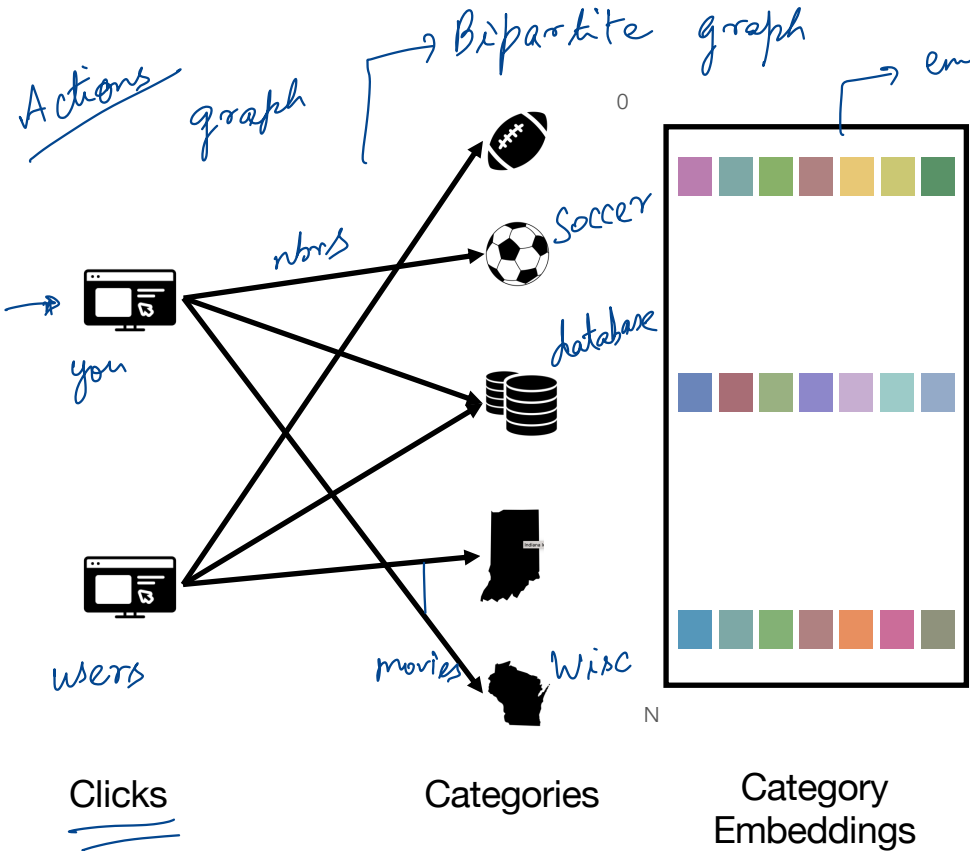
*~ 21.2 average*



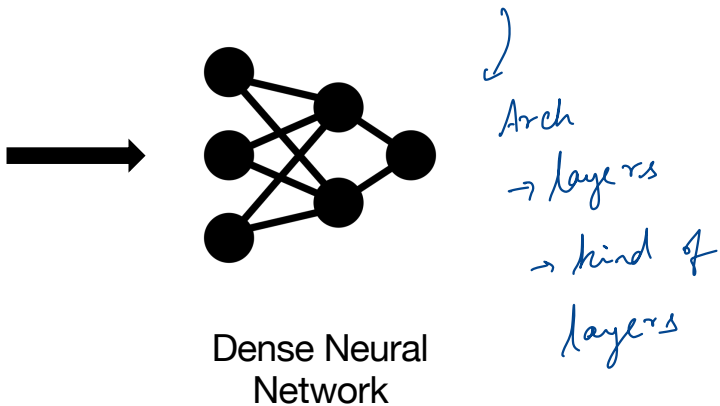
Powergraph  
↳ Page Rank  
↳ Single mc  
Distributed  
↳ ML on graphs  
↳ Single GPU  
Marius  
ML x Distributed

# RECOMMENDATION MODELS = DLRM

Bipartite graph  
 large num of nbrs  
 sparse / structured



Examples: DLRM, DeepFM,  
Wide & Deep

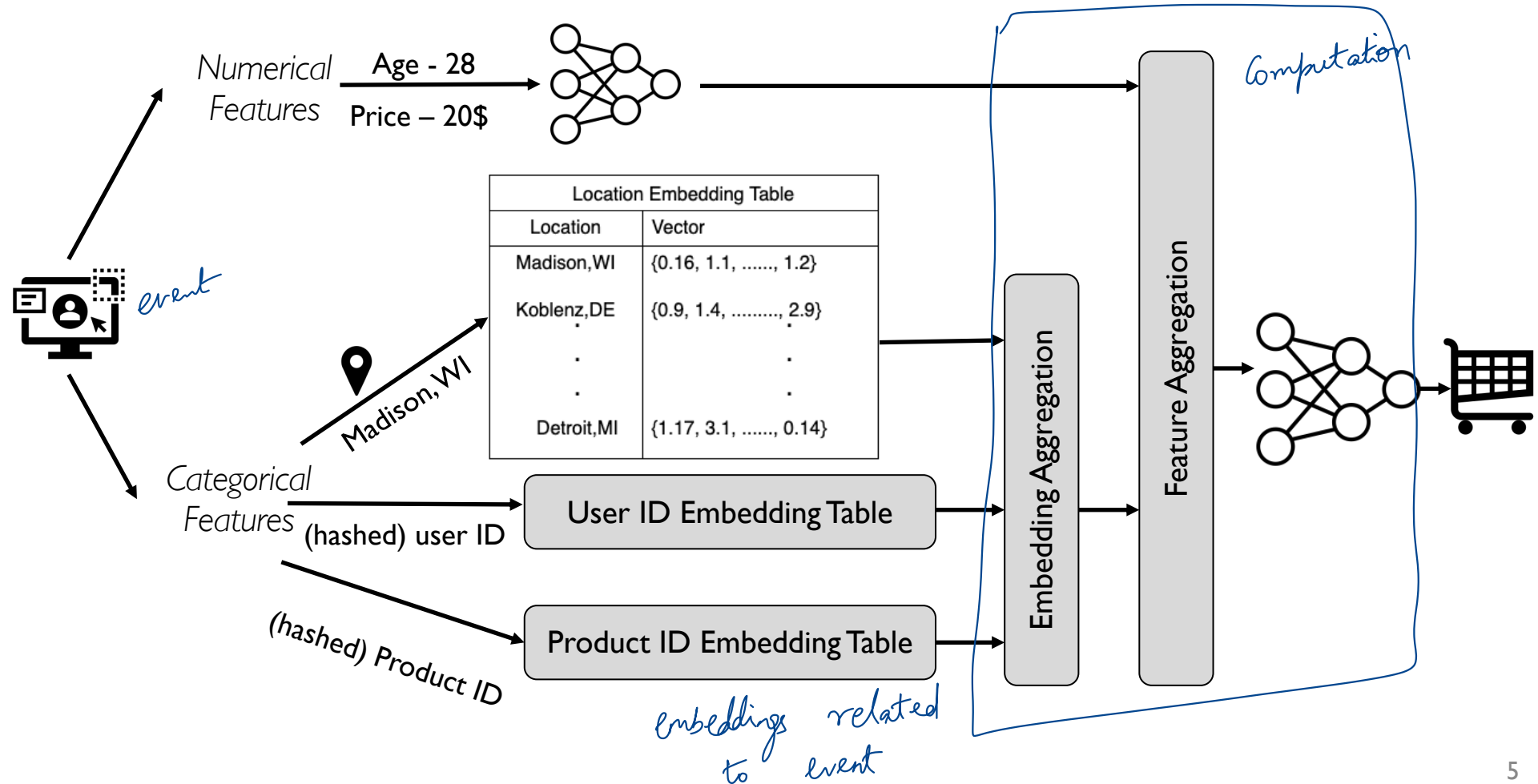


events → batch of event . For every event  
 ↳ Fetch Nbrs  
 → Fwd / Bwd 4

Clicks

Categories

Category Embeddings



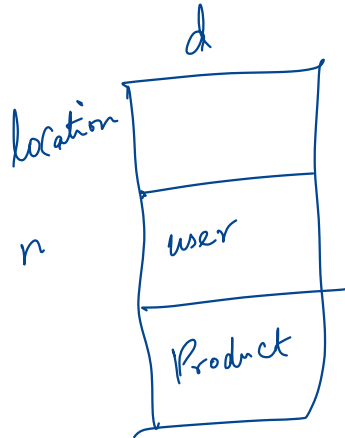
# EMBEDDING TABLES

Convert categorical features to numerical features

Example: Geographic Location to a vector

Extremely memory intensive, could be up to TBs

Have sparse access pattern

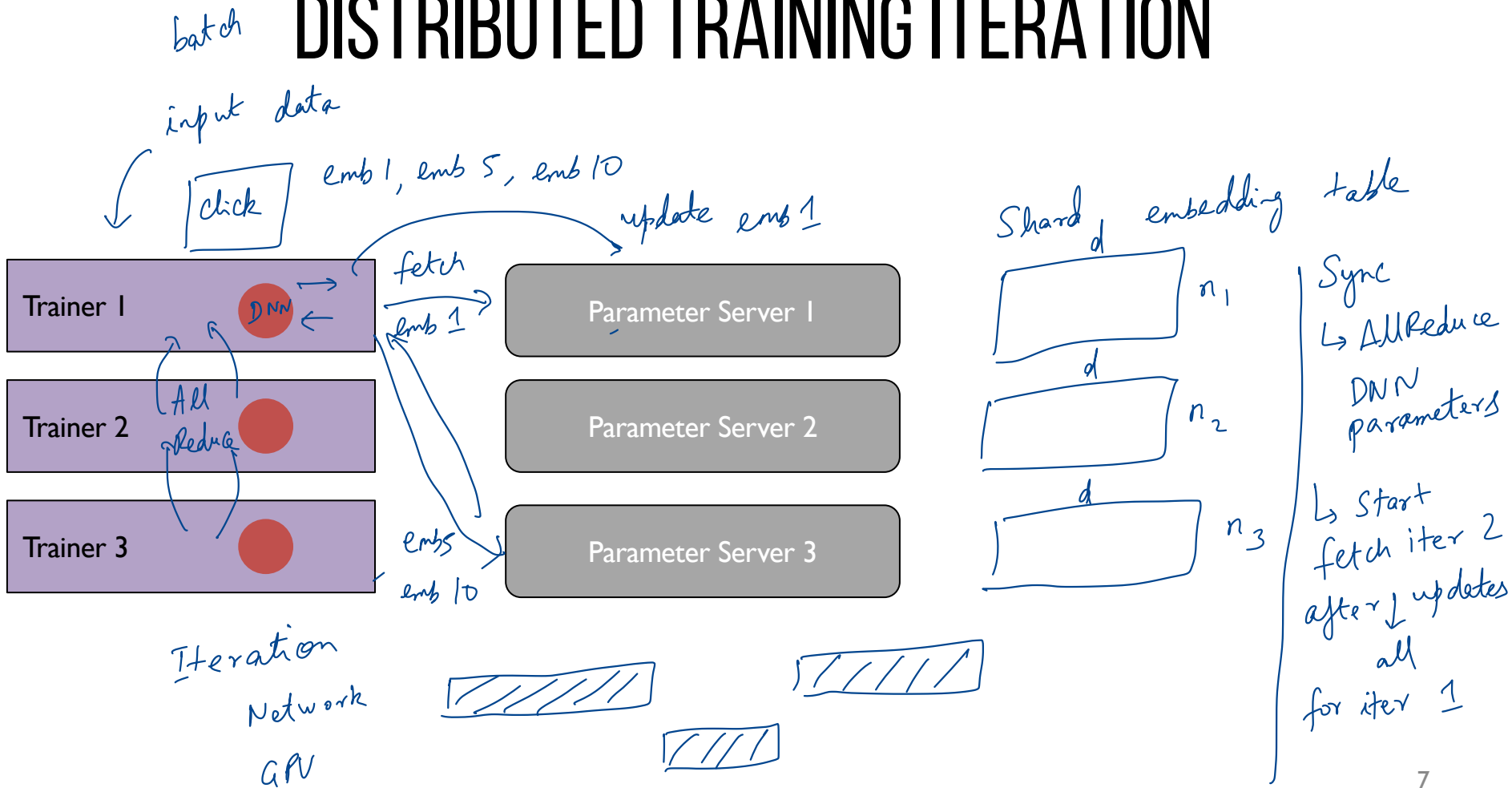


Geographical Location  
Embedding Table

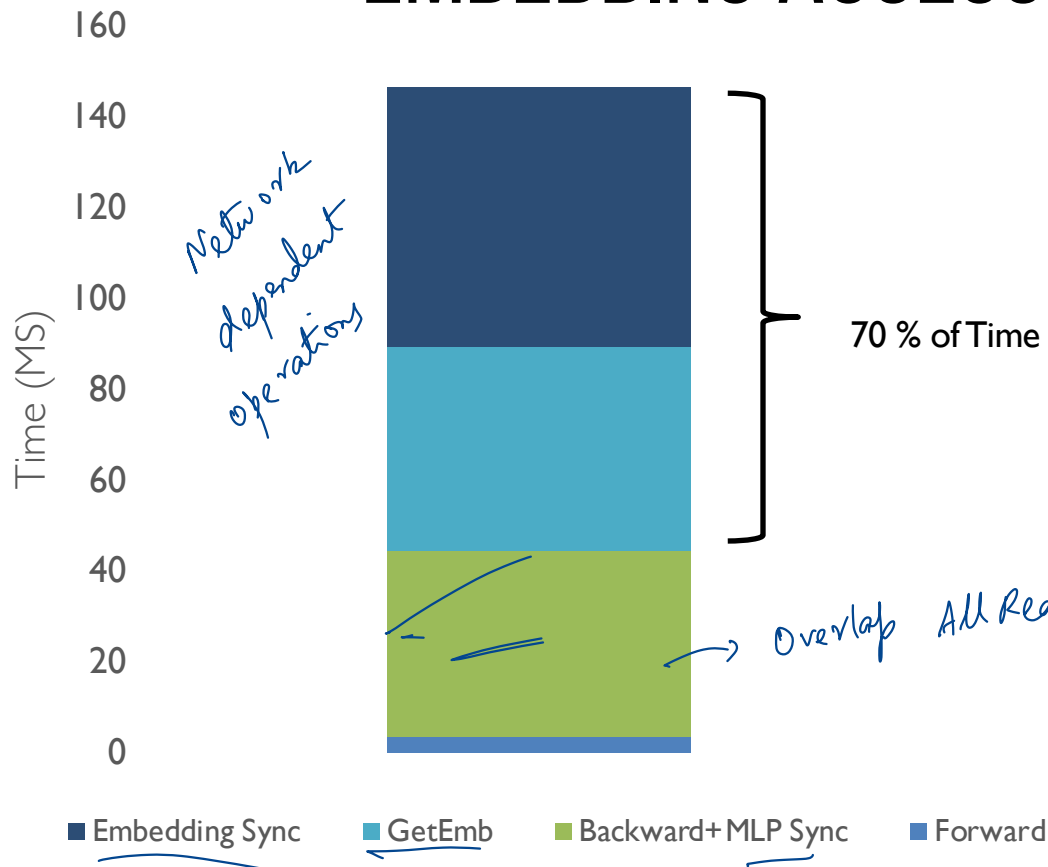
UserID Embedding Table

Product ID Embedding Table

# DISTRIBUTED TRAINING ITERATION



# EMBEDDING ACCESS OVERHEADS

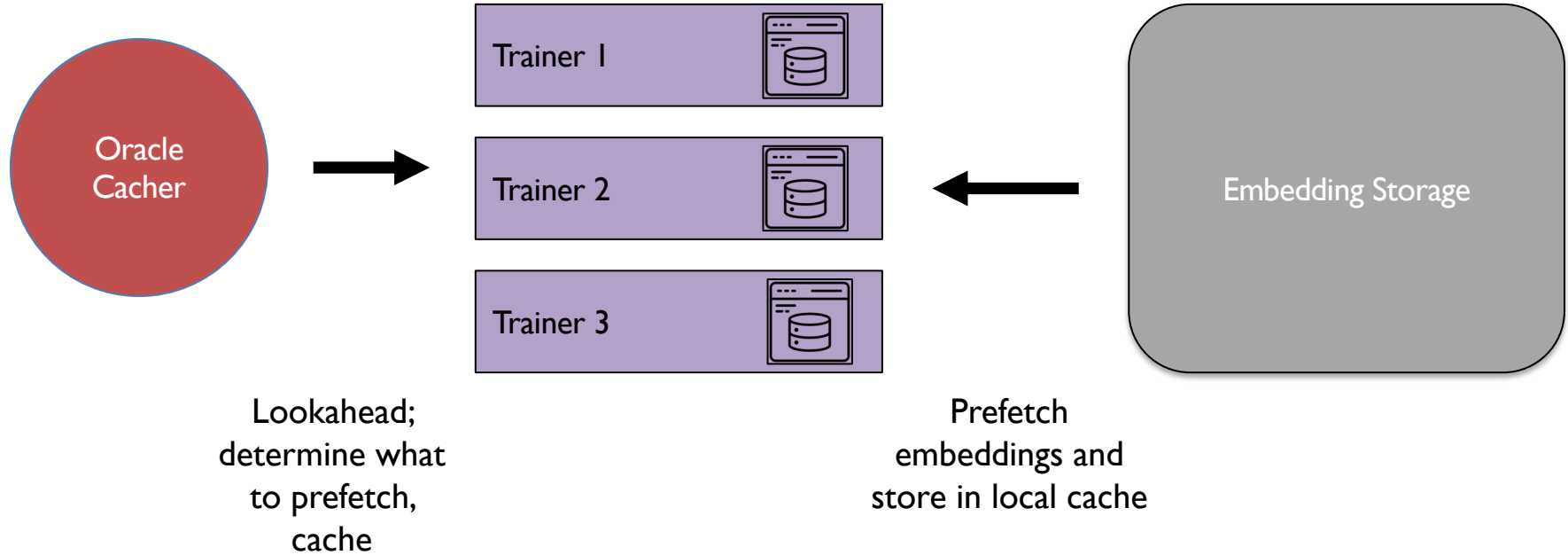


## Setup:

- DLRM model
- 8 trainers (p3.2xlarge EC2 instances | V100 each node).
- Batch 2048 per machine.
- Criteo Terabyte Dataset

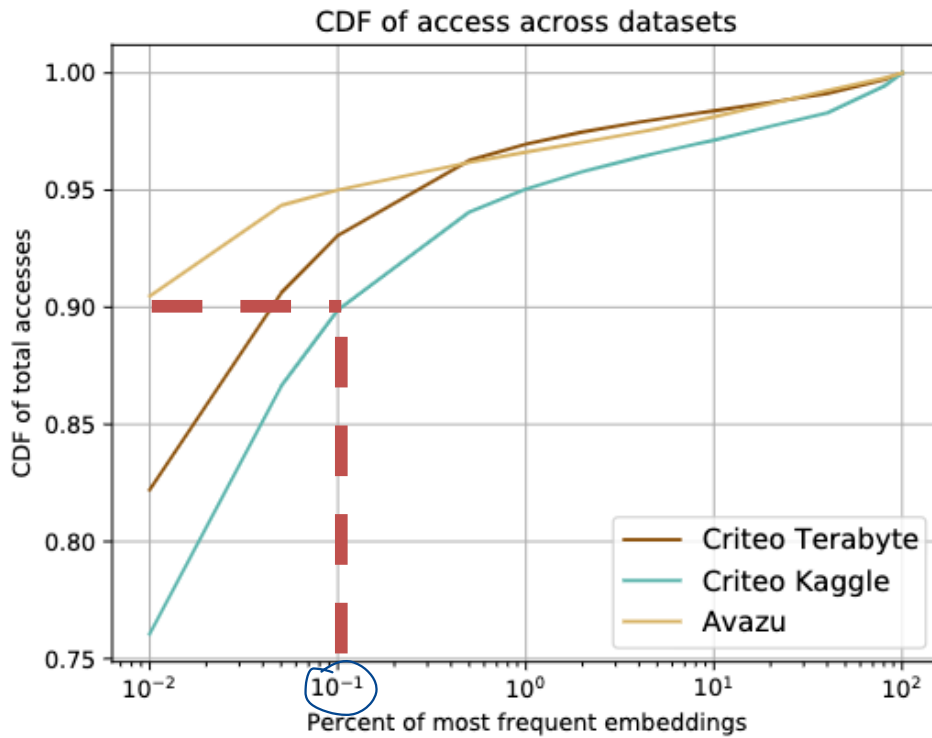


# BAGPIPE DESIGN



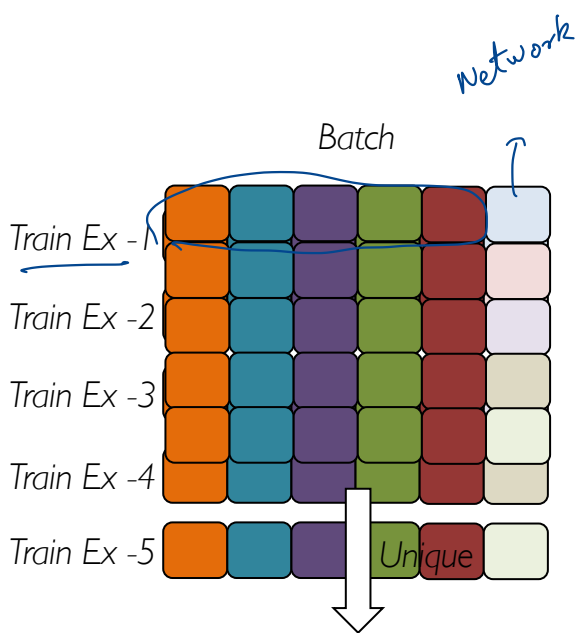
# EMBEDDING ACCESS PATTERNS

Power-law  
in terms of  
emb. popularity



0.1% of  
embeddings  
account for  
90% of  
access

# LONG TAIL ACCESSSES



Network

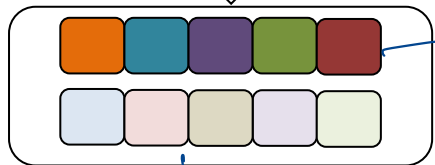
6 emb  
each  
example  
fetches

5 of those  
are popular  
embeddings

Models are trained with a batch of examples.

- For a batch only fetch unique embeddings

- Since hot embeddings are replicated, unique embeddings are comprised of long-tail accesses.



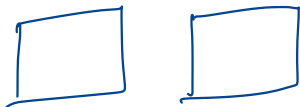
popular embeddings  
are fetched  
only once

50% of  
accesses miss  
your cache!

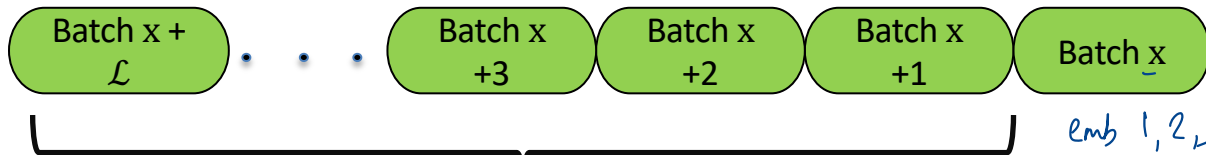
# LOOKAHEAD ALGORITHM

Look at “ $L$ ” next batches ahead of current batch to extract access pattern of embeddings by future batches

fixed dataset



→ form batches of input before processed



Perfect Ideal Caching

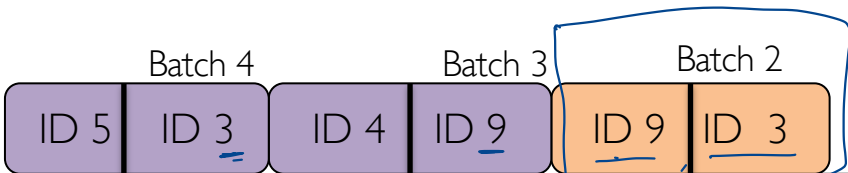
Caching

look into future to see what emb. will be accessed

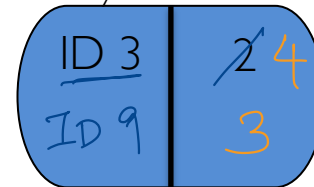
- ① If you will re-use emb future. Keep it cache
- ② Long tail prefetch it before batch starts

$L=2$

# Look-ahead Value of 2, Batch size of 2

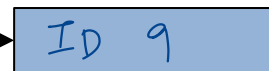


Cache Key      Cache Expiry

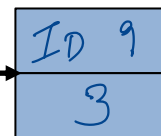


Current Cache State

Lookahead Algorithm

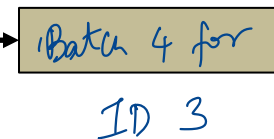


EmbID to Prefetch



Cache EMB

Expiry



Expiration Update

Perfect 1 cache hits for trainer

# LOOKAHEAD GUARANTEES

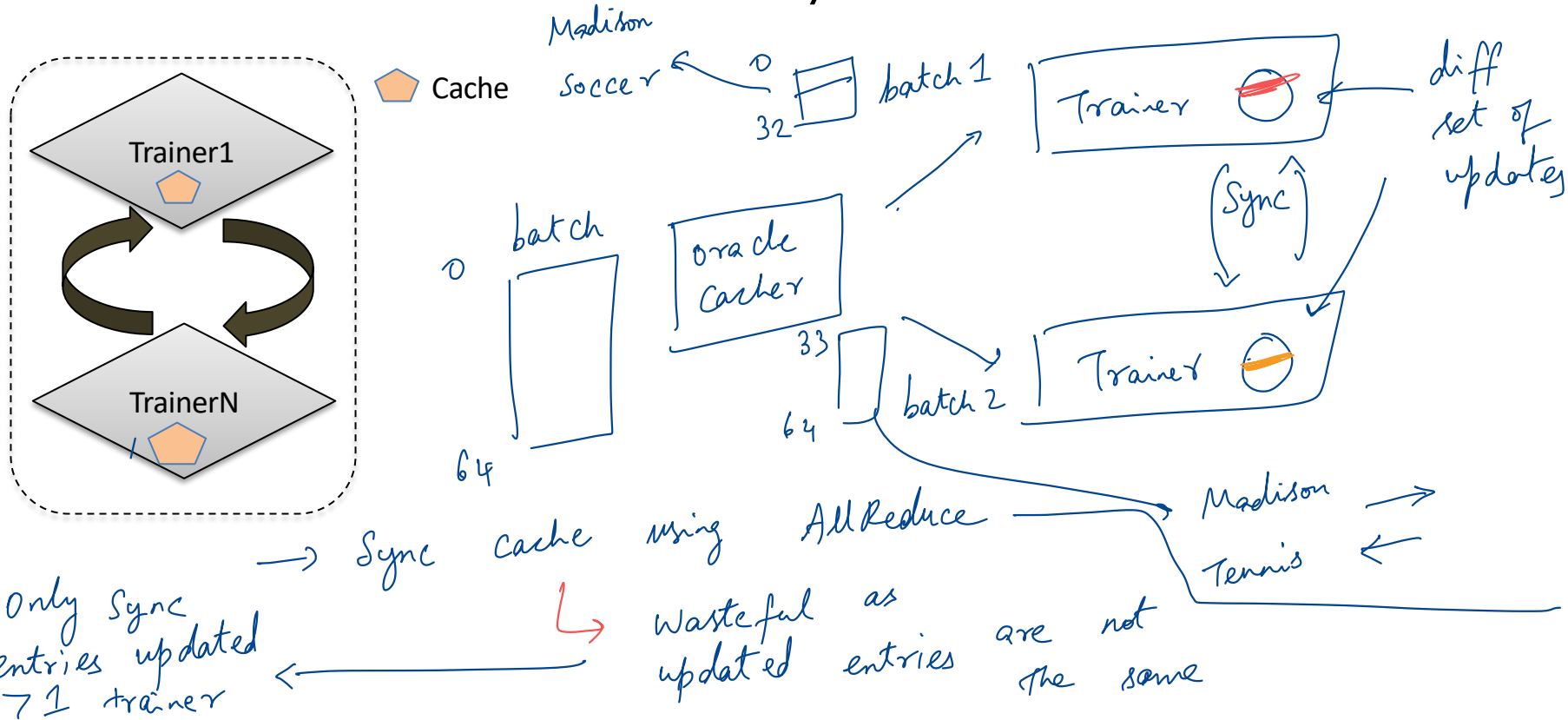
An embedding used by batch  $x$ , will either be available in cache, or no preceding batch in range  $[x - \mathcal{L}, x)$  has accessed it.

Consequently, we can prefetch embeddings used by batch  $x$ , once embeddings for batch  $x - \mathcal{L}$  have been updated

When can I prefetch an embedding 10,  $\mathcal{L} = 5$   
- batch 10, I wait till batch 5 is flushed  
and after that it's safe to prefetch

# CACHE SYNCHRONIZATION

At the end of each iteration, each trainer synchronizes caches



# SUMMARY

Recommendation models: Embeddings access overheads

BagPipe: Efficient distributed training

- Lookahead to pre-fetch and cache embeddings

- Cache synchronization across trainers

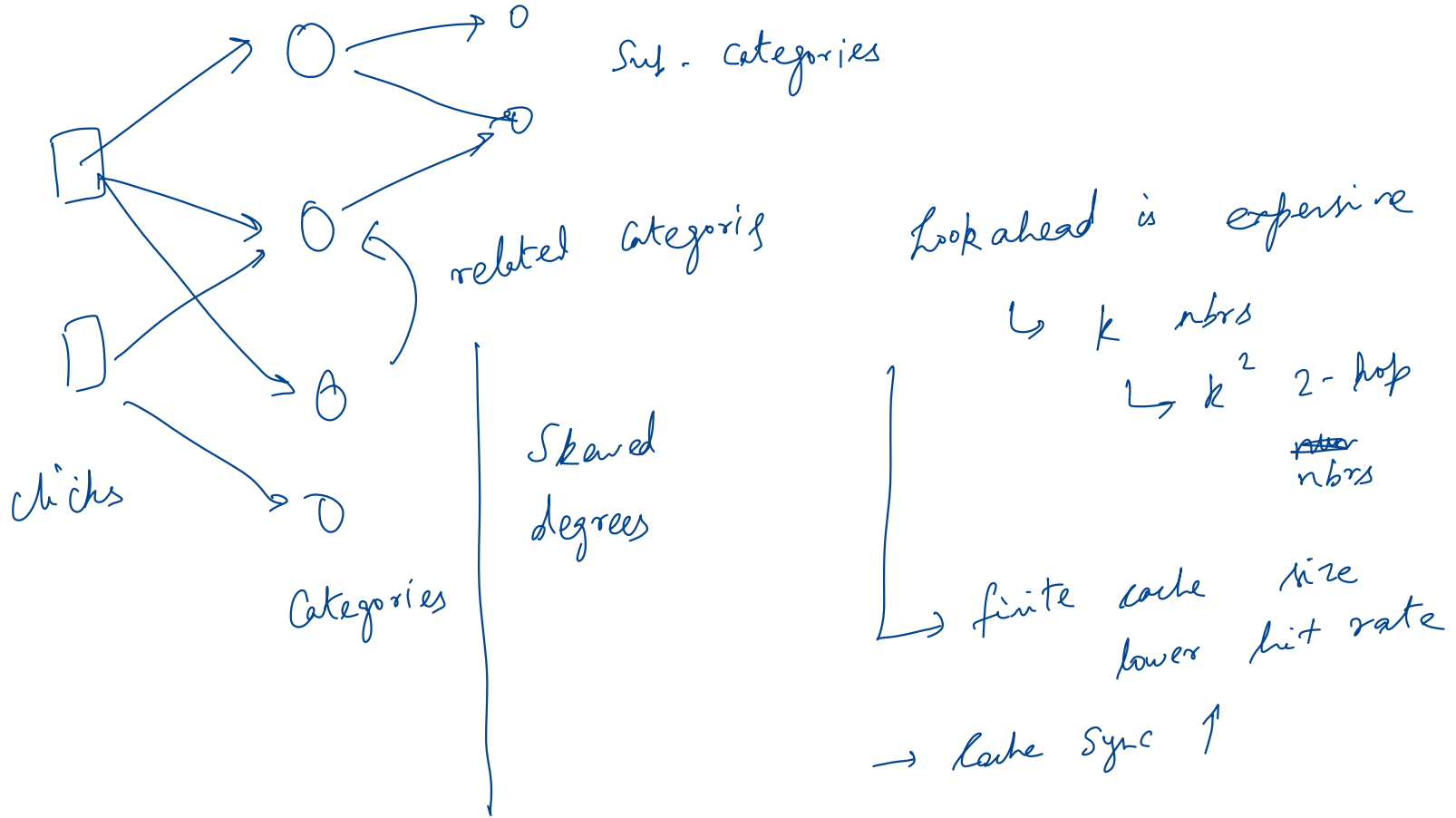




# DISCUSSION

<https://forms.gle/xfTAHiQ5bNENZk7m9>

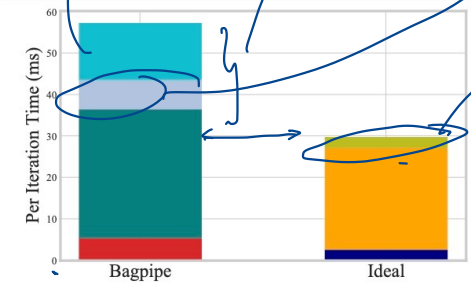
Consider a recommendation model trained on a graph where we use 2-hop neighbors. What are some challenges in using BagPipe-style ideas for such a workload?



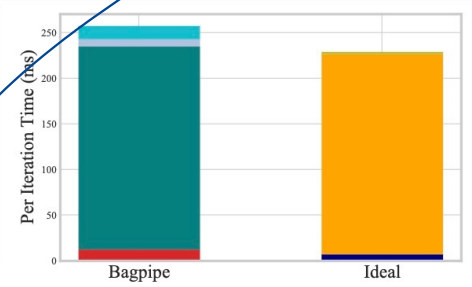
Cache Synch. adds overhead

gap is similar across models

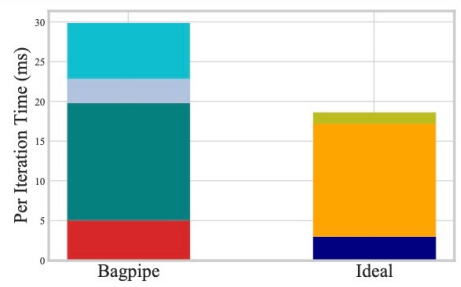
network contention between sync & prefetching etc.



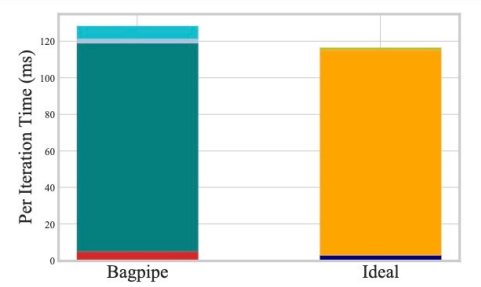
(a) DLRM: p3.2xlarge



(b) DeepFM: p3.2xlarge



(c) DLRM: g5.8xlarge



(d) DeepFM: g5.8xlarge

# NEXT STEPS

Next class: Serverless computing

Project check-ins next week