# CS 744: BAGPIPE

Shivaram Venkataraman

Spring 2024

# ADMINISTRIVIA

- Midterm grades on Gradescope!

   - Submit regrade requests through Gradescope

- Course Project: Check in by April 16th

# RECOMMENDATION MODELS

Examples: DLRM, DeepFM, Wide & Deep

0                    d

N

Clicks

Categories

Category Embeddings

Dense Neural Network

Numerical Features

Age - 28
Price – 20$

Location Embedding Table

| Location | Vector |
|---|---|
| Madison,WI | {0.16, 1.1, ......, 1.2} |
| Koblenz,DE | {0.9, 1.4, ........., 2.9} |
| . | . |
| . | . |
| . | . |
| Detroit,MI | {1.17, 3.1, ......, 0.14} |

Madison,WI

Categorical Features

(hashed) user ID

User ID Embedding Table

(hashed) Product ID

Product ID Embedding Table

Embedding Aggregation

Feature Aggregation

# EMBEDDING TABLES

Convert categorical features to numerical features
Example: Geographic Location to a vector

| Geographical Location Embedding Table |
|---|

Extremely memory intensive, could be up to TBs

| UserID Embedding Table |
|---|

Have sparse access pattern

| Product ID Embedding Table |
|---|

# DISTRIBUTED TRAINING ITERATION

Trainer 1

Trainer 2

Trainer 3

Parameter Server 1

Parameter Server 2

Parameter Server 3

# EMBEDDING ACCESS OVERHEADS



70 % of Time
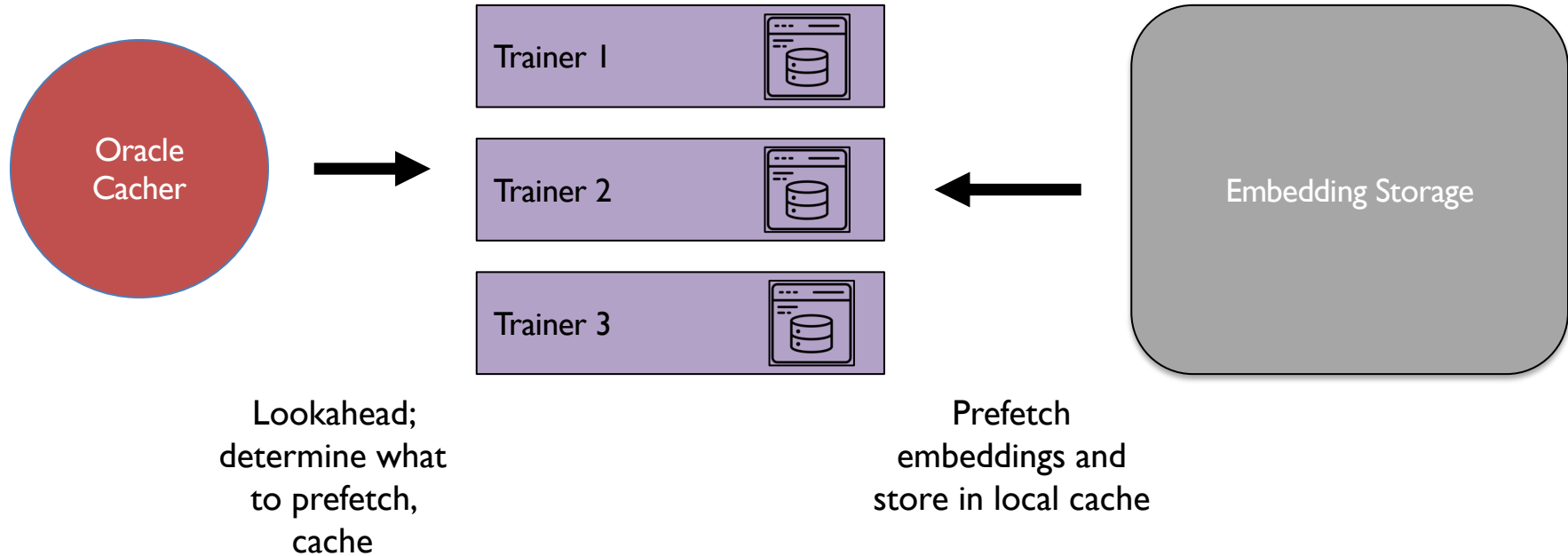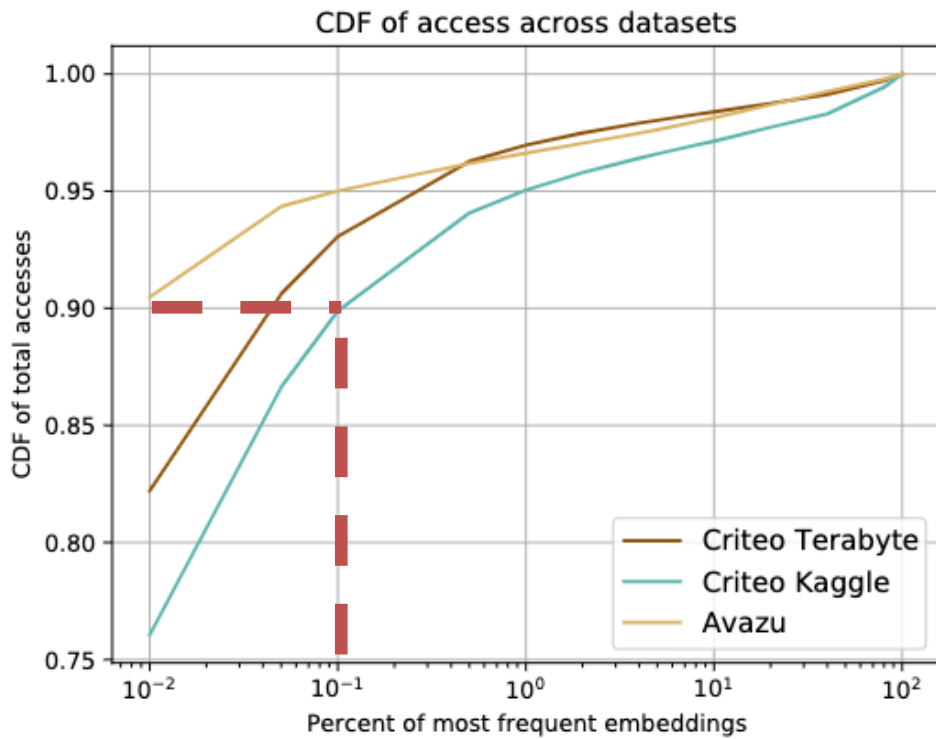
**Setup:**
- DLRM model
- *8 trainers (p3.2xlarge* EC2 instances 1 V100 each node).
- Batch 2048 per machine.
- Criteo Terabyte Dataset

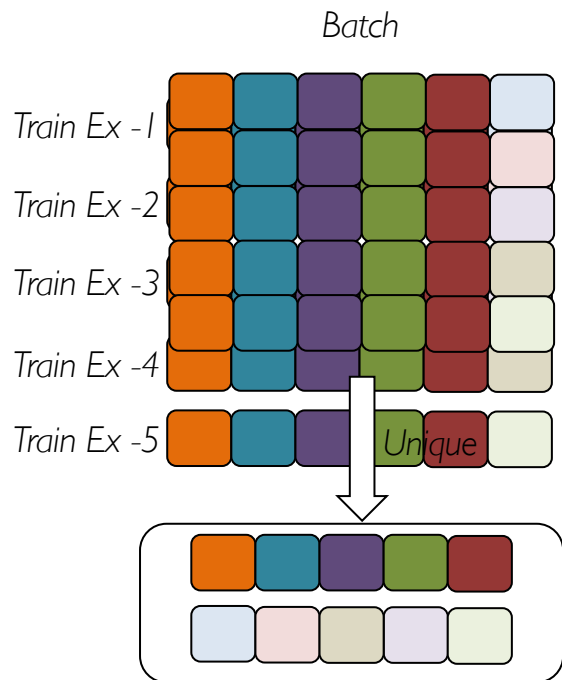Legend: ■ Embedding Sync ■ GetEmb ■ Backward+ MLP Sync ■ Forward

Y-axis: Time (MS) — 0, 20, 40, 60, 80, 100, 120, 140, 160

# BAGPIPE DESIGN



Oracle Cacher

Trainer 1

Trainer 2

Trainer 3

Embedding Storage

Lookahead; determine what to prefetch, cache

Prefetch embeddings and store in local cache

# EMBEDDING ACCESS PATTERNS



CDF of access across datasets

# LONG TAIL ACCESSES
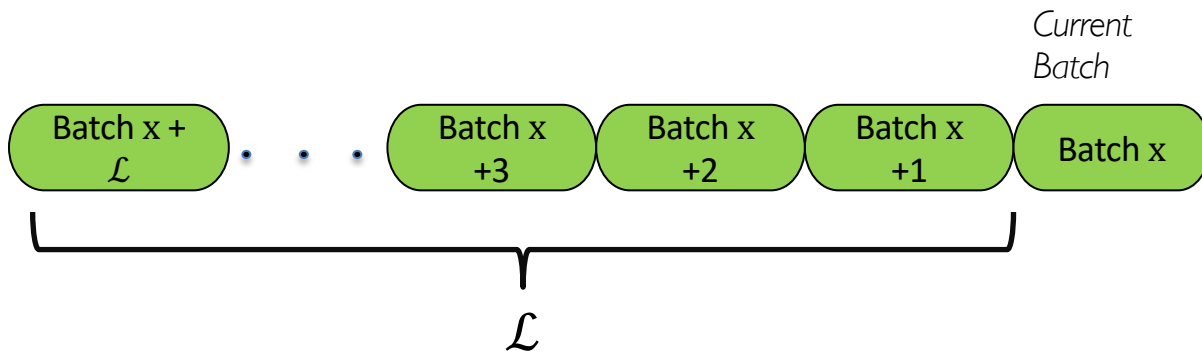


Batch

Train Ex -1
Train Ex -2
Train Ex -3
Train Ex -4
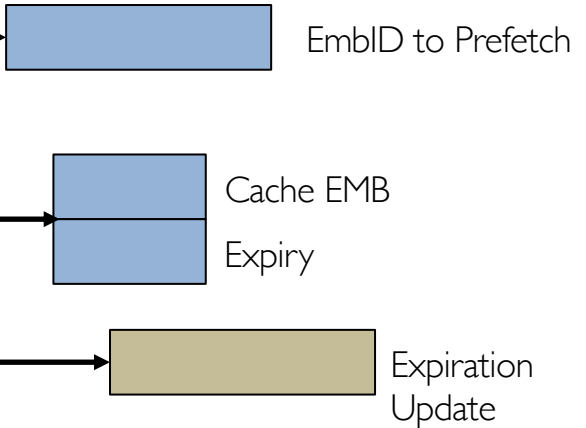Train Ex -5

Unique

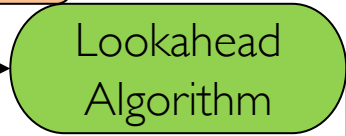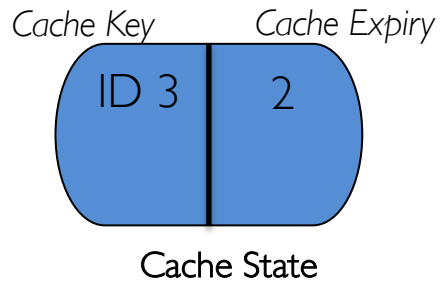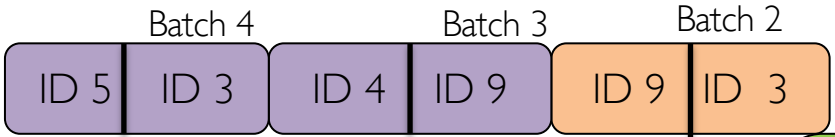Models are trained with a batch of examples.

*- For a batch only fetch unique embeddings*

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

# LOOKAHEAD ALGORITHM

Look at "$\mathcal{L}$" next batches ahead of current batch to extract access pattern of embeddings by future batches

*Current Batch*

| Batch x + $\mathcal{L}$ | . . . | Batch x +3 | Batch x +2 | Batch x +1 | Batch x |

$\mathcal{L}$

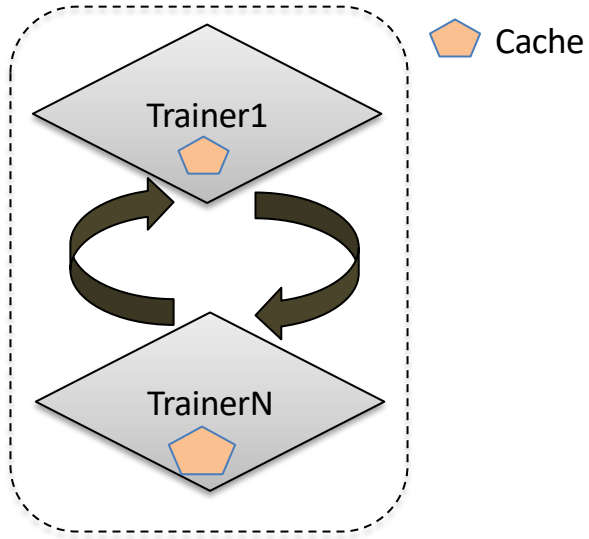Look-ahead Value of 2, Batch size of 2

# 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

# CACHE SYNCHRONIZATION

At the end of each iteration, each trainer synchronizes caches

# SUMMARY

Recommendation models: Embeddings access overheads

BagPipe: Efficient distributed training
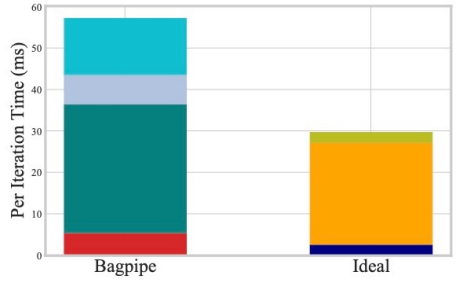
    Lookahead to pre-fetch and cache embeddings
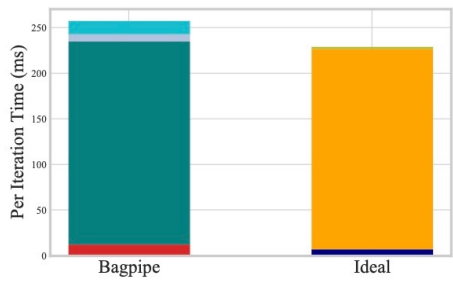
    Cache synchronization across trainers

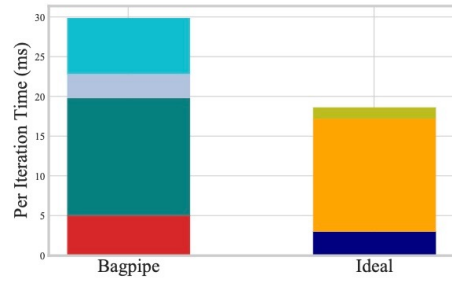# DISCUSSION

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?
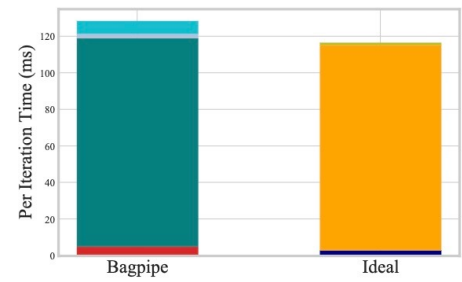
**(a)** DLRM: p3.2xlarge     **(b)** DeepFM: p3.2xlarge     **(c)** DLRM: g5.8xlarge     **(d)** DeepFM: g5.8xlarge

# NEXT STEPS

Next class: Serverless computing

Project check-ins next week