# Leveraging AI for Faster Storage Access: a Graph-Neural-Network-Based Prefetcher

**Faradawn Yang**

Advisor:          Haryadi S. Gunawi
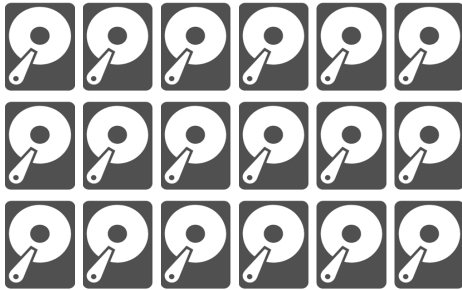
THE UNIVERSITY OF
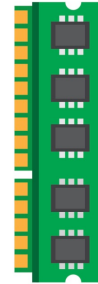## CHICAGO

# Backend storage is slow
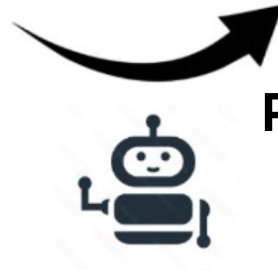
**Bulk Storage**                              **Memory**          **CPU**

🥶
**Slow**          **Fast**

**1.5 day**          **1 second**

**Prefetch!**

# Prefetching improves access speed

**No prefetch**

| Run | Load | | |

**With prefetch**

| Run | |
| Load | |

**Faster!**

# Previous work

['92] Stride

2 4 6 → Stride → 8

['20] Leap

Find majority

2 2 4 2 → 2

['21] LSTM

Neural net

2 4 6 → 8

['24] Baleen

Limit prefetch →

# Current problem

## Modern workload

**User 1**

**User 2**

Address

Time

100  102  104  106

10

Average hit rate on Alibaba, Tencent, and Microsoft traces

## Previous method



Bar chart — Hit Rate vs Algorithm:
- Leap: 8.71
- LSTM: 9.89
- Stride: 10.1
- Markov Chain: 27.05
- Spectral GNN: 31.9

1) Not prone against short-term irregularities.

2) Can't see a global picture.

# Methodology

# Step 1: Convert to deltas

**Address**

| 100 | 102 | 104 | 10 | 106 |

→ Time

**Past 32 accesses.**
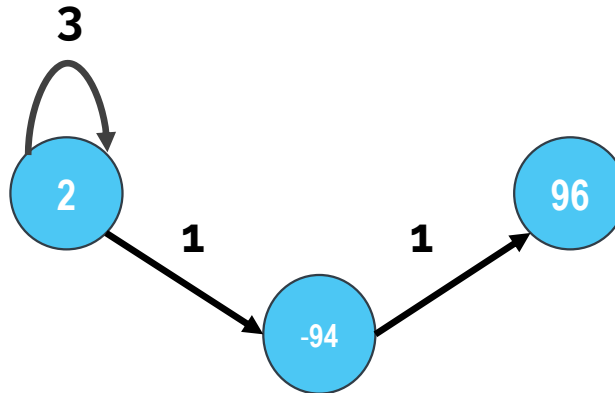
**Delta**

( 2 ) ( 2 ) ( 2 ) ( -94 ) ( 96 )

→ Time

# Step 2: Build a graph
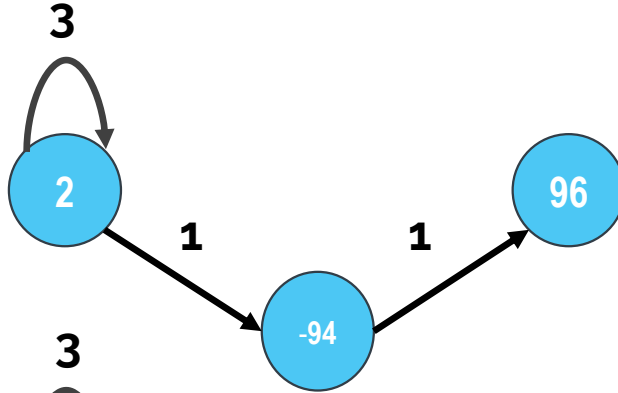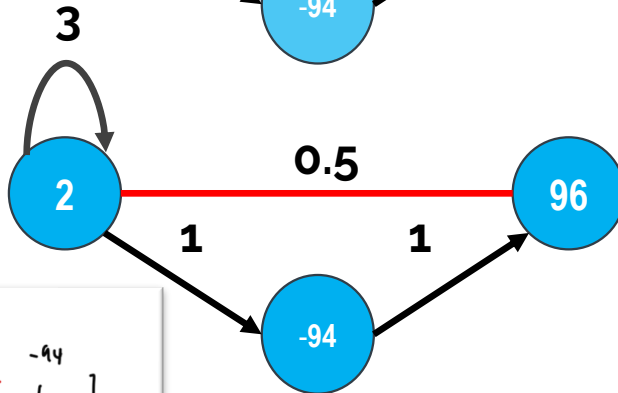
**Delta**



Time

**Graph**

# Step 3: Add full connections

**Graph**

**Directed edges capture temporal information.**

**Fully-connected Graph**

**Fully-connected edges reveal spatial pattern.**



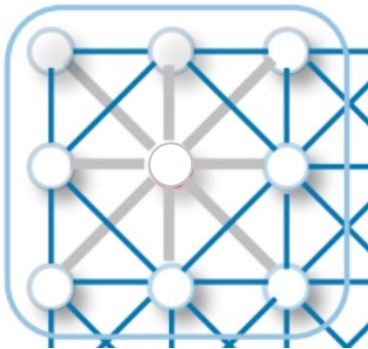**Another example**

# Step 4: Feed into ML



1000 delta addresses.

# What ML model

**[SGDP '23]**

**Message passing**

Local patterns.
Many iterations.
😞

**[Ours]**
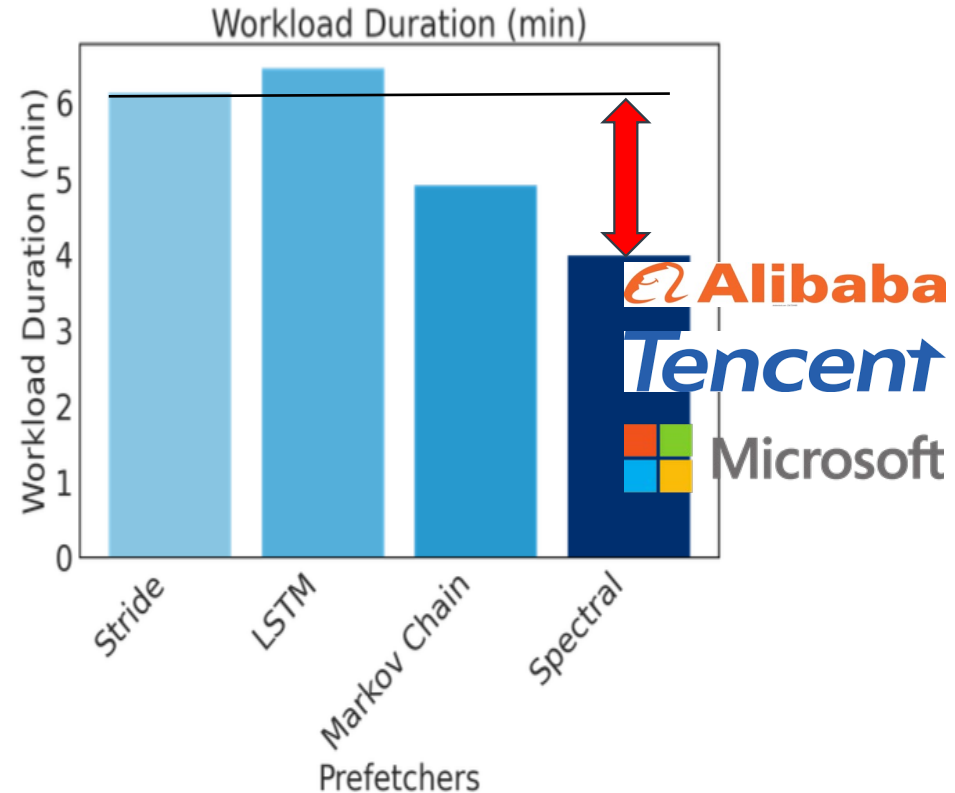
**Spectral network**

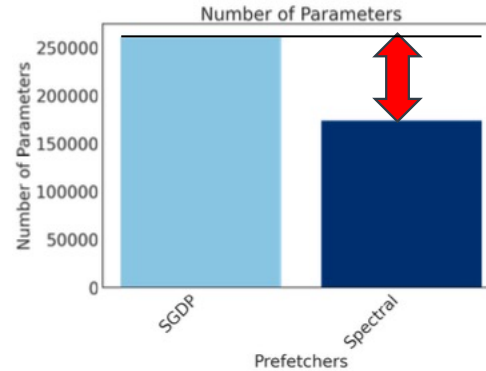Global structure.
One pass.
😄

# Evaluations

# Faster access speed

**Saves 33% access time!**

- ❑ Setup: python simulator.

- ❑ Respected arrival time of requests



Workload Duration (min)

# Saves memory and time

❏ Similar hit rate, but much smaller model and faster training and inference than state-of-the-art, [SGDP '23].



**Saves 33.4% model size.**

**Saves 79% training time.**

# Conclusion

# Entire pipeline



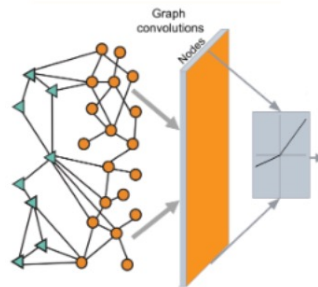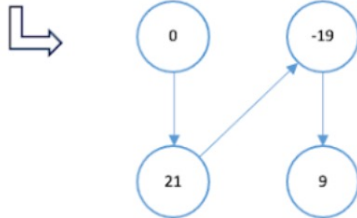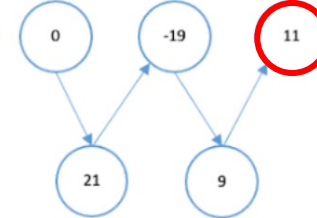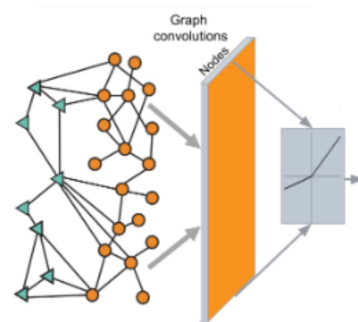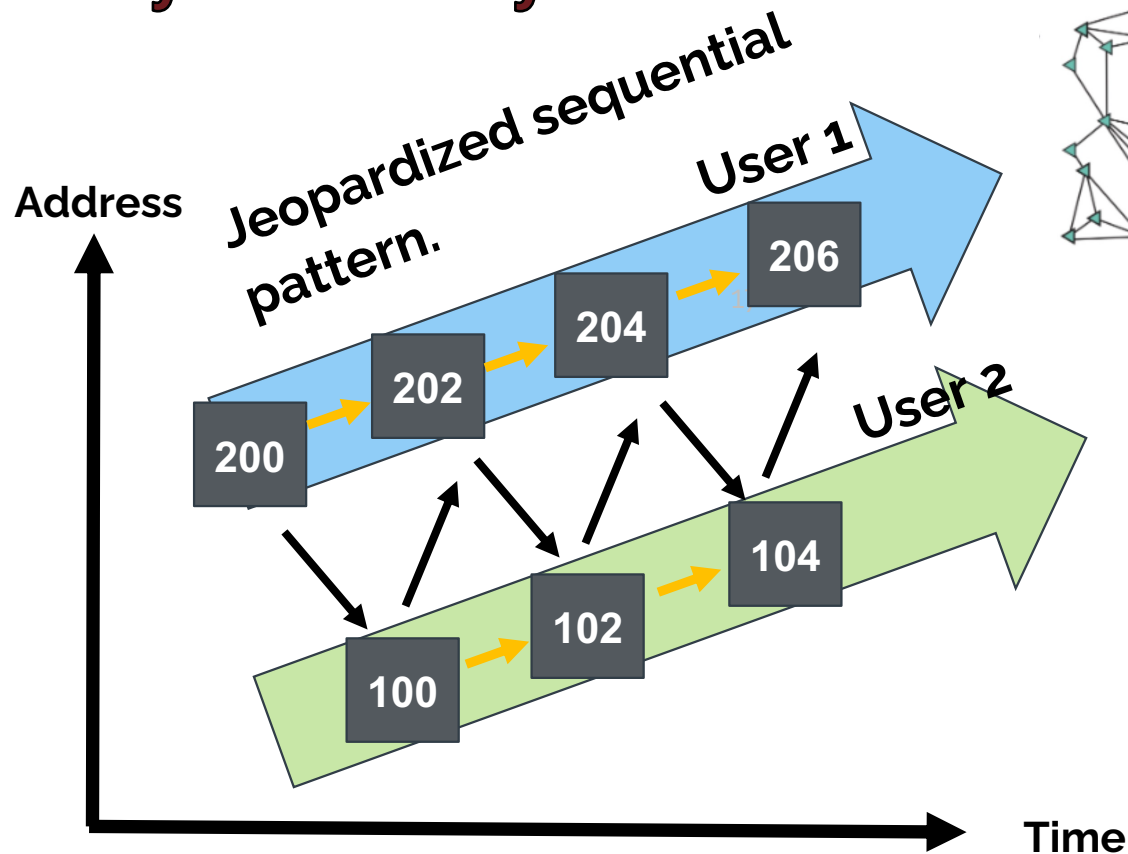**Build a distilled graph.**

**Predict a new node.**

# Key takeaway

Use **graph** to detect a global pattern.

# Thank you!

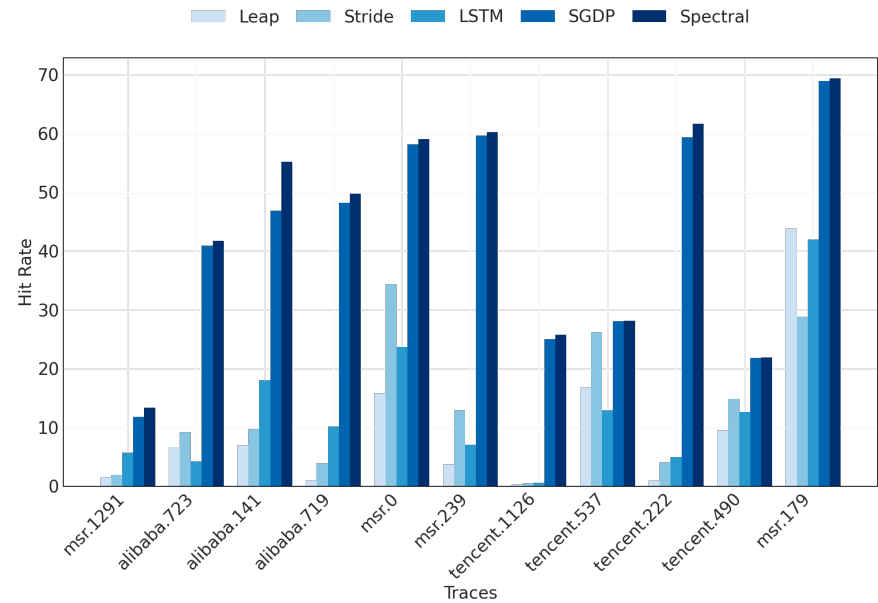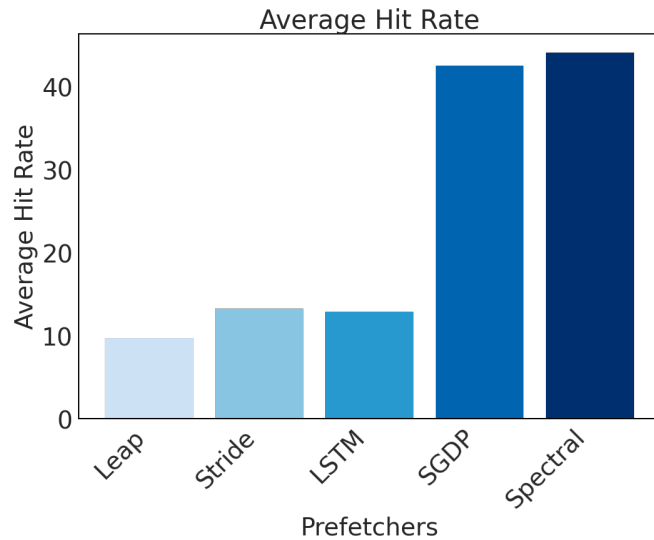Leveraging AI for Faster Storage Access: a Graph-Neural-Network-Based Prefetcher

# Backup Slides

# Is graph building too costly

❑ Only putting 32 integers into a matrix.

❑ Main Bottleneck is disk access (20ms), then inference (2ms), then graph building (1ms).

❑ Spectral < SGDP < LSTM < transformer.

$$\begin{pmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}$$
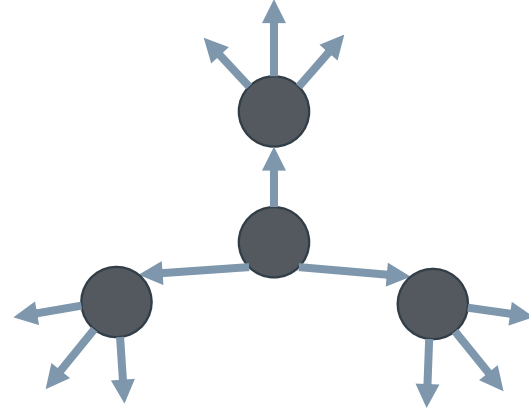
# Higher hit rate

# Why spectral method

## Message passing

$$m_v^{(t+1)} = \sum_{u \in \mathcal{N}(v)} M^{(t)} \left( h_v^{(t)}, h_u^{(t)}, e_{uv} \right)$$

$$h_v^{(t+1)} = U^{(t)} \left( h_v^{(t)}, m_v^{(t+1)} \right)$$

Many iterations

## Spectral method

$$H^{(2)} = \sigma_2 \left( \hat{A} \sigma_1 \left( \hat{A} X W^{(0)} \right) W^{(1)} \right)$$

One go