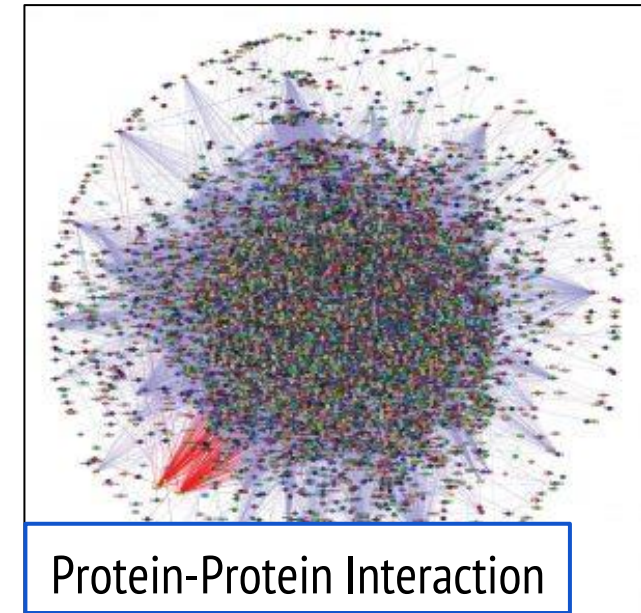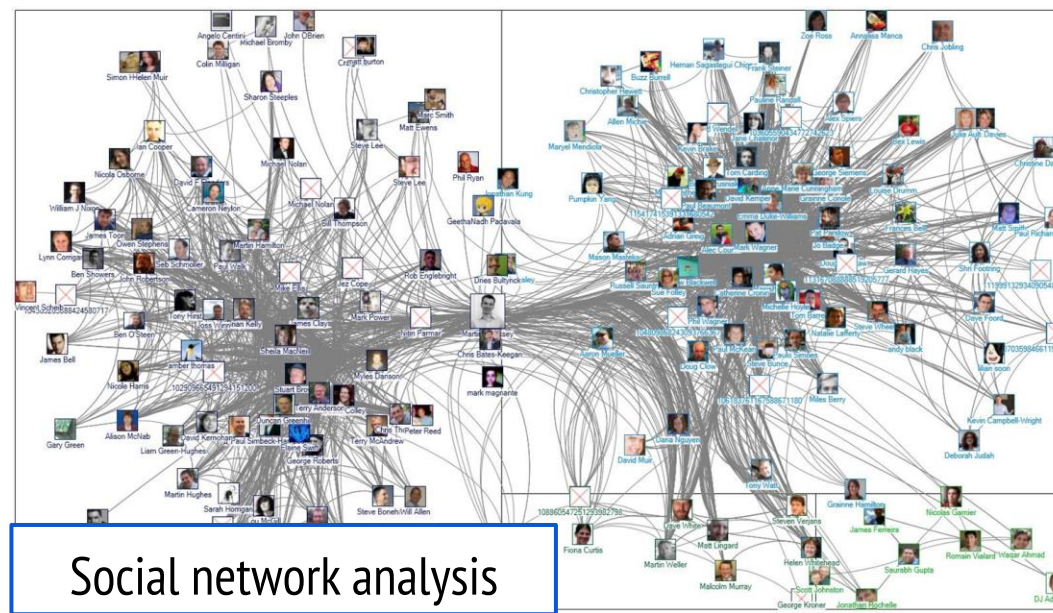# CMU 18-344: Computer Systems and the Hardware/Software Interface

Fall 2024, Prof. Brandon Lucia
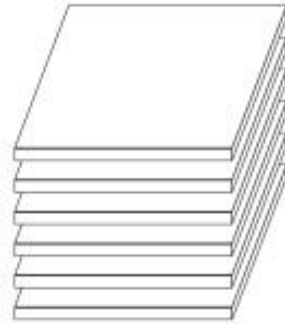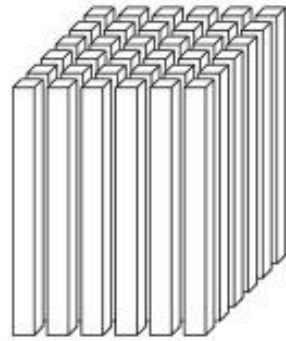
# Today: Sparse Problems
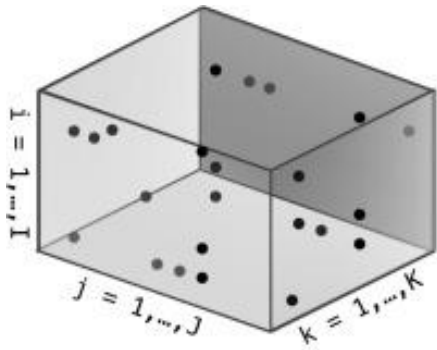
- What is a sparse problem?  Why are they called "sparse"?

- What makes sparse problems hard?

- Roofline performance modeling

- Hardware and software strategies for optimizing sparse problems

# Graph Processing Problems are Sparse Problems



Path Planning

Social network analysis

Protein-Protein Interaction

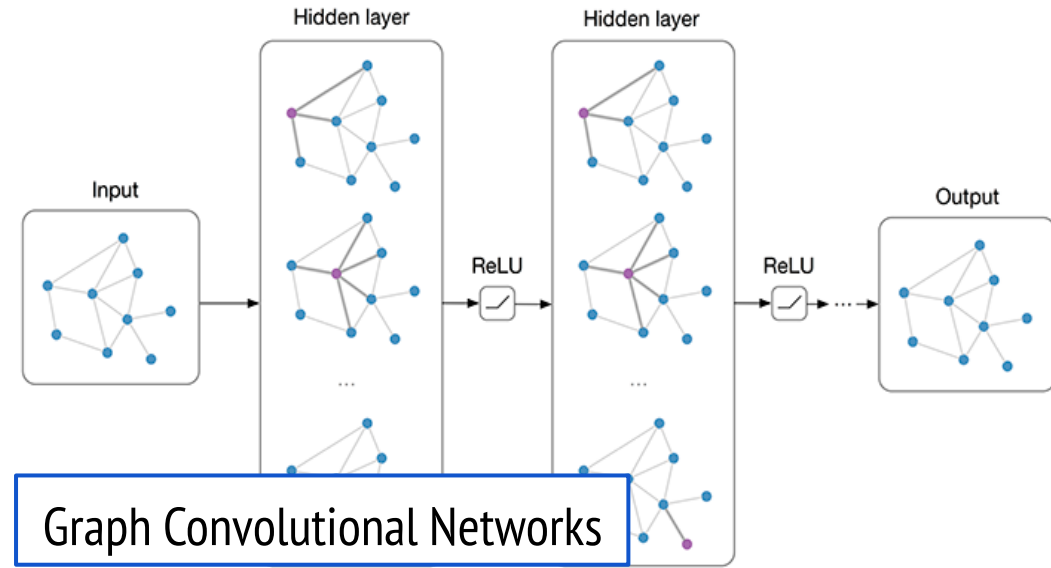The canonical examples of sparse problems are graph processing applications.
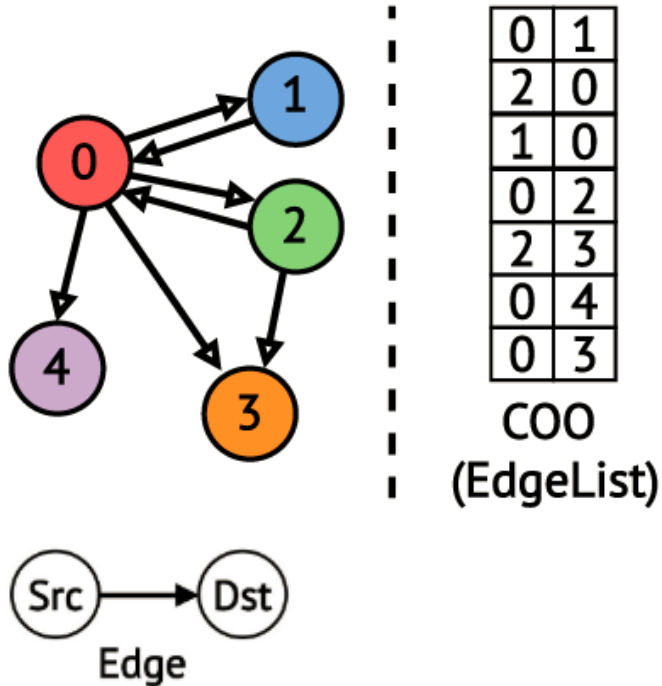
# Machine Learning Problems are Sparse Problems



(b) Mode-1 fibers:
$$\mathbf{f}_{:jk} = \mathcal{X}(:, j, k)$$

(c) Slices:
$$\mathbf{S}_{::k} = \mathcal{X}(:, :, k)$$

Data Mining

Graph Convolutional Networks

# What does a graph processing program look like?



COO
(EdgeList)

Edge

```
for e in EL:
    dstData[e.dst] =
        f(srcData[e.src],dstData[e.dst])
```
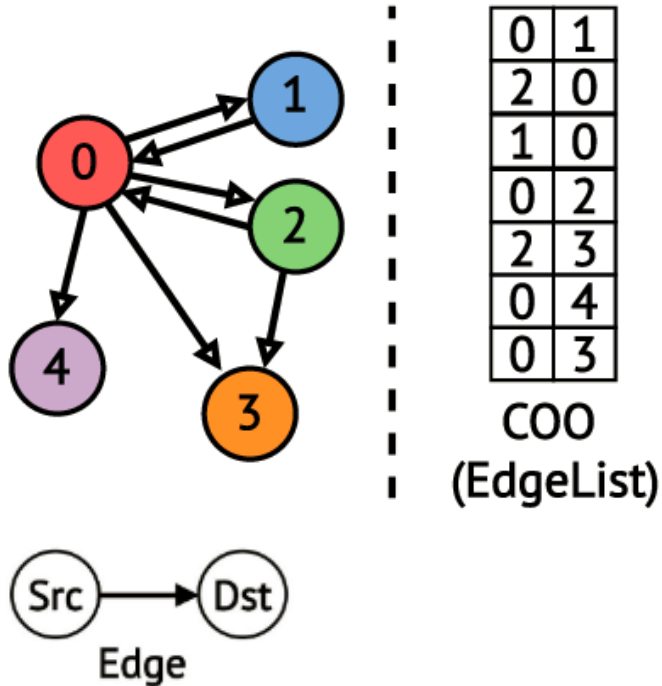
dstData

srcData

**stores vertex property information**
**if srcData == dstData, updating in-place;**
**often "swap" srcData & dstData from 1 iteration to the next iteration**

# What does a graph processing program look like?

| | |
|---|---|
| 0 | 1 |
| 2 | 0 |
| 1 | 0 |
| 0 | 2 |
| 2 | 3 |
| 0 | 4 |
| 0 | 3 |

COO
(EdgeList)

Src —→ Dst

Edge

```
PageRank(-ish){
    for e in EL:
        rank_n[e.dst] =
            rank_nminus1[e.src] + rank_n[e.dst]
}
```
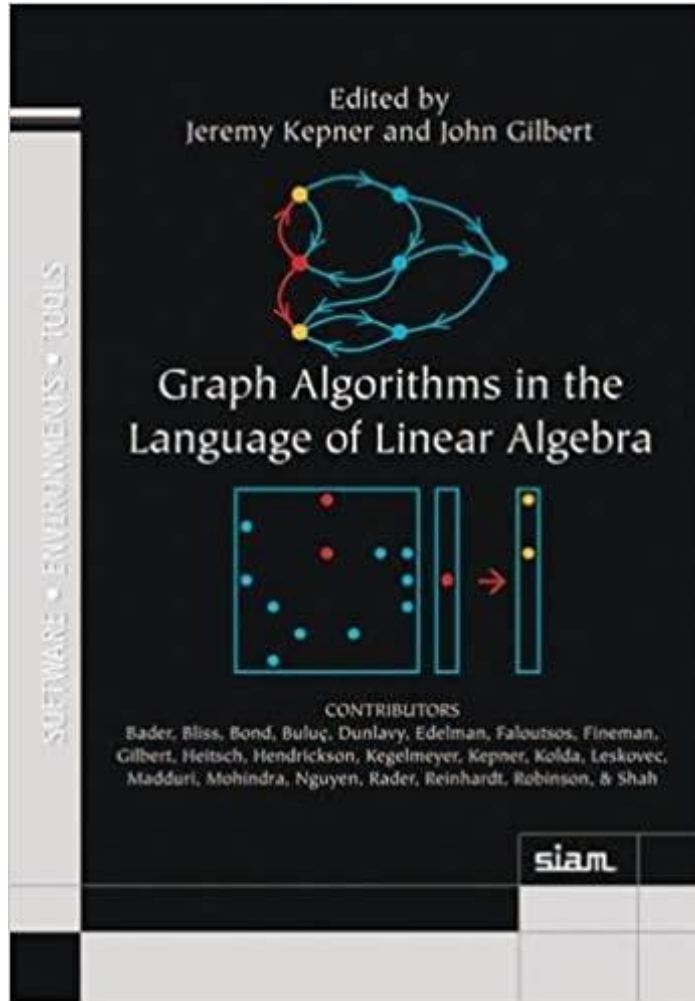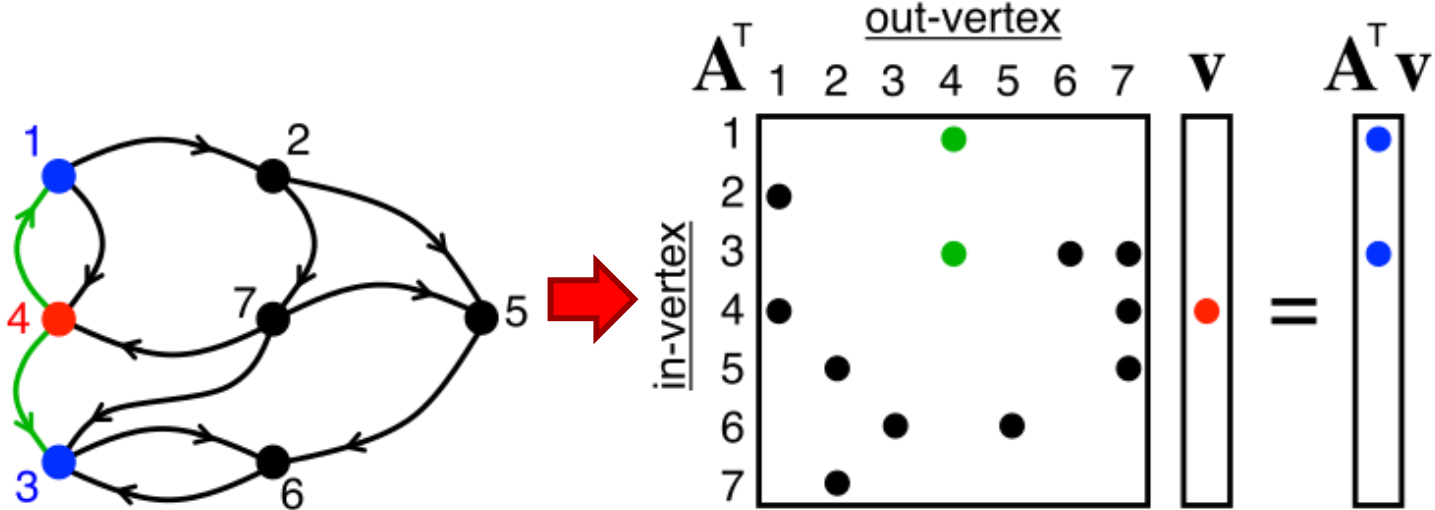
dstData

srcData

**rank_n is a webpage's rank in this iteration,
rank_nminus1 is rank_n from the last iteration**

# Graph Analytics can be mapped to Sparse Linear Algebra

# Graph Analytics can be mapped to Sparse Linear Algebra

# Graph Analytics can be mapped to Sparse Linear Algebra

# How do graph applications correspond to linear algebra?

# How do graph applications correspond to linear algebra?



**Matrix-transpose-vector product is one BFS iteration**

$$A^T x_i = x_{i+1}$$

Initial $x_i$ vector is starting vertex for BFS.

# How do graph applications correspond to linear algebra?



**Matrix-transpose-vector product is one BFS iteration**

$$A^T x_i = x_{i+1}$$

**A** Transpose

Initial $x_i$ vector is starting vertex for BFS.

Initial $x_{i+1}$ is vertices reachable from $x_i$

# How do graph applications correspond to linear algebra?



**Matrix transpose vector product is one BFS iteration**

$$A^T x_i = x_{i+1}$$

**The next iteration is computed by performing the next matrix transpose vector product**

# How do graph applications correspond to linear algebra?



**Matrix transpose vector product is one BFS iteration**

$$A^T x_i = x_{i+1}$$

**The next iteration is computed by performing the next matrix transpose vector product**

# How do graph applications correspond to linear algebra?

**Matrix transpose vector product is one BFS iteration**

$$A^T x_i = x_{i+1}$$

**The next iteration is computed by performing the next matrix transpose vector product**
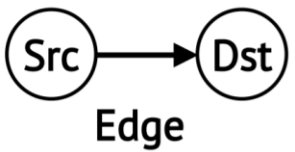
**Search done when no new vertices added (or all visited)**

# How do graph applications correspond to linear algebra?

Turns out that other graph applications also correspond to roughly this formulation if you change the operations you use (min/+ instead of +/*) or consider weighted edges

$$A^T x_i = x_{i+1}$$

SSSP, BFS, PageRank, Connected-Components, Betweenness-Centrality, triangle counting... BFS is a representative sparse problem.

Search done when no new vertices added (or all visited)

# Nobody EVER uses the adjacency matrix!



Why would the Adjacency Matrix not be used?

# Nobody EVER uses the adjacency matrix!

Src → Dst
**Edge**

Dst →

|     | D_0 | D_1 | D_2 | D_3 | D_4 |
|-----|-----|-----|-----|-----|-----|
| S_0 |     |     | 1   |     |     |
| S_1 | 1   |     |     |     | 1   |
| S_2 | 1   | 1   |     | 1   |     |
| S_3 |     | 1   |     |     | 1   |
| S_4 | 1   |     | 1   |     |     |

Src ↓

Reasons Adjacency Matrix is never used:
- **Sparsity:** % of Non-Zero Entries ~ $10^{-5}$
- **Total Size:** 32M nodes => (32M * 32M) = 1PB

# Compressed Sparse Data Structures for Feasible Memory Size



Offsets Array (OA)

Neighbors Array (NA)

**Compressed Sparse Row (CSR)**
*Outgoing Neighbors*

Vertex Property Array
i.e., srcData / dstData

**Often we will leave the vertex property array implicitly defined when we talk about sparse structures, but it is always there**

# Compressed Sparse Data Structures for Feasible Memory Size

Offsets Array (OA)

| 0 | 1 | 3 | 6 | 8 |

**index is src id**

Neighbors Array (NA)

| 2 | 0 | 4 | 0 | 1 | 3 | 1 | 4 | 0 | 2 |

EdgeList sorted by SrcIDs

**Compressed Sparse Row (CSR)**
*Outgoing Neighbors*

Src → Dst
**Edge**

Dst →

|         | $D_0$ | $D_1$ | $D_2$ | $D_3$ | $D_4$ |
|---------|-------|-------|-------|-------|-------|
| $S_0$   |       |       | 1     |       |       |
| $S_1$   | 1     |       |       |       | 1     |
| $S_2$   | 1     | 1     |       | 1     |       |
| $S_3$   |       | 1     |       |       | 1     |
| $S_4$   | 1     |       | 1     |       |       |

Src ↓

**OA indexed by vertex ID of src of edge**
**Value in OA is *offset* into NA**

**start index for edges w/ src == vertex i** = OA[i]
**#edges with src == vertex i** = OA[i+1] – OA[i]

**Dense** encoding of **sparse** structure

20

# Compressed Sparse Data Structures for Feasible Memory Size



Offsets Array (OA)
`0  1  3  6  8`

Neighbors Array (NA)
`2  0  4  0  1  3  1  4  0  2`

EdgeList sorted by SrcIDs

**Compressed Sparse Row (CSR)**
*Outgoing Neighbors*

**The CSC is the *transpose* of the CSR**

Offsets Array (OA)
`0  3  5  7  8`

Neighbors Array (NA)
`1  2  4  2  3  0  4  2  1  3`

EdgeList sorted by DstIDs

**Compressed Sparse Column (CSC)**
*Incoming Neighbors*

21

# Building the CSR / CSC from a Graph's Edge List



COO
(EdgeList)

| 0 | 1 |
| 2 | 0 |
| 1 | 0 |
| 0 | 2 |
| 2 | 3 |
| 0 | 4 |
| 0 | 3 |

```
for e in EL:
    neigh_count[e.dst]++; /*e.src*/
```

| 2 | 1 | 1 | 2 | 1 |

neigh_count

Src → Dst
Edge

# Building the CSR / CSC from a Graph's Edge List

| 0 | 1 |
|---|---|
| 2 | 0 |
| 1 | 0 |
| 0 | 2 |
| 2 | 3 |
| 0 | 4 |
| 0 | 3 |

COO
(EdgeList)

Src → Dst

Edge

```
for e in EL:
    neigh_count[e.dst]++; /*e.src*/
```

| 2 | 1 | 1 | 2 | 1 |
|---|---|---|---|---|

neigh_count

| 2 | 1 | 1 | 2 | 1 |
|---|---|---|---|---|

neigh_count_dup

```
sum = 0
for i in 0 .. |V|:
    tmp = neigh_count[i]
    neigh_count[i] = sum;
    neigh_count_dup[i] = sum;
    sum += tmp
```

# Building the CSR / CSC from a Graph's Edge List



```
for e in EL:
    neigh_count[e.dst]++; /*e.src*/
```

| | | | | |
|---|---|---|---|---|
| 2 | 1 | 1 | 2 | 1 |

**neigh_count**

```
sum = 0
for i in 0 .. |V|:
    tmp = neigh_count[i]
    neigh_count[i] = sum; //OA
    neigh_count_dup[i] = sum;
    sum += tmp
```

| | | | | |
|---|---|---|---|---|
| 0 | 2 | 3 | 4 | 6 |

**OA (also OA_dup)**

COO
(EdgeList)

| 0 | 1 |
|---|---|
| 2 | 0 |
| 1 | 0 |
| 0 | 2 |
| 2 | 3 |
| 0 | 4 |
| 0 | 3 |

Src ⟶ Dst
Edge

# Building the CSR / CSC from a Graph's Edge List

COO
(EdgeList)

| | |
|---|---|
| 0 | 1 |
| 2 | 0 |
| 1 | 0 |
| 0 | 2 |
| 2 | 3 |
| 0 | 4 |
| 0 | 3 |

Src → Dst
Edge

OA (also OA_dup)

| 0 | 2 | 3 | 4 | 6 |
|---|---|---|---|---|

```
for e in EL:
  neigh_ind = OA[e.src]
  NA[neigh_ind] = e.dst
  OA[e.src]++ /*sacrificial OA*/
//i.e., NA[ OA[e.src]++ ] = e.dst
```

OA_dup

| 0 | 2 | 3 | 4 | 6 |
|---|---|---|---|---|

NA

| 1 | 2 | 0 | 0 | 0 | 2 | 0 |
|---|---|---|---|---|---|---|

**Completed CSC**

# Compressed Representations ⇒ Irregular Memory Accesses



**Pull (CSC Traversal)**

```
for dst in G:
    for src in in_neighs(dst):
        dstData[dst] += srcData[src]
```

**Pull traversal performs *irregular read operations* that lack locality**

**i.e., $x_{i+1}$**

**e.g., current rank of page I,
e.g., current shortest path
from source vertex**

**CSC**

# Compressed Representations ⇒ Irregular Memory Accesses

**Push (CSR Traversal)**

```
for src in G:
    for dst in out_neighs(src):
        dstData[dst] += srcData[src]
```

**Push traversal performs *irregular write operations* that lack locality**

**i.e., $x_{i+1}$**

**dstData**

0   1   2   3   4

**srcData**  | 5 | 20 | 10 | 2 | 1 |

**e.g., current rank of page I,
e.g., current shortest path
from source vertex**

**OA**  | 0 | 1 | 3 | 6 | 8 |

**NA**  | 2 | 0 | 4 | 0 | 1 | 3 | 1 | 4 | 0 | 2 |

**CSR**

# Compressed Representations ⇒ Irregular Memory Accesses

Dst →

**Push (CSR Traversal)**

```
for src in G:
    for dst in out_neighs(src):
        dstData[dst] += srcData[src]
```

**Push traversal performs *irregular write operations* that lack locality**

**i.e., $x_{i+1}$**

Irregular Data Footprint >> LLC Size

Size of srcData ~ **256MB** (32M vertices * 8B)

e.g., current rank of page i,
**e.g., current shortest path
from source vertex**

**CSR**

# Irregular Accesses Lead to Poor Locality

## LLC Miss Rate (%)

**Why such bleak cache performance?**
**Consequence of bleak cache performance?**



**Running on RMAT27**
**Graph w/ 35MB LLC**

*Right Figure from "Optimizing Cache Performance for Graph Analytics" ArXiv v1;*

# Irregular Accesses Lead to Poor Locality

**LLC Miss Rate (%)**



**Running on RMAT27 Graph w/ 35MB LLC**

**Dst coordinate of edge is index in dstData: totally input dependent & random!!!**



COO
(EdgeList)

**dstData**
**Remember: dstData[e.dst] ++**
**and e.dst is random, from edge list**

*Right Figure from "Optimizing Cache Performance for Graph Analytics" ArXiv v1;*

# Irregular Accesses Lead to Poor Locality

**LLC Miss Rate (%)**



**Running on RMAT27
Graph w/ 35MB LLC**

**Dst coordinate of edge is index in dstData:
totally input dependent & random!!!**

COO
(EdgeList)

miss

**dstData**
**Remember: dstData[e.dst] ++
and e.dst is random, from edge list**

*Right Figure from "Optimizing Cache Performance for Graph Analytics" ArXiv v1;*

# Irregular Accesses Lead to Poor Locality

## LLC Miss Rate (%)



Running on RMAT27
Graph w/ 35MB LLC

**Dst coordinate of edge is index in dstData: totally input dependent & random!!!**



COO
(EdgeList)

miss

**dstData**
**Remember: dstData[e.dst] ++**
**and e.dst is random, from edge list**

*Right Figure from "Optimizing Cache Performance for Graph Analytics" ArXiv v1;*

# Irregular Accesses Lead to Poor Locality

**LLC Miss Rate (%)**



**Running on RMAT27 Graph w/ 35MB LLC**

**Dst coordinate of edge is index in dstData: totally input dependent & random!!!**

COO (EdgeList)

(You get lucky sometimes)

hit

**dstData**

**Remember: dstData[e.dst] ++ and e.dst is random, from edge list**

*Right Figure from "Optimizing Cache Performance for Graph Analytics" ArXiv v1;*

# Irregular Accesses Lead to Poor Locality

## LLC Miss Rate (%)



Running on RMAT27
Graph w/ 35MB LLC

**Dst coordinate of edge is index in dstData: totally input dependent & random!!!**

| | |
|---|---|
| 0 | 1 |
| 2 | 0 |
| 1 | 0 |
| 0 | 2 |
| 2 | 3 |
| 0 | 4 |
| 0 | 1 |

COO
(EdgeList)

miss

| | 0 | | | 0 | |
|---|---|---|---|---|---|

| 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|

**dstData**
**Remember: dstData[e.dst] ++ and e.dst is random, from edge list**

*Right Figure from "Optimizing Cache Performance for Graph Analytics" ArXiv v1;*

# Irregular Accesses Lead to Poor Locality

**LLC Miss Rate (%)**



**Running on RMAT27 Graph w/ 35MB LLC**

**Dst coordinate of edge is index in dstData: totally input dependent & random!!!**



COO (EdgeList)

miss

**dstData**

**Remember: dstData[e.dst] ++ and e.dst is random, from edge list**

*Right Figure from "Optimizing Cache Performance for Graph Analytics" ArXiv v1;*

# Irregular Accesses Lead to Poor Locality

**LLC Miss Rate (%)**

**Dst coordinate of edge is index in dstData: totally input dependent & random!!!**

miss

COO
(EdgeList)

**dstData**

**Running on RMAT27 Graph w/ 35MB LLC**

**Remember: dstData[e.dst] ++ and e.dst is random, from edge list**

*Right Figure from "Optimizing Cache Performance for Graph Analytics" ArXiv v1;*

# Irregular Accesses Lead to Poor Locality

**LLC Miss Rate (%)**



**Running on RMAT27
Graph w/ 35MB LLC**

**Dst coordinate of edge is index in dstData:
totally input dependent & random!!!**

COO
(EdgeList)

miss

**dstData**
**Remember: dstData[e.dst] ++
and e.dst is random, from edge list**

*Right Figure from "Optimizing Cache Performance for Graph Analytics" ArXiv v1;*

# Irregular Accesses Lead to Poor Locality

**LLC Miss Rate (%)**



**Cycles stalled on DRAM / Total Cycles**



**Cache miss latency *cannot be hidden by anything else in the program*. Each miss incurs DRAM latency!**

*Right Figure from "Optimizing Cache Performance for Graph Analytics" ArXiv v1;*

# Irregular Accesses Lead to Poor Locality

**LLC Miss Rate (%)**



**Cycles stalled on DRAM / Total Cycles**



**Problem: Sparse representations make processing large graphs feasible, but graph processing still entails a large working set with poor locality**

*Right Figure from "Optimizing Cache Performance for Graph Analytics" ArXiv v1;*

# Even Building the CSR / CSC is an Irregular Access Pattern!

```
for e in EL:
    neigh_count[e.dst]++;
```

| 0 | 1 |
| 2 | 0 |
| 1 | 0 |
| 0 | 2 |
| 2 | 3 |
| 0 | 4 |
| 0 | 3 |

COO
(EdgeList)

| 2 | 1 | 1 | 2 | 1 | neigh_count

Src → Dst

Edge

## Why is this irregular?

# Even Building the CSR / CSC is an Irregular Access Pattern!

```
for e in EL:
    neigh_count[e.dst]++; /*e.src*/
```

neigh_count

Updates to the neigh_count array are to random elements determined by order of edges in edge list

COO
(EdgeList)

| 0 | 1 |
| 2 | 0 |
| 1 | 0 |
| 0 | 2 |
| 2 | 3 |
| 0 | 4 |
| 0 | 3 |

Src → Dst
Edge

| 2 | 1 | 1 | 2 | 1 |

# Even Building the CSR / CSC is an Irregular Access Pattern!



```
for e in EL:
    NA[ OA[e.src]++ ] = e.dst
```

COO (EdgeList):

| 0 | 1 |
|---|---|
| 2 | 0 |
| 1 | 0 |
| 0 | 2 |
| 2 | 3 |
| 0 | 4 |
| 0 | 3 |

OA:

| 0 | 2 | 3 | 4 | 6 |
|---|---|---|---|---|

NA:

| 1 | 2 | 0 | 0 | 0 | 2 | 0 |
|---|---|---|---|---|---|---|

**Completed CSC**

Src → Dst
Edge

Why is the NA update part irregular?

# Even Building the CSR / CSC is an Irregular Access Pattern!



```
for e in EL:
    NA[ OA[e.src]++ ] = e.dst
```

COO
(EdgeList)

Src → Dst
Edge

OA

NA

Completed CSC

Updates to NA based on EL order & OA[e.src]
NA[ OA[e.src]++ ] = e.dst

# Roofline Performance Analysis of Graph Applications

# The Roofline Model

CPU
(compute, flop/s)

DRAM Bandwidth
(GB/s)

DRAM
(data, GB)

*GFLOPS = Giga-Floating Point Operations Per Second*

*Yes, this is not a proper acronym*

Memory-
Bound

Compute-
Bound

*Peak ops/s*

Throughput
*(operations per second)*

*Peak Mem BW*

Operational Intensity
*(operations per byte)*

# The Roofline Model



CPU (compute, flop/s)

DRAM Bandwidth (GB/s)

DRAM (data, GB)

What does Roofline help us understand about a program?

Memory-Bound ⟵ ⟶ Compute-Bound

Throughput *(GFLOP/s)*

*Peak FLOPS*

Peak Mem BW

● App 2

● App 1

Operational Intensity *(FLOPS/Byte)*

46

# The Roofline Model



CPU
(compute, flop/s)

DRAM Bandwidth
(GB/s)

DRAM
(data, GB)

What does Roofline help us understand about a program?
**Tell us what limits performance & how close to peak an app is.**

Memory-Bound

Compute-Bound

*Peak FLOPS*

Throughput
*(GFLOP/s)*

Peak Mem BW

App 2

App 1

Operational Intensity
*(FLOPS/Byte)*

# The Roofline Model



CPU
(compute, flop/s)

DRAM Bandwidth
(GB/s)

DRAM
(data, GB)

Memory-Bound

Compute-Bound

*Peak FLOPS*

Throughput
*(GFLOP/s)*

Peak Mem BW

● App 2

● App 1

"Ridge point" is a property of a particular machine

What does Roofline help us understand about a program?
**Tell us what limits performance & how close to peak an app is.**

Operational Intensity
*(FLOPS/Byte)*

48

# The Roofline Model



CPU (compute, flop/s)

DRAM Bandwidth (GB/s)

DRAM (data, GB)

Memory-Bound ⬅ ➡ Compute-Bound

*Peak FLOPS*

Throughput *(GFLOP/s)*

*Peak Mem BW*

● App 2

● App 1

"Ridge point" is a property of a particular machine

As a program does more operations per byte, memory has more time to deliver next byte, **relieving Mem BW pressure & increasing compute pressure**

Operational Intensity *(FLOPS/Byte)*

# The Roofline Model



Memory-
Bound

Compute-
Bound

CPU
(compute, flop/s)

DRAM Bandwidth
(GB/s)

DRAM
(data, GB)

Throughput
*(GFLOP/s)*

*Peak Mem BW*

*Peak FLOPS*

App 2

App 1

What is this point?

"Ridge point" is a property of a particular machine

As a program does more operations per byte, memory has more time to deliver next byte, **relieving Mem BW pressure & increasing compute pressure**

Operational Intensity
*(FLOPS/Byte)*

50

# The Roofline Model

CPU (compute, flop/s)

DRAM Bandwidth (GB/s)

DRAM (data, GB)

As a program does more operations per byte, memory has more time to deliver next byte, **relieving Mem BW pressure & increasing compute pressure**

Throughput *(GFLOP/s)*

Memory-Bound

Compute-Bound

*Peak FLOPS*

*Peak Mem BW*

● App 2

**Compare App1 and App2.** What are they doing differently from one another?

● App 1

**What is this point?**

Operational Intensity *(FLOPS/Byte)*

51

# Operational Intensity of Irregular Graph Applications

| 0 | 1 |
|---|---|
| 2 | 0 |
| 1 | 0 |
| 0 | 2 |
| 2 | 3 |
| 0 | 4 |
| 0 | 3 |

COO
(EdgeList)

```
for e in EL:
    dstData[e.dst] += srcData[e.src]
```

What is the operational intensity of a
random update kernel like this one?

# Operational Intensity of Irregular Graph Applications



COO
(EdgeList)

```
for e in EL:
    dstData[e.dst] += srcData[e.src]
```

What is the operational intensity of a
random update kernel like this one?
**Operations per byte:**

# Operational Intensity of Irregular Graph Applications



COO
(EdgeList)

```
for e in EL:
    dstData[e.dst] += srcData[e.src]
```

What is the operational intensity of a random update kernel like this one?

**Operations per byte:**

**Operations:** 1 addition

**Bytes to Load:** 8B for edge, 4B srcData, 4B dstData

**Operational Intensity** = 1 / (8+4+4) = **1/16**

# Graph Applications are Memory-Bound



Memory-Bound ⟵

Compute-Bound ⟶

Peak FLOPS

Throughput (GFLOP/s)

Peak Mem BW

1/16

250

Operational Intensity (FLOPS/Byte)

# Graph Applications are Memory-Bound

Memory-
Bound

Compute-
Bound

*Peak*
*FLOPS*

Throughput
*(GFLOP/s)*

Peak Mem BW

DRAM BW utilization in graph apps is ~50%

**Why would we have spare BW capacity to go to memory and not use it?**

*1/16*

*250*

Operational Intensity
*(FLOPS/Byte)*

# Graph Applications are Memory-Bound

DRAM BW utilization in graph apps is ~50%

**Why would we have spare BW capacity to go to memory and not use it?**

<span style="color:red">**Don't know what to fetch next (no temporal locality), can't use extra stuff we fetch (no spatial locality). Limited ability to send more memory requests (limited mem. parallelism).**</span>

Memory-Bound ⇐   ⇒ Compute-Bound

*Peak FLOPS*

Throughput *(GFLOP/s)*

*Peak Mem BW*

*1/16*

*250*

Operational Intensity *(FLOPS/Byte)*

# Graph Applications are Memory-Bound

Memory-
Bound

Compute-
Bound

*Peak FLOPS*

Throughput
*(GFLOP/s)*

*Peak Mem BW*

How to improve BW utilization?

**Option #1**: Improve Locality → Reduce Bytes moved → Improve OI

*1/16*

*250*

Operational Intensity
*(FLOPS/Byte)*

# Graph Applications are Memory-Bound

Throughput *(GFLOP/s)*

Memory-Bound ← → Compute-Bound

*Peak Mem BW*

*Peak FLOPS*

**Option #2**: Improve Memory to handle more parallel requests

**Option #1**: Improve Locality → Reduce Mem → Improve OI

How to improve BW utilization?

1/16          250

Operational Intensity *(FLOPS/Byte)*

# Operational Intensity of Irregular Graph Applications



COO
(EdgeList)

```
for e in EL:
    dstData[e.dst] += srcData[e.src]
```

Ideal Best Possible Operational Intensity?

**Operations per byte:**

**Operations:** 1 addition

**Bytes to Load:**

**Operational Intensity =**

# Ideal Operational Intensity of Irregular Graph Applications



COO
(EdgeList)

```
for e in EL:
   dstData[e.dst] += srcData[e.src]
```

Ideal Best Possible Operational Intensity?

**Operations per byte:**

**Operations:** 1 addition

**Bytes to Load:** 8B for edge, 0B srcData, 0B dstData

**Operational Intensity =** 1 / (8+0+0) **= 1/8**

# Improving Performance by Improving Locality

# Improving Performance by Improving Locality



**CPU**
(compute, flop/s)

DRAM Bandwidth
(GB/s)

**DRAM**
(data, GB)

Throughput
*(GFLOP/s)*

Memory-
Bound

Compute-
Bound

*Peak
FLOPS*

*Peak Cache BW*

*Peak Mem BW*

**Locality wins:** If we can operate out of cache, higher ceiling.

Why is cache BW > DRAM BW?

*1/16*    *1/8*         *250*

Operational Intensity
*(FLOPS/Byte)*

# Improving Performance by Improving Locality

# Improving Performance by Improving Locality



CPU
(compute, flop/s)

DRAM Bandwidth
(GB/s)

DRAM
(data, GB)

Throughput
*(GFLOP/s)*

Memory-
Bound

Compute-
Bound

*Peak
FLOPS*

Peak Cache BW

Peak Mem BW

**Key Question:** How to improve locality, to reduce data movement, for peak performance?

**Locality wins:** If we can operate out of cache, higher ceiling & more leftward ridge point.

Why is cache BW > DRAM BW?
**Smaller SRAM caches much faster.**

*1/16*  *1/8*  *250*

Operational Intensity
*(FLOPS/Byte)*

# Propagation Blocking: Optimizing Sparse Irregular Writes to Improve Cache Locality

# Propagation Blocking: Reorganize Input to Make Memory Being Randomly Written Fit in Cache

**Bad for the cache:** the size of the *domain* of vertex data array entries is |V|, but the cache holds only |C| << |V| entries

|Domain| = |V| = 5 vertices

| 0 | 0 | 0 | 0 | 0 |

|Cache| = 2 vertices

| 0 | 1 |
|---|---|
| 2 | 0 |
| 1 | 0 |
| 0 | 2 |
| 2 | 3 |
| 0 | 4 |
| 0 | 3 |

COO
(EdgeList)

**Recall:** irregular accesses into vertex data array based on e.dst *which are essentially random*

# Propagation Blocking: Reorganize Input to Make Memory Being Randomly Written Fit in Cache

**Bad for the cache:** the size of the *domain* of vertex data array entries is |V|, but the cache holds only |C| << |V| entries

|Domain| = |V| = 5 vertices

| 0 | 1 |
|---|---|
| 2 | 0 |
| 1 | 0 |
| 0 | 2 |
| 2 | 3 |
| 0 | 4 |
| 0 | 3 |

COO
(EdgeList)

0 0 0 0 0

|Cache| = 2 vertices

**Recall:** irregular accesses into vertex data array based on e.dst *which are essentially random*

**Key idea in propagation blocking:** Limit the domain of updates to a *sub-space* of vertices, **V\***, so that |V\*| <= |C| and do multiple sub-spaces of V\*s, so that all V\*s together = V

# Propagation Blocking: Reorganize Input to Make Memory Being Randomly Written Fit in Cache

Create "Bins" that hold input elements (edges from the edge list)

COO
(EdgeList)

| 0 | 1 |
| 2 | 0 |
| 1 | 0 |
| 0 | 2 |
| 2 | 3 |
| 0 | 4 |
| 0 | 3 |

**Bin 0:**
**dst 0-1**

| 0 | 1 |
| 2 | 0 |
| 1 | 0 |

**Bin 1:**
**dst 2-3**

| 0 | 2 |
| 2 | 3 |
| 0 | 3 |

**Bin 2:**
**dst 4-5**

| 0 | 4 |

| 0 | 0 | 0 | 0 | 0 |

**dstData**
**Remember: dstData[e.dst] ++**
**and e.dst is random, from edge list**

# Propagation Blocking: Reorganize Input to Make Memory Being Randomly Written Fit in Cache

COO
(EdgeList)

Bin 0:
dst 0-1

Bin 1:
dst 2-3

Bin 2:
dst 4-5

**Execute the kernel for one bin at a time**

**dstData**
**Remember: dstData[e.dst] ++**
**and e.dst is random, from edge list**

# Propagation Blocking: Reorganize Input to Make Memory Being Randomly Written Fit in Cache



COO
(EdgeList)

Bin 0:
dst 0-1

Bin 1:
dst 2-3

Bin 2:
dst 4-5

**Execute the kernel for one bin at a time**

**dstData**
**Remember: dstData[e.dst] ++**
**and e.dst is random, from edge list**

# Propagation Blocking: Reorganize Input to Make Memory Being Randomly Written Fit in Cache



COO
(EdgeList)

Bin 0:
dst 0-1

Bin 1:
dst 2-3

Bin 2:
dst 4-5

**Execute the kernel for one bin at a time**

**dstData**
**Remember: dstData[e.dst] ++**
**and e.dst is random, from edge list**

# Propagation Blocking: Reorganize Input to Make Memory Being Randomly Written Fit in Cache



COO
(EdgeList)

Bin 0:
dst 0-1

Bin 1:
dst 2-3

Bin 2:
dst 4-5

hit

dstData

**Execute the kernel for one bin at a time**

**Remember: dstData[e.dst] ++
and e.dst is random, from edge list**

# Propagation Blocking: Reorganize Input to Make Memory Being Randomly Written Fit in Cache

| 0 | 1 |
|---|---|
| 2 | 0 |
| 1 | 0 |
| 0 | 2 |
| 2 | 3 |
| 0 | 4 |
| 0 | 3 |

COO
(EdgeList)

**Bin 0:**
**dst 0-1**

| 0 | 1 |
| 2 | 0 |
| 1 | 0 |

**Bin 1:**
**dst 2-3**

| 0 | 2 |
| 2 | 3 |
| 0 | 3 |

**Bin 2:**
**dst 4-5**

| 0 | 4 |

**Execute the kernel for one bin at a time**

| 0 | | 0 |
|---|---|---|

| 0 | 0 | 0 | 0 | 0 |

**dstData**
**Remember: dstData[e.dst] ++**
**and e.dst is random, from edge list**

# Propagation Blocking: Reorganize Input to Make Memory Being Randomly Written Fit in Cache



COO
(EdgeList)

Bin 0:
dst 0-1

Bin 1:
dst 2-3

Bin 2:
dst 4-5

**Execute the kernel for one bin at a time**

**dstData**
**Remember: dstData[e.dst] ++**
**and e.dst is random, from edge list**

# Propagation Blocking: Reorganize Input to Make Memory Being Randomly Written Fit in Cache



COO
(EdgeList)

Bin 0:
dst 0-1

Bin 1:
dst 2-3

Bin 2:
dst 4-5

Execute the kernel for one bin at a time

hit

dstData
Remember: dstData[e.dst] ++
and e.dst is random, from edge list

76

# Propagation Blocking: Reorganize Input to Make Memory Being Randomly Written Fit in Cache

COO
(EdgeList)

Bin 0:
dst 0-1

Bin 1:
dst 2-3

Bin 2:
dst 4-5

**How to decide how many vertices go in each of your Propagation Blocker's bins?**

hit

**dstData**
**Remember: dstData[e.dst] ++**
**and e.dst is random, from edge list**

# Propagation Blocking: Reorganize Input to Make Memory Being Randomly Written Fit in Cache

| | |
|---|---|
| 0 | 1 |
| 2 | 0 |
| 1 | 0 |
| 0 | 2 |
| 2 | 3 |
| 0 | 4 |
| 0 | 3 |

COO
(EdgeList)

**Bin 0:**
**dst 0-2**

| 0 1 | 0 2 |
|---|---|
| 2 0 | |
| 1 0 | |

**Bin 1:**
**dst 3-5**

| 2 3 |
|---|
| 0 3 |
| 0 4 |

**Match destinations per bin to number of vertices worth of dstData that can fit in cache at one time**

hit

| | 0 | | 0 | | 0 | |
|---|---|---|---|---|---|---|

| 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|

**dstData**
**Remember: dstData[e.dst] ++**
**and e.dst is random, from edge list**

78

# Propagation Blocking: Performance Analysis

**Traverse the edge list twice instead of once**

**Binning**

| 0 | 1 |
|---|---|
| 2 | 0 |
| 1 | 0 |
| 0 | 2 |
| 2 | 3 |
| 0 | 4 |
| 0 | 3 |

COO
(EdgeList)

**Bin Read**

| 0 1 | 0 2 |
|---|---|
| 2 0 | |
| 1 0 | |

**Bin 0:**
**dst 0-2**

| 2 3 |
|---|
| 0 3 |
| 0 4 |

**Bin 1:**
**dst 3-5**

**All locations written fit in cache! Compulsory misses on dstData[] only: all the rest hit.**

hit

| | 0 | | 0 | | 0 | |
|---|---|---|---|---|---|---|

| 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|

**dstData**

**Remember: dstData[e.dst] ++**
**and e.dst is random, from edge list**

# Propagation Blocking: Performance Analysis

**Traverse the edge list twice instead of once**

**Binning**

| 0 | 1 |
| 2 | 0 |
| 1 | 0 |
| 0 | 2 |
| 2 | 3 |
| 0 | 4 |
| 0 | 3 |

COO
(EdgeList)

**Bin Read**

| 0 1 | 0 2 |
| 2 0 | |
| 1 0 | |

Bin 0:
dst 0-2

| 2 3 |
| 0 3 |
| 0 4 |

Bin 1:
dst 3-5

**What about the performance of reading the edge list during binning?**

**All locations written fit in cache! Compulsory misses on dstData[] only: all the rest hit.**

hit

| | 0 | | 0 | | 0 | |

| 0 | 0 | 0 | 0 | 0 |

**dstData**
**Remember: dstData[e.dst] ++**
**and e.dst is random, from edge list**

# Propagation Blocking: Performance Analysis

**Traverse the edge list twice instead of once**

Usually save a little space in cache for *streaming edge list* data.  Easy to cache.



**Binning**

COO
(EdgeList)

**Bin Read**

Bin 0:
dst 0-2

Bin 1:
dst 3-5

**Streaming**

**Random Access, but always in cache**

**dstData**

**Remember: dstData[e.dst] ++**
**and e.dst is random, from edge list**

**What about propagation blocking for irregular reads?**

81

# Propagation Blocking

```
PropagationBlocking_EdgeCount(EdgeList E){

  Bins B[];
  for edge in E{
    add_to_bin( find_bin(edge) )
  }


  for bin in B{
    for e in bin{
      dstData[e.dst]++
    }
  }

}
```

**Reducing Pagerank Communication via Propagation Blocking**

Scott Beamer*
Computational Research Division
Lawrence Berkeley National Laboratory
Berkeley, California
sbeamer@lbl.gov

Krste Asanović    David Patterson
Electrical Engineering & Computer Sciences Department
University of California
Berkeley, California
{krste,pattrsn}@eecs.berkeley.edu

Application of Propagation Blocking for Graph Applications (Page Rank only, at first) discovered in 2017
(Prior work on "radix partitioning" applied the idea to other domains, but not graphs)

# Cache Locality determines Overall Performance
# What about better replacement policies?

# Existing Replacement Policies Are Insufficient

# Existing Replacement Policies Are Insufficient

**App - PageRank**

Lower is Better

LLC MPKI

State-of-the-Art Policies

Cache Replacement Policies

LRU   DRRIP   SHIP-PC   SHIP-MEM-INF   HAWKEYE

# Existing Replacement Policies Are Insufficient



App - PageRank

Lower is Better

Marginal Benefit over LRU

State-of-the-Art Policies

LLC MPKI

Cache Replacement Policies

LRU    DRRIP    SHIP-PC    SHIP-MEM-INF    HAWKEYE

86

# Existing Replacement Policies Are Insufficient



App - PageRank

Marginal Benefit over LRU

**Problem:** **Heuristics used by SOTA policies fail to capture the complex reuse patterns of graph data**

Cache Replacement Policies

LRU    DRRIP    SHIP-PC    SHIP-MEM-INF    HAWKEYE

# Is It Possible To Do Better Cache Replacement?

Element To Be
Inserted $\rightarrow$ $E_{new}$

**Belady's MIN Replacement Policy**

Elements
in Cache
$E_1$
$E_2$
$\vdots$
$E_k$

# Is It Possible To Do Better Cache Replacement?



Element To Be Inserted → $E_{new}$

Belady's MIN Replacement Policy

Elements in Cache: $E_1$, $E_2$, ..., $E_k$

Time

Next References

# Is It Possible To Do Better Cache Replacement?

Element To Be Inserted → $E_{new}$

Belady's MIN Replacement Policy

Elements in Cache

$E_1$

$E_2$

$\vdots$

$E_k$

Evict the element accessed furthest in the **future**

Time

**Next References**

# Is It Possible To Do Better Cache Replacement?

Element To Be Inserted → **E**<sub>new</sub>

**Belady's MIN Replacement Policy**

Elements in Cache → E₁, E₂, ⋮, Eₖ

Evict the element accessed furthest in the **future**

Recall: Prescience required. Problem?

**Next References**

Time

# Is It Possible To Do Better Cache Replacement? YES!



Element To Be Inserted → $E_{new}$

Belady's MIN Replacement Policy

Elements in Cache → $E_1$, $E_2$, ..., $E_k$

Evict the element accessed furthest in the **future**

Prescience is a flawed requirement: any hope?

Time

**Key Observation:** The Graph's Transpose Efficiently Encodes Future Accesses

# Key Graph Application Property That Enables Belady's OPT

# Key Graph Application Property That Enables Belady's OPT



Dst →

|  | $D_0$ | $D_1$ | $D_2$ | $D_3$ | $D_4$ |
|---|---|---|---|---|---|
| $S_0$ |  |  | 1 |  |  |
| $S_1$ | 1 |  |  |  | 1 |
| $S_2$ | 1 | 1 |  | 1 |  |
| $S_3$ |  | 1 |  |  | 1 |
| $S_4$ | 1 |  | 1 |  |  |

Src ↓

Pull Traversal Pattern

**Pull Execution (*CSC Traversal*)**

```
for dst in G:
    for src in in_neighs(dst):
        dstData[dst] += srcData[src]
```

OA

| 0 | 3 | 5 | 7 | 8 |
|---|---|---|---|---|

NA

| 1 | 2 | 4 | 2 | 3 | 0 | 4 | 2 | 1 | 3 |
|---|---|---|---|---|---|---|---|---|---|

**CSC**

# Key Graph Application Property That Enables Belady's OPT

Dst $\longrightarrow$

|       | $D_0$ | $D_1$ | $D_2$ | $D_3$ | $D_4$ |
|-------|-------|-------|-------|-------|-------|
| $S_0$ |       |       | 1     |       |       |
| $S_1$ | 1     |       |       |       | 1     |
| $S_2$ | 1     | 1     |       | 1     |       |
| $S_3$ |       | 1     |       |       | 1     |
| $S_4$ | 1     |       | 1     |       |       |

Src $\downarrow$

Pull Traversal Pattern

**Pull Execution (*CSC Traversal*)**

```
for dst in G:
    for src in in_neighs(dst):
        dstData[dst] += srcData[src]
```

OA: 0 3 5 7 8

NA: 1 2 4 2 3 0 4 2 1 3

**CSC**

**CurrDs    Irregular Data Stream**

| CurrDs | Irregular Data Stream |
|--------|-----------------------|
| $t\,D_0$ | srcData[$S_1$] |
| $D_0$ | srcData[$S_2$] |
| $D_0$ | srcData[$S_4$] |
| $D_1$ | srcData[$S_2$] |
| $D_1$ | srcData[$S_3$] |

Time $\downarrow$

$\vdots$

95

# Key Graph Application Property That Enables Belady's OPT

Dst →

D_0  D_1  D_2  D_3  D_4

|       | $D_0$ | $D_1$ | $D_2$ | $D_3$ | $D_4$ |
|-------|-------|-------|-------|-------|-------|
| $S_0$ |       |       | 1     |       |       |
| $S_1$ | 1     |       |       |       | 1     |
| $S_2$ | 1     | 1     |       | 1     |       |
| $S_3$ |       | 1     |       |       | 1     |
| $S_4$ | 1     |       | 1     |       |       |

Src ↓

Pull Traversal Pattern

**Pull Execution (*CSC Traversal*)**

```
for dst in G:
    for src in in_neighs(dst):
        dstData[dst] += srcData[src]
```

OA | 0 | 3 | 5 | 7 | 8 |

NA | 1 | 2 | 4 | 2 | 3 | 0 | 4 | 2 | 1 | 3 |

**CSC**

**Key Property:** Dst-IDs are like timestamps for irregular accesses

**CurrDs  Irregular Data Stream**

| CurrDs | Irregular Data Stream |
|--------|----------------------|
| $t\,D_0$ | srcData[$S_1$] |
| $D_0$ | srcData[$S_2$] |
| $D_0$ | srcData[$S_4$] |
| $D_1$ | srcData[$S_2$] |
| $D_1$ | srcData[$S_3$] |

Time ↓

⋮

96

# Key Graph Application Property That Enables Belady's OPT



**Pull Execution (*CSC Traversal*)**

```
for dst in G:
    for src in in_neighs(dst):
        dstData[dst] += srcData[src]
```

**Key Property:** Dst-IDs are like timestamps for irregular accesses

**CurrDs**  **Irregular Data Stream**

$t_{D_0}$  srcData[$S_1$]

$D_0$  **srcData[$S_2$]**

$D_0$  srcData[$S_4$]

$D_1$  **srcData[$S_2$]**

$D_1$  srcData[$S_3$]

Time

**CSC**

# Key Graph Application Property That Enables Belady's OPT

**Pull Execution (*CSC Traversal*)**

```
for dst in G:
    for src in in_neighs(dst):
        dstData[dst] += srcData[src]
```

**Key Property:** Dst-IDs are like timestamps for irregular accesses

**CurrDs**   **Irregular Data Stream**

$t D_0$   srcData[$S_1$]

$D_0$   **srcData[$S_2$]**

$D_0$   srcData[$S_4$]

$D_1$   **srcData[$S_2$]**

$D_1$   srcData[$S_3$]

Time

Dst →

| | $D_0$ | $D_1$ | $D_2$ | $D_3$ | $D_4$ |
|---|---|---|---|---|---|
| $S_0$ | | | 1 | | |
| $S_1$ | 1 | | | | 1 |
| $S_2$ | 1 | 1 | | 1 | |
| $S_3$ | | 1 | | | 1 |
| $S_4$ | 1 | | 1 | | |

Src ↓

Pull Traversal Pattern

OA  | 0 | 3 | 5 | 7 | 8 |

NA  | 1 | 2 | 4 | 2 | 3 | 0 | 4 | 2 | 1 | 3 |

**CSC**

# Key Graph Application Property That Enables Belady's OPT

Dst →

|     | $D_0$ | $D_1$ | $D_2$ | $D_3$ | $D_4$ |
|-----|-------|-------|-------|-------|-------|
| $S_0$ |     |     | 1   |     |     |
| $S_1$ | 1   |     |     |     | 1   |
| $S_2$ | 1   | 1   |     | 1   |     |
| $S_3$ |     | 1   |     |     | 1   |
| $S_4$ | 1   |     | 1   |     |     |

Src ↓

Pull Traversal Pattern

**Pull Execution (*CSC Traversal*)**

```
for dst in G:
    for src in in_neighs(dst):
        dstData[dst] += srcData[src]
```

OA: 0 3 5 7 8

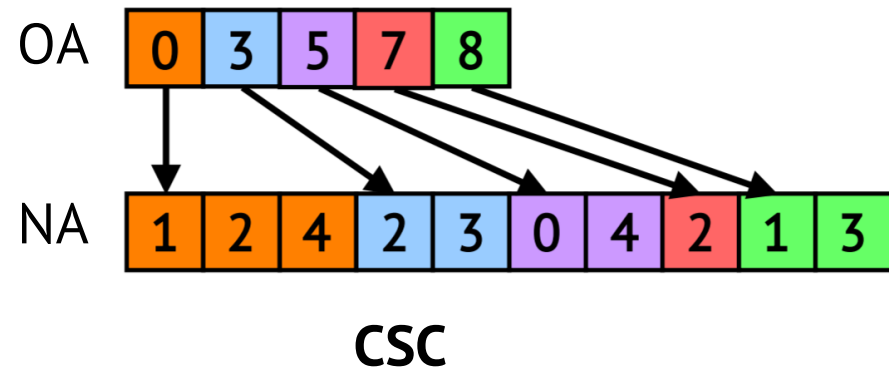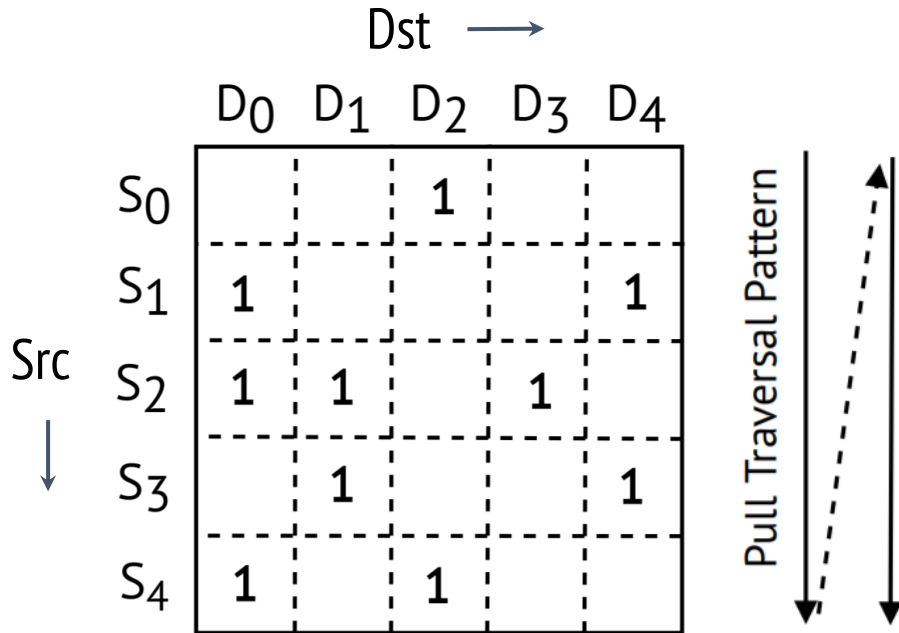NA: 1 2 4 2 3 0 4 2 1 3

**CSC**

**Key Property:** Dst-IDs are like timestamps for irregular accesses

srcData[$S_2$] is accessed at $D_0 \Rightarrow D_1 \Rightarrow D_3$

**CurrDs   Irregular Data Stream**

| | |
|---|---|
| $t D_0$ | srcData[$S_1$] |
| $D_0$ | **srcData[$S_2$]** |
| $D_0$ | srcData[$S_4$] |
| $D_1$ | **srcData[$S_2$]** |
| $D_1$ | srcData[$S_3$] |

Time ↓

⋮

# Using The Graph's Transpose For Optimal Replacement

# Using The Graph's Transpose For Optimal Replacement

Dst →

|  | $D_0$ | $D_1$ | $D_2$ | $D_3$ | $D_4$ |
|---|---|---|---|---|---|
| $S_0$ |  |  | 1 |  |  |
| $S_1$ | 1 |  |  |  | 1 |
| $S_2$ | 1 | 1 |  | 1 |  |
| $S_3$ |  | 1 |  |  | 1 |
| $S_4$ | 1 |  | 1 |  |  |

Src ↓

Pull Traversal Pattern

**Pull Execution (*CSC Traversal*)**

```
for dst in G:
    for src in in_neighs(dst):
        dstData[dst] += srcData[src]
```

**CurrDs**  **Irregular Data Stream**

| CurrDs | Irregular Data Stream |
|---|---|
| **t** $D_0$ | srcData[$S_1$] |
| $D_0$ | srcData[$S_2$] |
| $D_0$ | srcData[$S_4$] |
| $D_1$ | srcData[$S_2$] |
| $D_1$ | srcData[$S_3$] |

Time ↓

2-way Set-Associative

Assumptions:
1. One srcData elem per line
2. Only irregular data enters the cache

# Using The Graph's Transpose For Optimal Replacement

Dst →

| | $D_0$ | $D_1$ | $D_2$ | $D_3$ | $D_4$ |
|---|---|---|---|---|---|
| $S_0$ | | | 1 | | |
| $S_1$ | 1 | | | | 1 |
| $S_2$ | 1 | 1 | | 1 | |
| $S_3$ | | 1 | | | 1 |
| $S_4$ | 1 | | 1 | | |

Src ↓

Pull Traversal Pattern

**Pull Execution (*CSC Traversal*)**

```
for dst in G:
    for src in in_neighs(dst):
        dstData[dst] += srcData[src]
```

**CurrDs  Irregular Data Stream**

| CurrDs | Irregular Data Stream |
|---|---|
| $t_{D_0}$ | **srcData[$S_1$]** ← |
| $D_0$ | srcData[$S_2$] |
| $D_0$ | srcData[$S_4$] |
| $D_1$ | srcData[$S_2$] |
| $D_1$ | srcData[$S_3$] |

Time ↓

2-way Set-Associative

102

# Using The Graph's Transpose For Optimal Replacement

Dst ⟶

D$_0$  D$_1$  D$_2$  D$_3$  D$_4$

Src ↓

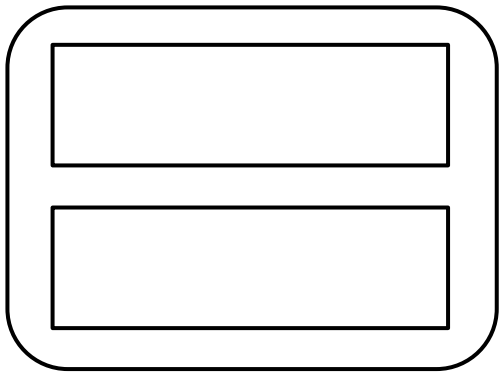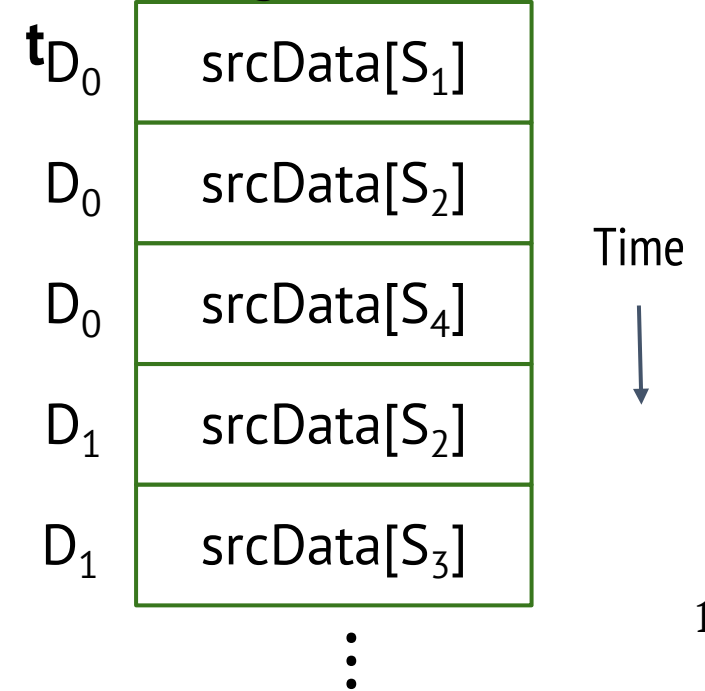|      | D$_0$ | D$_1$ | D$_2$ | D$_3$ | D$_4$ |
|------|-------|-------|-------|-------|-------|
| S$_0$|       |       | 1     |       |       |
| S$_1$| 1     |       |       |       | 1     |
| S$_2$| 1     | 1     |       | 1     |       |
| S$_3$|       | 1     |       |       | 1     |
| S$_4$| 1     |       | 1     |       |       |

Pull Traversal Pattern

**Pull Execution (*CSC Traversal*)**

```
for dst in G:
    for src in in_neighs(dst):
        dstData[dst] += srcData[src]
```

**CurrDs**   **Irregular Data Stream**

| **t**D$_0$ | srcData[S$_1$] |
| **D$_0$** | **srcData[S$_2$]** ⬅ |
| D$_0$ | srcData[S$_4$] |
| D$_1$ | srcData[S$_2$] |
| D$_1$ | srcData[S$_3$] |

Time ↓

| srcData[S$_1$] |
|                |

2-way Set-Associative

103

# Using The Graph's Transpose For Optimal Replacement

Dst $\longrightarrow$

|       | $D_0$ | $D_1$ | $D_2$ | $D_3$ | $D_4$ |
|-------|-------|-------|-------|-------|-------|
| $S_0$ |       |       | 1     |       |       |
| $S_1$ | 1     |       |       |       | 1     |
| $S_2$ | 1     | 1     |       | 1     |       |
| $S_3$ |       | 1     |       |       | 1     |
| $S_4$ | 1     |       | 1     |       |       |

Src $\downarrow$

Pull Traversal Pattern

**Pull Execution (*CSC Traversal*)**

```
for dst in G:
    for src in in_neighs(dst):
        dstData[dst] += srcData[src]
```

**CurrDs   Irregular Data Stream**
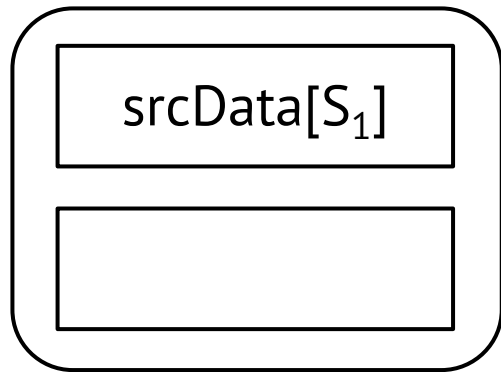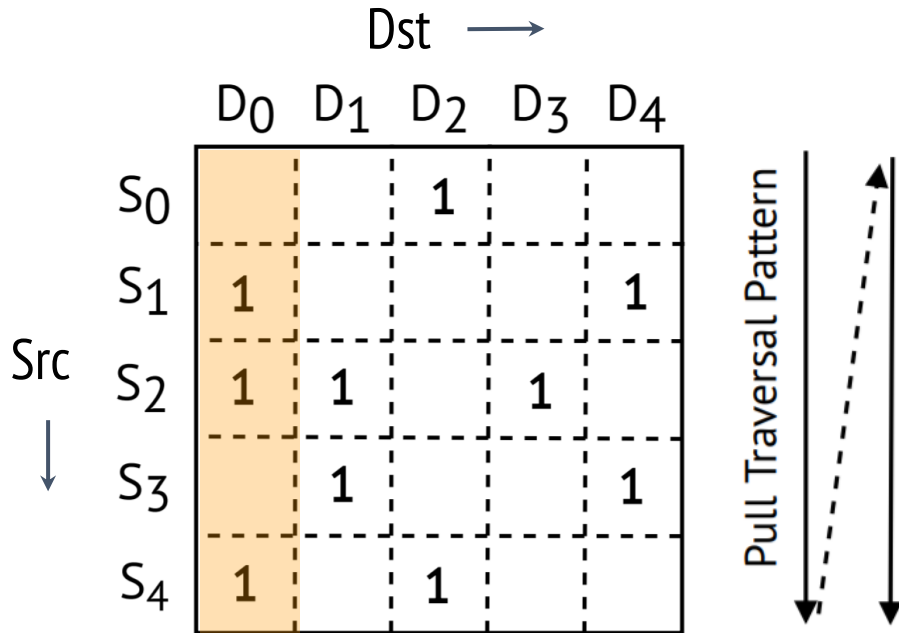
| CurrDs | Irregular Data Stream |
|--------|-----------------------|
| $\mathbf{t}D_0$ | srcData[$S_1$] |
| $D_0$ | srcData[$S_2$] |
| $\mathbf{D_0}$ | **srcData[$S_4$]** |
| $D_1$ | srcData[$S_2$] |
| $D_1$ | srcData[$S_3$] |

Time

srcData[$S_1$]

srcData[$S_2$]

2-way Set-Associative

104

# Using The Graph's Transpose For Optimal Replacement

Dst →

|     | $D_0$ | $D_1$ | $D_2$ | $D_3$ | $D_4$ |
|-----|-------|-------|-------|-------|-------|
| $S_0$ |     |       | 1     |       |       |
| $S_1$ | 1   |       |       |       | 1     |
| $S_2$ | 1   | 1     |       | 1     |       |
| $S_3$ |     | 1     |       |       | 1     |
| $S_4$ | 1   |       | 1     |       |       |

Src ↓

Pull Traversal Pattern

**Pull Execution (*CSC Traversal*)**

```
for dst in G:
        for src in in_neighs(dst):
                dstData[dst] += srcData[src]
```

**CurrDs  Irregular Data Stream**
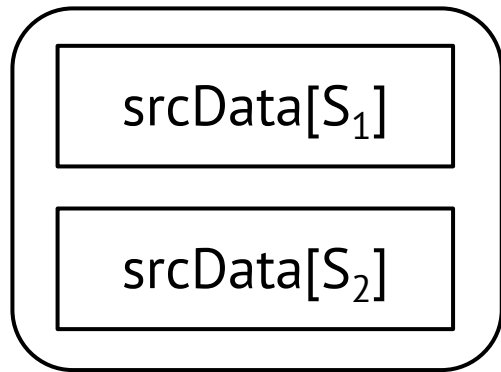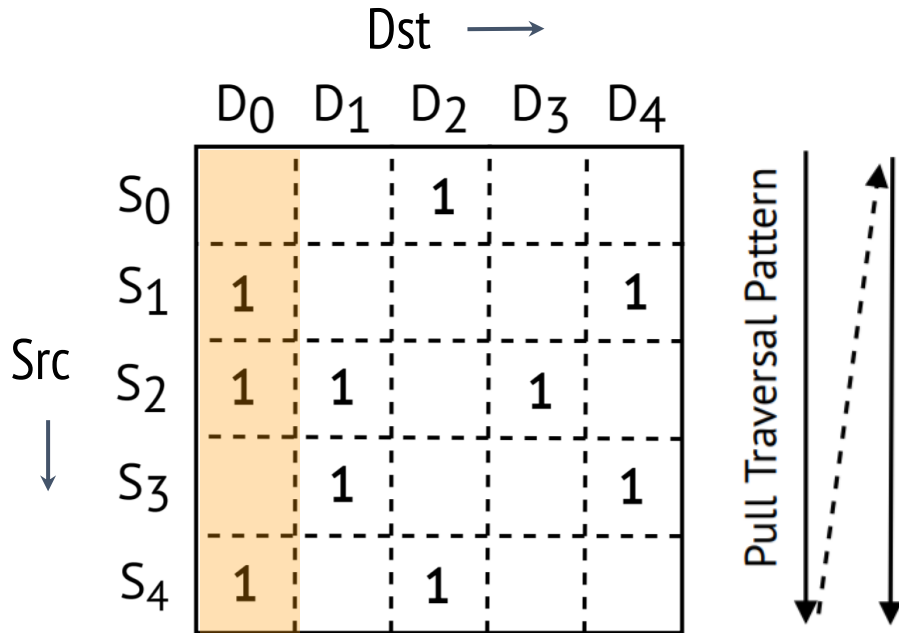
| | |
|---|---|
| **t**$D_0$ | srcData[$S_1$] |
| $D_0$ | srcData[$S_2$] |
| **$D_0$** | **srcData[$S_4$]** ← Time ↓ |
| $D_1$ | srcData[$S_2$] |
| $D_1$ | srcData[$S_3$] |

Which line should we evict?:
- srcData[$S_1$]
- srcData[$S_2$]

srcData[$S_1$]

srcData[$S_2$]

2-way Set-Associative

105

# Using The Graph's Transpose For Optimal Replacement

Dst →

D_0 D_1 D_2 D_3 D_4

S_0 | | | 1 | |
S_1 | 1 | | | | 1
S_2 | 1 | 1 | | 1 |
S_3 | | 1 | | | 1
S_4 | 1 | | 1 | |

Src ↓

Pull Traversal Pattern

**Pull Execution (*CSC Traversal*)**

```
for dst in G:
    for src in in_neighs(dst):
        dstData[dst] += srcData[src]
```

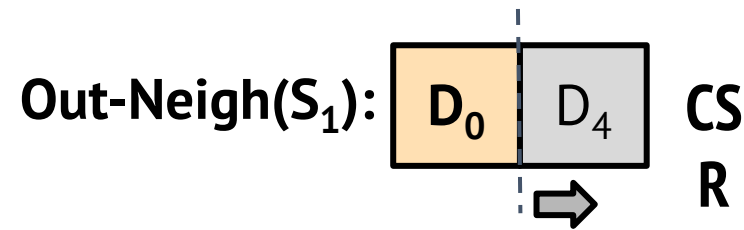**CurrDs** **Irregular Data Stream**

| | |
|---|---|
| **t** $D_0$ | srcData[$S_1$] |
| $D_0$ | srcData[$S_2$] |
| **$D_0$** | **srcData[$S_4$]** |
| $D_1$ | srcData[$S_2$] |
| $D_1$ | srcData[$S_3$] |

Time →

Which line should we evict?:
- srcData[$S_1$]  **(nextRef @ $D_4$)**
- srcData[$S_2$]

srcData[$S_1$]

srcData[$S_2$]

2-way Set-Associative

**Out-Neigh($S_1$):** | $D_0$ | $D_4$ | **CSR**

106

# Using The Graph's Transpose For Optimal Replacement

Dst →

|     | $D_0$ | $D_1$ | $D_2$ | $D_3$ | $D_4$ |
|-----|-------|-------|-------|-------|-------|
| $S_0$ |     |     | 1 |     |     |
| $S_1$ | 1 |     |     |     | 1 |
| $S_2$ | 1 | 1 |     | 1 |     |
| $S_3$ |     | 1 |     |     | 1 |
| $S_4$ | 1 |     | 1 |     |     |

Src ↓

Pull Traversal Pattern

**Pull Execution (*CSC Traversal*)**

```
for dst in G:
    for src in in_neighs(dst):
        dstData[dst] += srcData[src]
```

**CurrDs**    **Irregular Data Stream**

| CurrDs | Irregular Data Stream |
|--------|----------------------|
| **t**$D_0$ | srcData[$S_1$] |
| $D_0$ | srcData[$S_2$] |
| **$D_0$** | **srcData[$S_4$]** ← Time |
| $D_1$ | srcData[$S_2$] |
| $D_1$ | srcData[$S_3$] |

Which line should we evict?:
- srcData[$S_1$] **(nextRef @ $D_4$)**
- srcData[$S_2$] **(nextRef @ $D_1$)**

srcData[$S_1$]

srcData[$S_2$]

2-way Set-Associative

**Out-Neigh($S_2$):** | **$D_0$** | $D_1$ | $D_3$ | **CSR**

# Using The Graph's Transpose For Optimal Replacement

Dst $\longrightarrow$

|  | $D_0$ | $D_1$ | $D_2$ | $D_3$ | $D_4$ |
|---|---|---|---|---|---|
| $S_0$ |  |  | 1 |  |  |
| $S_1$ | 1 |  |  |  | 1 |
| $S_2$ | 1 | 1 |  | 1 |  |
| $S_3$ |  | 1 |  |  | 1 |
| $S_4$ | 1 |  | 1 |  |  |

Src $\downarrow$

Pull Traversal Pattern

**Pull Execution (*CSC Traversal*)**

```
for dst in G:
    for src in in_neighs(dst):
        dstData[dst] += srcData[src]
```
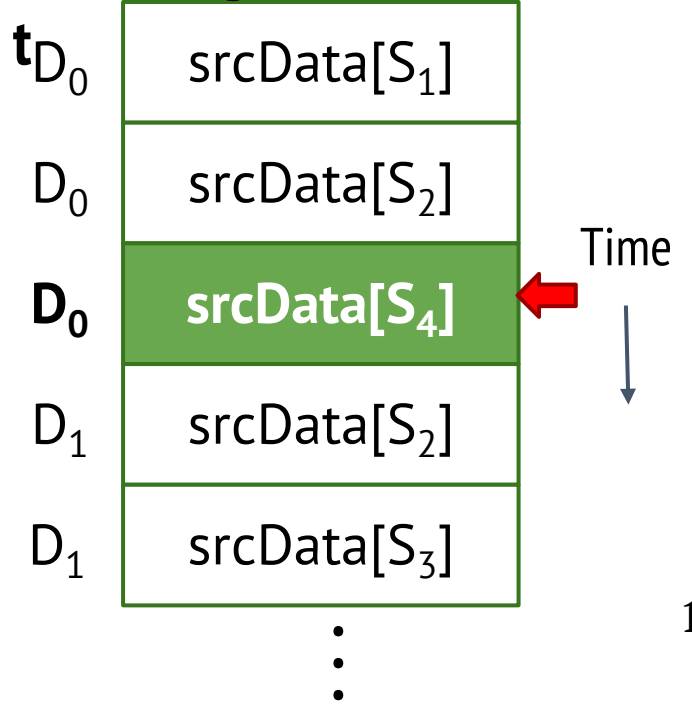
**CurrDs**  **Irregular Data Stream**

| | |
|---|---|
| **t** $D_0$ | srcData[$S_1$] |
| $D_0$ | srcData[$S_2$] |
| **$D_0$** | **srcData[$S_4$]** |
| $D_1$ | srcData[$S_2$] |
| $D_1$ | srcData[$S_3$] |

Time

Which line should we evict?:
- srcData[$S_1$]  **(nextRef @ $D_4$)** ✔
- srcData[$S_2$]  **(nextRef @ $D_1$)**

**srcData[$S_1$]**

srcData[$S_2$]

2-way Set-Associative

108

# Using The Graph's Transpose For Optimal Replacement

Dst →

| | $D_0$ | $D_1$ | $D_2$ | $D_3$ | $D_4$ |
|---|---|---|---|---|---|
| $S_0$ | | | 1 | | |
| $S_1$ | 1 | | | | 1 |
| $S_2$ | 1 | 1 | | 1 | |
| $S_3$ | | 1 | | | 1 |
| $S_4$ | 1 | | 1 | | |

Src ↓

Pull Traversal Pattern

**Pull Execution (*CSC Traversal*)**

```
for dst in G:
    for src in in_neighs(dst):
        dstData[dst] += srcData[src]
```
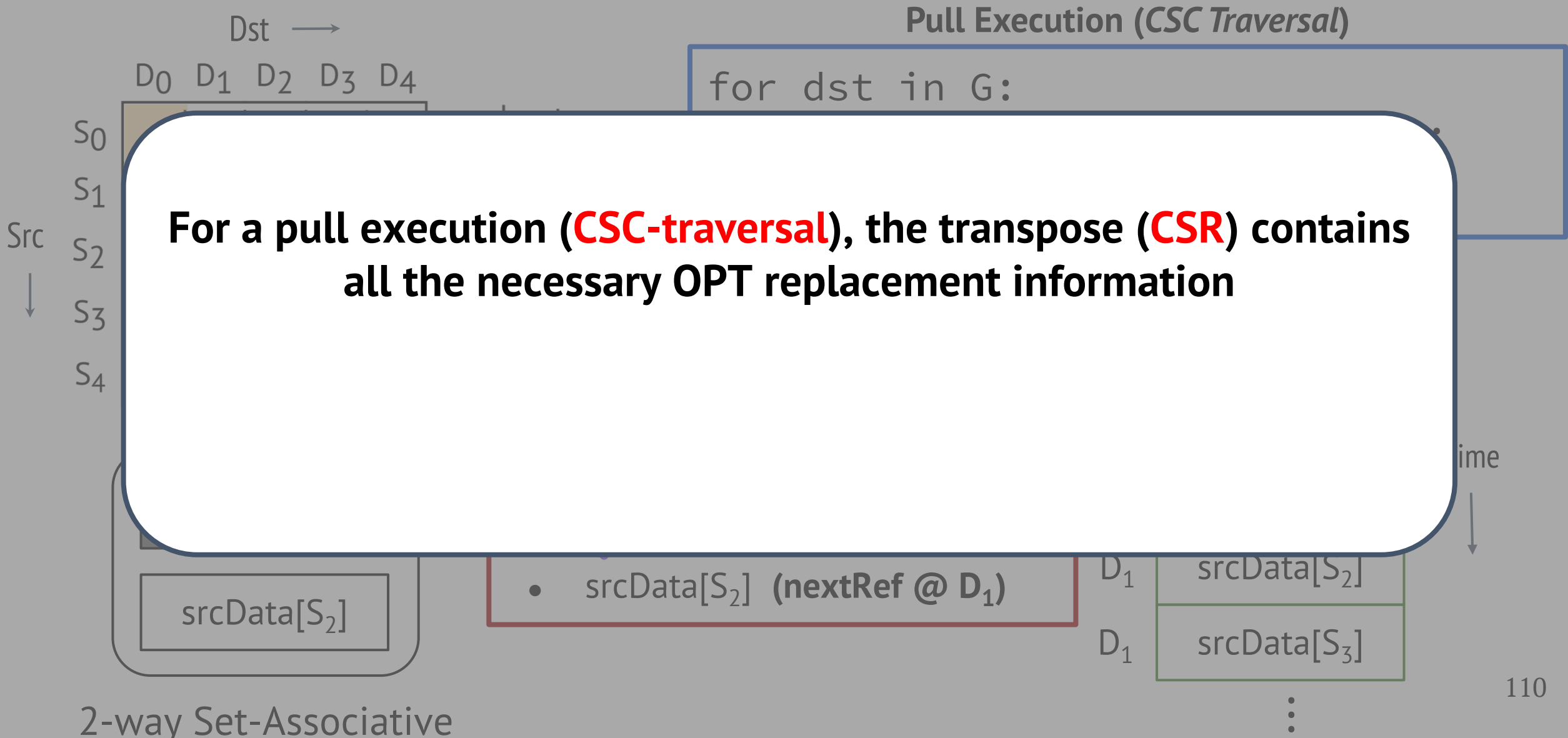
**CurrDs    Irregular Data Stream**

| CurrDs | Irregular Data Stream |
|---|---|
| $t_{D_0}$ | srcData[$S_1$] |
| $D_0$ | srcData[$S_2$] |
| **$D_0$** | **srcData[$S_4$]** ← Time |
| $D_1$ | srcData[$S_2$] |
| $D_1$ | srcData[$S_3$] |

**Which line should we evict?:**
- srcData[$S_1$]  **(nextRef @ $D_4$)** ✔
- srcData[$S_2$]  **(nextRef @ $D_1$)**

**srcData[$S_1$]**

srcData[$S_2$]

2-way Set-Associative

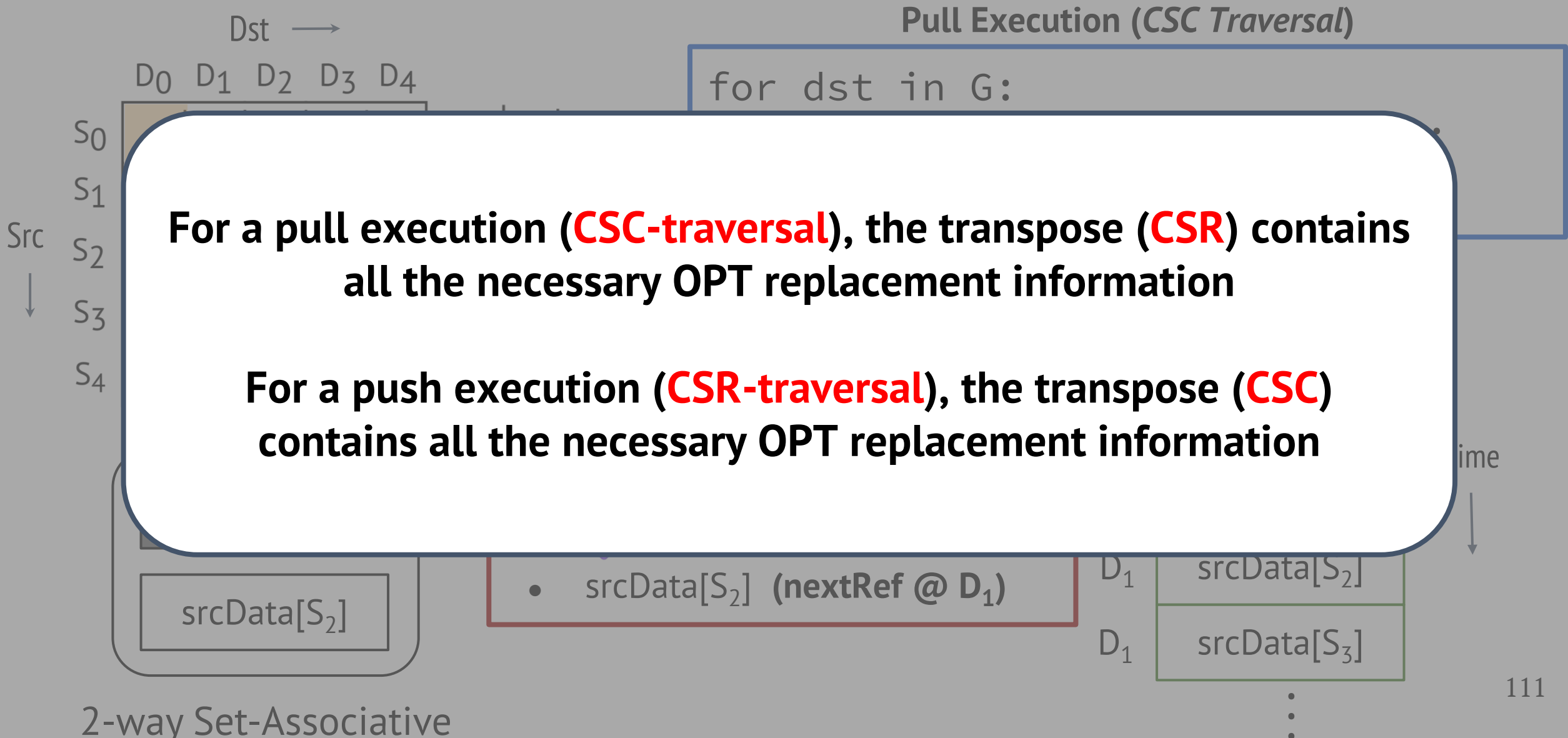**Key Question:** how to query next reference while running the program?
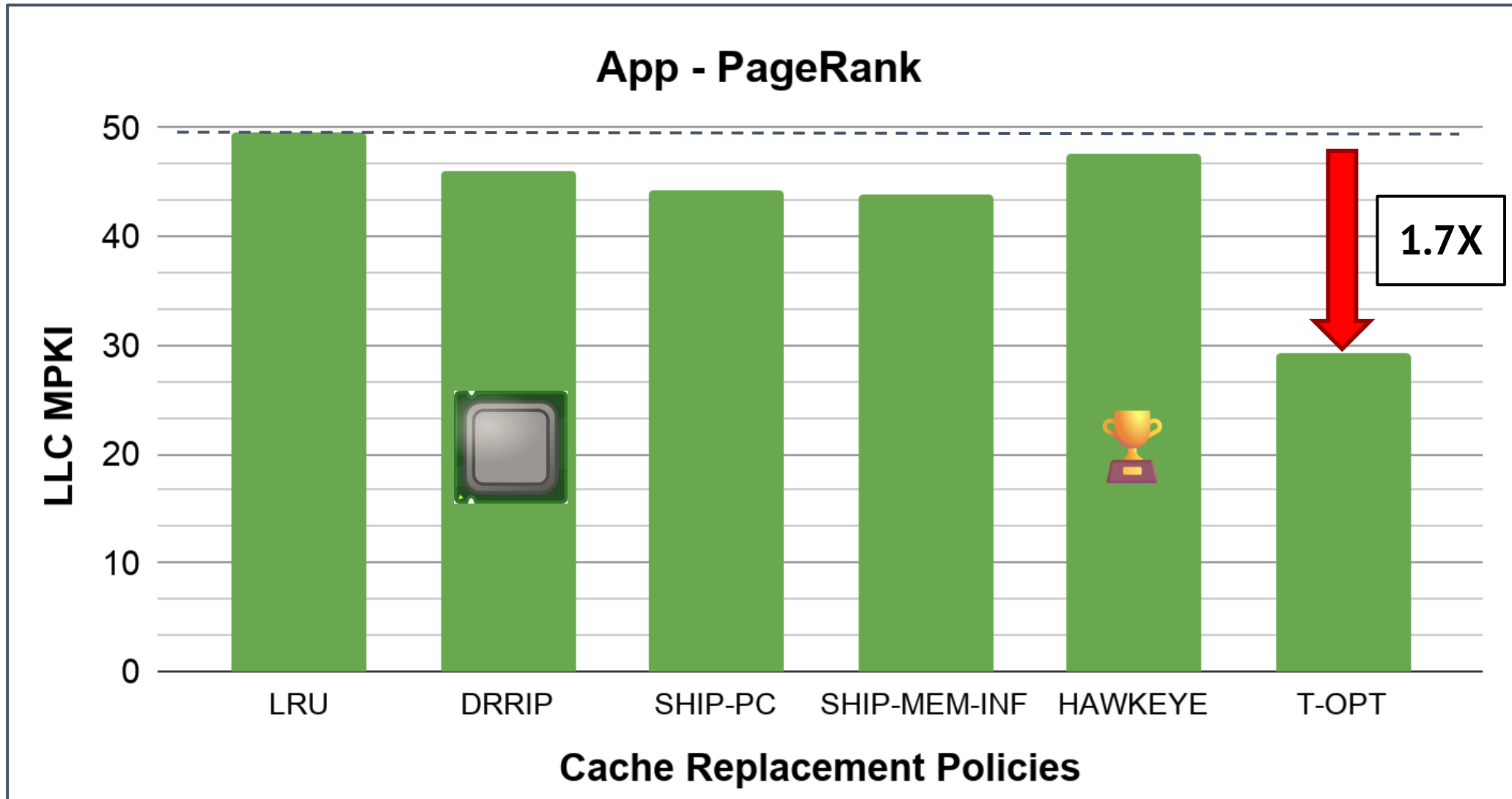
109

# Using The Graph's Transpose For Optimal Replacement

Dst →

D₀ D₁ D₂ D₃ D₄

**Pull Execution (*CSC Traversal*)**

```
for dst in G:
```

S₀

S₁

Src

S₂

S₃

S₄

**For a pull execution (CSC-traversal), the transpose (CSR) contains
all the necessary OPT replacement information**

ime

srcData[S₂] **(nextRef @ D₁)**

D₁   srcData[S₂]

srcData[S₂]

D₁   srcData[S₃]

2-way Set-Associative

# Using The Graph's Transpose For Optimal Replacement

Dst ⟶

D₀ D₁ D₂ D₃ D₄

**Pull Execution (*CSC Traversal*)**

```
for dst in G:
```

Src

S₀
S₁
S₂
S₃
S₄

**For a pull execution (CSC-traversal), the transpose (CSR) contains all the necessary OPT replacement information**

**For a push execution (CSR-traversal), the transpose (CSC) contains all the necessary OPT replacement information**

ime

srcData[S₂]

• srcData[S₂] (**nextRef @ D₁**)

D₁    srcData[S₂]

D₁    srcData[S₃]

2-way Set-Associative

# Transpose-based OPT (T-OPT) Provides Large Gains

# Transpose-based OPT (T-OPT) Provides Large Gains



App - PageRank

# Transpose-based OPT (T-OPT) Incurs Overheads



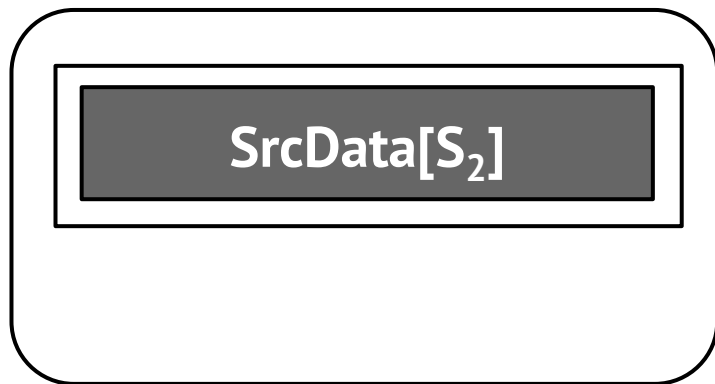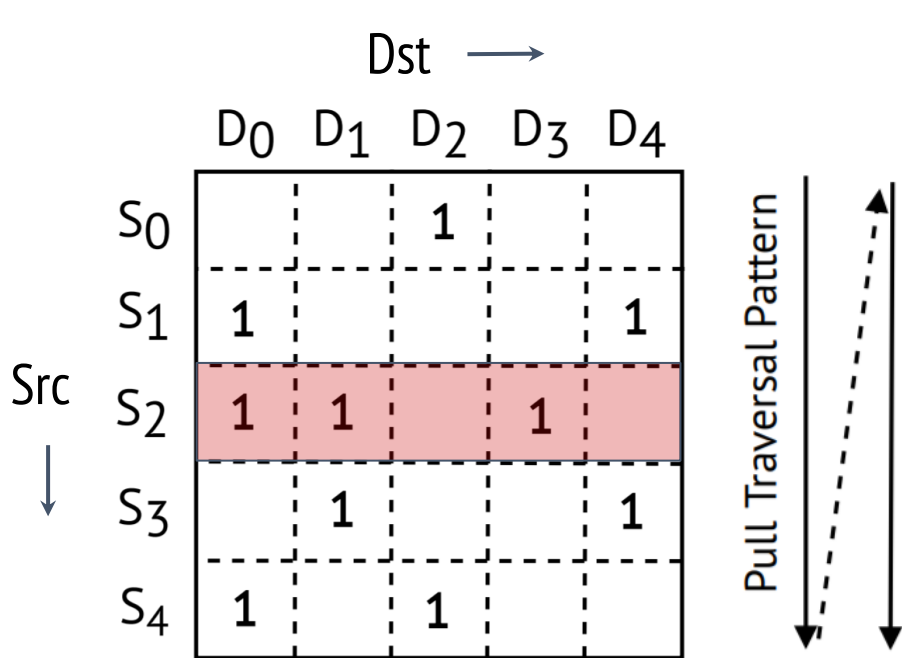**Finding Next References Using The Transpose**

# Transpose-based OPT (T-OPT) Incurs Overheads



Set-Associative Cache

CSR *(Transpose)*

# Transpose-based OPT (T-OPT) Incurs Overheads



Dst →

Src ↓

Pull Traversal Pattern

**Finding Next References Using The Transpose**

DRAM Access overhead

Runtime Traversal overhead

**SrcData[S_2]**

Set-Associative Cache

OA: 0 1 3 6 8

NA: 2 0 4 0 1 3 1 4 0 2

**CSR** *(Transpose)*

**Out-Neigh(S_2):** D_0 D_1 D_3

Need to scan neighbors

# Transpose-based OPT (T-OPT) Incurs Overheads
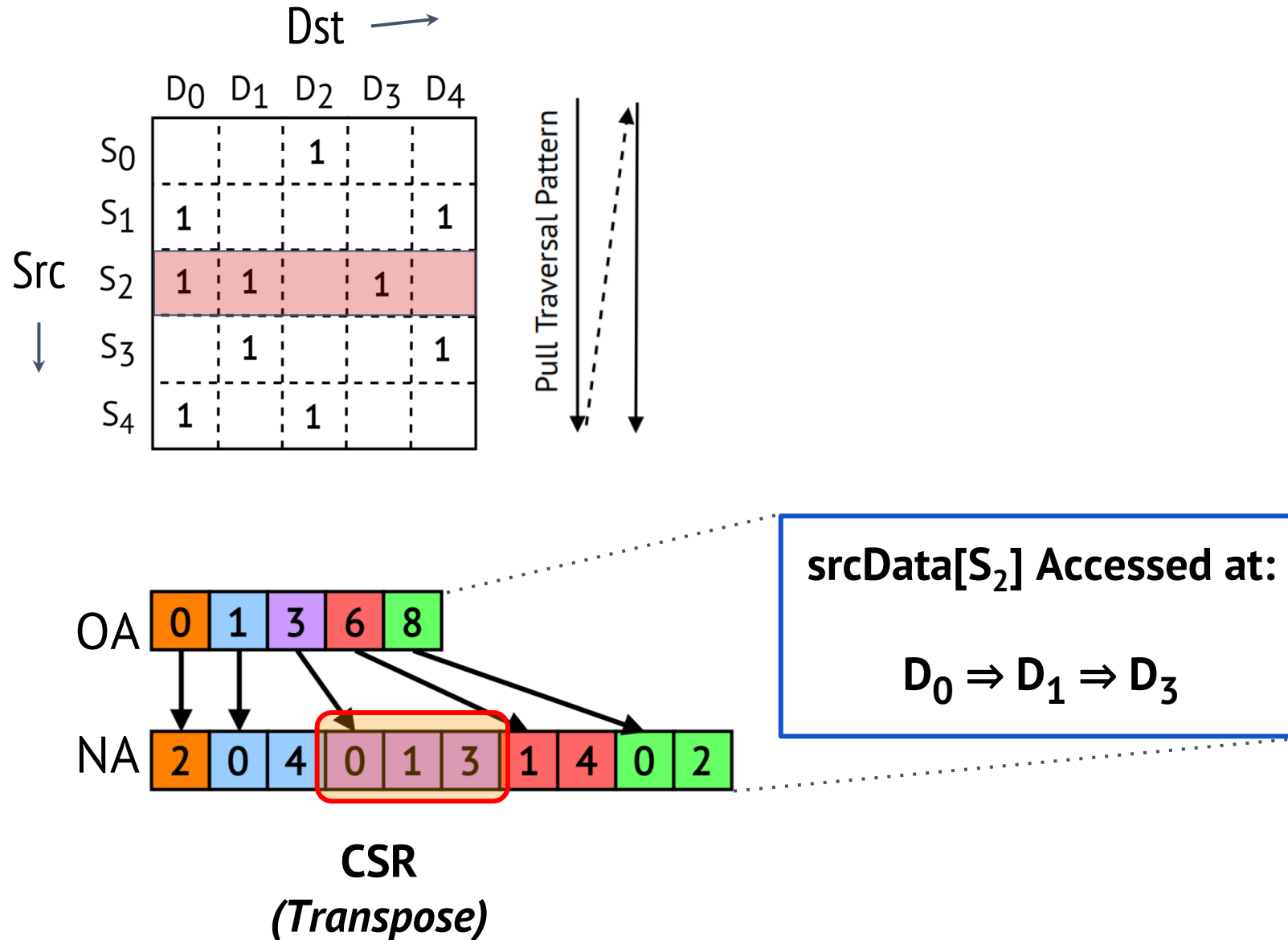


Set-Associative Cache

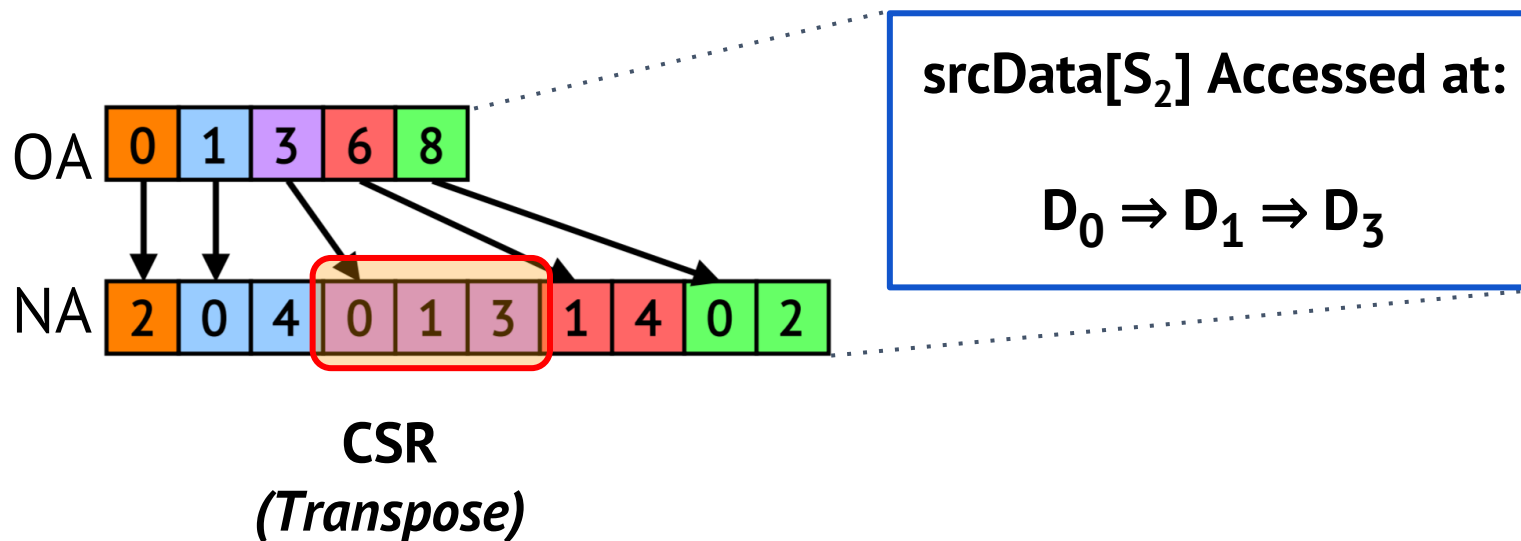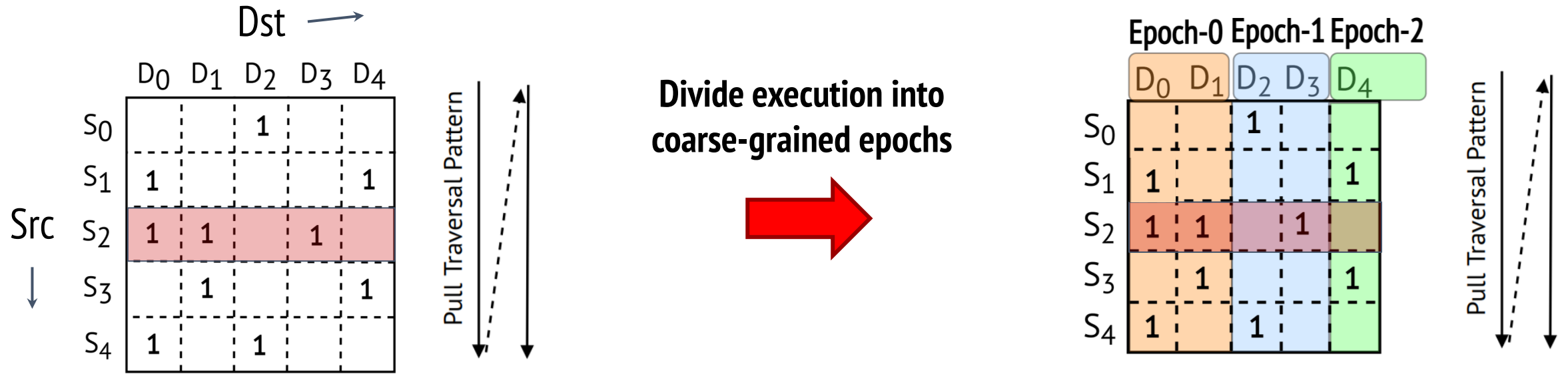**Finding Next References Using The Transpose**

DRAM Access overhead

Runtime Traversal overhead

**Question:** How do we retrieve the next reference information from the graph's transpose **without all the cost of traversing the graph?**
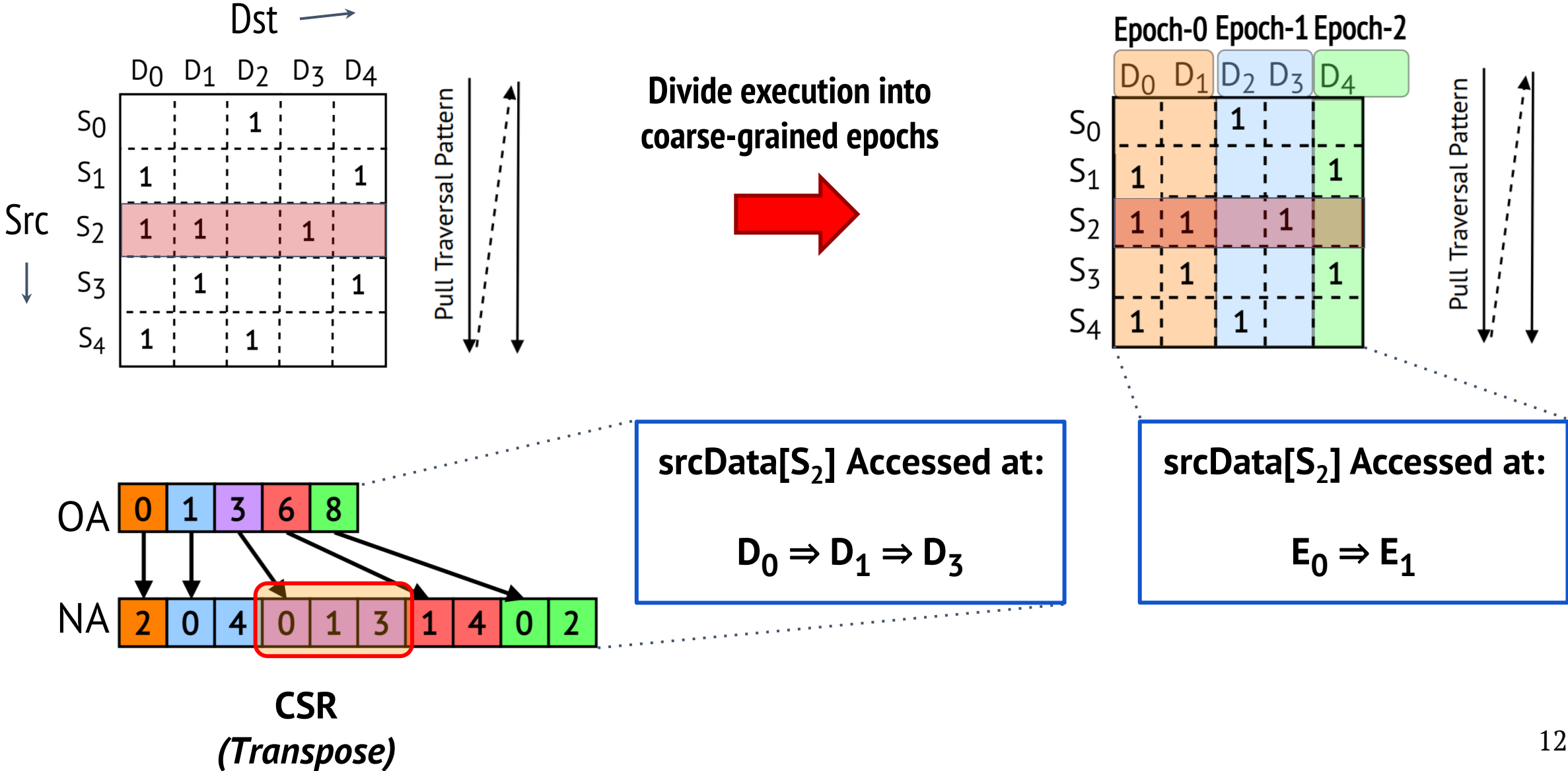
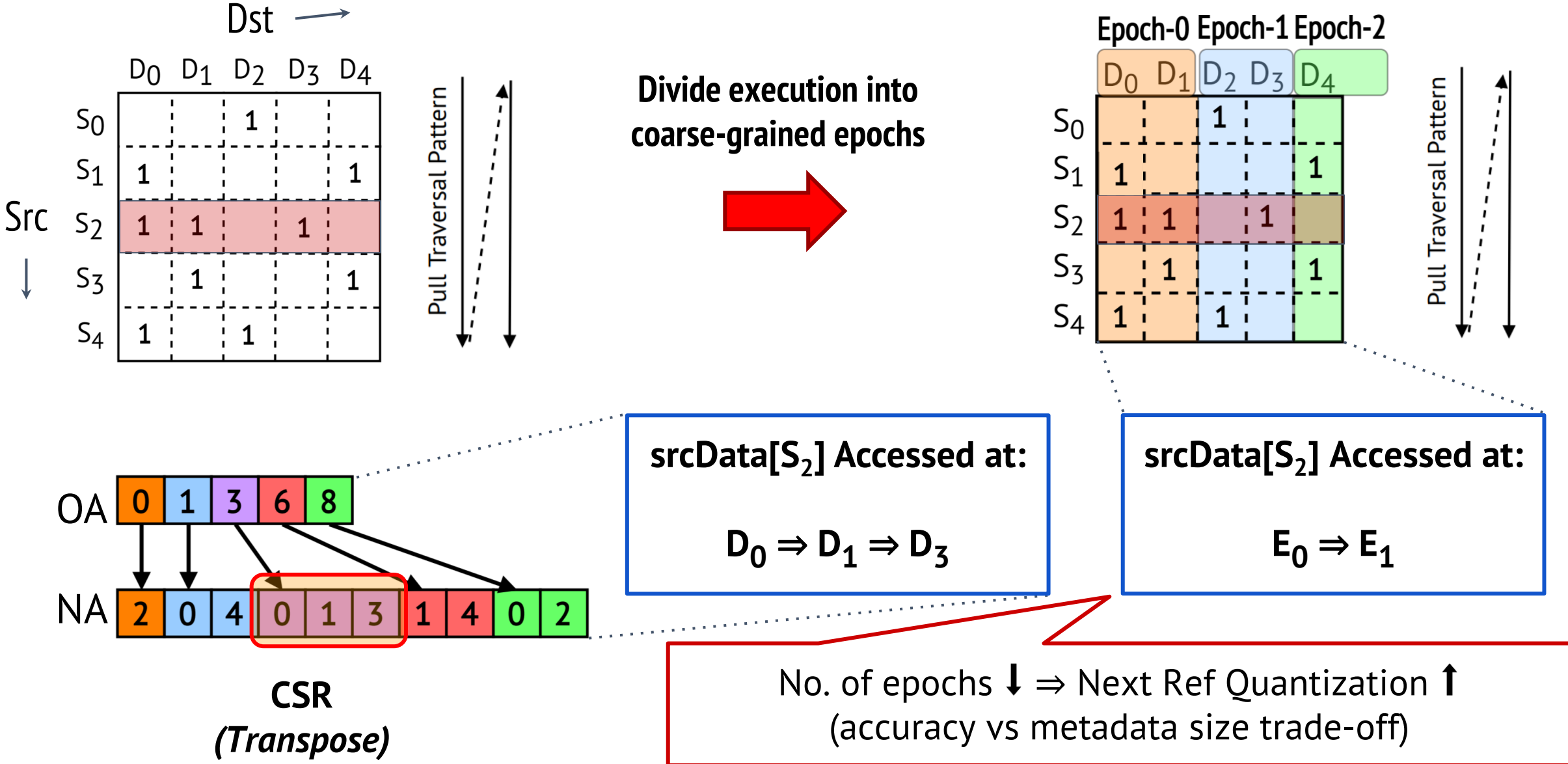# Main Technique: Use Quantization To Compress The Transpose



srcData[$S_2$] Accessed at:

$$D_0 \Rightarrow D_1 \Rightarrow D_3$$

CSR
*(Transpose)*

# Main Technique: Use Quantization To Compress The Transpose



Divide execution into coarse-grained epochs

srcData[$S_2$] Accessed at:

$$D_0 \Rightarrow D_1 \Rightarrow D_3$$

CSR
*(Transpose)*

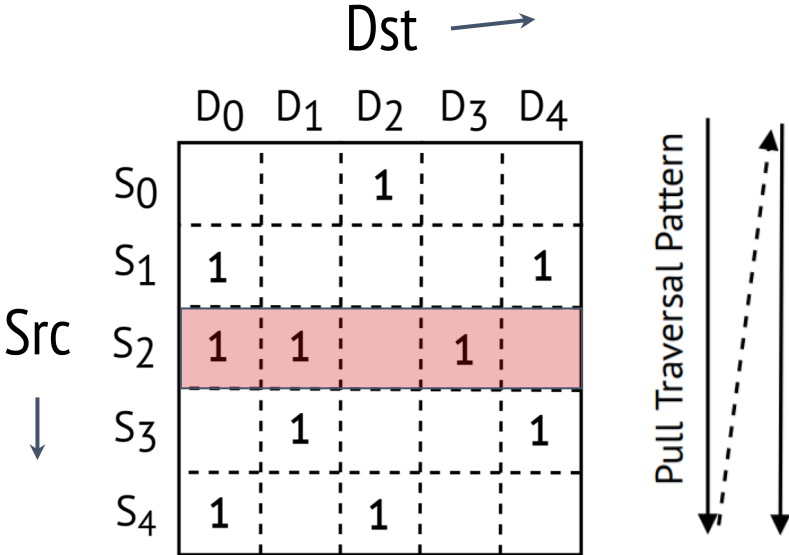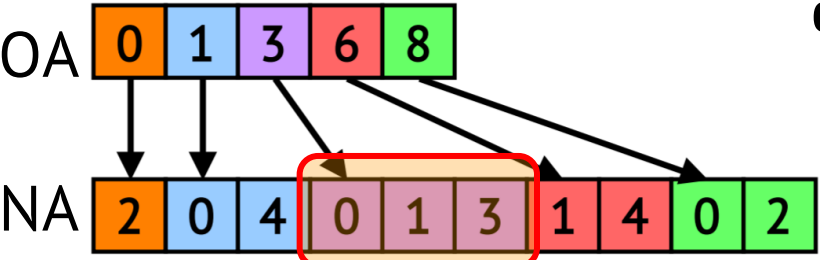# Main Technique: Use Quantization To Compress The Transpose



**Divide execution into coarse-grained epochs**

**srcData[$S_2$] Accessed at:**

$$D_0 \Rightarrow D_1 \Rightarrow D_3$$

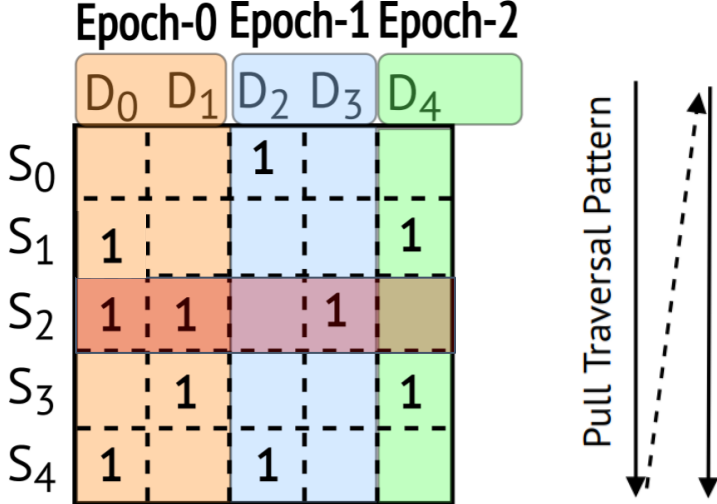**srcData[$S_2$] Accessed at:**

$$E_0 \Rightarrow E_1$$

CSR
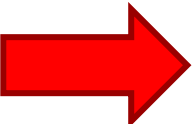*(Transpose)*

120

# Main Technique: Use Quantization To Compress The Transpose



Divide execution into coarse-grained epochs

OA → NA

**CSR**
*(Transpose)*

**srcData[$S_2$] Accessed at:**

$D_0 \Rightarrow D_1 \Rightarrow D_3$

**srcData[$S_2$] Accessed at:**

$E_0 \Rightarrow E_1$

No. of epochs ↓ ⇒ Next Ref Quantization ↑
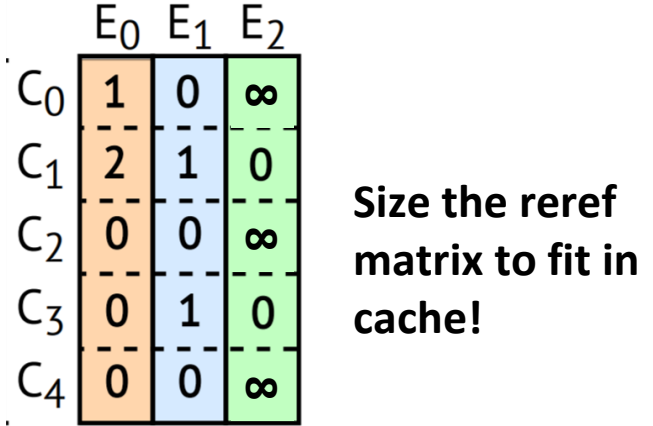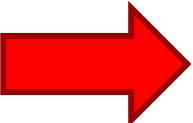(accuracy vs metadata size trade-off)

# Main Technique: Use Quantization To Compress The Transpose



Divide execution into coarse-grained epochs

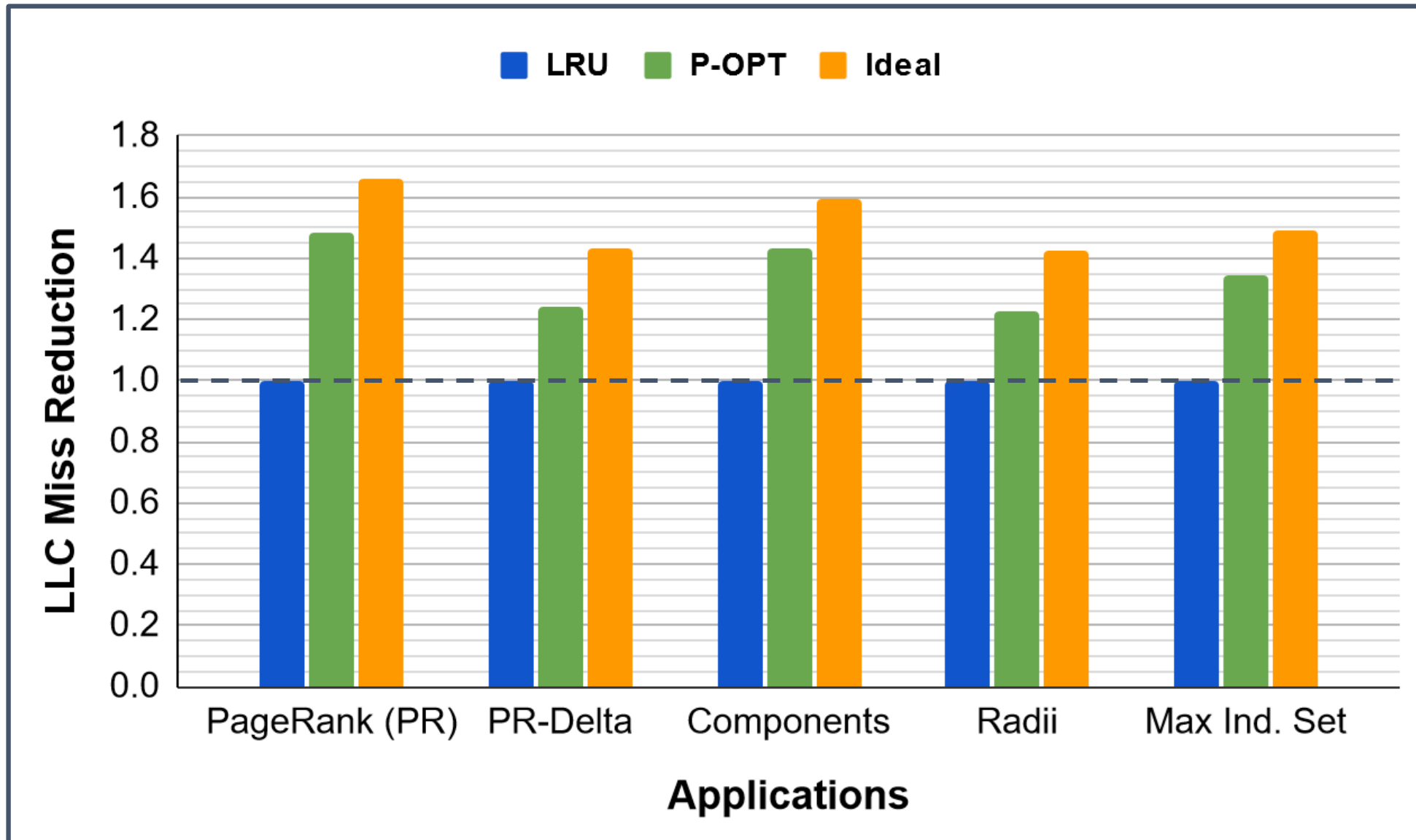Quantization enables compression of transpose data
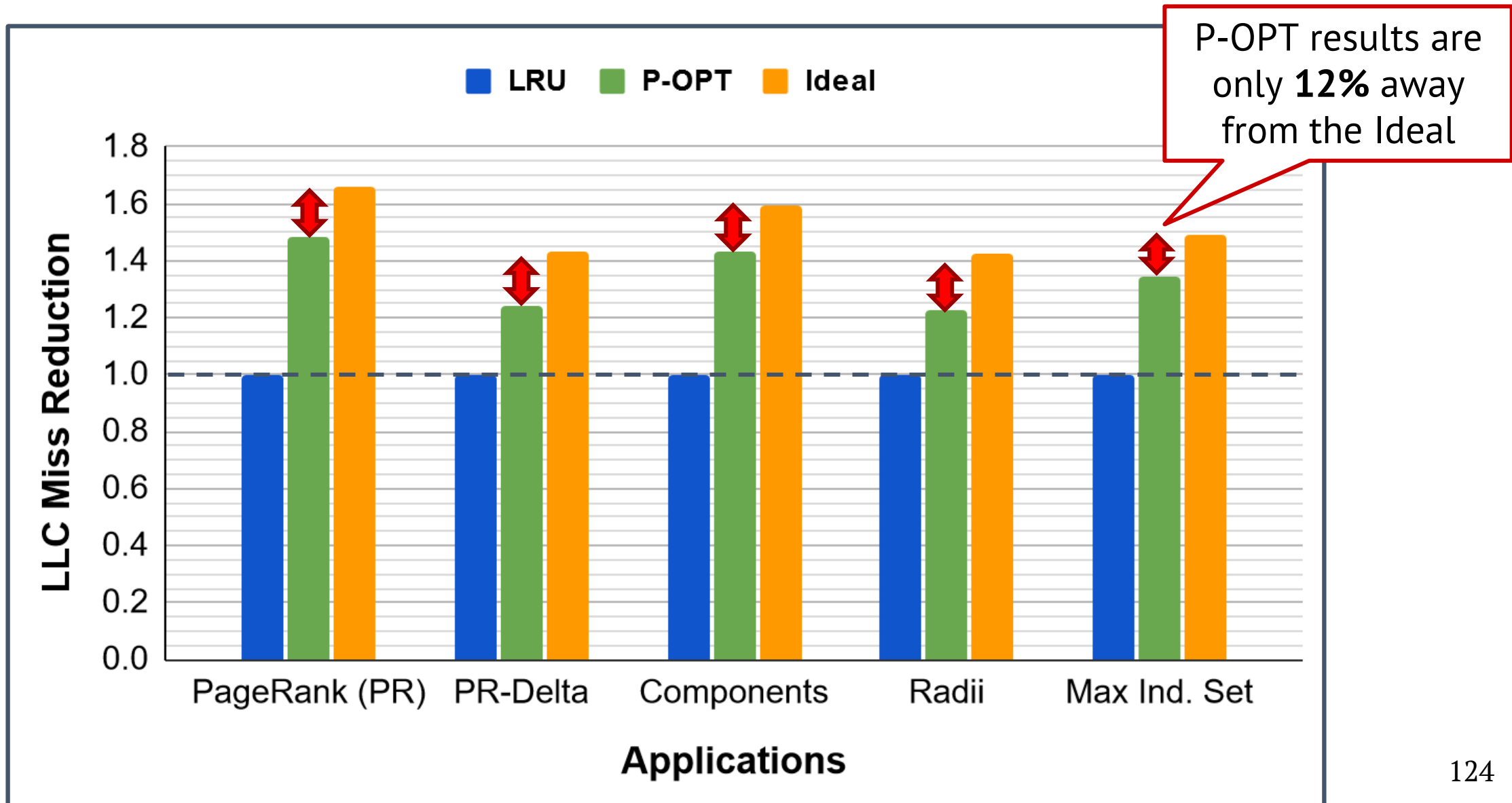
CSR
*(Transpose)*

Size the reref matrix to fit in cache!
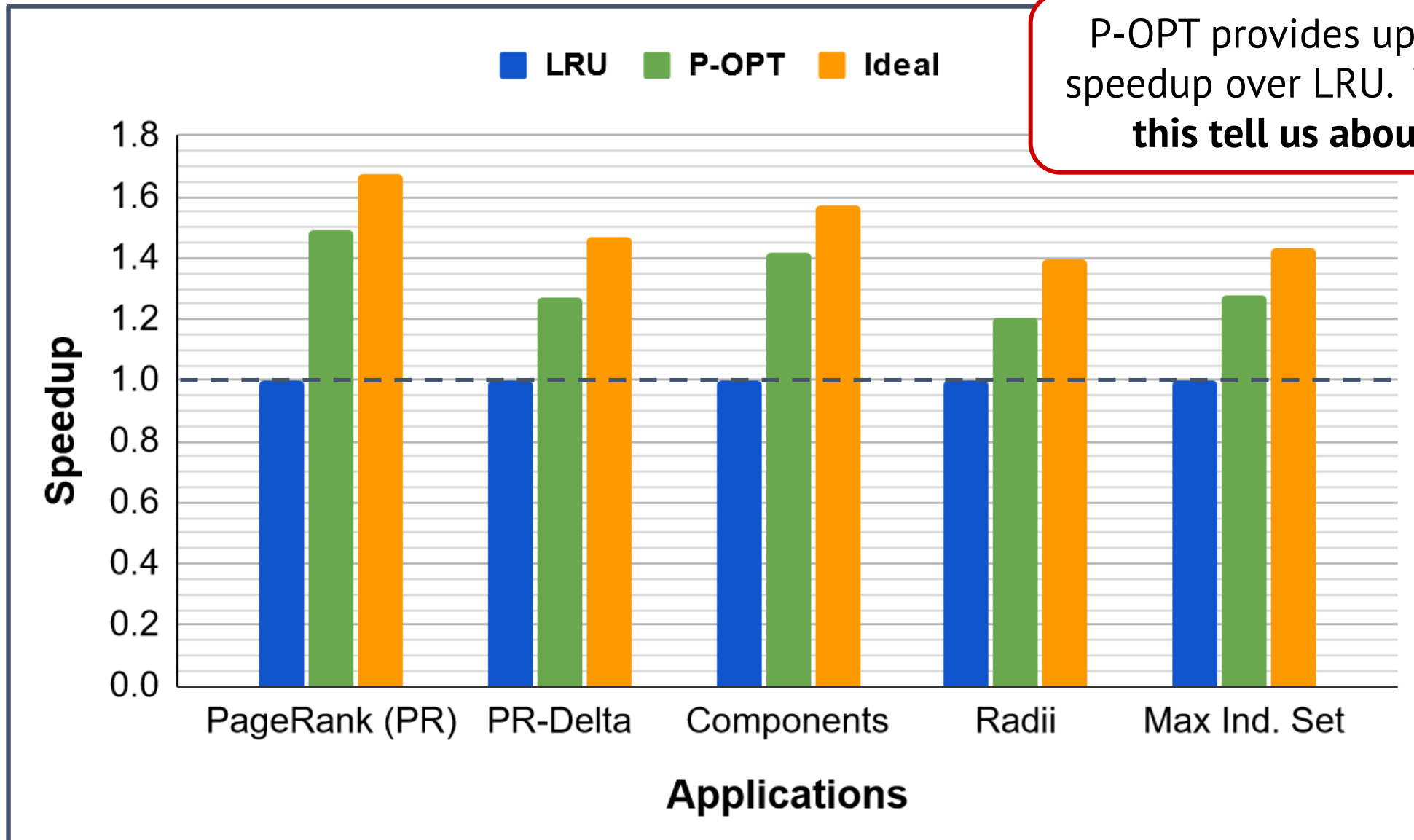
Rereference Matrix
*(Quantized Transpose)*

122

# P-OPT Improves Cache Locality

# P-OPT Improves Cache Locality



P-OPT results are only **12%** away from the Ideal

124

# P-OPT's LLC Miss Reductions Directly Translate To Speedups

P-OPT provides up to **1.56x** speedup over LRU. **What does this tell us about LRU?**

# What did we just learn?

- Sparse problems are ones that manipulate large, mostly-zero matrices

- Sparsity makes caching a useful part of the matrix hard

- Roofline model shows how close to peak perf. an app is

- Propagation blocking bins updates making irregular data fit in cache

- P-OPT is a *practical* implementation of Belady's OPT for graphs

# Takeaways

❖ Heuristic-based policies are ineffective for irregular memory access patterns

❖ The graph's transpose enables Belady's MIN replacement policy

❖ P-OPT achieves close to ideal performance (*quantization can be an effective tool in making a design practical*)