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

Fall 2023

## Course Description

**Lecture 15: Sparsity**

This course covers the design and implementation of computer systems from the perspective of the hardware software interface. The purpose of this course is for students to understand the relationship between the operating system, software, and computer architecture. Students that complete the course will have learned operating system fundamentals, computer architecture fundamentals, compilation to hardware abstractions, and how software actually executes from the perspective of the hardware software/boundary. The course will focus especially on understanding the relationships between software and hardware, and how those relationships influence the design of a computer system's software and hardware. The course will convey these topics through a series of practical, implementation-oriented lab assignments.
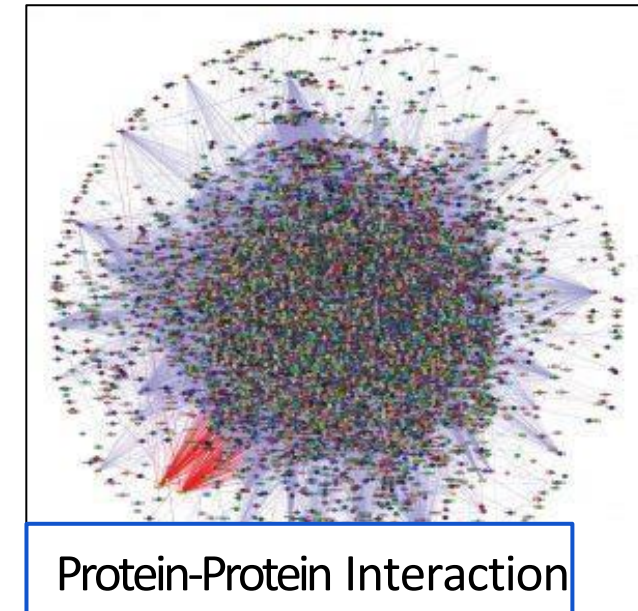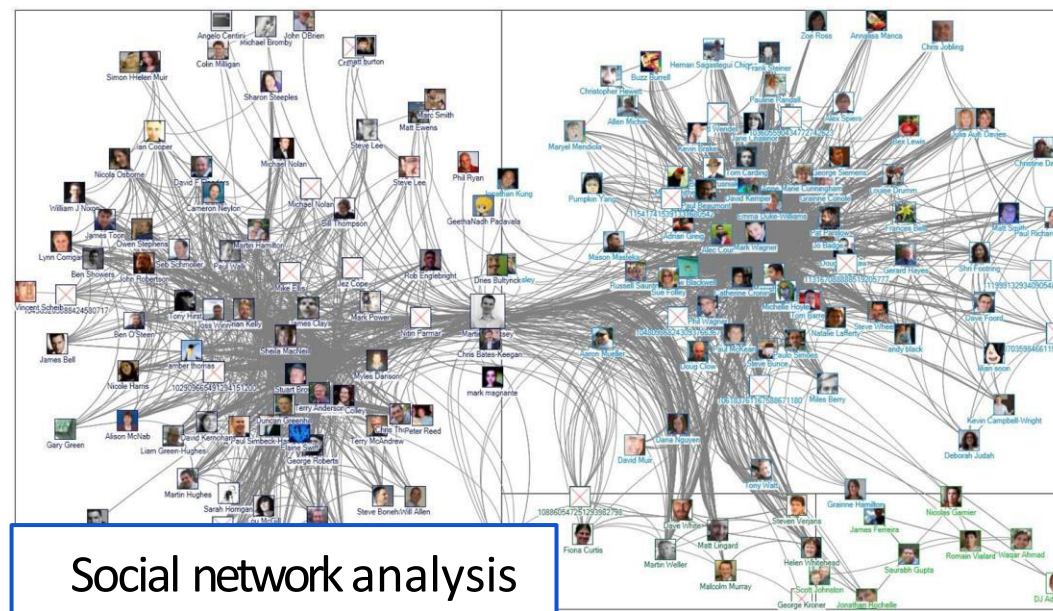
**Credit: Brandon Lucia**
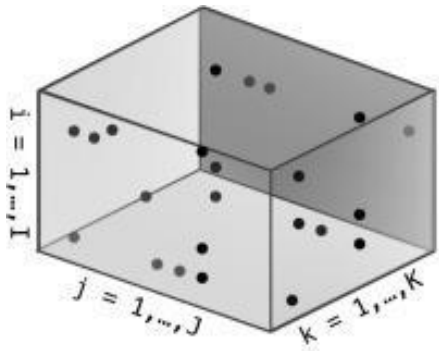
# Today: Sparse Problems

- 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

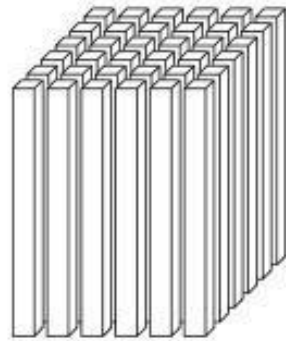# Graph Processing Problems are Sparse Problems



Path Planning



Social network analysis



Protein-Protein Interaction

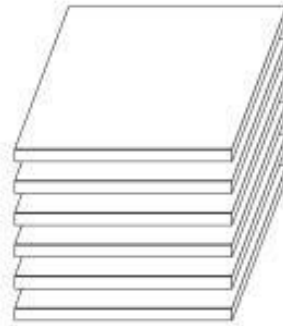The canonical examples of sparse problems are graph processing applications.

# Machine Learning Problems are Sparse Problems
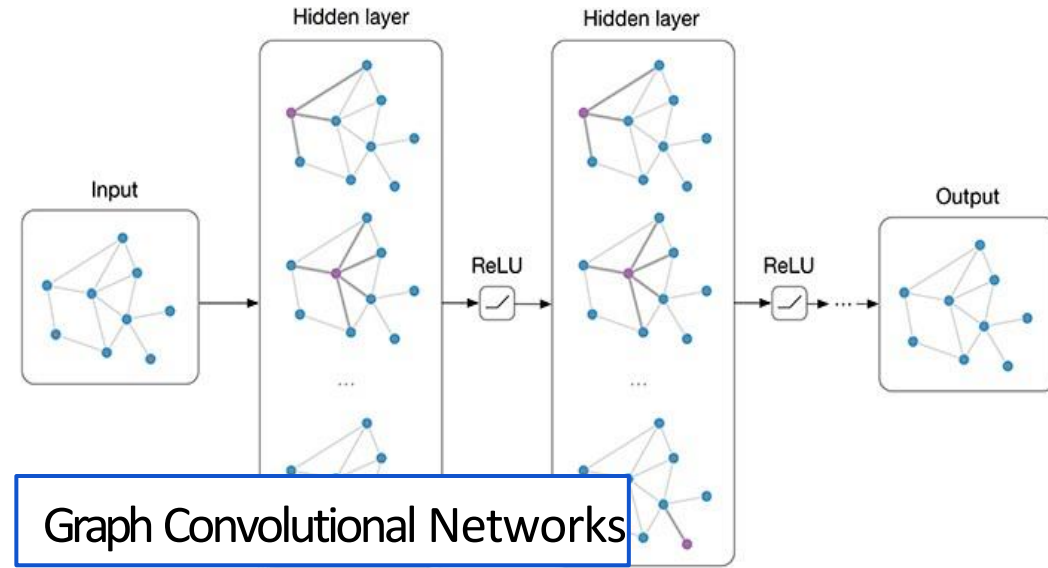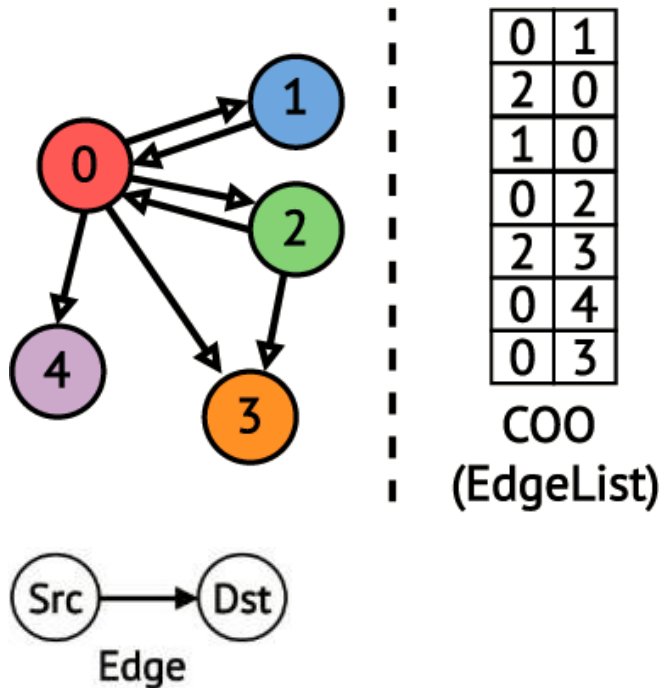


Data Mining

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

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

Graph Convolutional Networks

# What does a graph processing program look like?



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

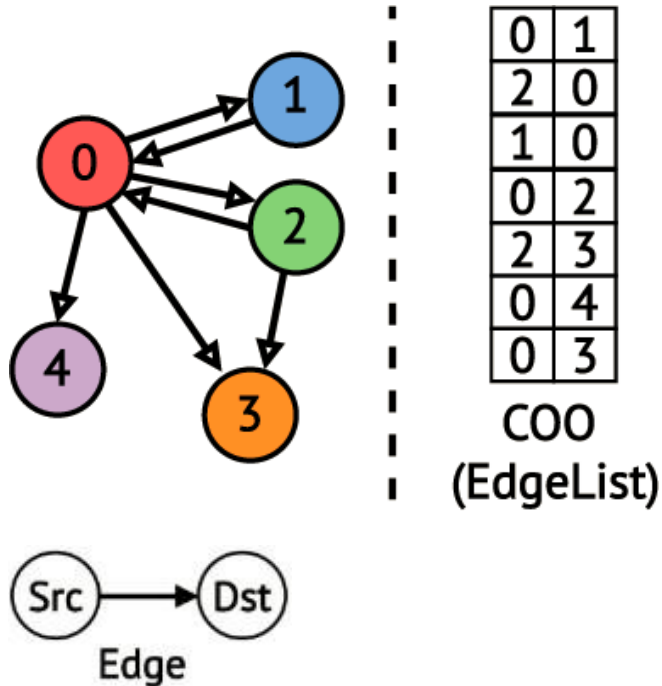COO
(EdgeList)

Src → Dst
Edge

```
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;**
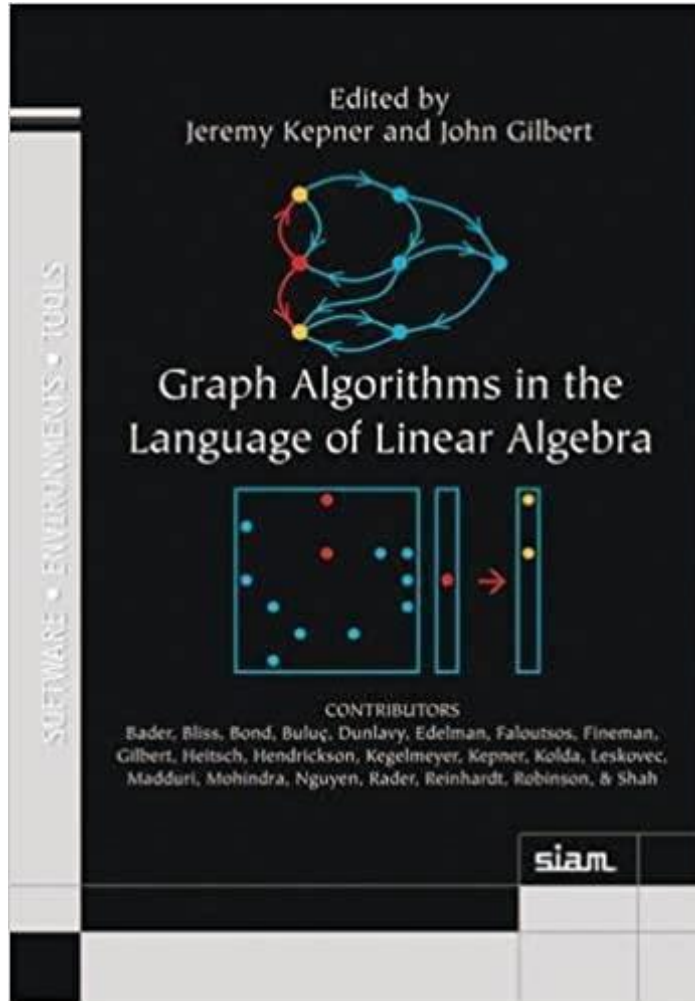**often "swap" srcData & dstData from 1 iteration to the next iteration**

# What does a graph processing program look like?

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

COO
(EdgeList)
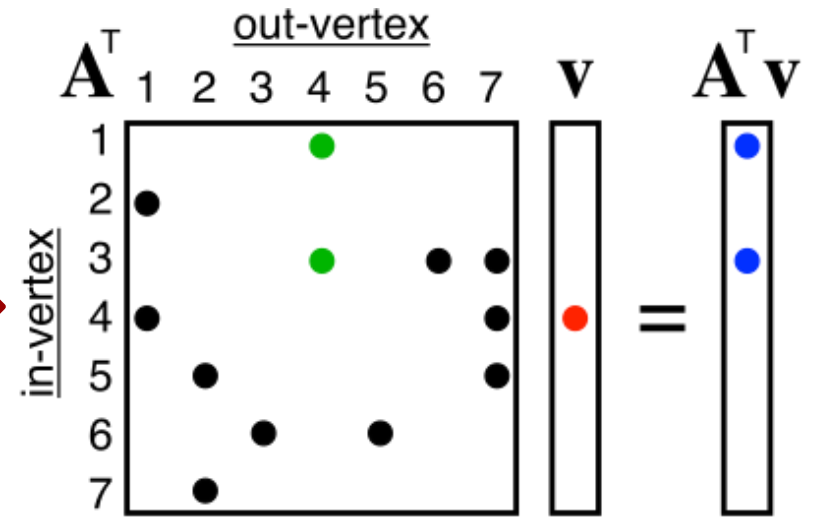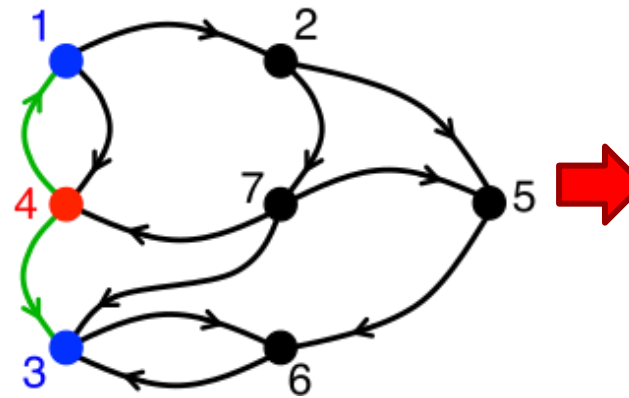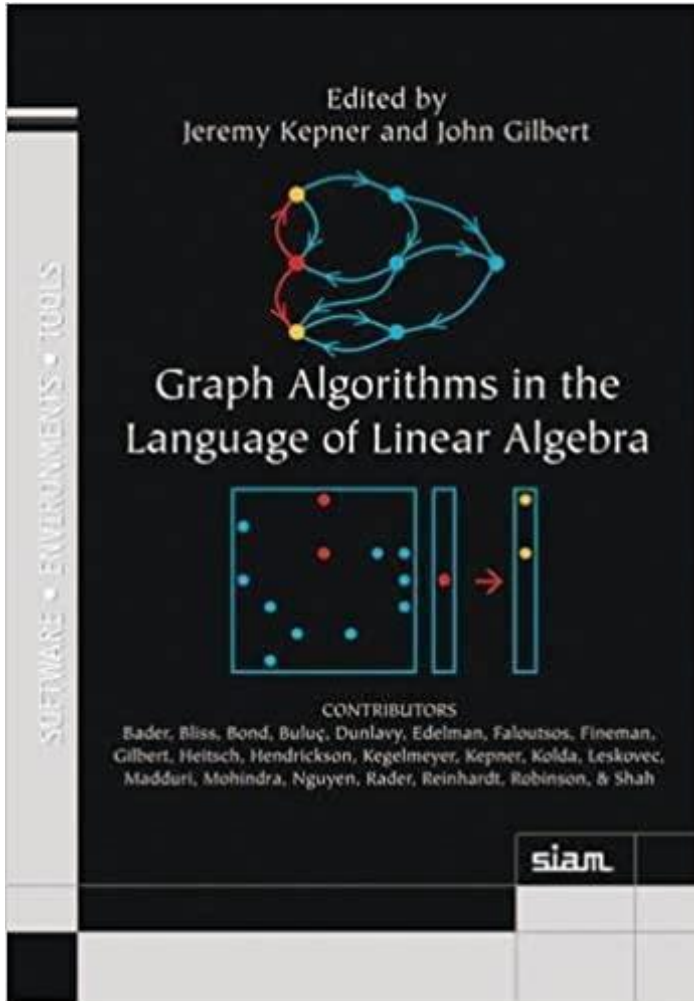
Src → Dst
Edge

dstData

srcData

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

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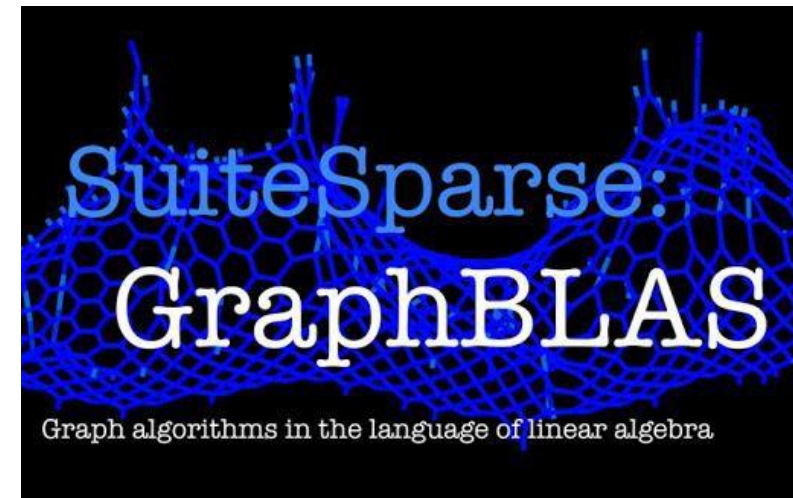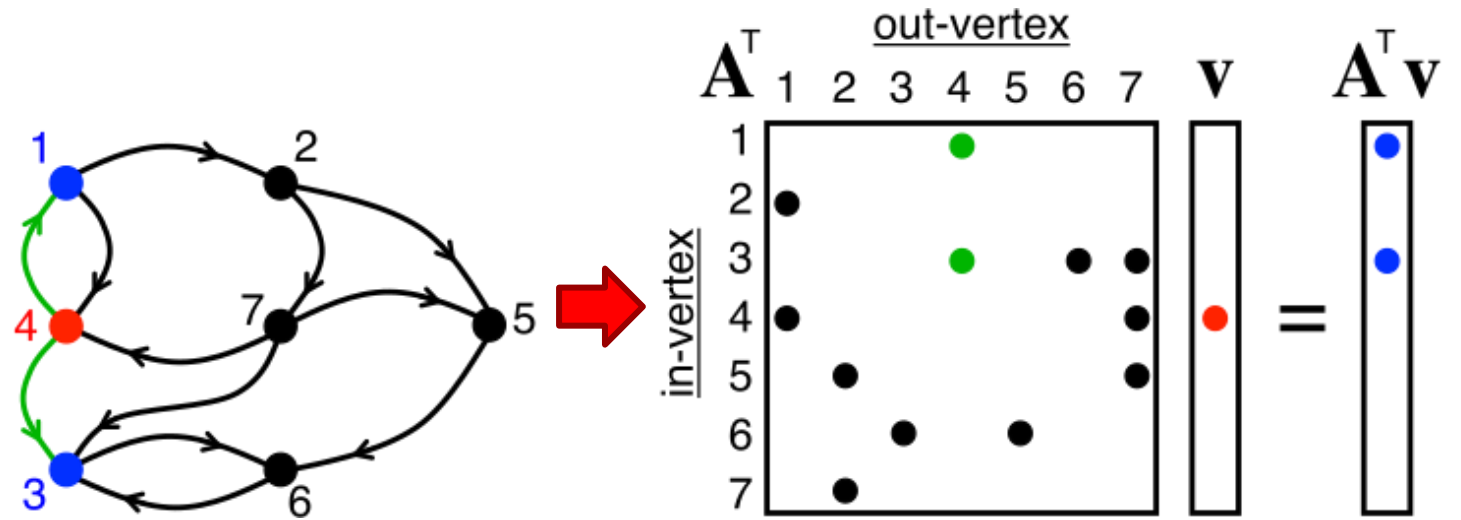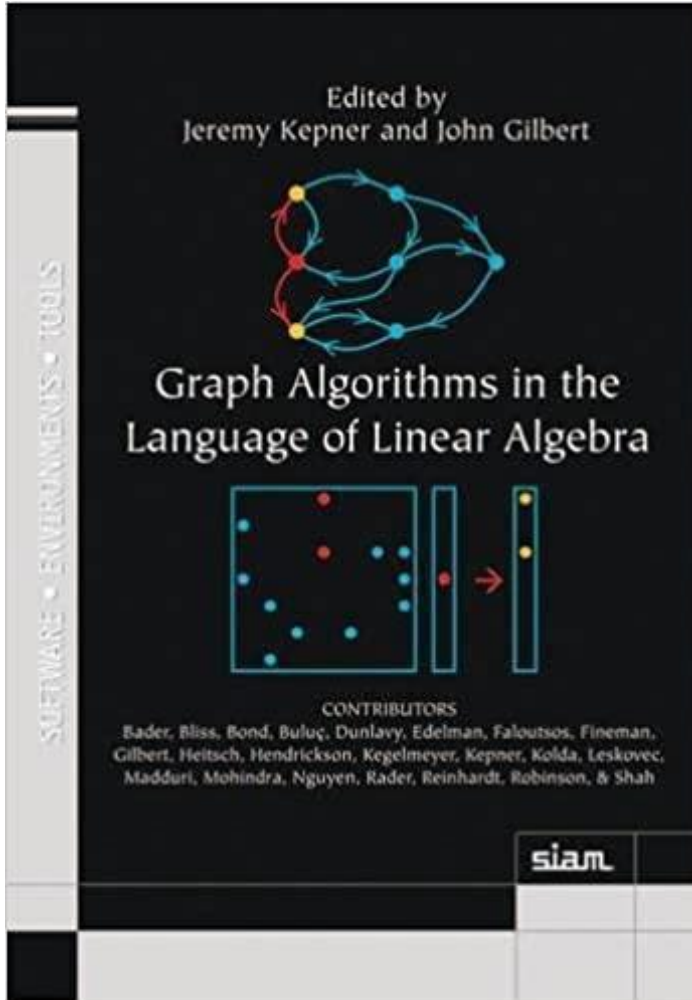rank_nminus1 is rank_n from the last iteration

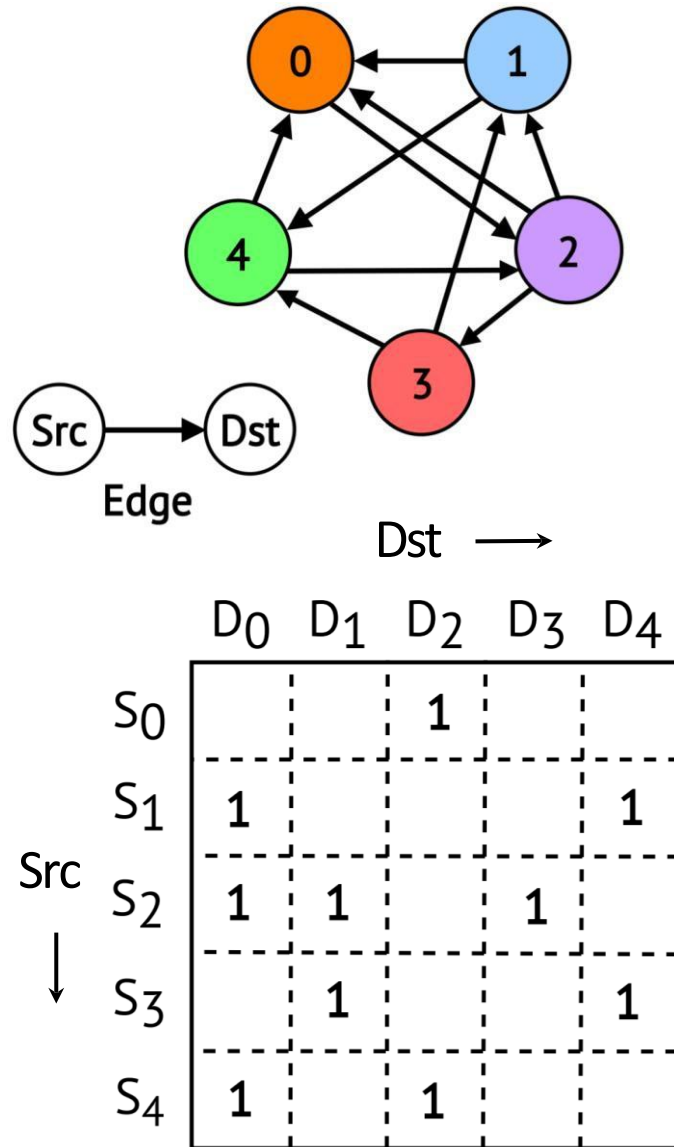# Graph Analytics can be mapped to Sparse Linear Algebra

# Graph Analytics can be mapped to Sparse Linear Algebra

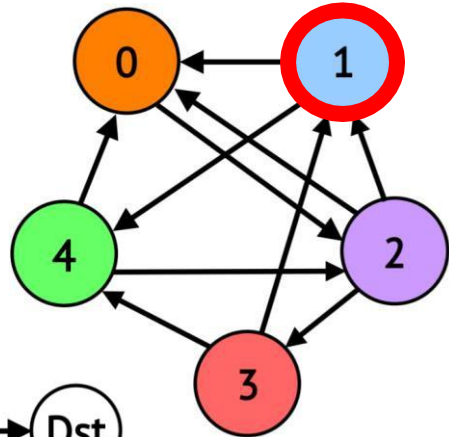# Graph Analytics can be mapped to Sparse Linear Algebra

# How do graph applications correspond to linear algebra?
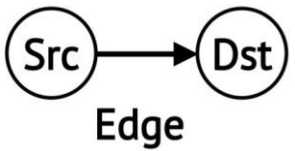
Src → Dst

**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 ↓

# How do graph applications correspond to linear algebra?



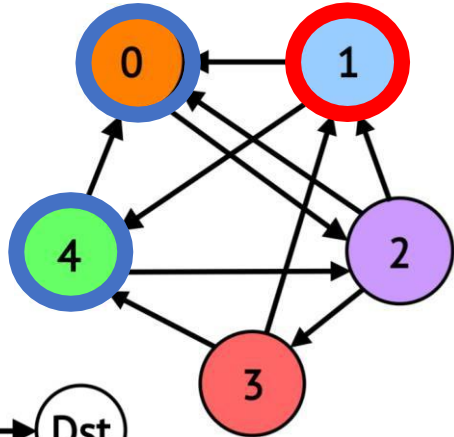**Matrix-transpose-vector product is one BFS iteration**

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

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.

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 ↓

$x_i$

1

$x_{i+1}$

1

1

=

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

12

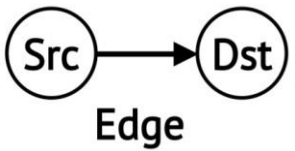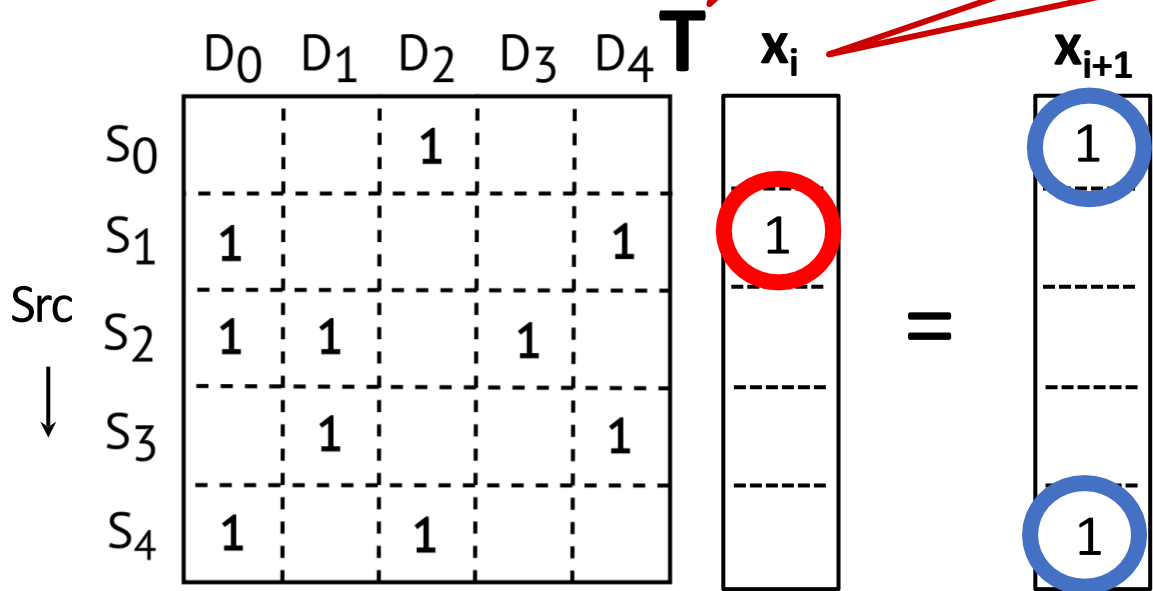# How do graph applications correspond to linear algebra?

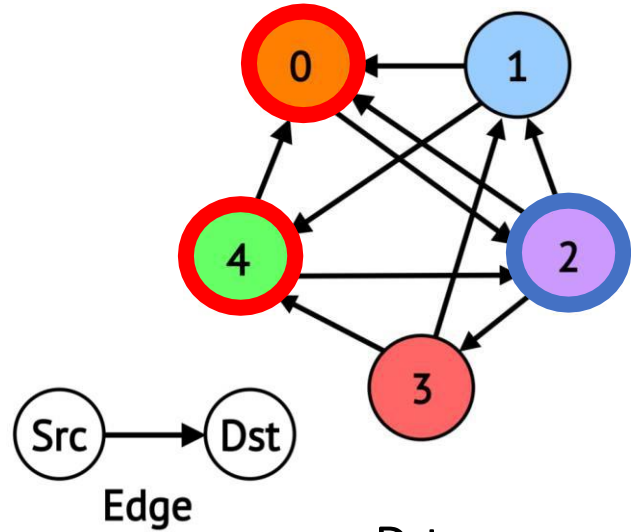**Matrix transpose vector product is one BFS iteration**

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

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

Src → Dst
Edge

Dst →

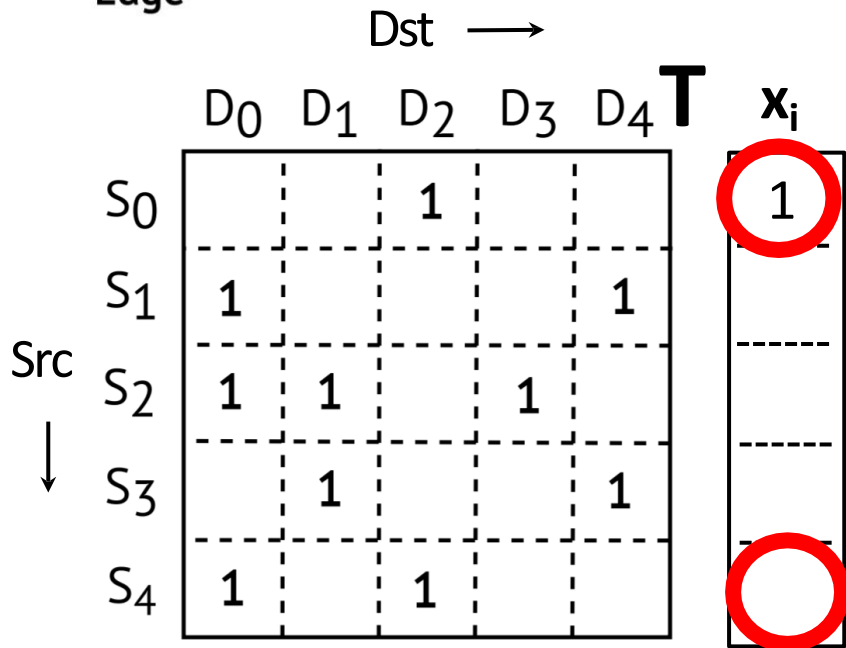|          | $D_0$ | $D_1$ | $D_2$ | $D_3$ | $D_4$ |
|----------|-------|-------|-------|-------|-------|
| $S_0$    |       |       | 1     |       |       |
| $S_1$    | 1     |       |       |       | 1     |
| $S_2$    | 1     | 1     |       | 1     |       |
| $S_3$    |       | 1     |       |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

Src ↓

**T**   $x_i$        $x_{i+1}$
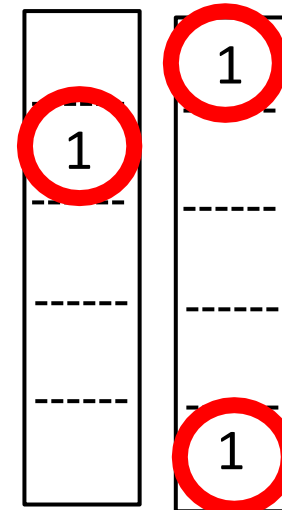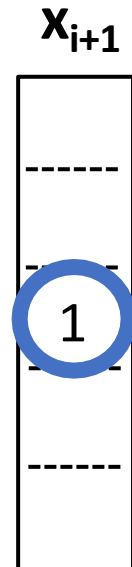
=

# How do graph applications correspond to linear algebra?

**Matrix transpose vector product is one BFS iteration**
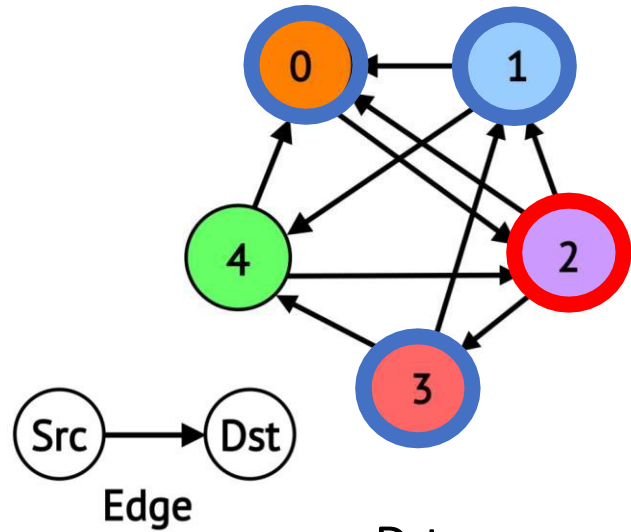
$$A^T x_i = x_{i+1}$$

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

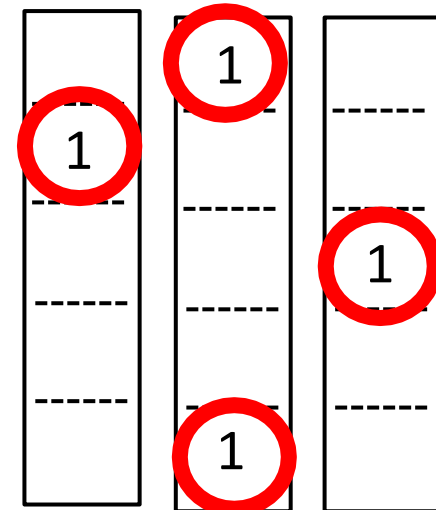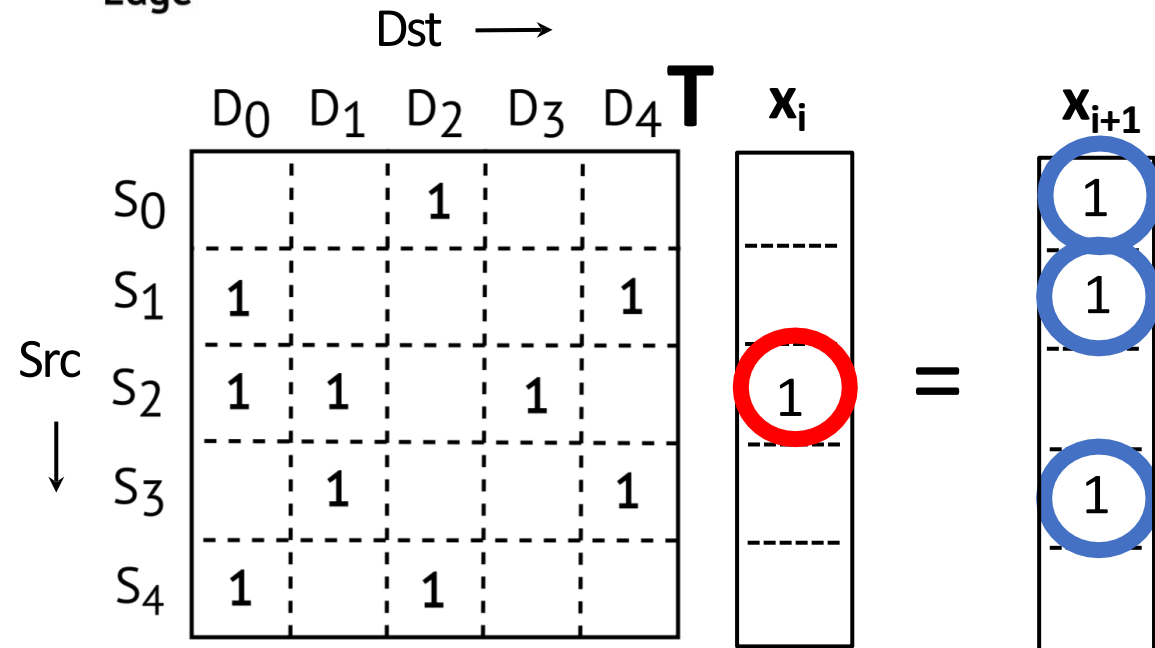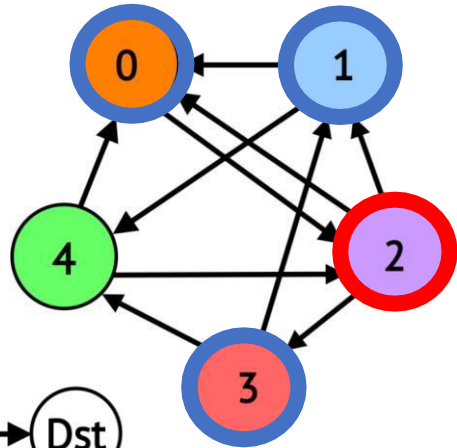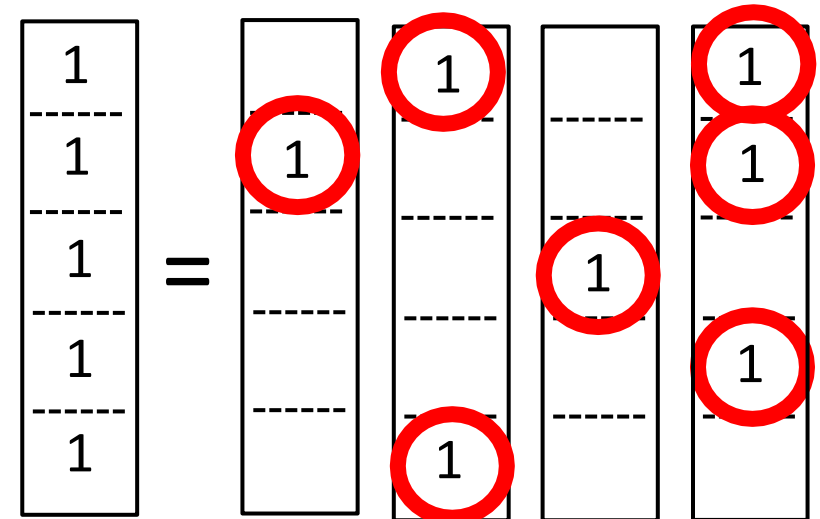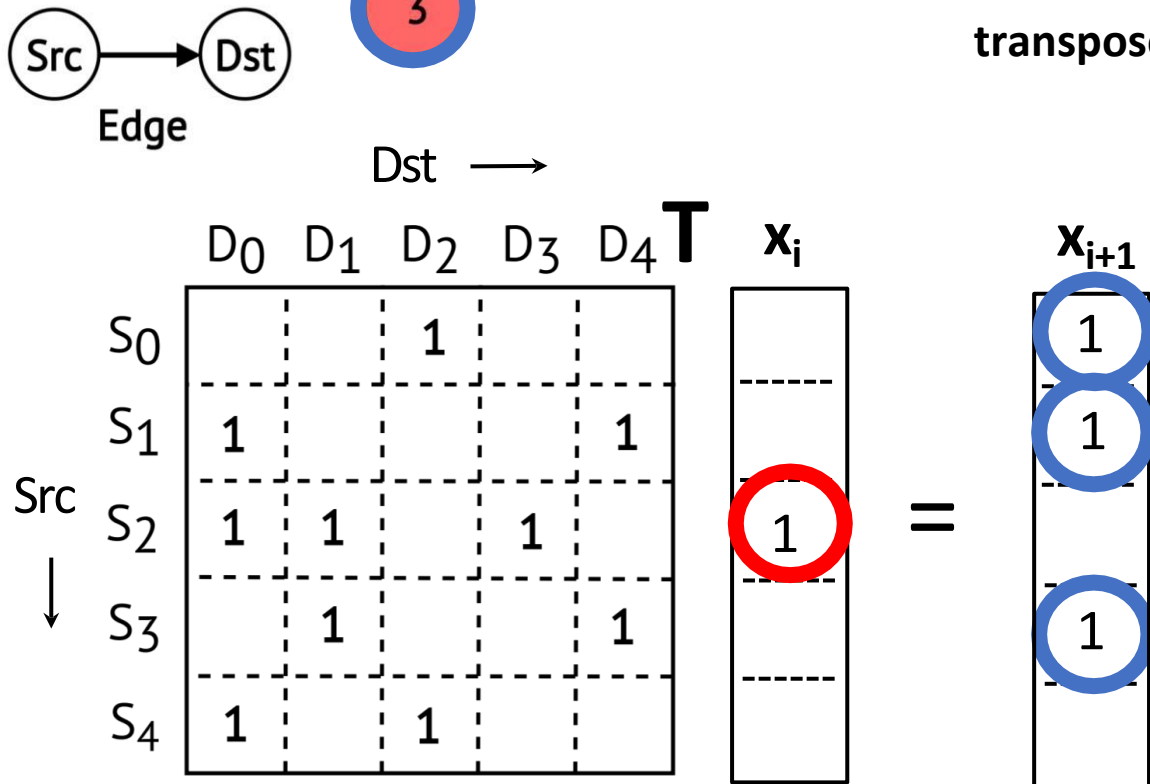# How do graph applications correspond to linear algebra?



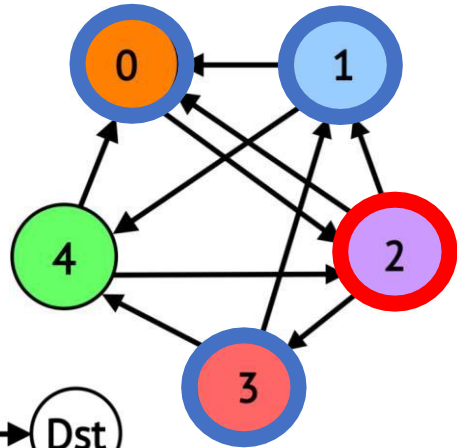**Matrix transpose vector product is one BFS iteration**

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

**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

18

# 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
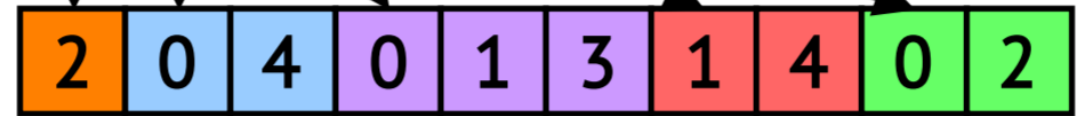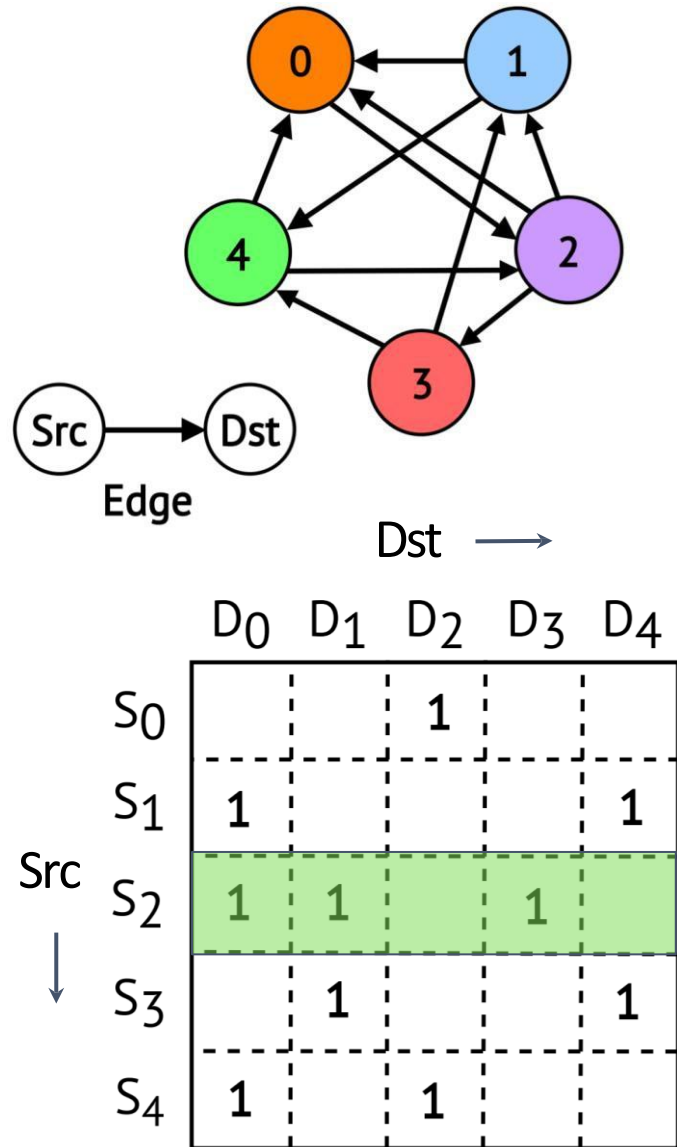


Offsets Array (OA)

**index is src id**

Neighbors Array (NA)

EdgeList sorted by SrcIDs

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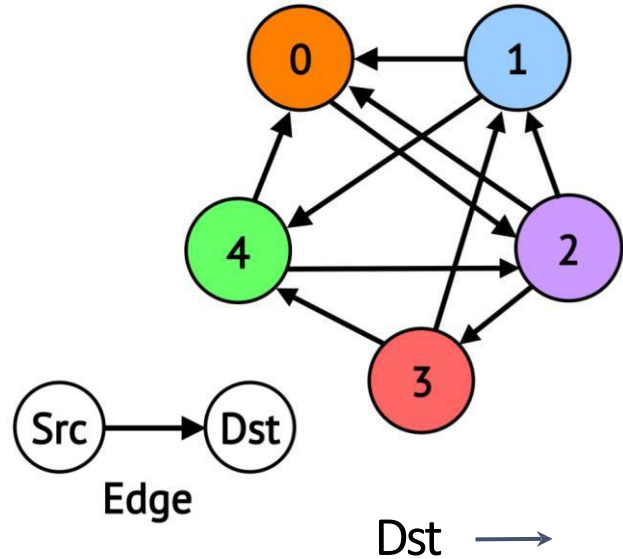*Outgoing Neighbors*

**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

# 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*

# 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

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

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

neigh_count

# 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

```
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



COO
(EdgeList)

```
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
```
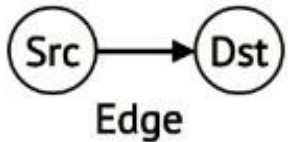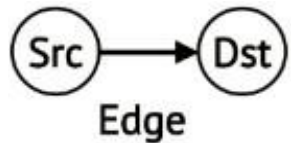
| 0 | 2 | 3 | 4 | 6 |   OA (also OA_dup)

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

OA (also OA_dup)

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

COO (EdgeList)

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

```
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**

Src → Dst

Edge

# Compressed Representations ⇒ Irregular Memory Accesses

Dst →

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

| 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

NA

**CSC**

# Compressed Representations ⇒ Irregular Memory Accesses

**Push (CSR Traversal)**

**Dst** ⟶

```
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
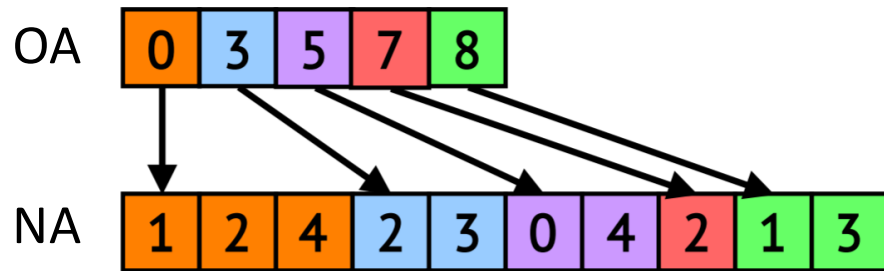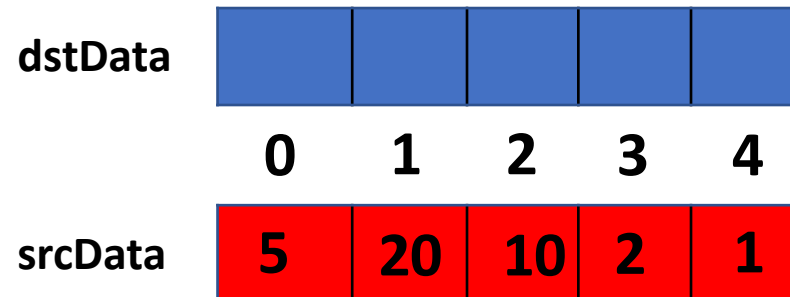
Dst →

**Push (CSR Traversal)**

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

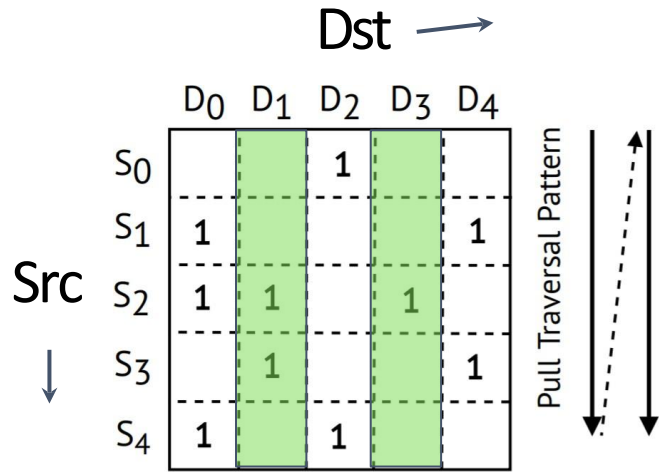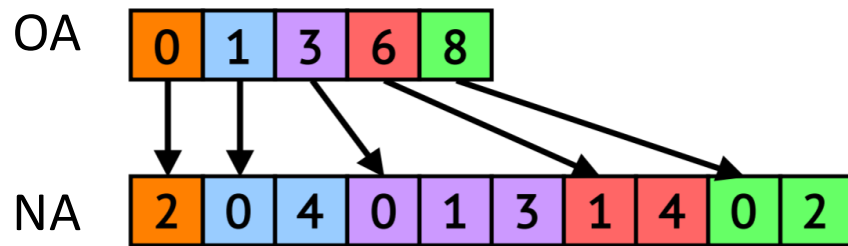**Push traversal performs *irregular write operations* that lack locality**

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

| | | | | | |
|---|---|---|---|---|---|
| **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

NA

**CSR**

Irregular Data Footprint >> LLC Size
Size of srcData ~ **256MB** (32M vertices * 8B)

# 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**

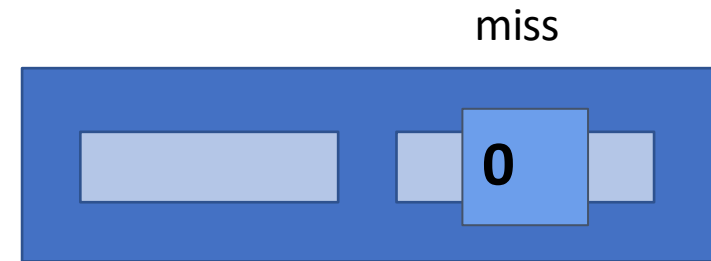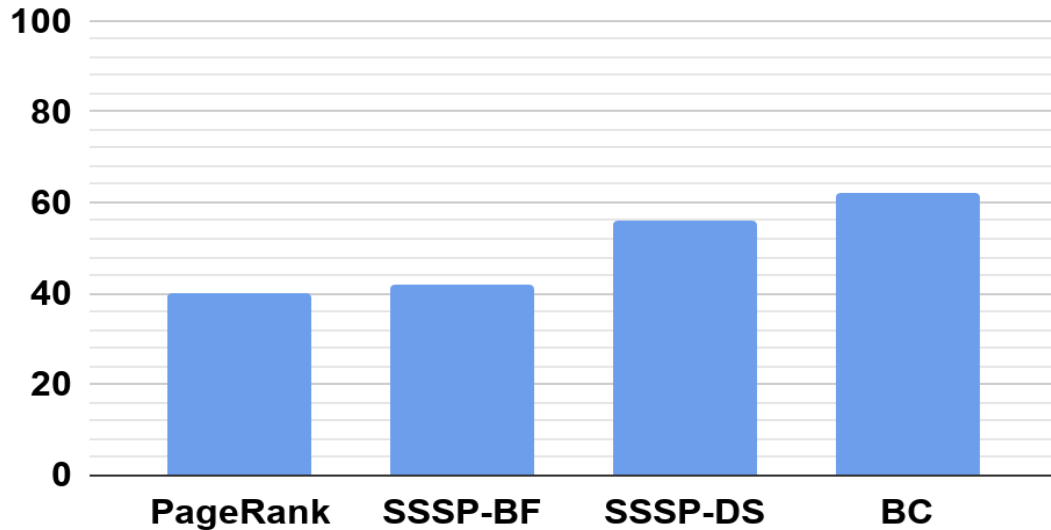# Irregular Accesses Lead to Poor Locality

**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;*
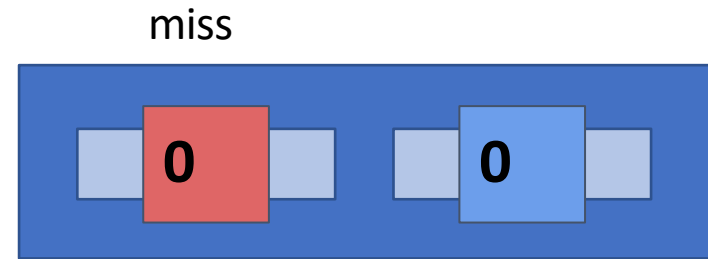
# Irregular Accesses Lead to Poor Locality

**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;*
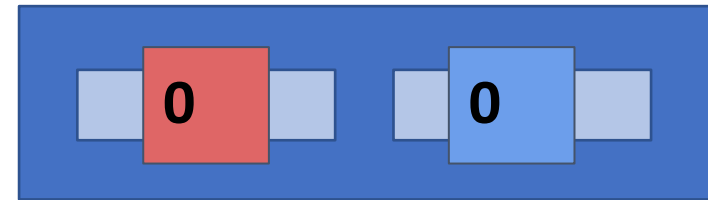
# Irregular Accesses Lead to Poor Locality

**LLC Miss Rate (%)**



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

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

COO
(EdgeList)

miss

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

34

*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!!!**



**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!!!**
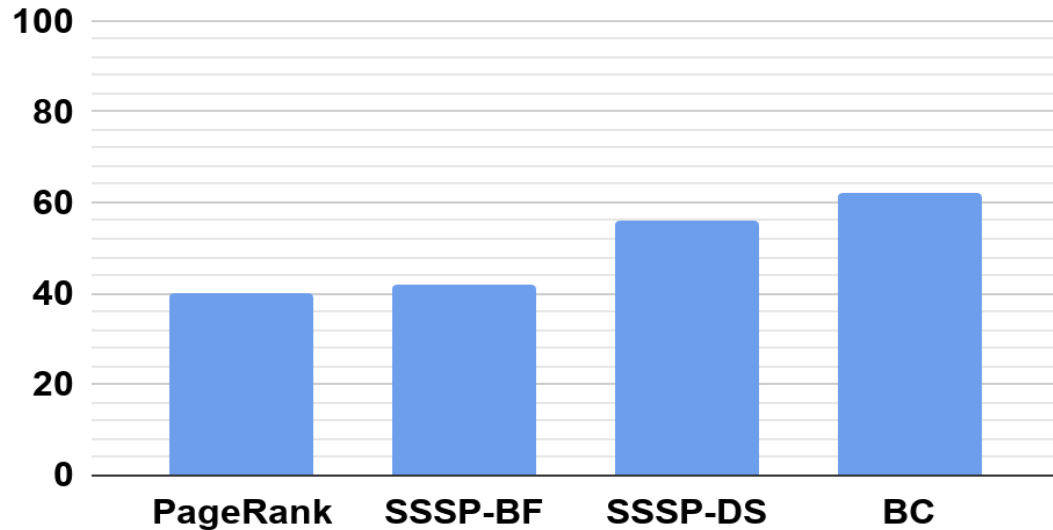


COO
(EdgeList)

miss

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

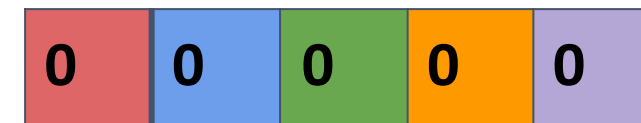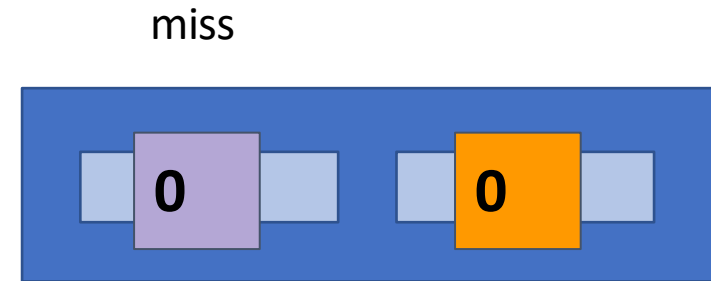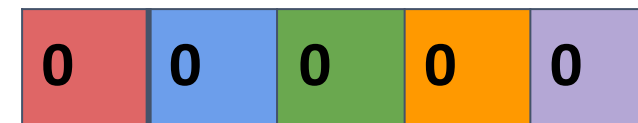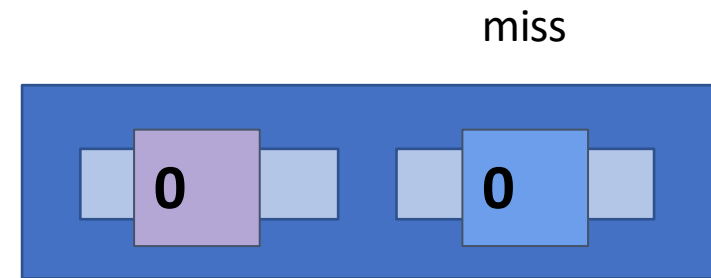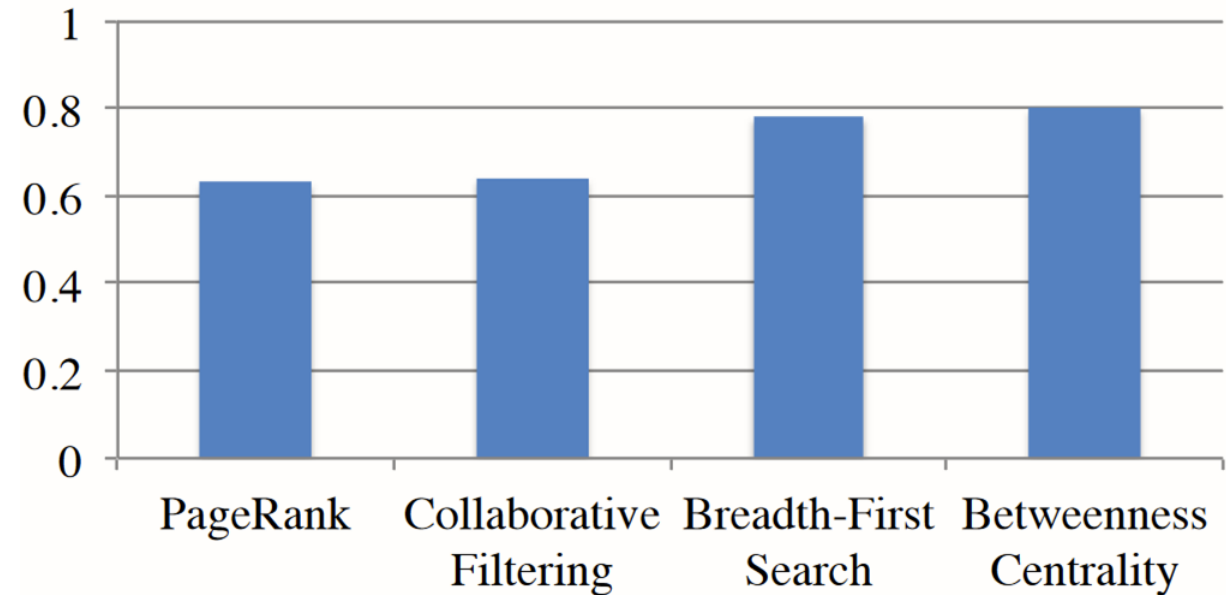# Irregular Accesses Lead to Poor Locality

### LLC Miss Rate (%)
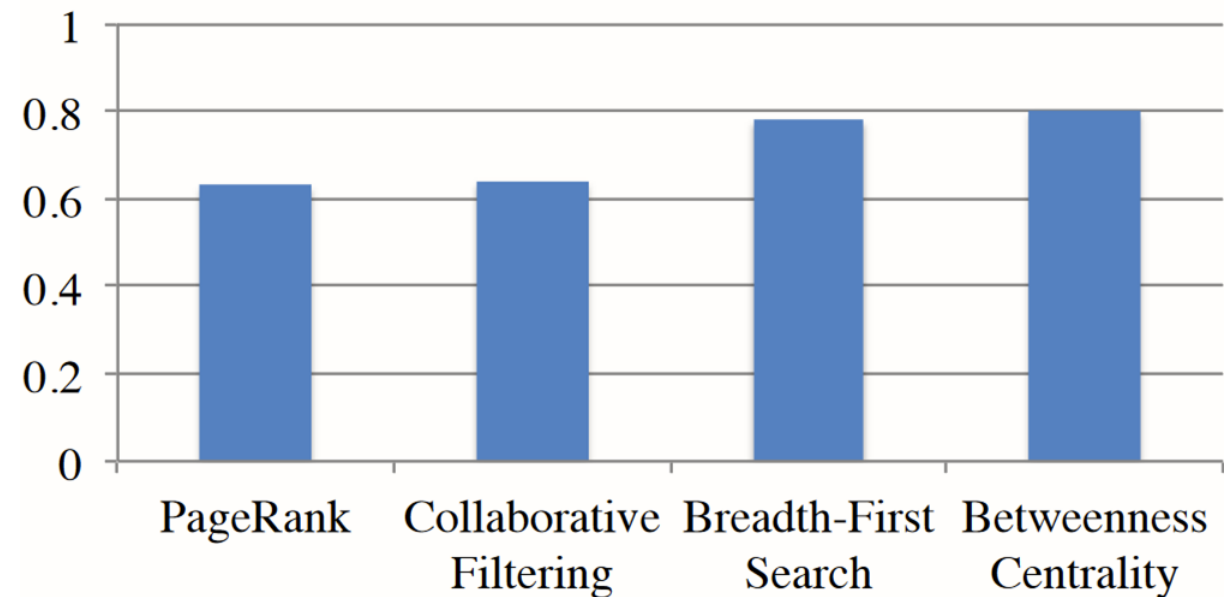


### Cycles stalled on DRAM / Total Cycles



**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
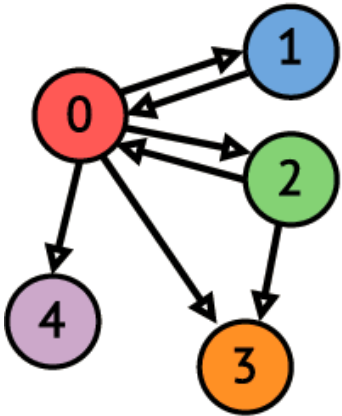
**LLC Miss Rate (%)**

**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!



COO
(EdgeList)

Src → Dst
Edge

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

| 2 | 1 | 1 | 2 | 1 |

neigh_count

## 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

# 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 |

Src → Dst
Edge

OA

| 0 | 2 | 3 | 4 | 6 |

NA

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

**Completed CSC**

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)

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)*

# The Roofline Model

# The Roofline Model

# 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



CPU (compute, flop/s)

DRAM Bandwidth (GB/s)

DRAM (data, GB)

Memory-Bound ← → Compute-Bound

Throughput *(GFLOP/s)*

*Peak FLOPS*

*Peak Mem BW*

● 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)*

# The Roofline Model



CPU
(compute, flop/s)

DRAM Bandwidth
(GB/s)

DRAM
(data, GB)

Memory-
Bound ⬅

Compute-
Bound ➡

*Peak FLOPS*

● App 2

● App 1

Throughput
*(GFLOP/s)*

*Peak Mem BW*

**What is this point?** ⬅

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

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)*

51

# Operational Intensity of Irregular Graph Applications

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

COO
(EdgeList)

Src → Dst

Edge

```
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)

Src → Dst
Edge

```
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)*

56

# 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?**

**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).**

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

Memory-Bound ⇐ ⇒ Compute-Bound

Throughput
*(GFLOP/s)*

*Peak*

*Peak Mem BW*

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

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

How to improve BW utilization?

1/16

250

Operational Intensity
*(FLOPS/Byte)*

59

# 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]
```
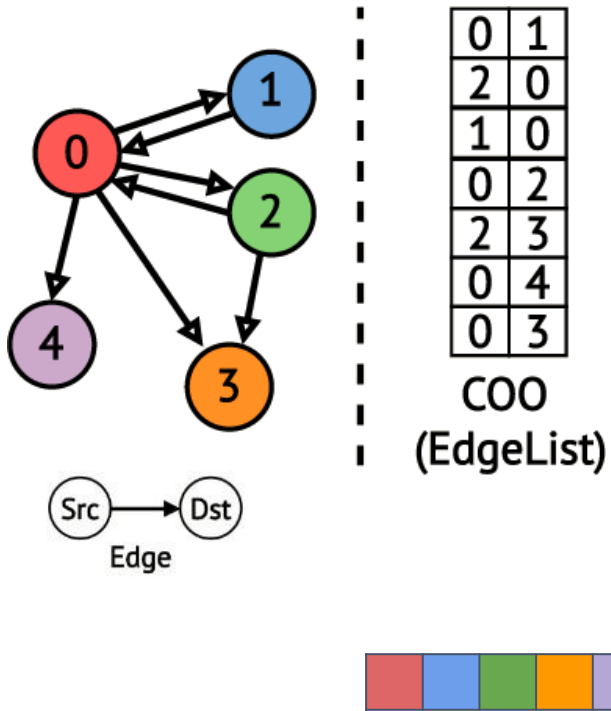
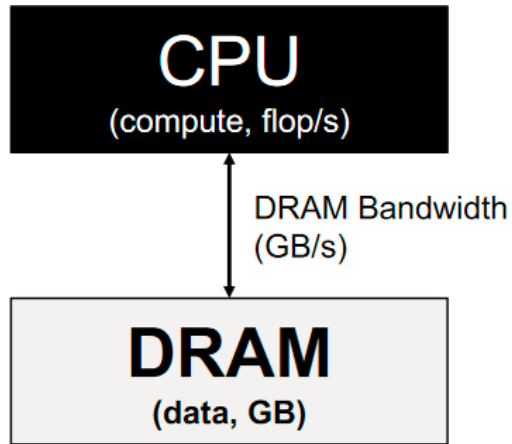Ideal Best Possible Operational Intensity?
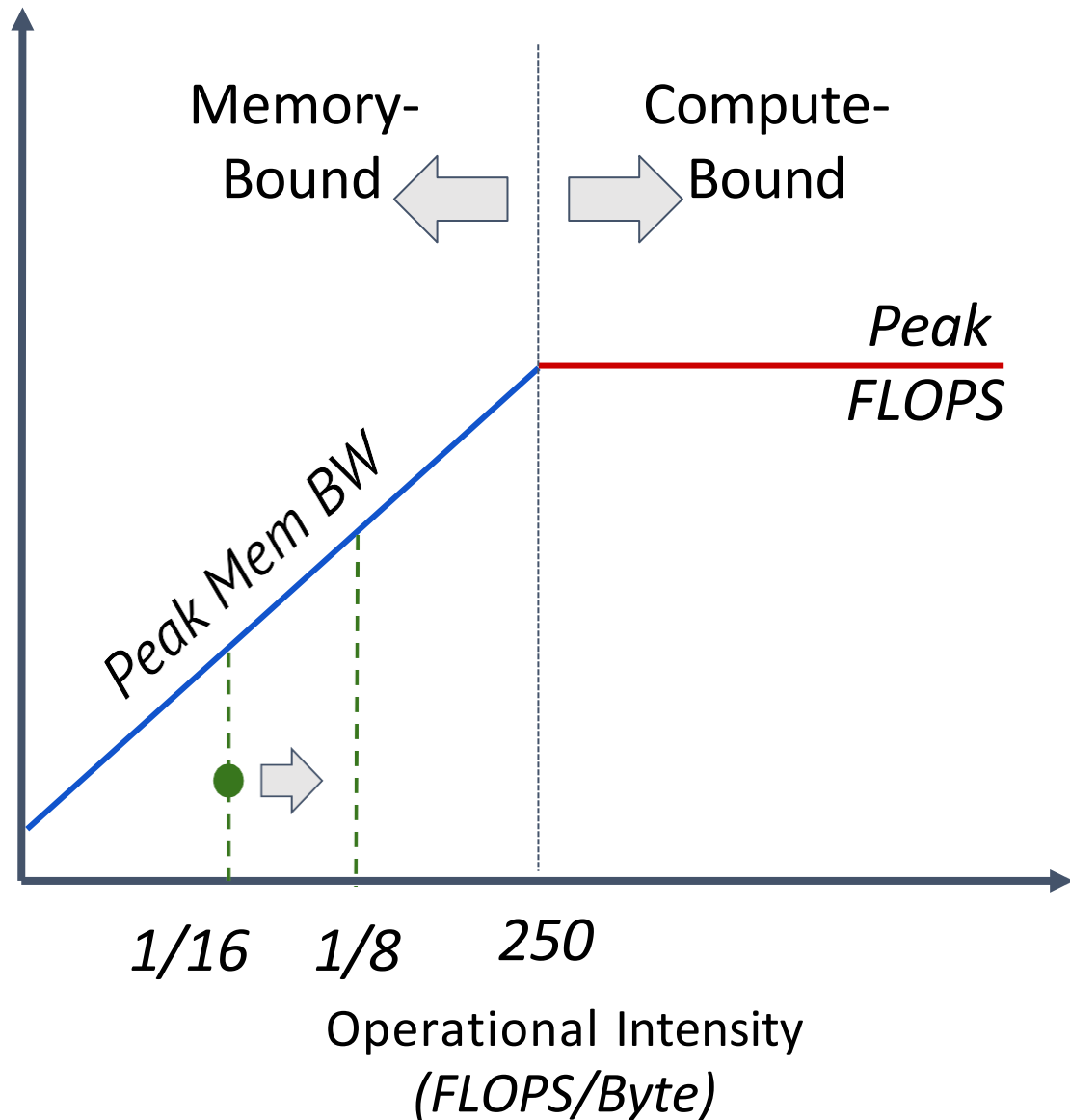
**Operations per byte:**

**Operations:** 1 addition

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

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

# Improving Operational Intensity (OI) by Improving Locality

# Improving Operational Intensity (OI) 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 & more leftward ridge point.

Why is cache BW > DRAM BW?

*1/16   1/8      250*

Operational Intensity
*(FLOPS/Byte)*

# Improving Operational Intensity (OI) 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 &  more
leftward ridge point.

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

*1/16    1/8        250*

Operational Intensity
*(FLOPS/Byte)*

# Improving Operational Intensity (OI) by Improving Locality

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

DRAM Bandwidth (GB/s)
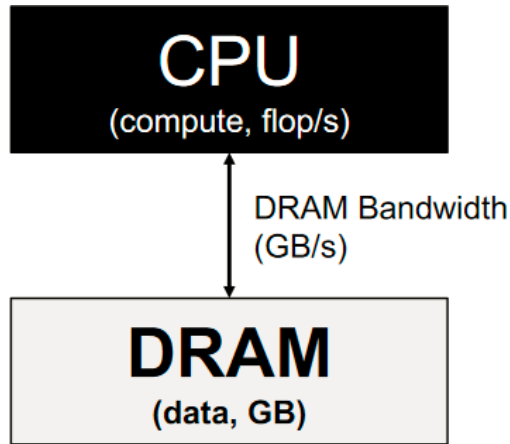
**DRAM** (data, GB)

Throughput *(GFLOP/s)*

Memory-Bound ⬅

Compute-Bound ➡

*Peak FLOPS*

Peak Cache BW

Peak Mem BW

**Key Question:** So how do we improve locality, to reduce data movement from memory, to increase OI, and move to peak performance?

**Locality wins:** If we can operate out of cache, higher ceiling & more leftward ridge point.
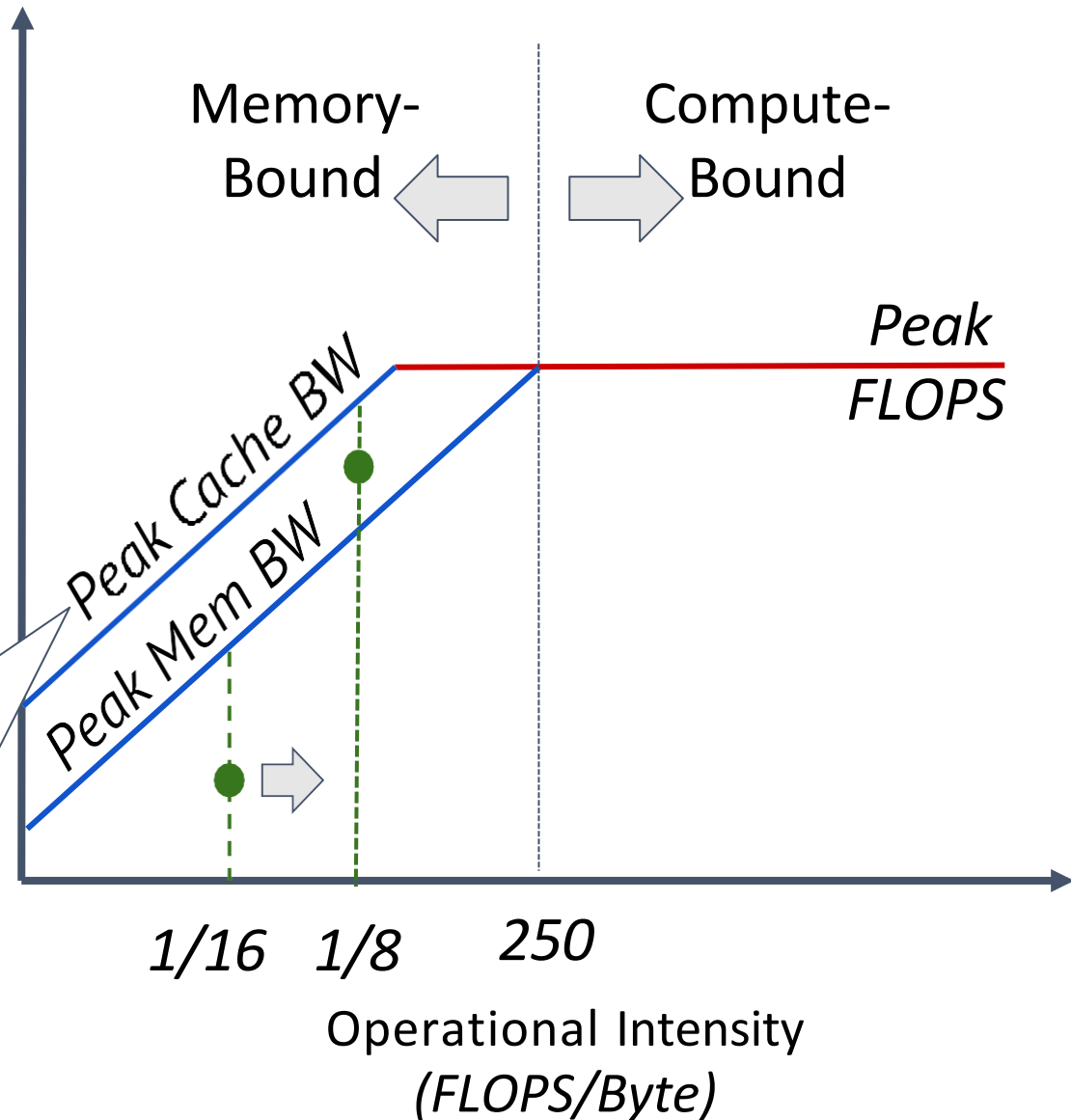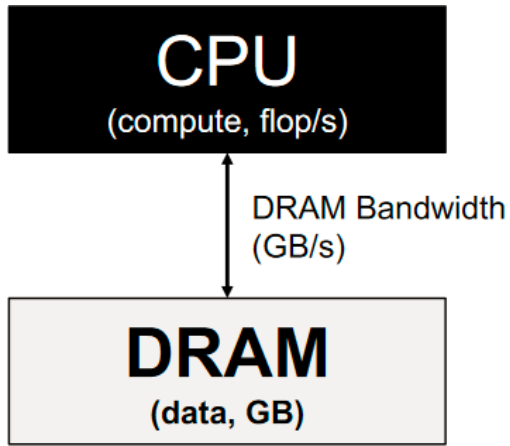
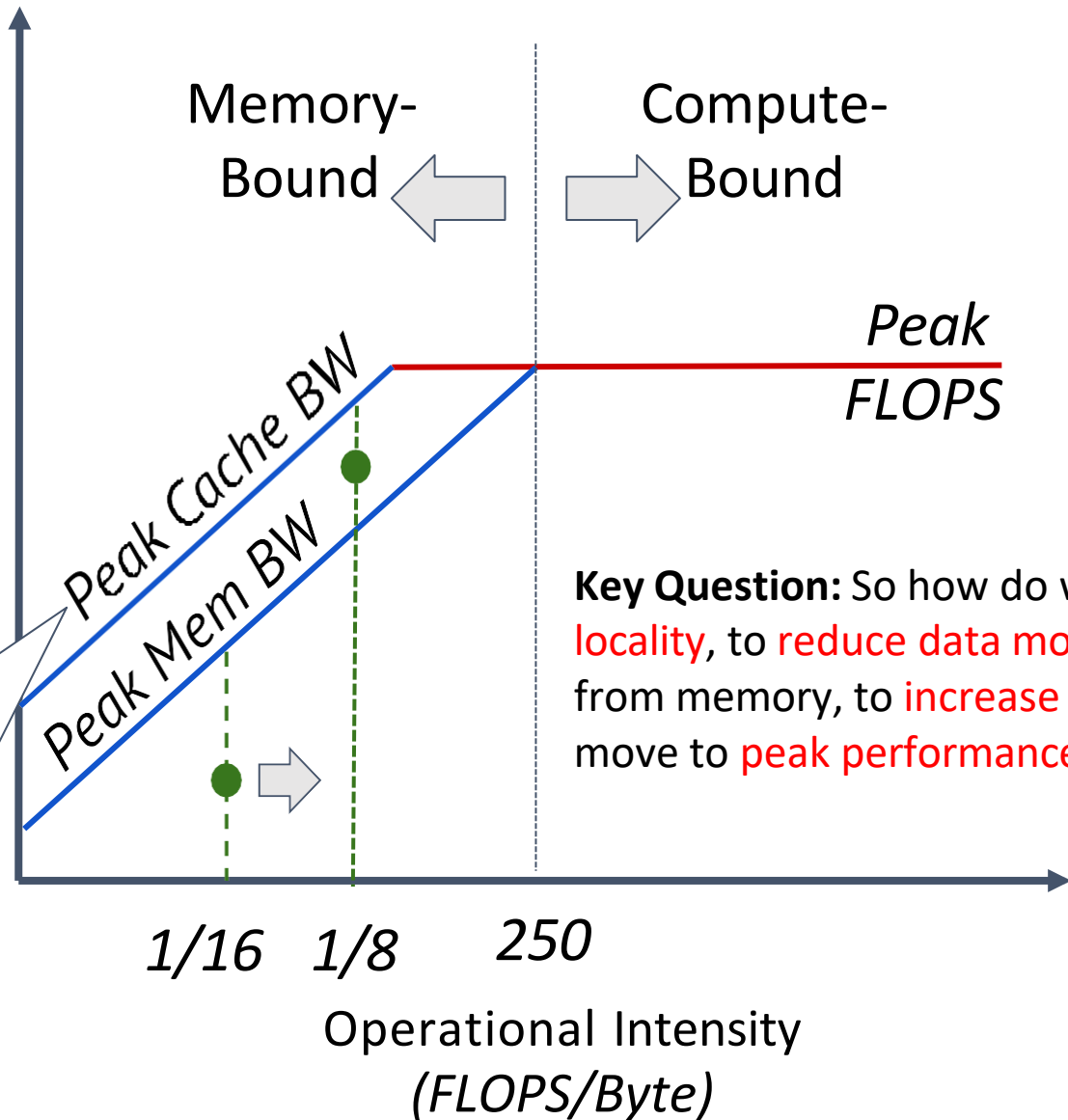Why is cache BW > DRAM BW?
**Smaller SRAM caches much faster.**

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

Operational Intensity *(FLOPS/Byte)*

# 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