

# Graph Processing in NVRAM and Streaming Settings

---

Laxman Dhulipala

MIT (Postdoc)

<https://ldhulipala.github.io/>

Based on joint work with

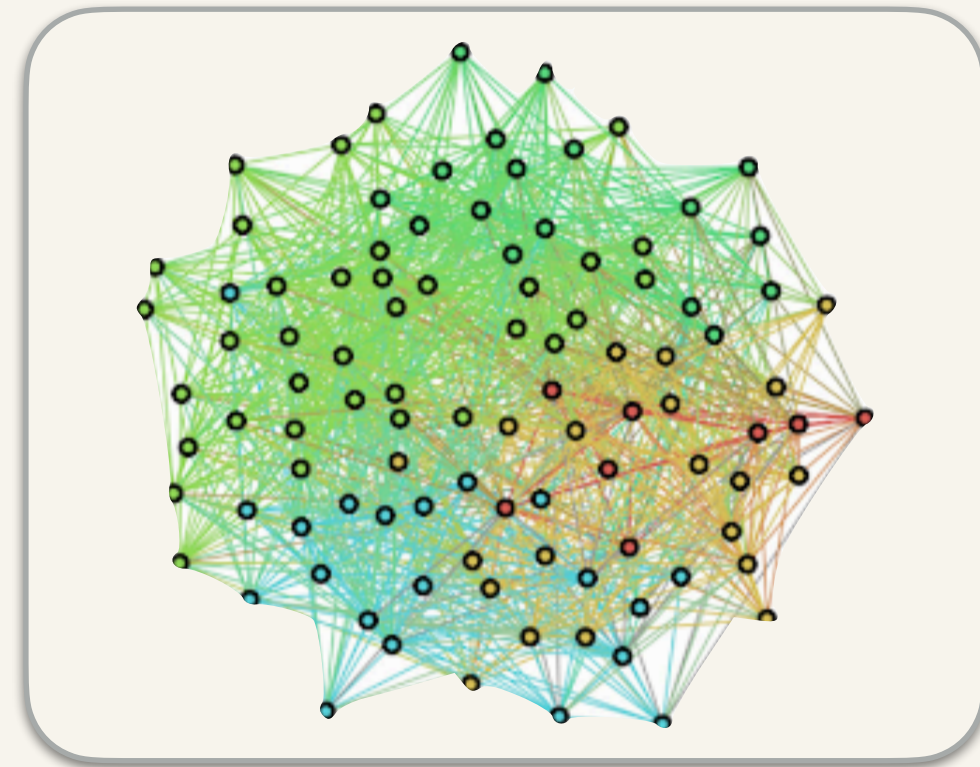
Guy Blelloch and Julian Shun (PLDI'19)

Charles McGuffey, Hong Kang, Yan Gu, Guy Blelloch, Phil Gibbons, and Julian Shun (VLDB'20)

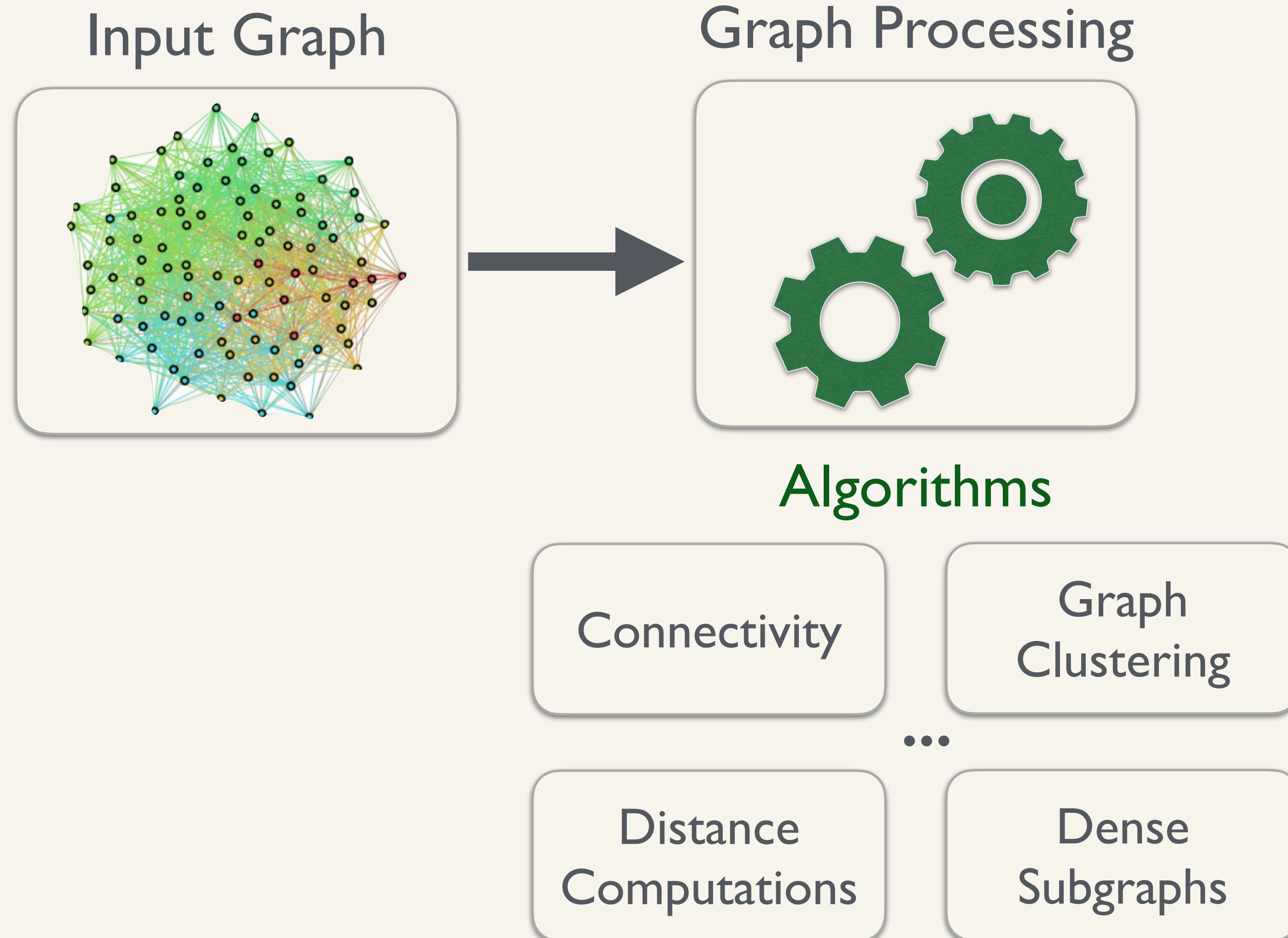
**Graph Processing:** algorithms and systems that enable us to analyze and understand graphs

# Graph Processing: algorithms and systems that enable us to analyze and understand graphs

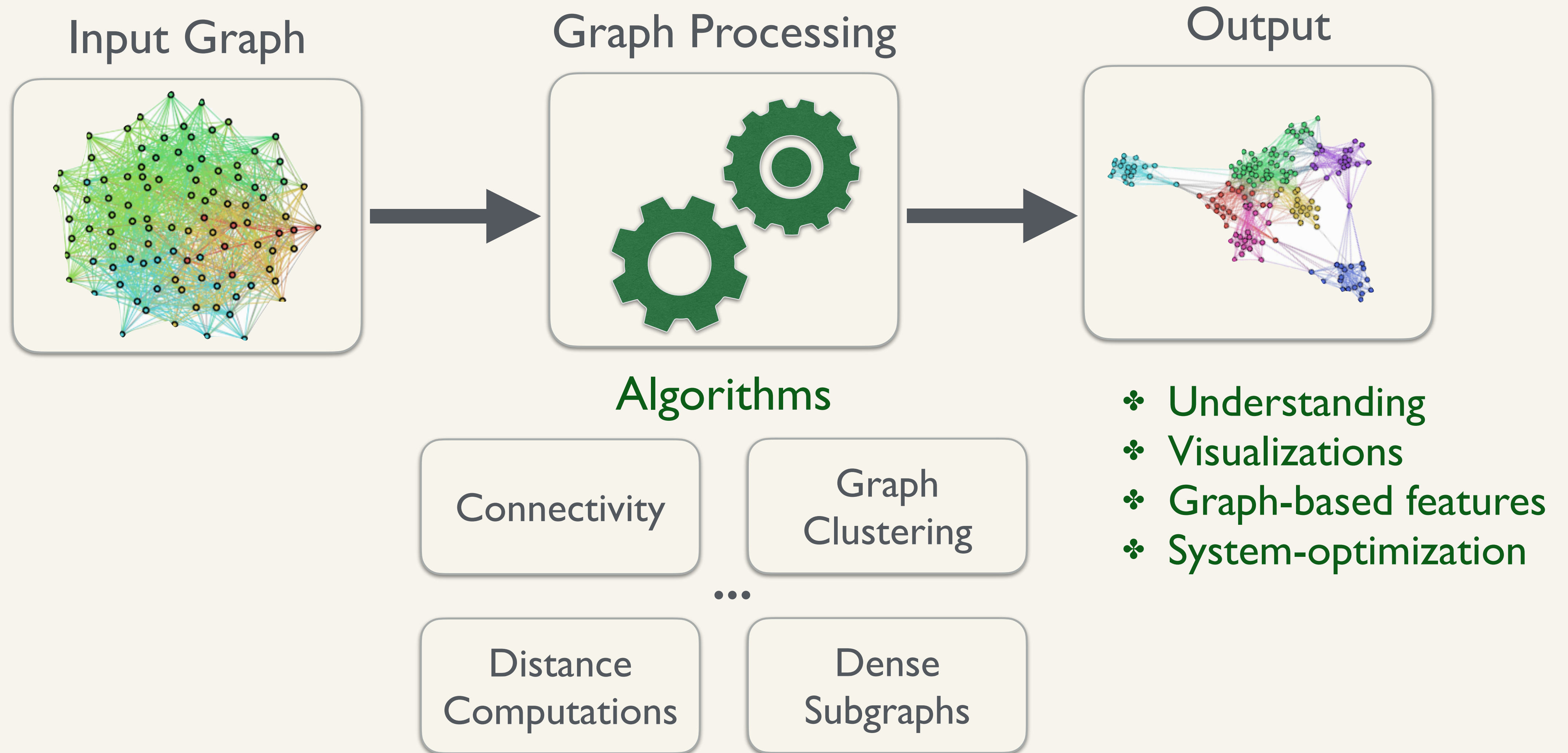
Input Graph



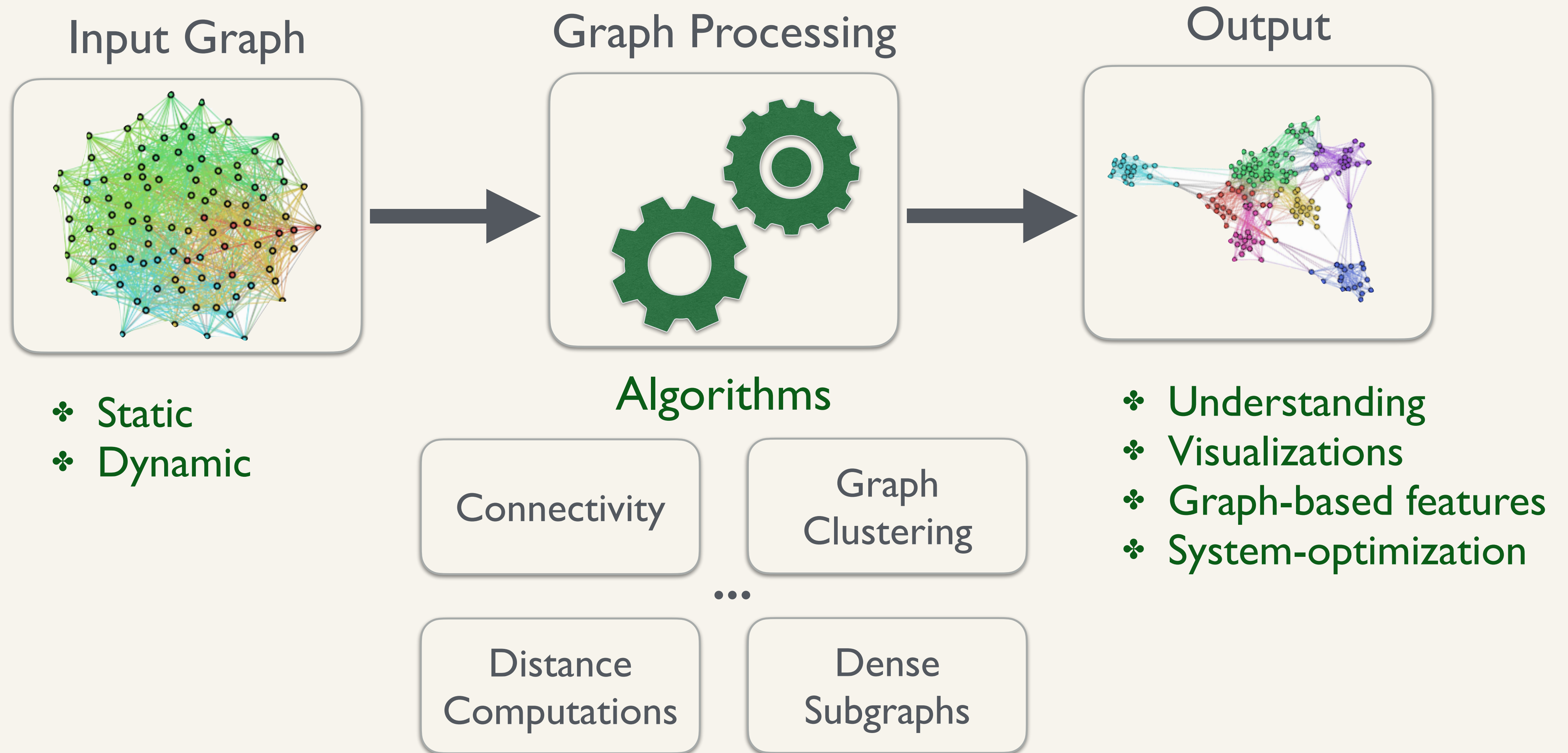
# Graph Processing: algorithms and systems that enable us to analyze and understand graphs



# Graph Processing: algorithms and systems that enable us to analyze and understand graphs



# Graph Processing: algorithms and systems that enable us to analyze and understand graphs

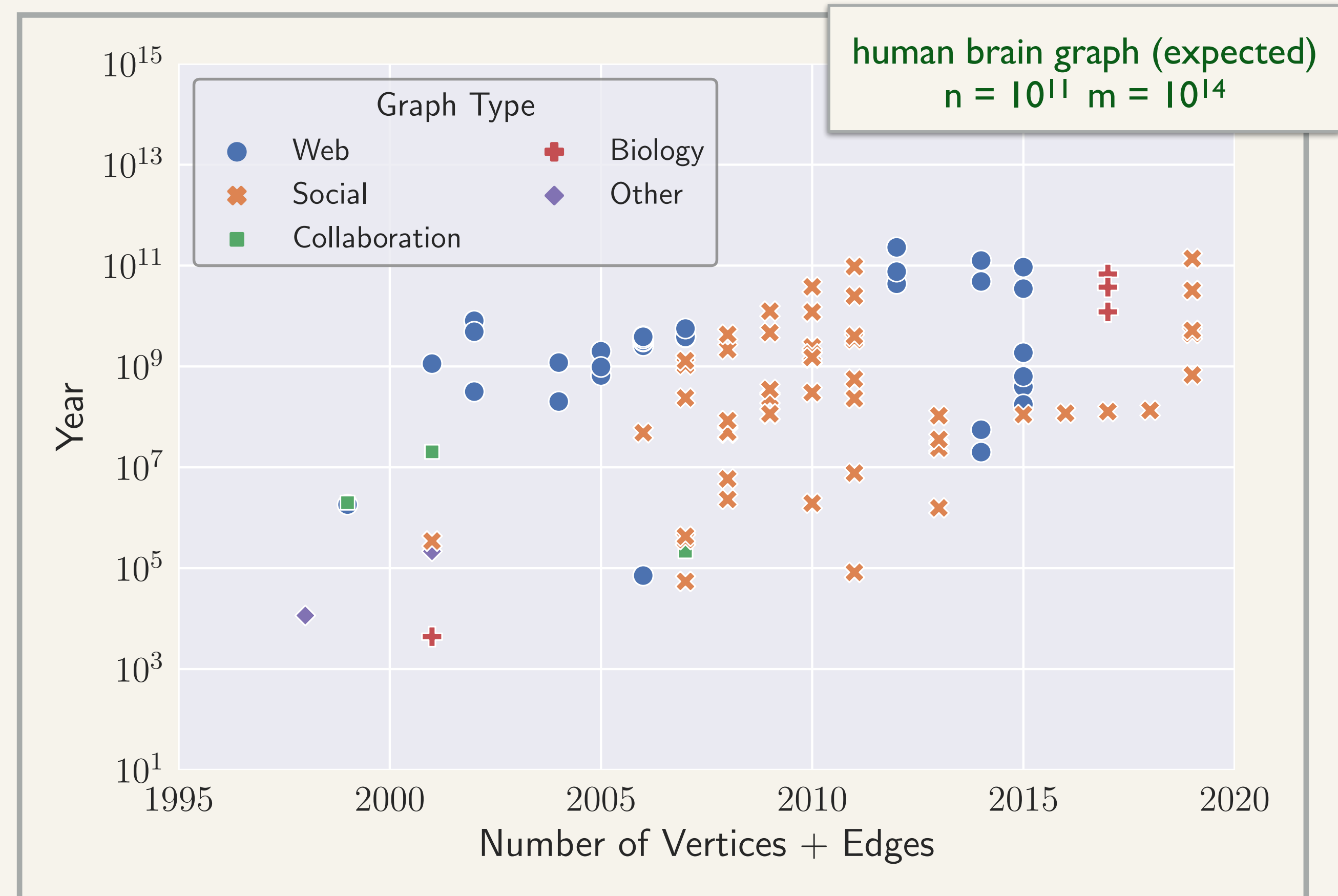


# Large-Scale Graph Processing

## WebDataCommons hyperlink graph

- ❖ 3.5 billion vertices and 128 billion edges
- ❖ ~1TB of memory to store
- ❖ Largest publicly available graph

*“...[the 2012 graph is the] largest hyperlink graph that is available to the public outside companies such as Google, Yahoo, and Microsoft.”*



Year of sourcing vs total number of vertices and edges for real-world graphs from the SNAP and LAW datasets

# Shared-Memory Parallelism

## Shared-Memory Machines

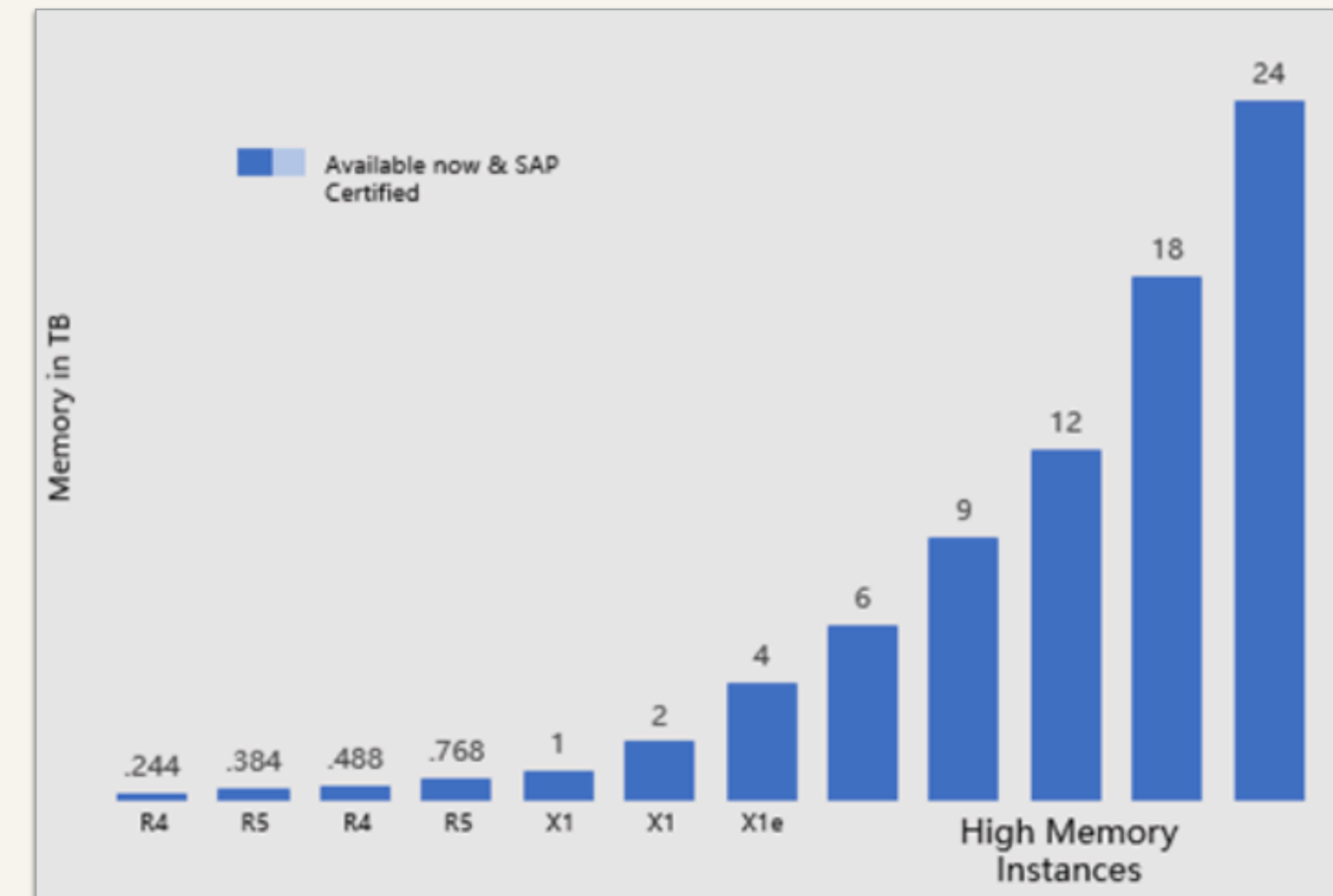
- Cost for a 1TB memory machine with 72 processors is about \$20,000.
- Can rent a similar machine (96 processors and 1.5TB memory) for \$11/hour on Google Cloud



## WebDataCommons Graph

- 3.5 billion vertices and 128 billion edges

A single shared-memory machine can already store the largest publicly available graph datasets, with plenty of room to spare





# Shared-Memory Parallelism

## Shared-Memory Machines

- Cost for a 1 TB memory machine with 72 processors is about \$20,000.
- Can rent a similar machine (96 processors and 1.5TB memory) for \$11/hour on Google Cloud

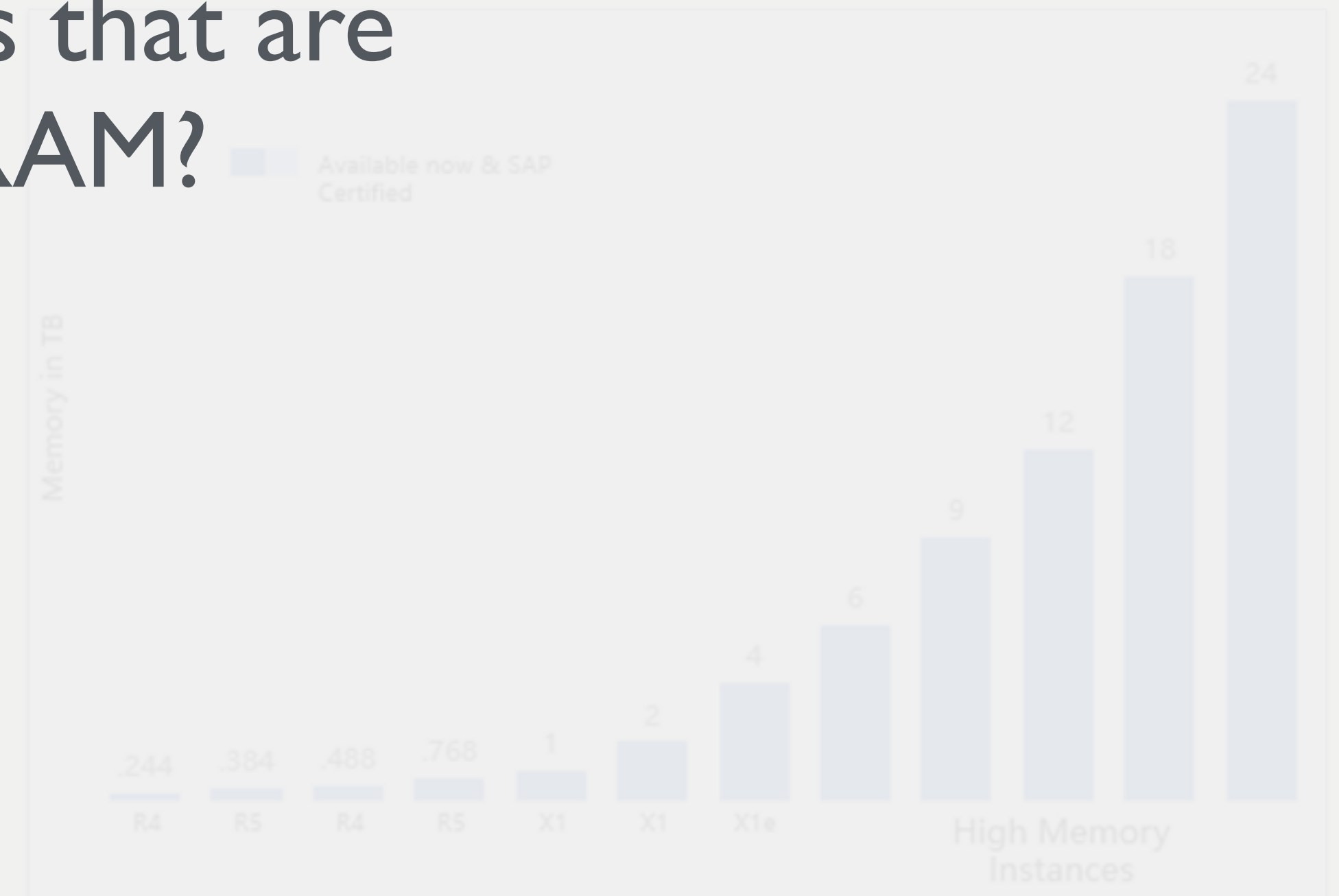


**What about graphs that are larger-than-DRAM?**

## WebDataCommons Graph

- 3.5 billion vertices and 128 billion edges

A single shared-memory machine can already store the largest publicly available graph datasets, with plenty of room to spare

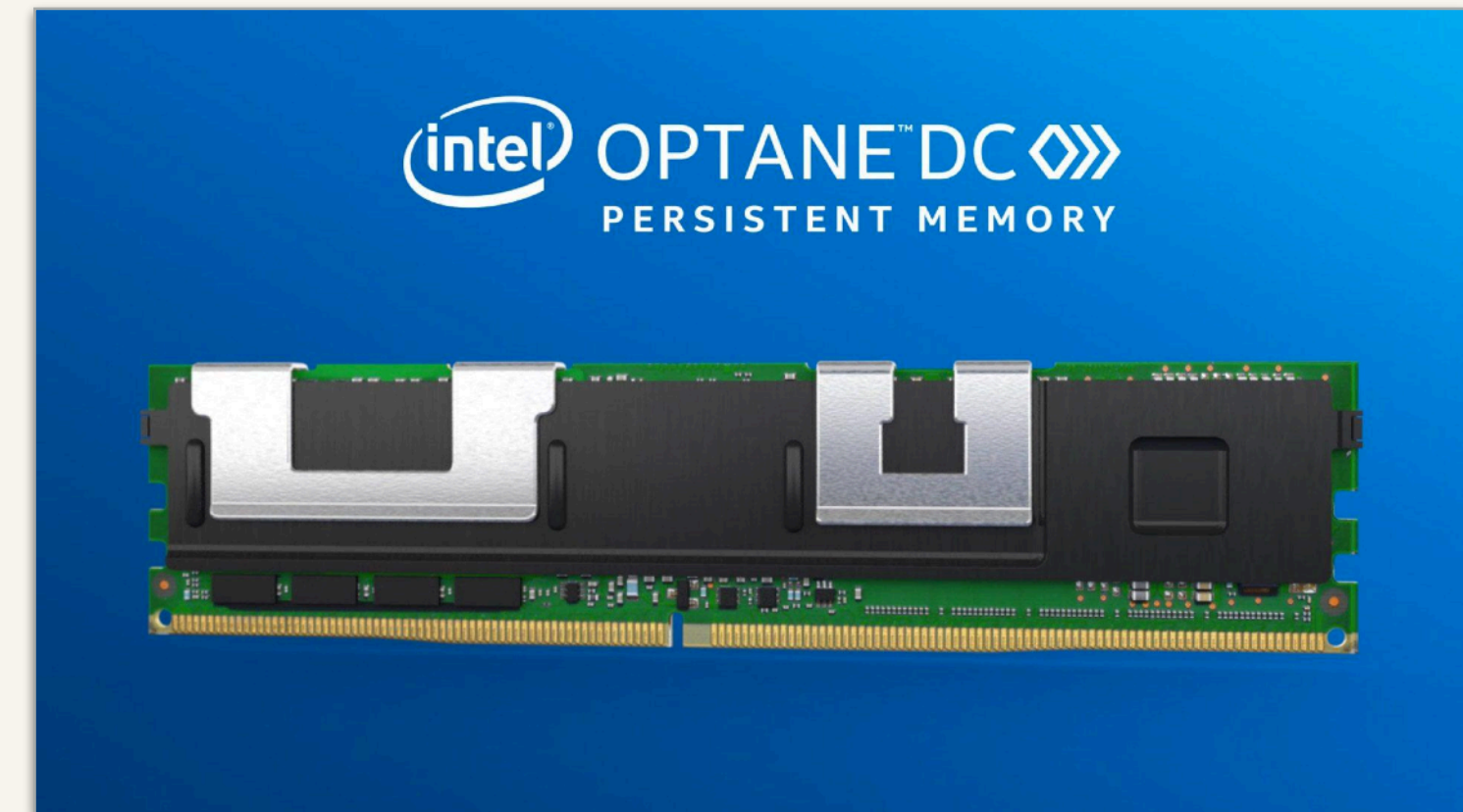


# NVRAM Graph Processing

# Non-Volatile Memory (NVRAM)

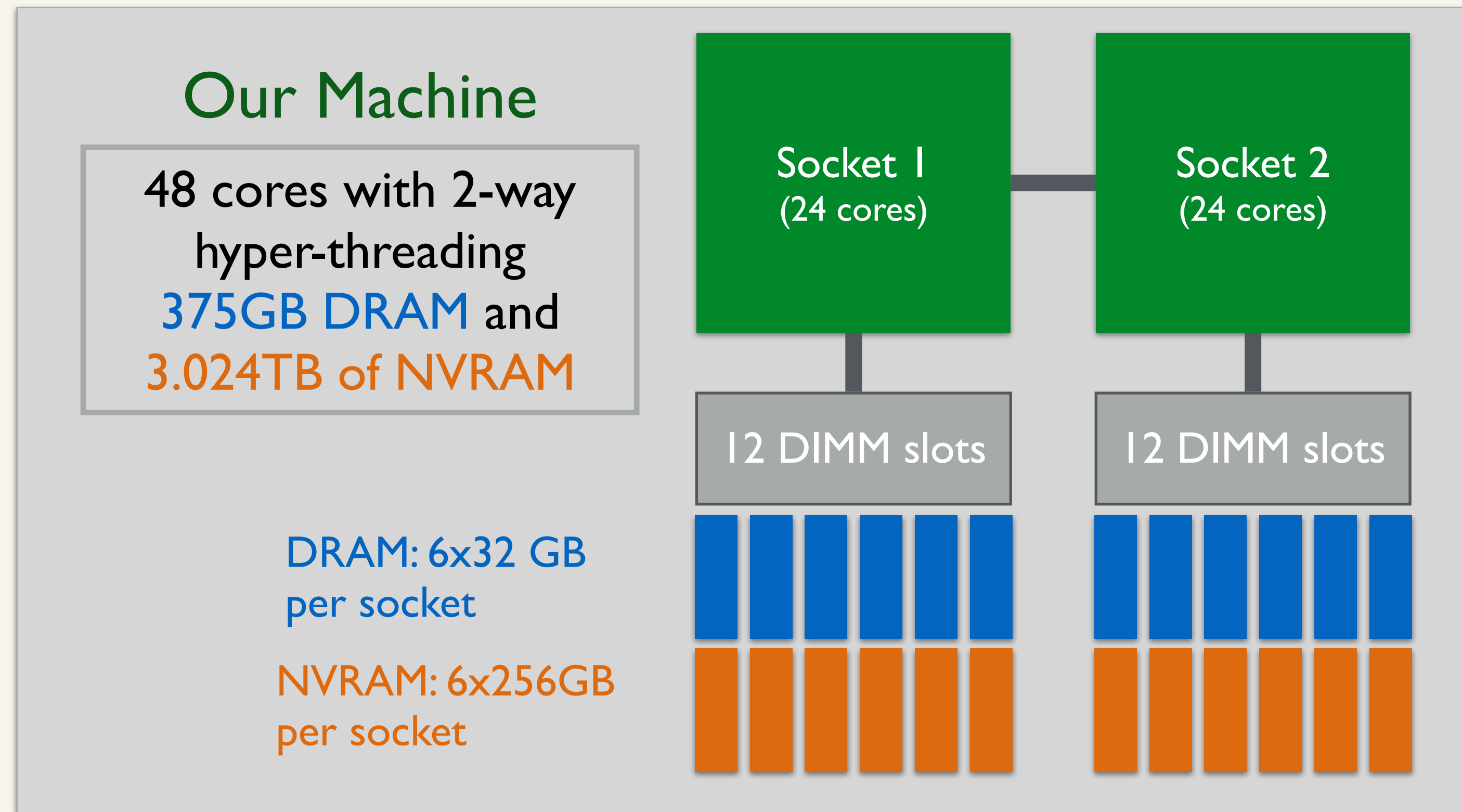
## Intel Optane DC Memory

- ❖ Cheaper than DRAM on a per-byte basis
- ❖ Order of magnitude more capacity
- ❖ Memory is *persistent* and *byte-addressable*



*Can we design algorithms that effectively use NVRAM as a higher-capacity memory while achieving DRAM-competitive performance?*

# NVRAM Characteristics



- ❖ 8x more NVRAM than DRAM
- ❖ NVRAM read throughput ~3x lower than DRAM read
- ❖ NVRAM write throughput further 4x lower

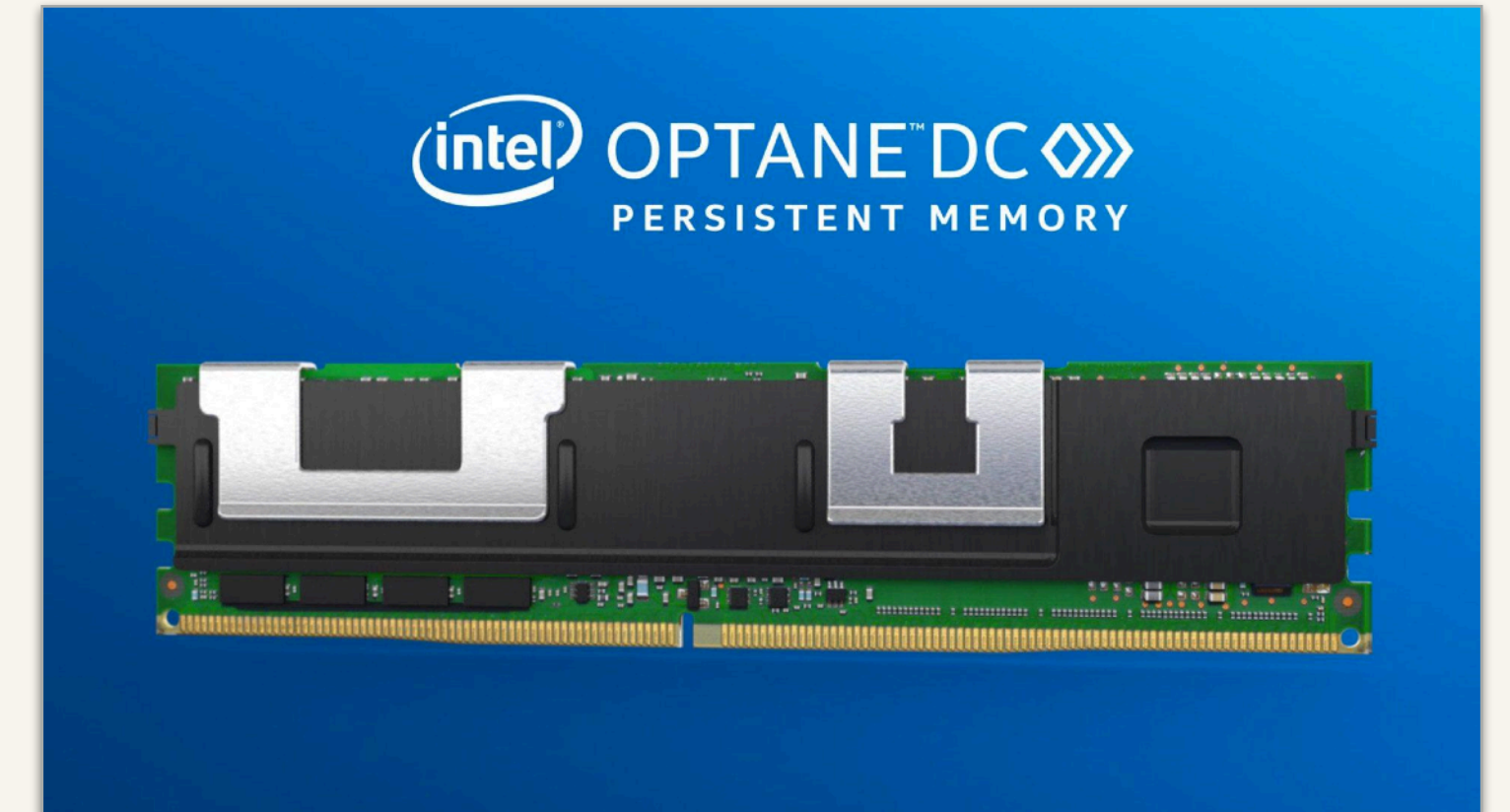
# Recent work on Asymmetry

## Benchmarking

- ❖ Two recent studies by Izraelevitz et al. [0] and van Renen et al. [1] perform careful benchmarking of Optane memory, and report similar asymmetries

## Algorithms and Systems for Asymmetric Settings

- ❖ Recent work explores how to minimize the number of NVRAM writes, e.g., [2 – 4], including many other papers
- ❖ Also significant work from systems, architecture, and database communities, e.g., [5 – 7], amongst many other papers



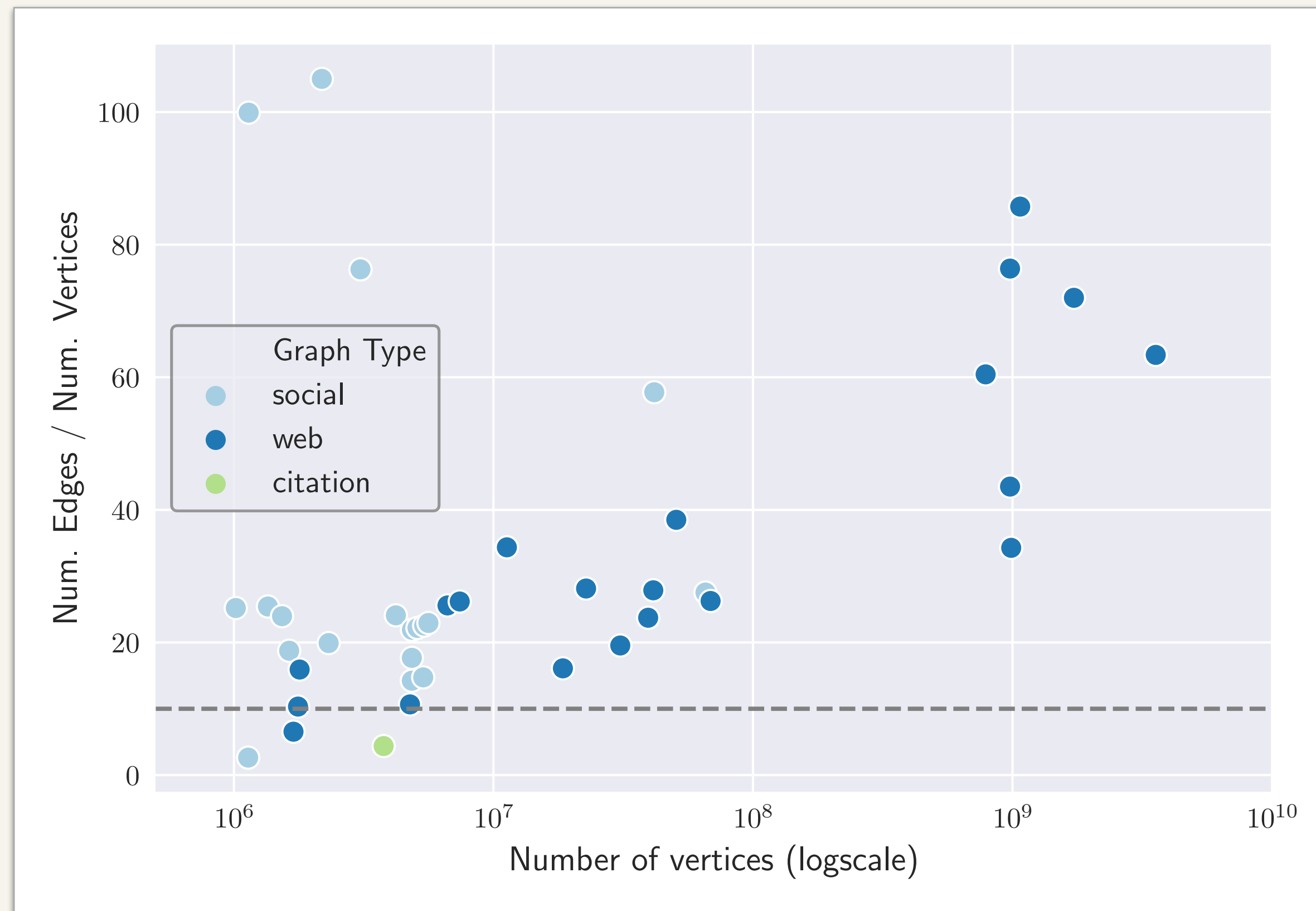
Sources:

- [0] Izraelevitz et al. [Basic performance measurements of the Intel Optane DC persistent memory module.](#) (2019)
- [1] van Renen et al. [Persistent Memory I/O Primitives](#) (2019)
- [2] Ben-David et al. [Parallel algorithms for asymmetric read-write costs](#) (2016)
- [3] Blelloch et al. [Efficient algorithms with asymmetric read and write costs](#) (2016)
- [4] Carson et al. [Write-avoiding algorithms](#) (2016)
- [5] Peng et al. [System Evaluation of the Intel Optane byte-addressable NVM](#) (2019)
- [6] Ni et al. [SSP: Eliminating Redundant Writes in Failure-Atomic NVRAMs via Shadow Sub-Paging](#) (2019)
- [7] Yang et al. [An Empirical Guide to the Behavior and Use of Scalable Persistent Memory](#) (2020)

# Semi-Asymmetric Parallel Graph Algorithms for NVRAMs [DMKGBGS'20]

Can we design practical and theoretically-sound techniques to overcome read/write asymmetry for graph problems on NVRAMs?

# Real World Graphs are not Ultra-Sparse



Over 90% of graphs with  $> 1M$  vertices from SNAP and LAW datasets have  $m/n \geq 10$

We expect that ratio of NVRAM/DRAM in future systems will be similar (our ratio is 8x)

Sources:

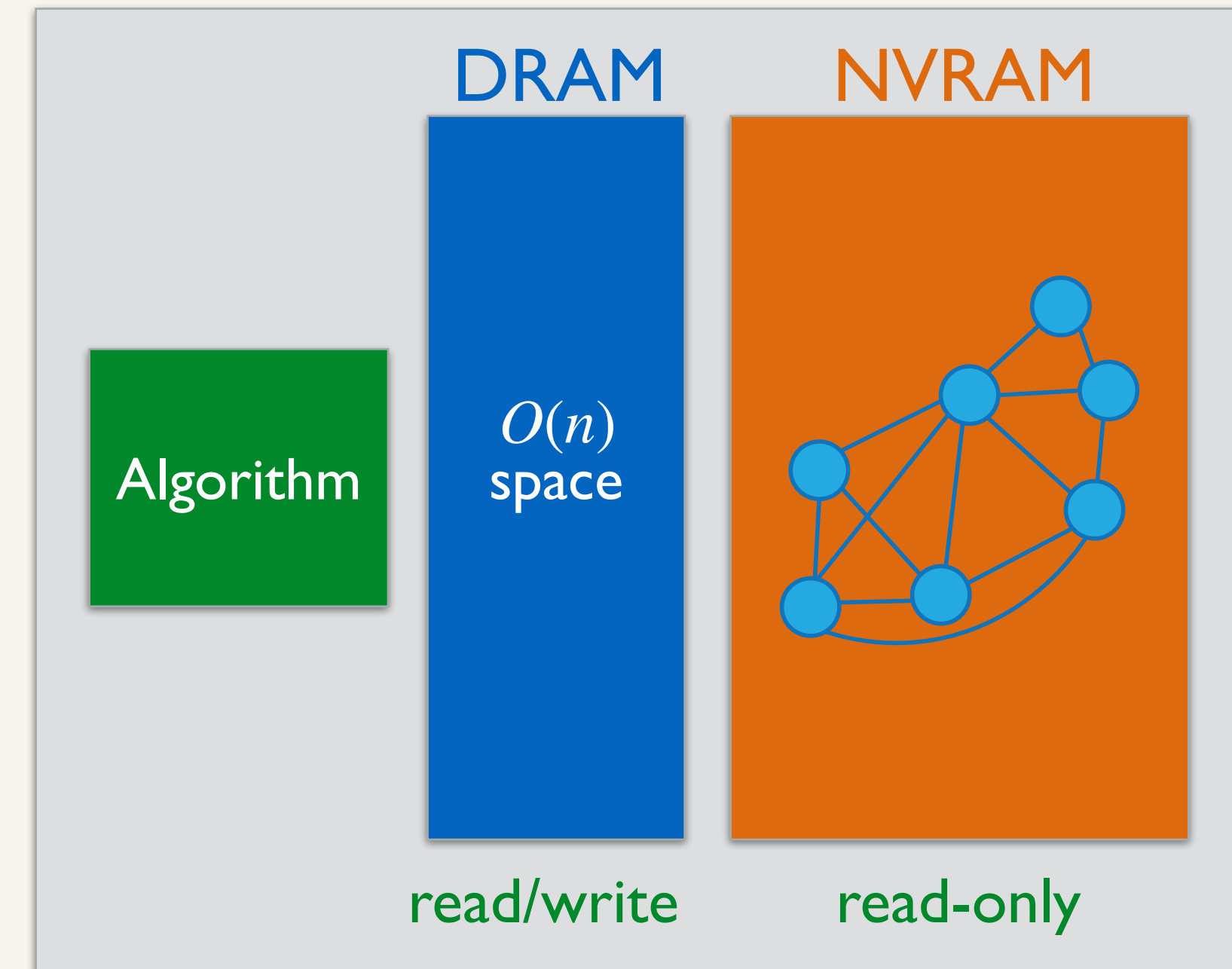
<https://snap.stanford.edu/data/>

<http://law.di.unimi.it/datasets.php>

# Our Approach

## Semi-Asymmetric Approach

- ❖ Graph stored in NVRAM and accessed in a *read-only mode*
- ❖ Amount of DRAM is proportional to the number of vertices





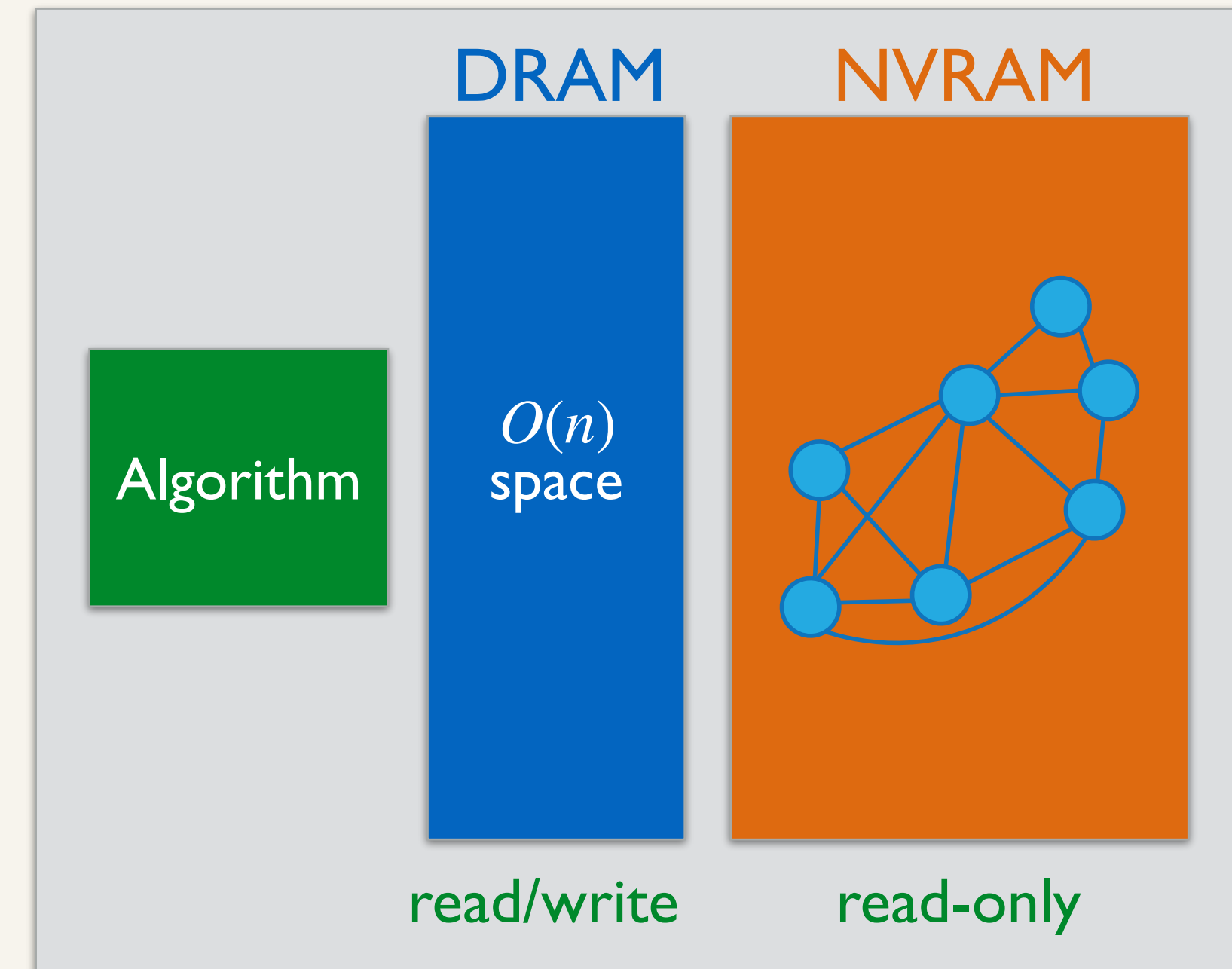
# Our Approach

## Semi-Asymmetric Approach

- ❖ Graph stored in NVRAM and accessed in a *read-only mode*
- ❖ Amount of DRAM is proportional to the number of vertices

## Benefits

- ❖ Algorithms avoid costly NVRAM writes, and algorithm design is independent of this cost
- ❖ Algorithms do not contribute to NVRAM wear-out



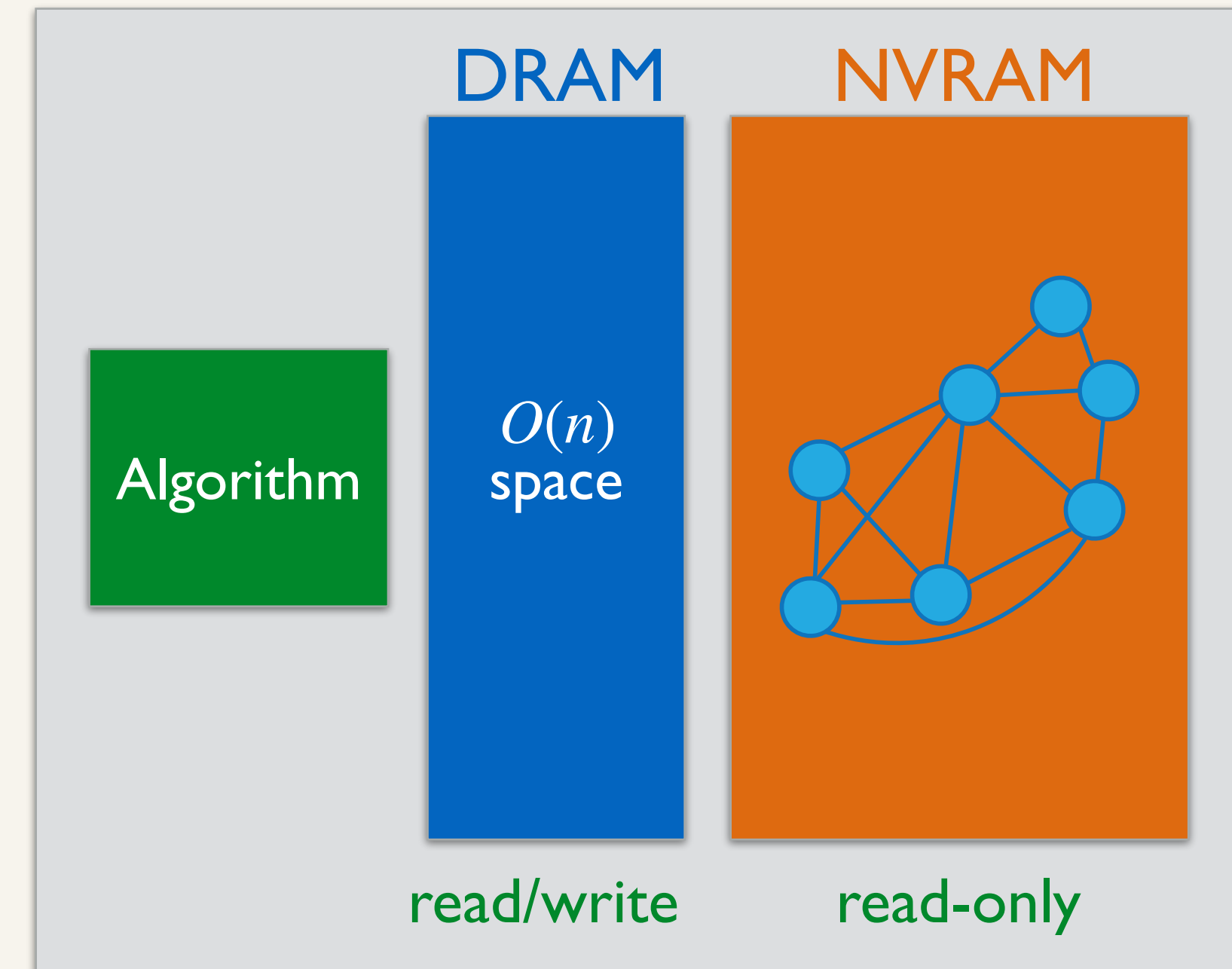
# Our Approach

## Semi-Asymmetric Approach

- ❖ Graph stored in NVRAM and accessed in a *read-only mode*
- ❖ Amount of DRAM is proportional to the number of vertices

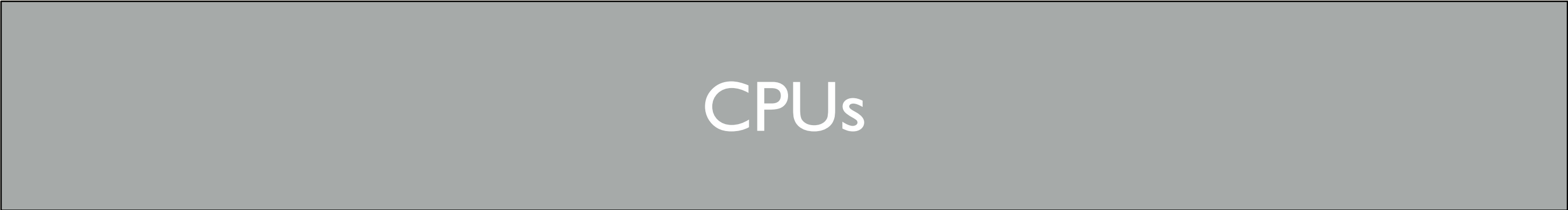
## Benefits

- ❖ Algorithms avoid costly NVRAM writes, and algorithm design is independent of this cost
- ❖ Algorithms do not contribute to NVRAM wear-out

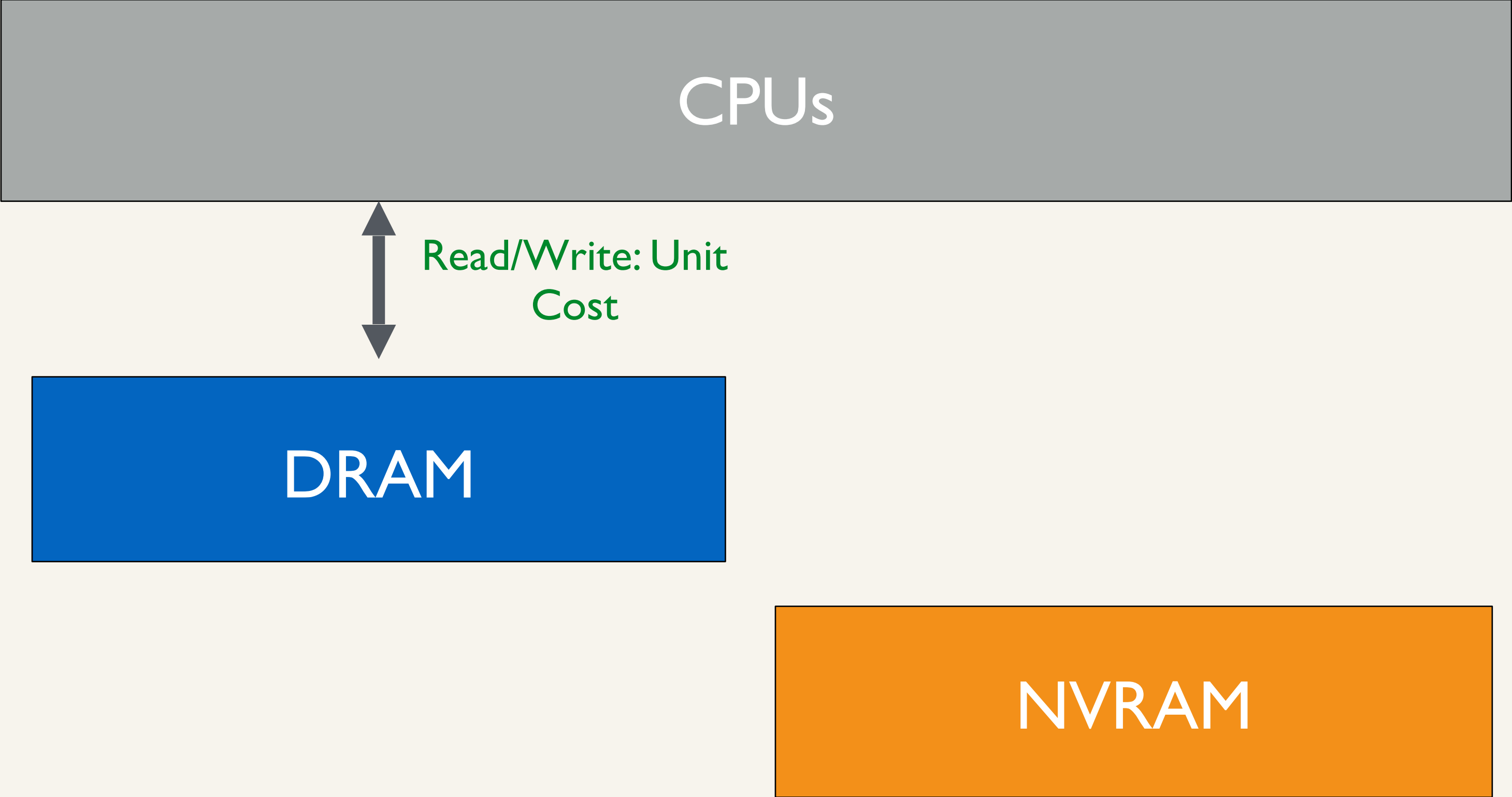


*Our contribution:*  
*This (restrictive) semi-asymmetric approach is effective for designing fast parallel graph algorithms*

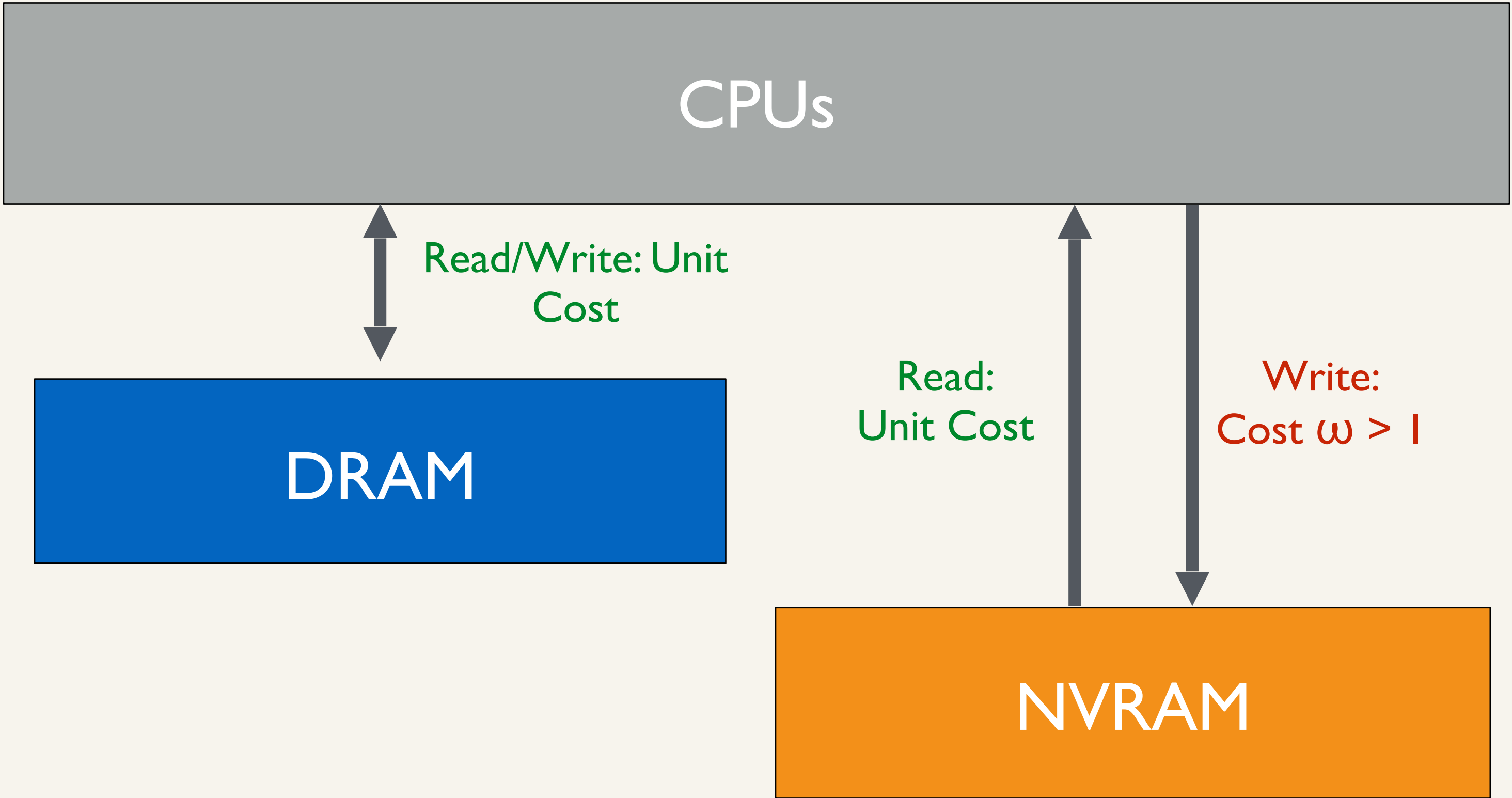
# Parallel Semi-Asymmetric Model (PSAM)



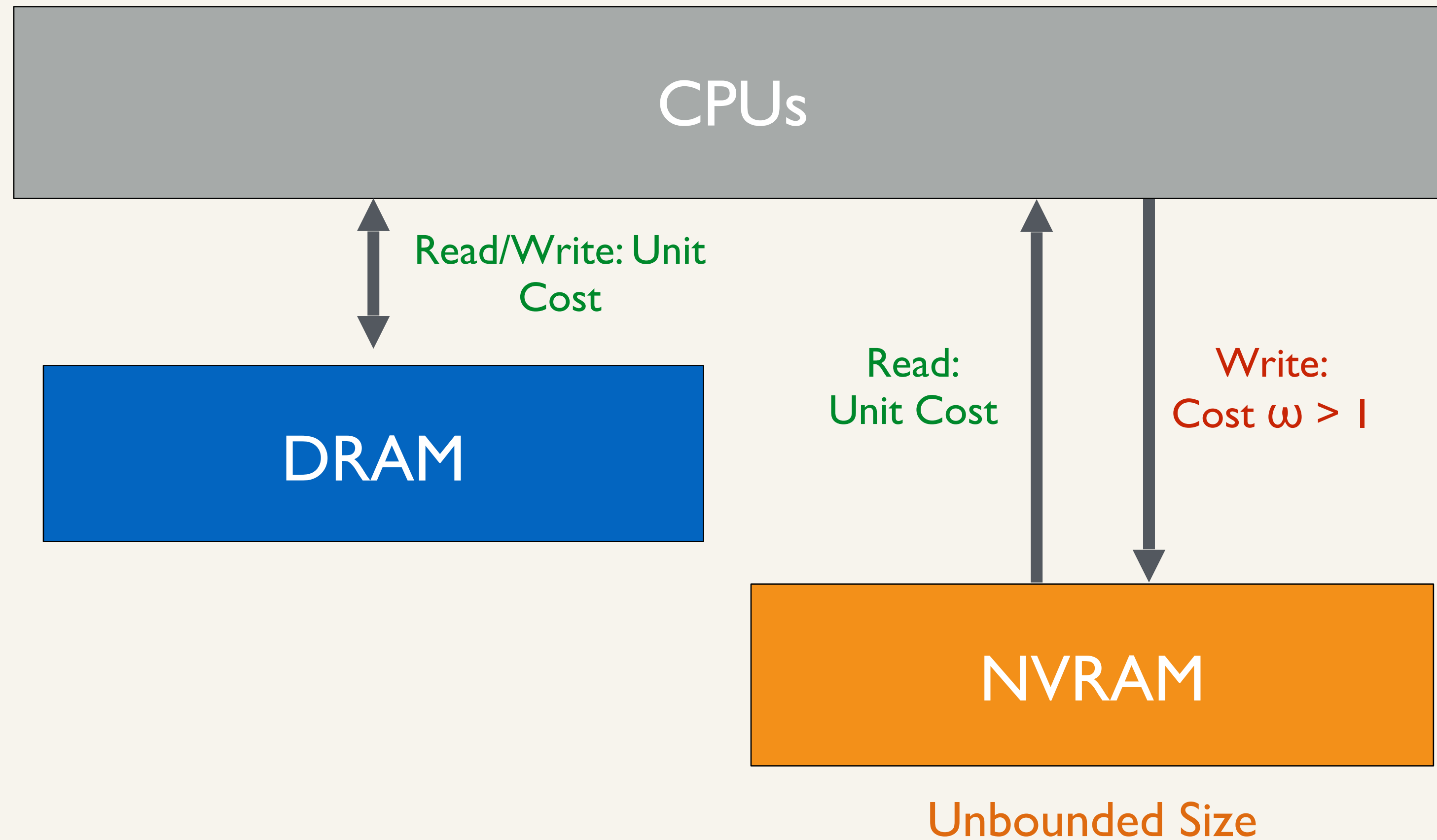
# Parallel Semi-Asymmetric Model (PSAM)



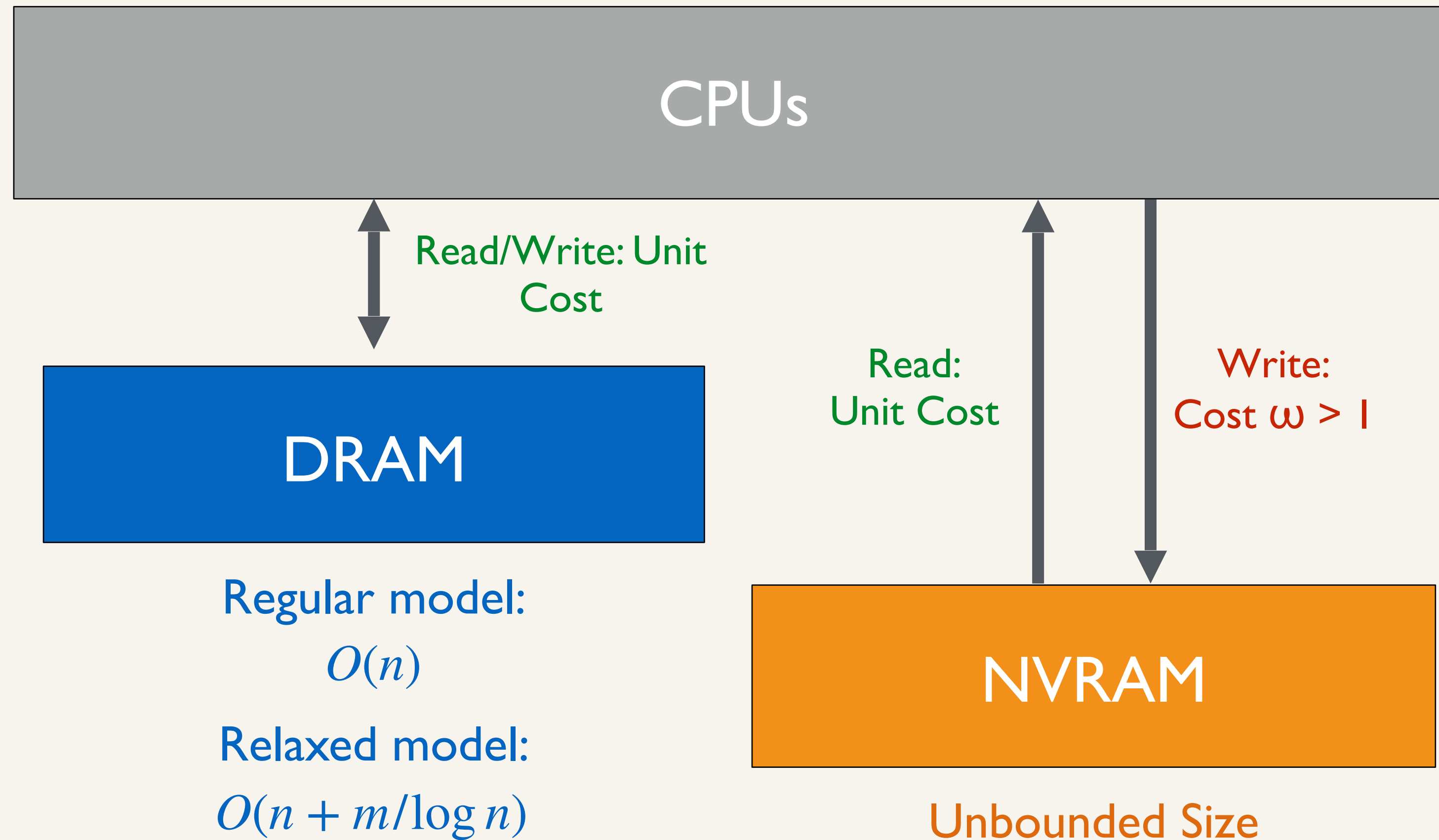
# Parallel Semi-Asymmetric Model (PSAM)



# Parallel Semi-Asymmetric Model (PSAM)



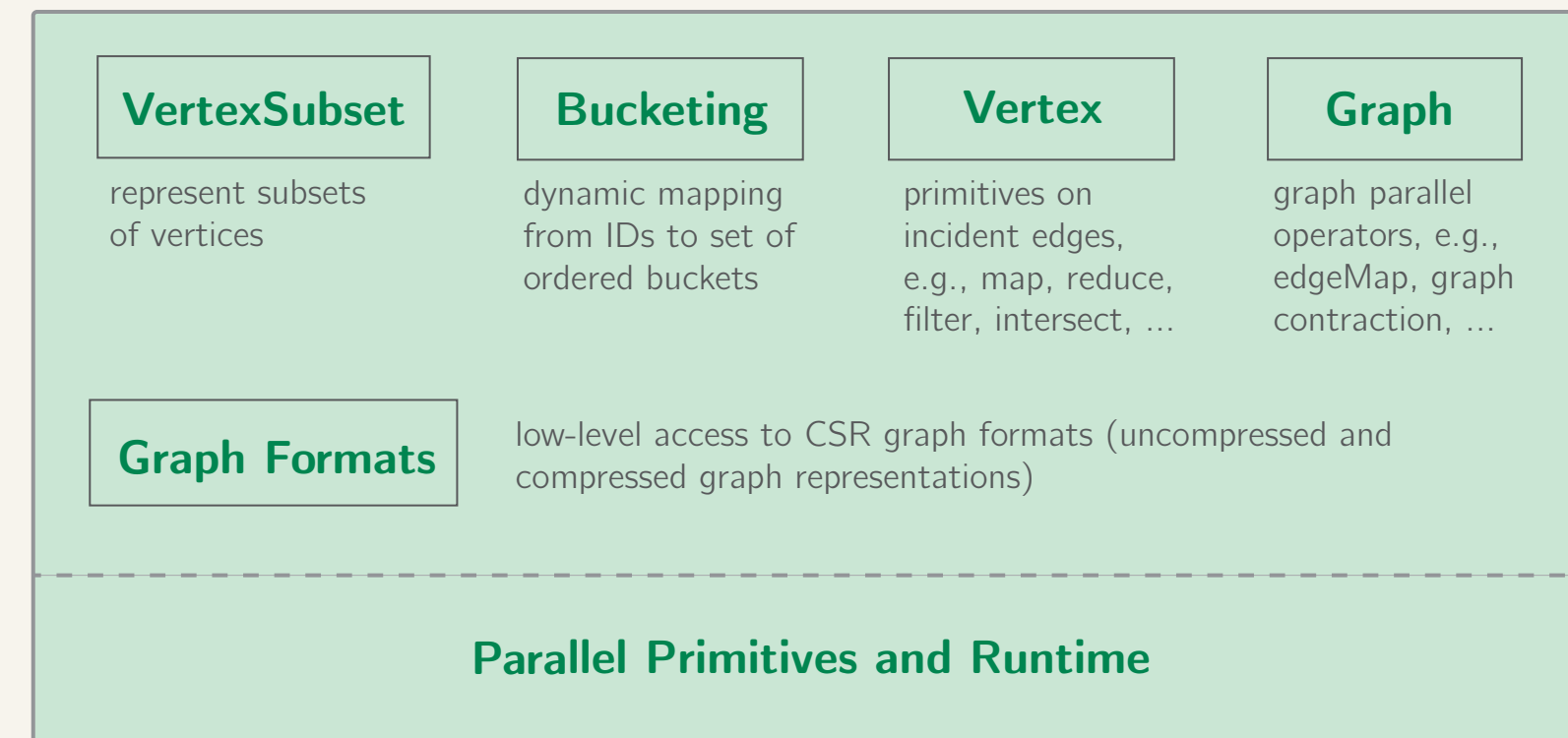
# Parallel Semi-Asymmetric Model (PSAM)



# Overview of Semi-Asymmetric Algorithms

- ❖ Start with work-efficient shared-memory algorithms from the Graph Based Benchmark Suite (GBBS)
- ❖ Implement interface primitives used by GBBS algorithms (*edgeMap* and *filtering*) efficiently in the PSAM

## GBBS Interface





# Overview of Semi-Asymmetric Algorithms

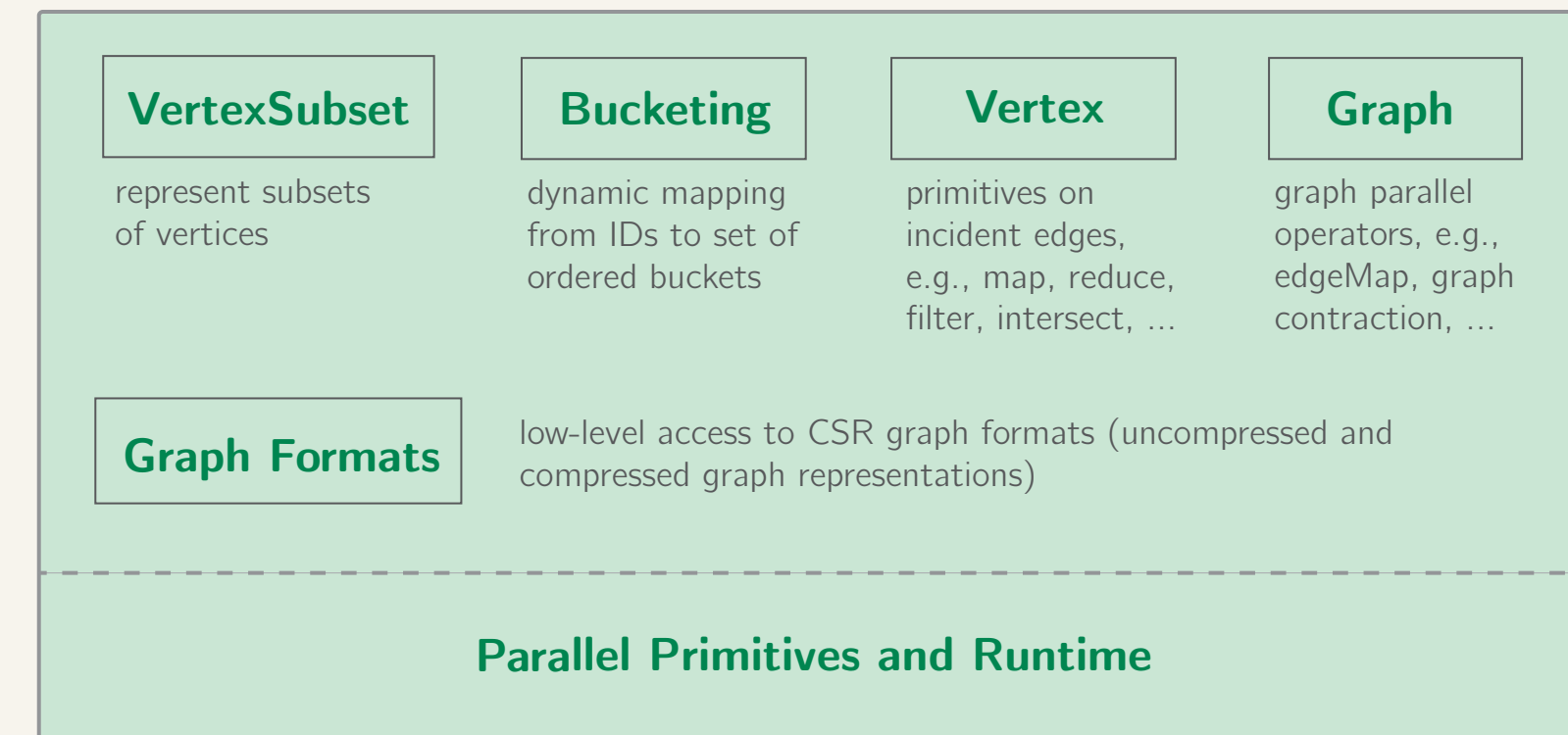
- ❖ Start with work-efficient shared-memory algorithms from the Graph Based Benchmark Suite (GBBS)
- ❖ Implement interface primitives used by GBBS algorithms (*edgeMap* and *filtering*) efficiently in the PSAM

## edgeMap

	Problem	GBBS Work	Sage Work	Sage Depth
EDGE MAP CHUNKED	<b>Breadth-First Search</b>	$O(\omega m)$	$O(m)$	$O(d_G \log n)$
	<b>Weighted BFS</b>	$O(\omega m)^*$	$O(m)^*$	$O(d_G \log n)^\ddagger$
	<b>Bellman-Ford</b>	$O(\omega d_G m)$	$O(d_G m)$	$O(d_G \log n)$
	<b>Single-Source Widest Path</b>	$O(\omega d_G m)$	$O(d_G m)$	$O(d_G \log n)$
	<b>Single-Source Betweenness</b>	$O(\omega m)$	$O(m)$	$O(d_G \log n)$
	<b><math>O(k)</math>-Spanner</b>	$O(\omega m)^*$	$O(m)^*$	$O(k \log n)^\ddagger$
	<b>LDD</b>	$O(\omega m)^*$	$O(m)^*$	$O(\log^2 n)^\ddagger$
	<b>Connectivity</b>	$O(\omega m)^*$	$O(m)^*$	$O(\log^3 n)^\ddagger$
	<b>Spanning Forest</b>	$O(\omega m)^*$	$O(m)^*$	$O(\log^3 n)^\ddagger$
	<b>Graph Coloring</b>	$O(\omega m)^*$	$O(m)^*$	$O(\log n + L \log \Delta)^*$
	<b>Maximal Independent Set</b>	$O(\omega m)^*$	$O(m)^*$	$O(\log^2 n)^\ddagger$

GBBS work indicates the work of naively converting existing shared-memory algorithms from GBBS to NVRAM algorithms

## GBBS Interface



## Filtering (relaxed model)

	Problem	GBBS Work	Sage Work	Sage Depth
Both	<b>Biconnectivity</b> <sup>†</sup>	$O(\omega m)^*$	$O(m)^*$	$O(d_G \log n + \log^3 n)^\ddagger$
	<b>Apx. Set Cover</b> <sup>†</sup>	$O(\omega m)^*$	$O(m)^*$	$O(\log^3 n)^\ddagger$
Filter	<b>Triangle Counting</b> <sup>†</sup>	$O(\omega(m+n) + m^{3/2})$	$O(m^{3/2})$	$O(\log n)$
	<b>Maximal Matching</b> <sup>†</sup>	$O(\omega m)^*$	$O(m)^*$	$O(\log^3 m)^\ddagger$

## Other Techniques

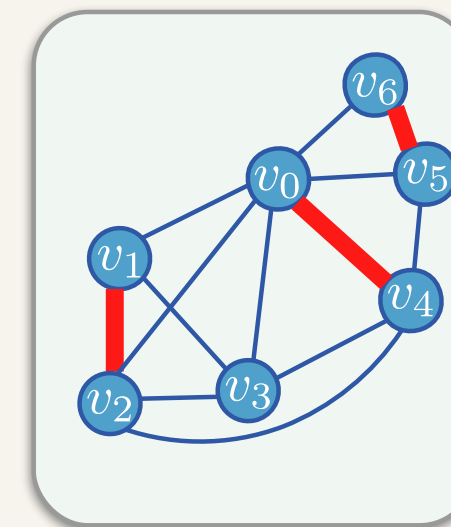
<b>PageRank Iteration</b>	$O(m + \omega n)$	$O(m)$	$O(\log n)$
<b>PageRank</b>	$O(P_{it}(m + \omega n))$	$O(P_{it}m)$	$O(P_{it} \log n)$
<b><math>k</math>-core</b>	$O(\omega m)^*$	$O(m)^*$	$O(\rho \log n)^\ddagger$
<b>Apx. Densest Subgraph</b>	$O(\omega m)$	$O(m)$	$O(\log^2 n)$

# Semi-Asymmetric Filtering

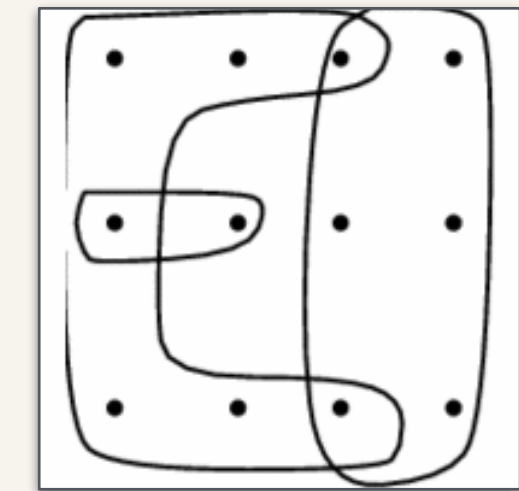
## Motivation

- ❖ Some algorithms *remove, or batch-delete edges* over the course of their operation for work-efficiency
- ❖ Modifying the graph directly requires writing to NVRAM

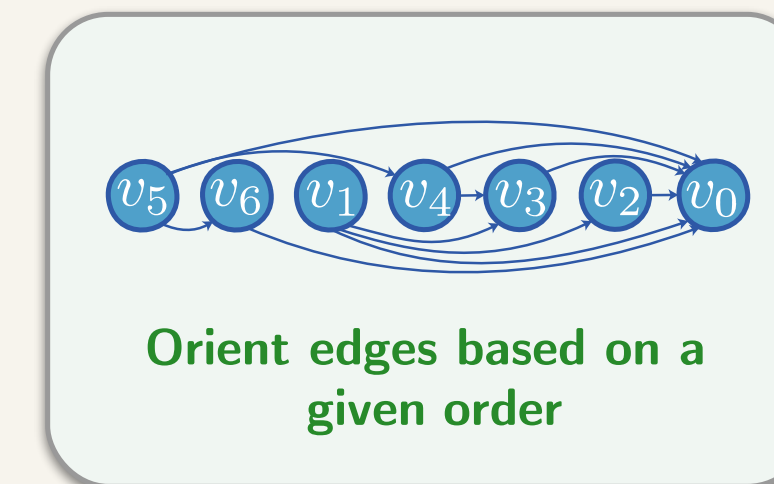
Maximal Matching



Parallel Approximate Set Cover



Triangle Counting



# Semi-Asymmetric Filtering

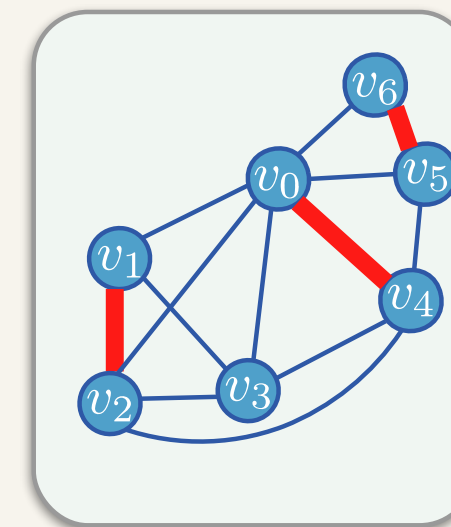
## Motivation

- ❖ Some algorithms *remove, or batch-delete edges* over the course of their operation for work-efficiency
- ❖ Modifying the graph directly requires writing to NVRAM

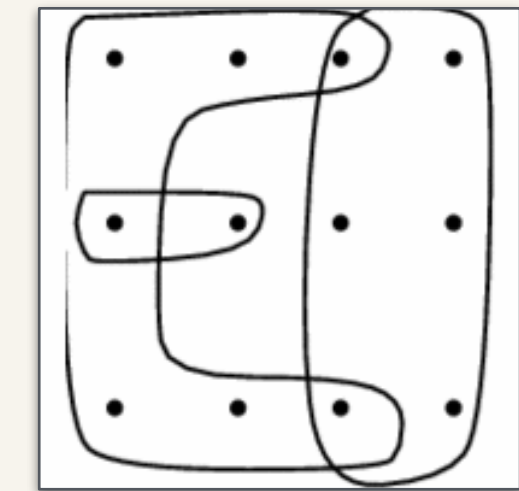
## Semi-Asymmetric Filtering

- ❖ Work in the relaxed model
- ❖ Use one bit per edge and mirror the CSR structure (in NVRAM) using a blocked approach in DRAM

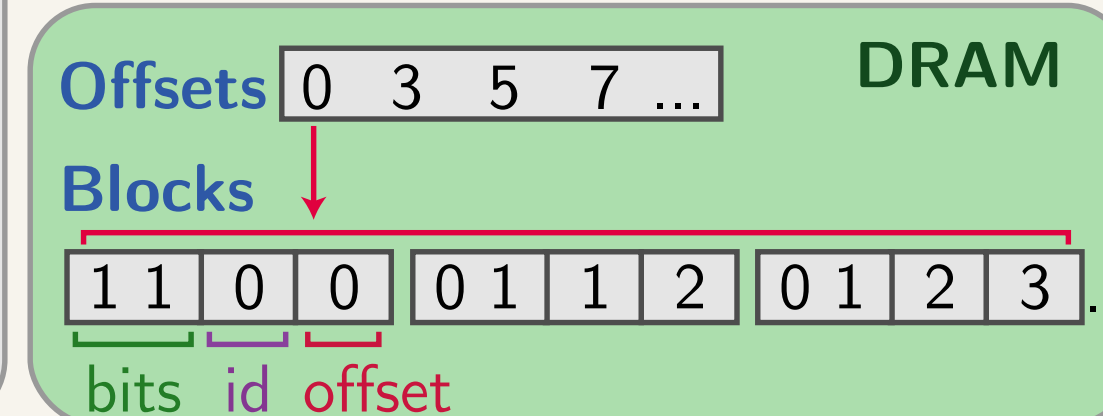
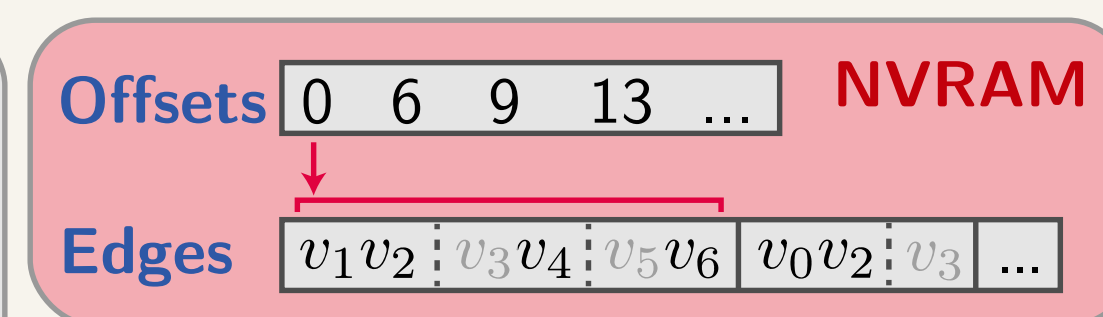
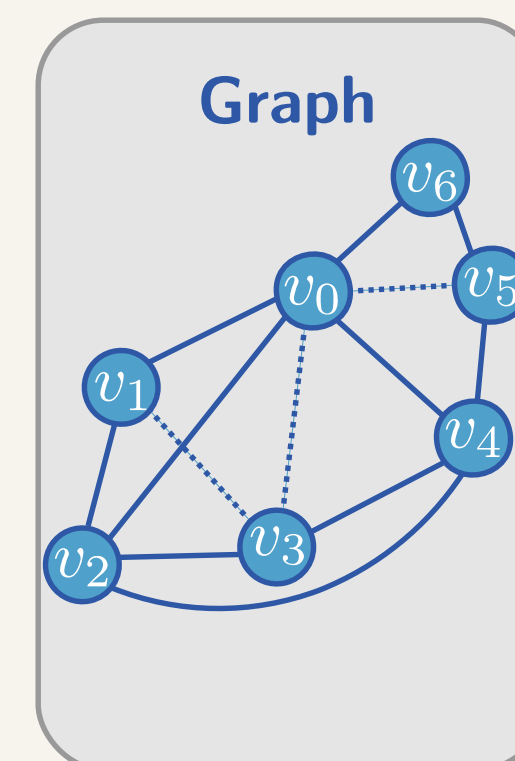
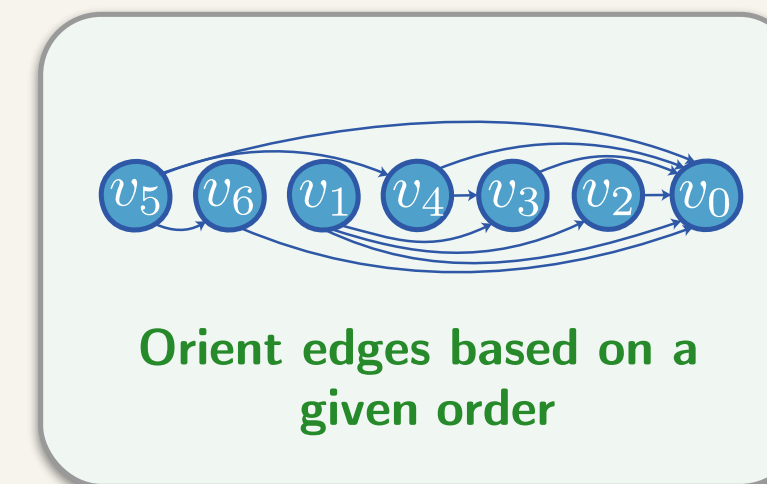
Maximal Matching



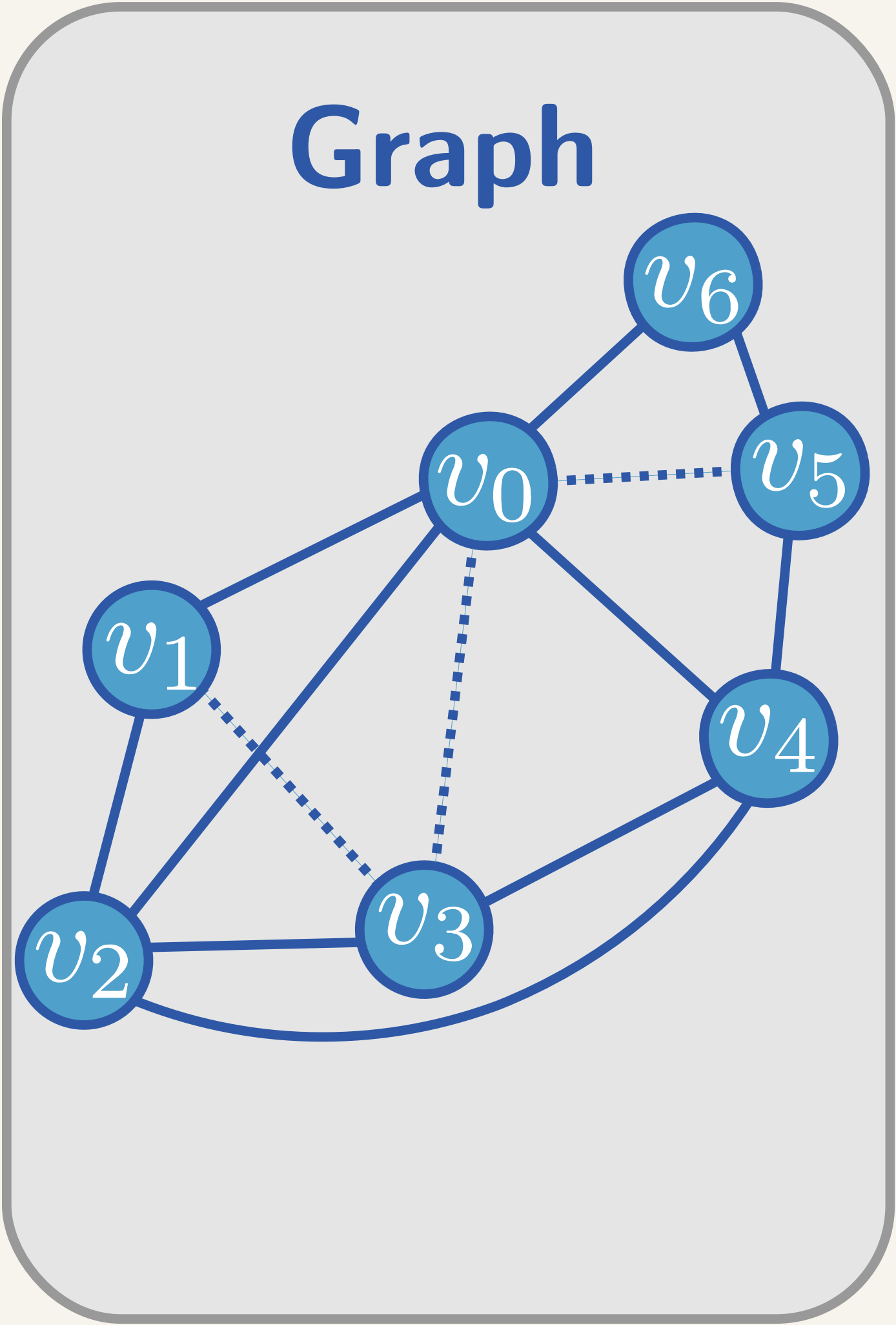
Parallel Approximate Set Cover



Triangle Counting



# Semi-Asymmetric Filtering

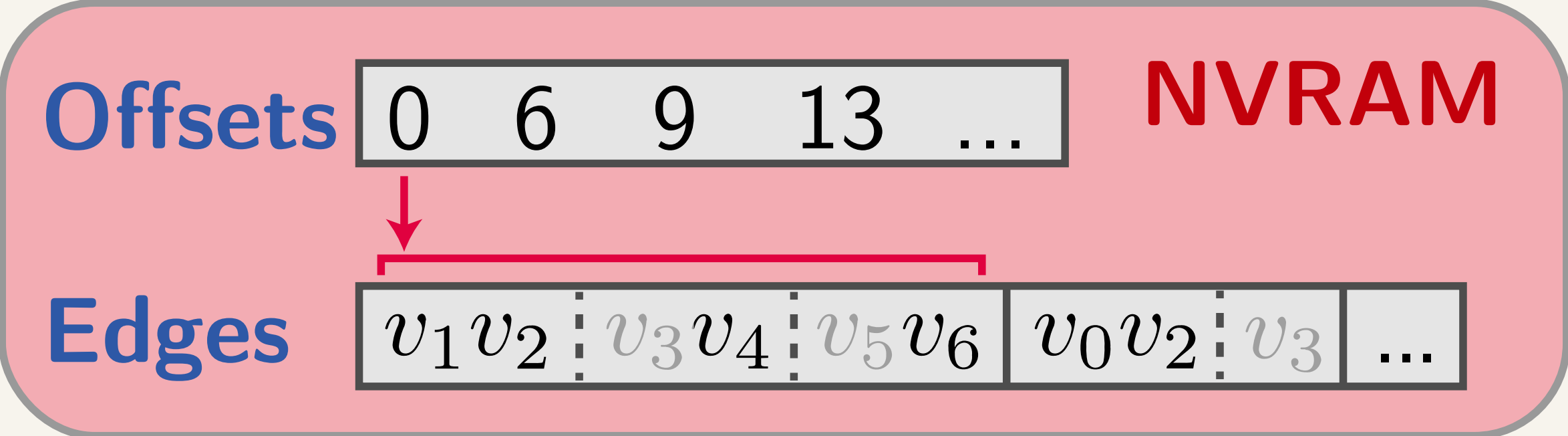


..... logically deleted  
—— present in graph

# Semi-Asymmetric Filtering

## High-level Approach

- (i) Set a filter block size, and logically chunk the CSR structure into chunks of this size

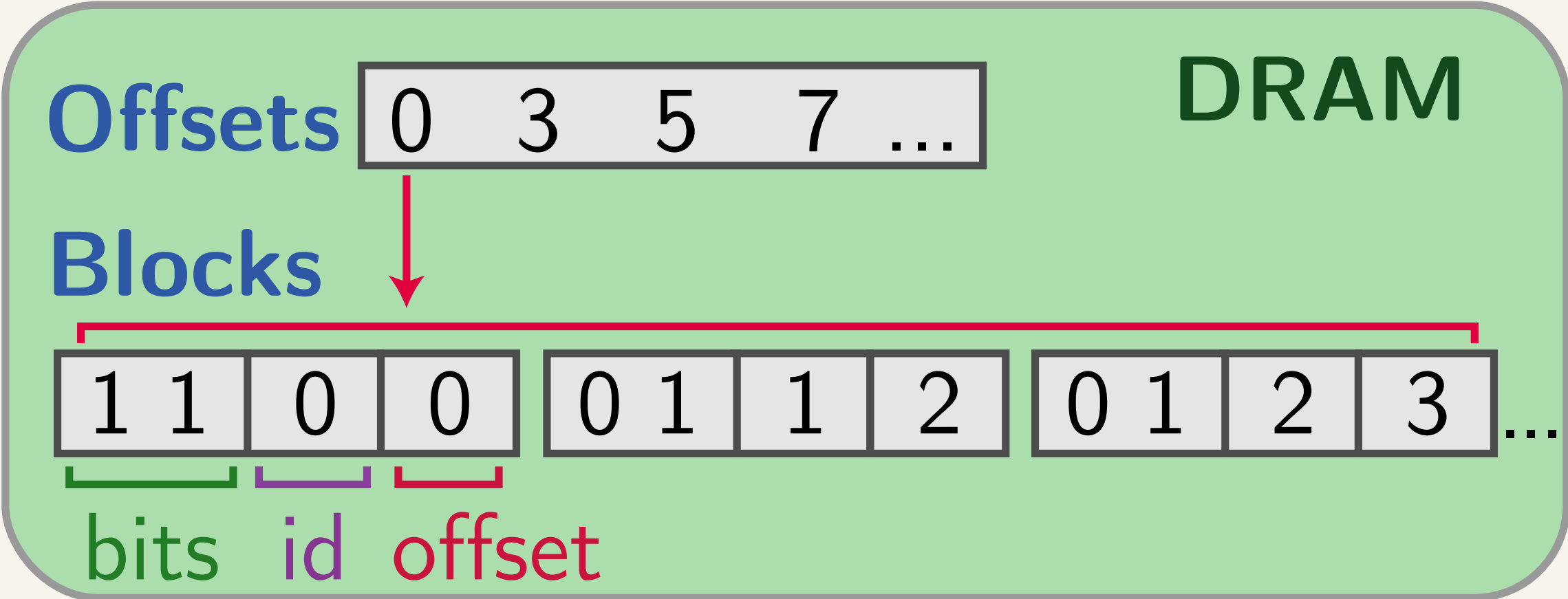


Graph in CSR format, stored in NVRAM ( $\mathcal{F}_B = 2$ )

# Semi-Asymmetric Filtering

## High-level Approach

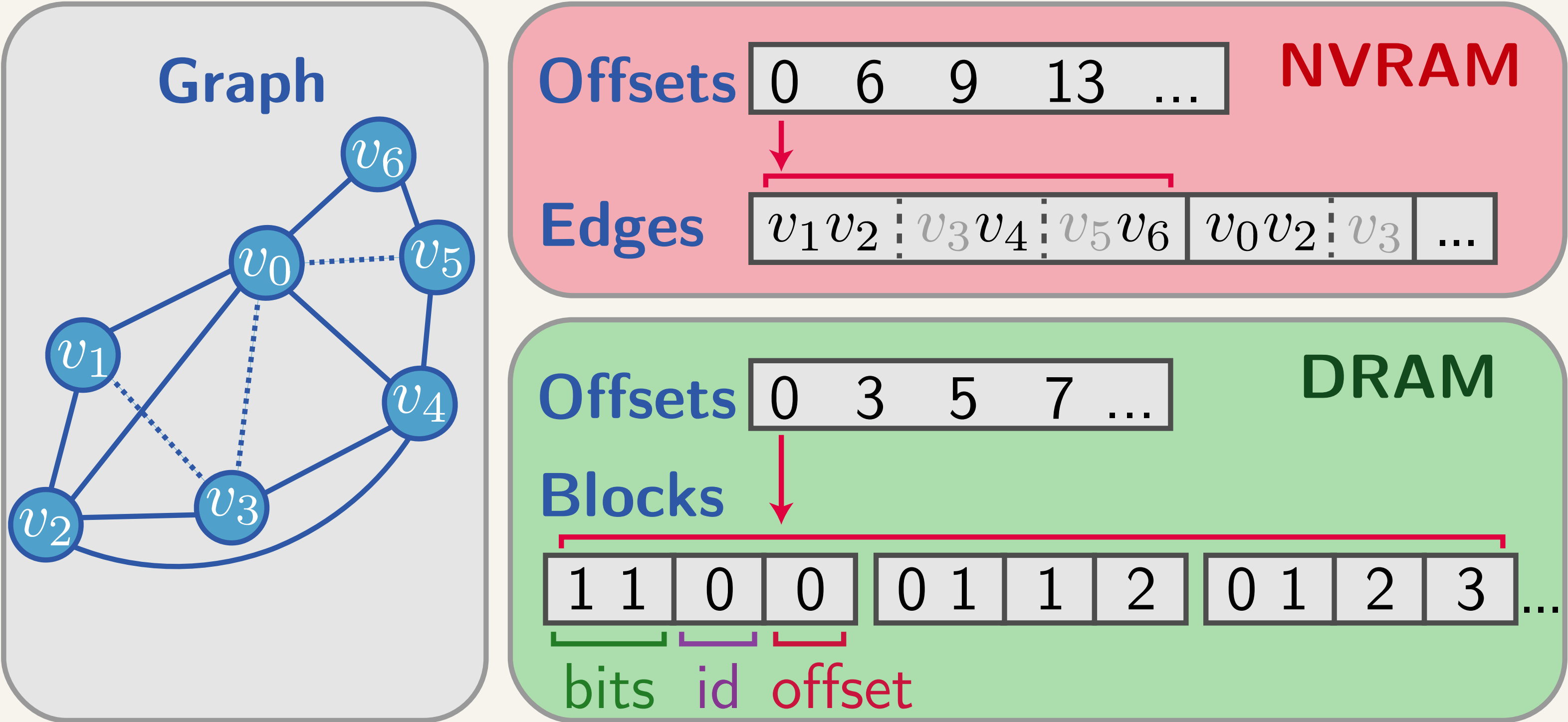
- (ii) Create a “mirrored” filter structure in DRAM, storing 1 bit per edge in NVRAM



GraphFilter in CSR format, stored in DRAM ( $\mathcal{F}_B = 2$ )

# Semi-Asymmetric Filtering

## Structure Overview



Note: Blocks with no "1" bits remaining are deleted

# Relationship to Other Models

## Semi-External Memory (SE) Model

- ❖ SE model performs block-transfers, with a focus on I/O cost [0, 1]
- ❖ Both PSAM and SE models provide the same amount of DRAM, but SE does not account for DRAM reads and writes

Sources:

[0] Abello et al. [A Functional Approach to External Graph Algorithms](#) (2002)

[1] Zheng et al. [FlashGraph: Processing Billion-Node Graphs on an Array of Commodity SSDs](#) (2015)

[2] Blelloch et al. [Efficient algorithms with asymmetric read and write costs](#) (2016)

[3] Ben-David et al. [Parallel algorithms for asymmetric read-write costs](#) (2016)



# Relationship to Other Models

## Semi-External Memory (SE) Model

- ❖ SE model performs block-transfers, with a focus on I/O cost [0, 1]
- ❖ Both PSAM and SE models provide the same amount of DRAM, but SE does not account for DRAM reads and writes

## Asymmetric RAM and Asymmetric Nested Parallel Models

- ❖ Both ARAM [2] and ANP [3] models capture asymmetry of writing to NVRAM
- ❖ Unlike ARAM/ANP models, the PSAM includes a fast memory, and is specialized for graph problems

Sources:

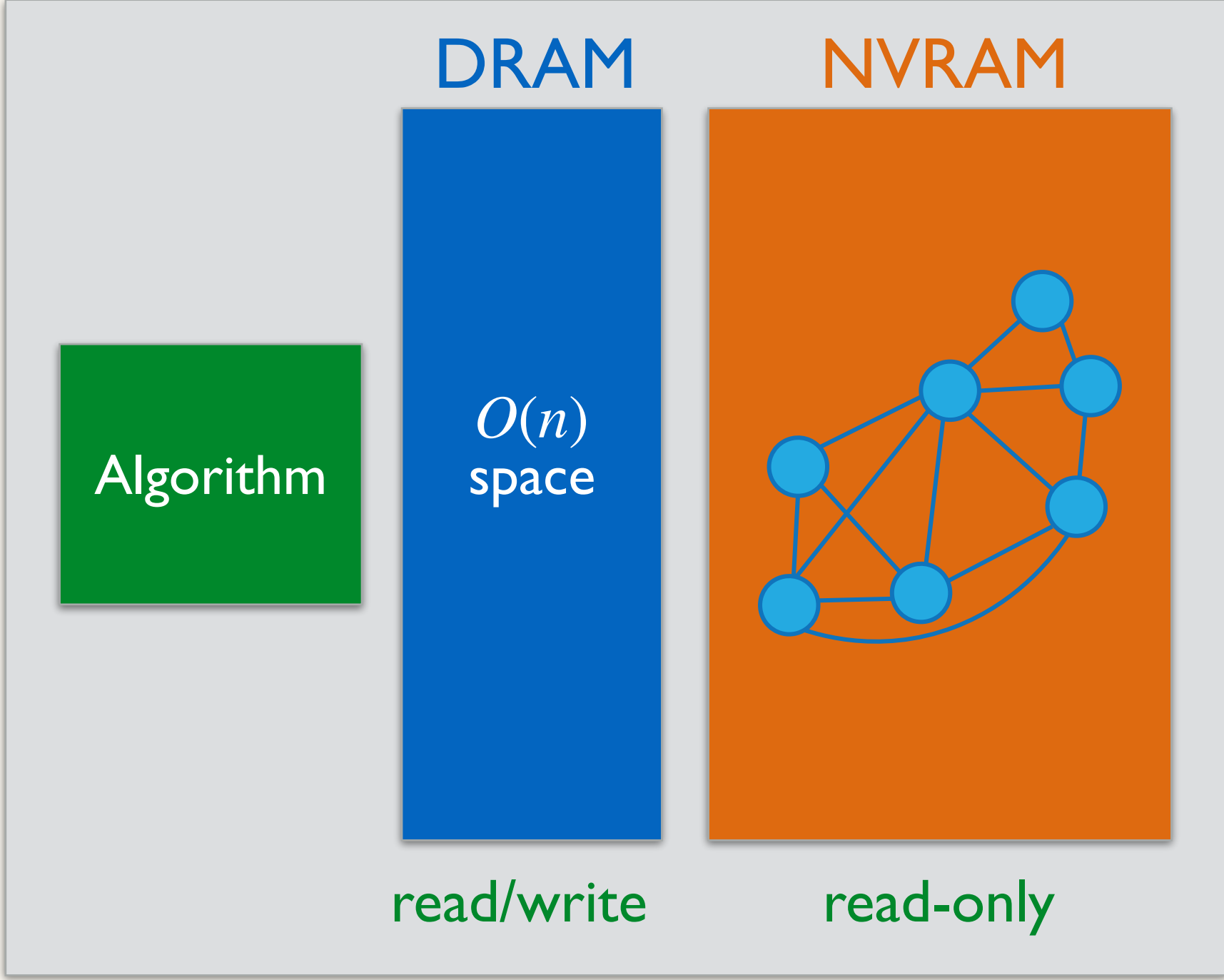
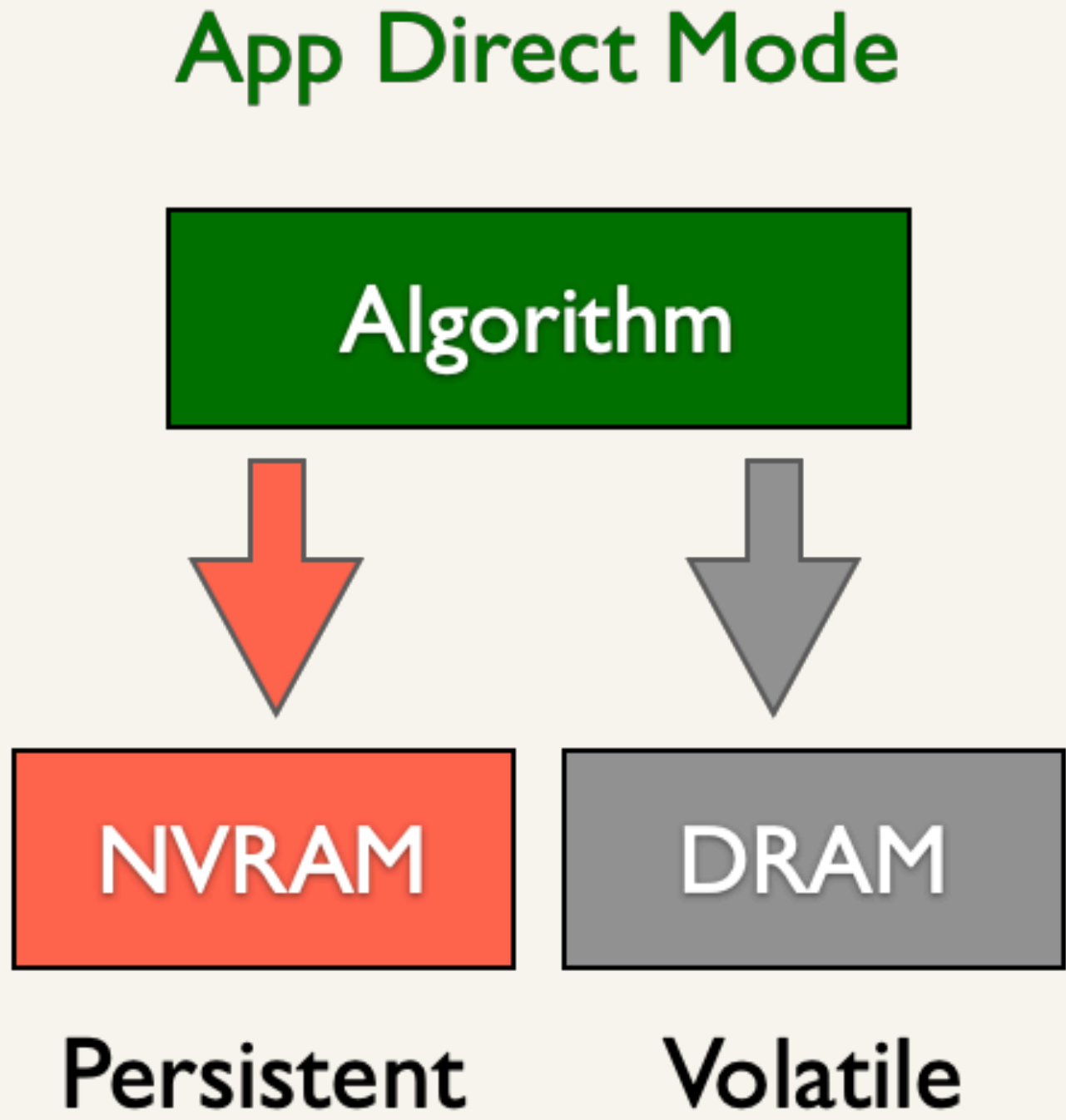
[0] Abello et al. [A Functional Approach to External Graph Algorithms](#) (2002)

[1] Zheng et al. [FlashGraph: Processing Billion-Node Graphs on an Array of Commodity SSDs](#) (2015)

[2] Blelloch et al. [Efficient algorithms with asymmetric read and write costs](#) (2016)

[3] Ben-David et al. [Parallel algorithms for asymmetric read-write costs](#) (2016)

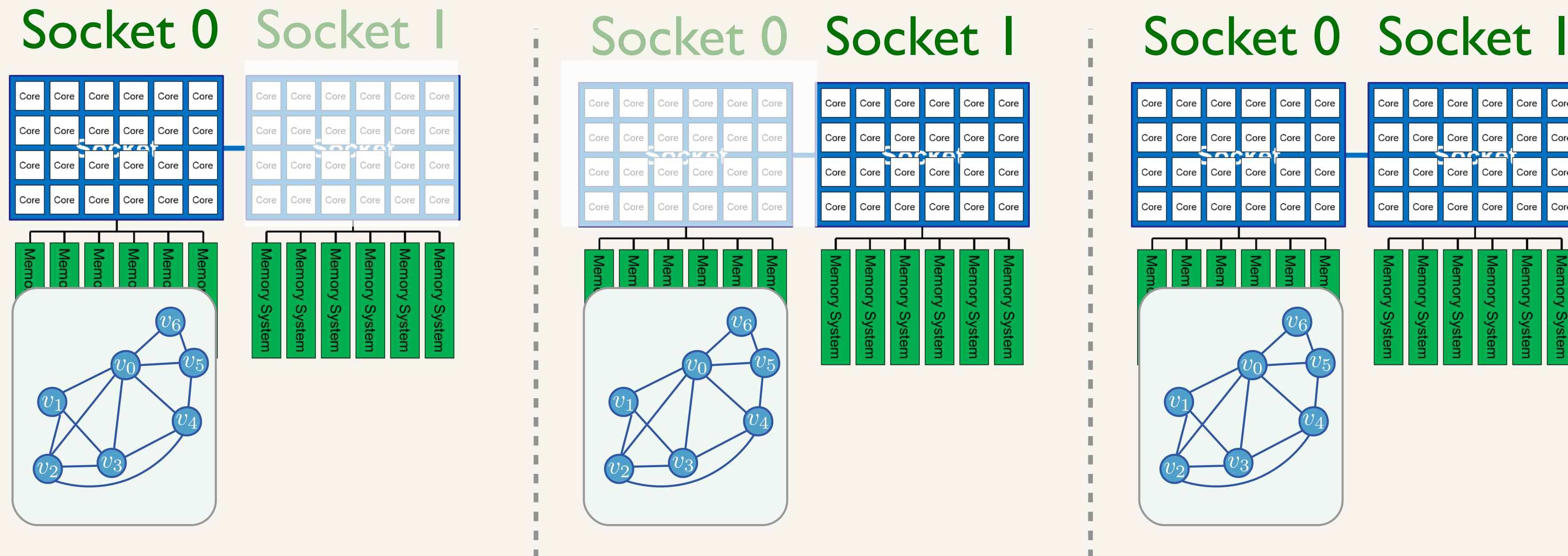
# Semi-Asymmetric Graph Engine (Sage) Approach



*AppDirect Mode enables a direct implementation of PSAM algorithms*

# NUMA Optimization in Sage

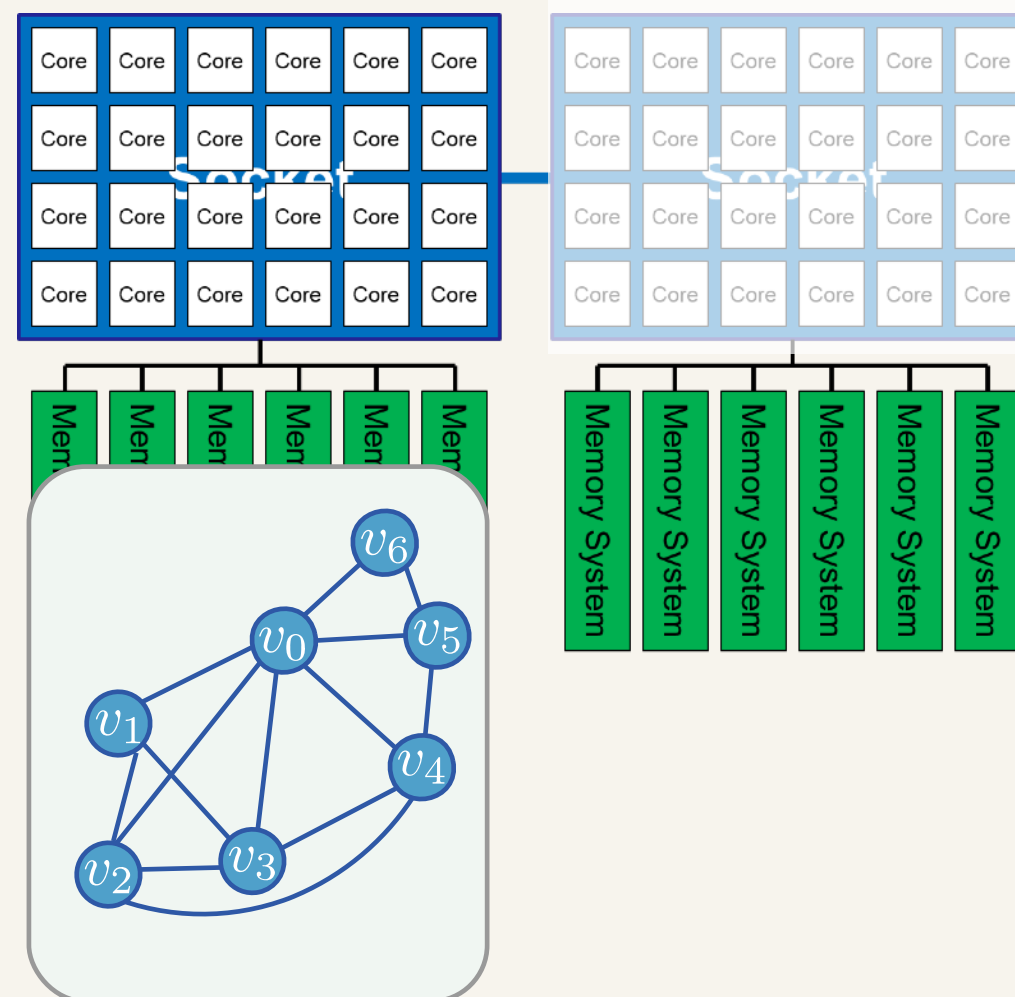
Consider an algorithm that maps over all vertices, and for each vertex performs a reduction over the neighbors of the vertex



Three experiments based on (threads, storage)

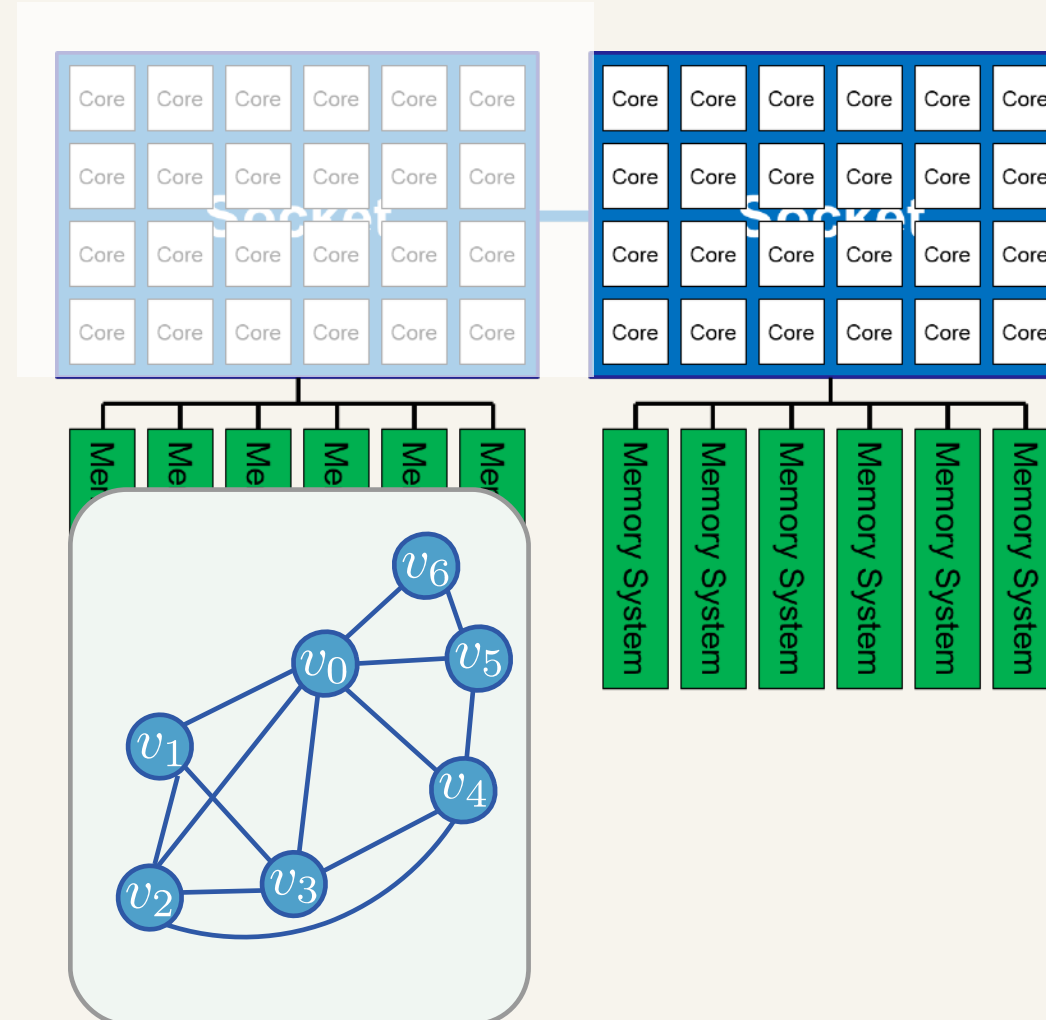
# NUMA Optimization in Sage

Socket 0 Socket 1



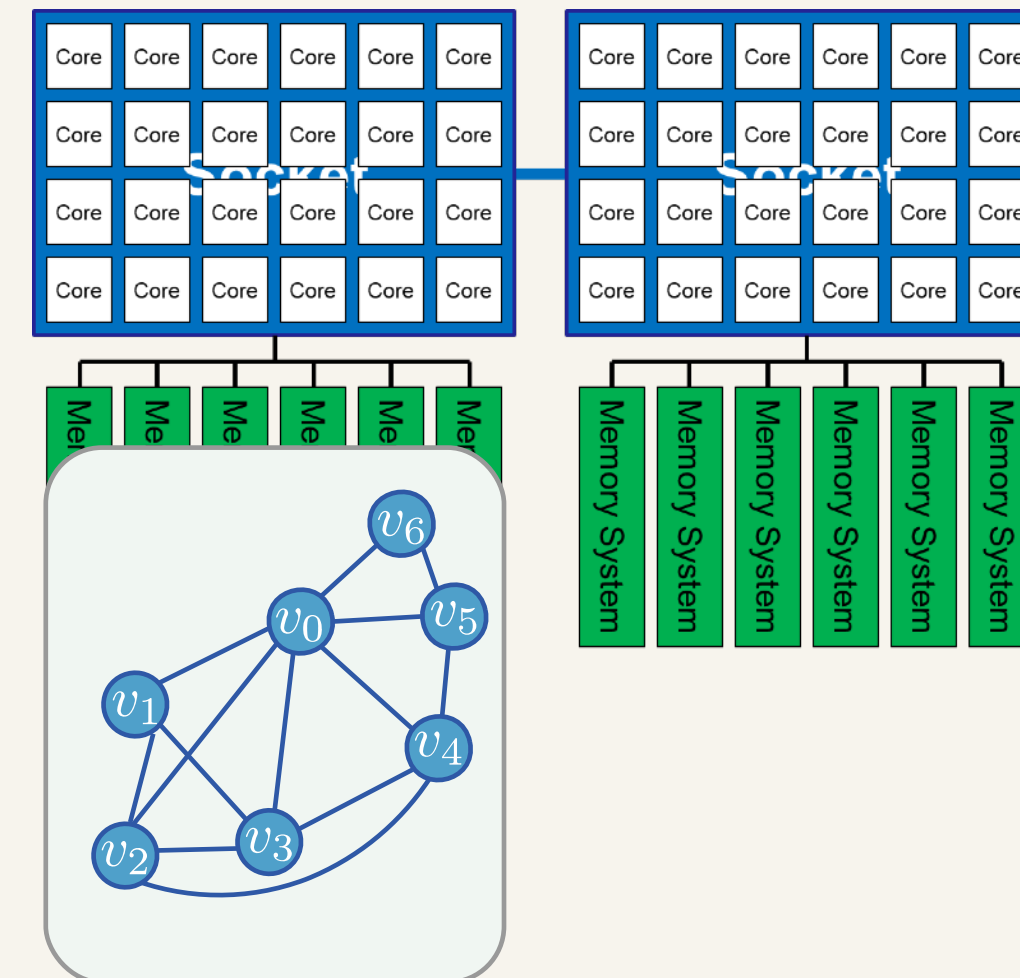
7 s

Socket 0 Socket 1



> 4x slower  
first run  
~7s subsequently

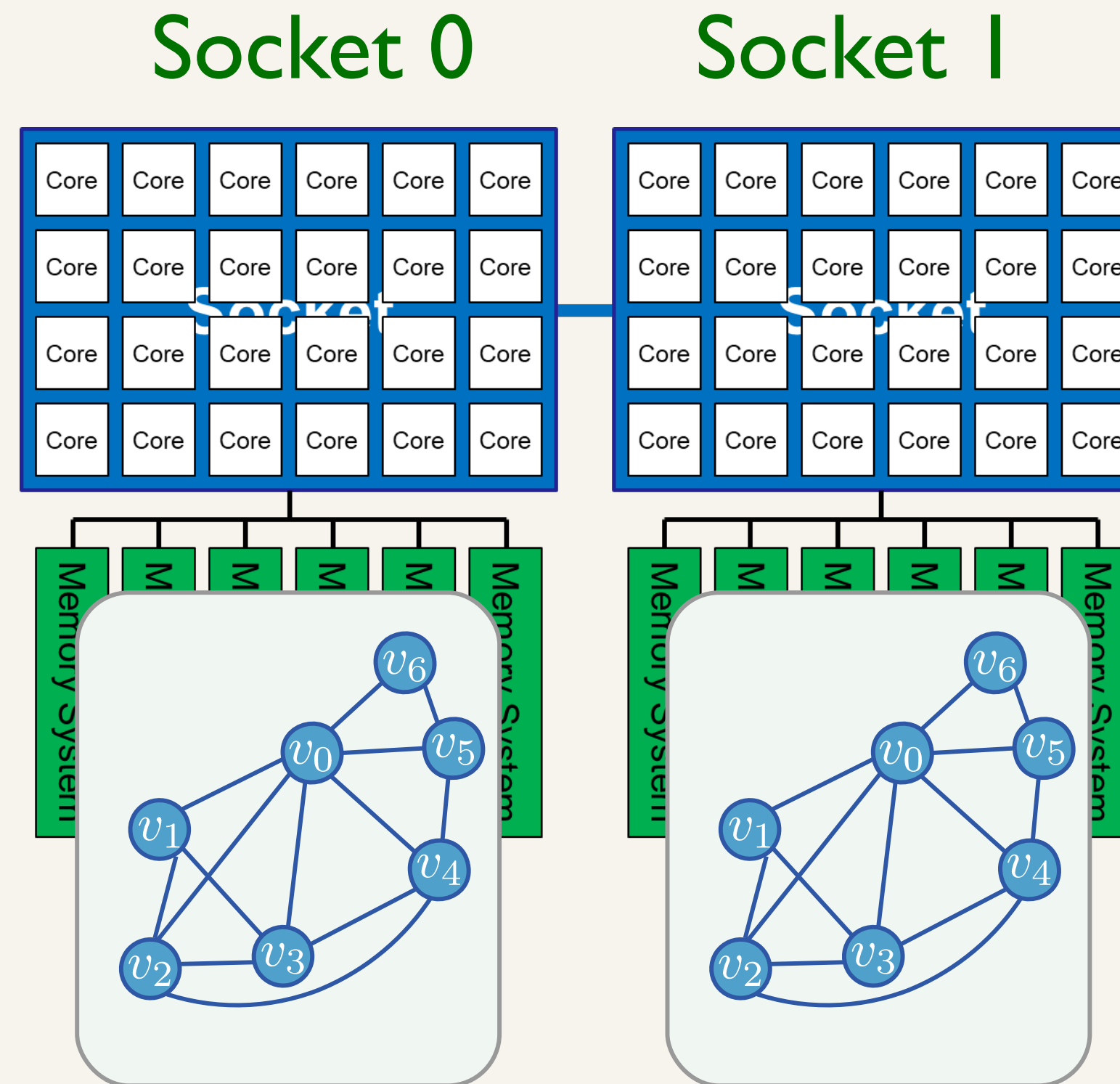
Socket 0 Socket 1



26 s

Cross-socket NVM reads should be avoided

# NUMA Optimization in Sage

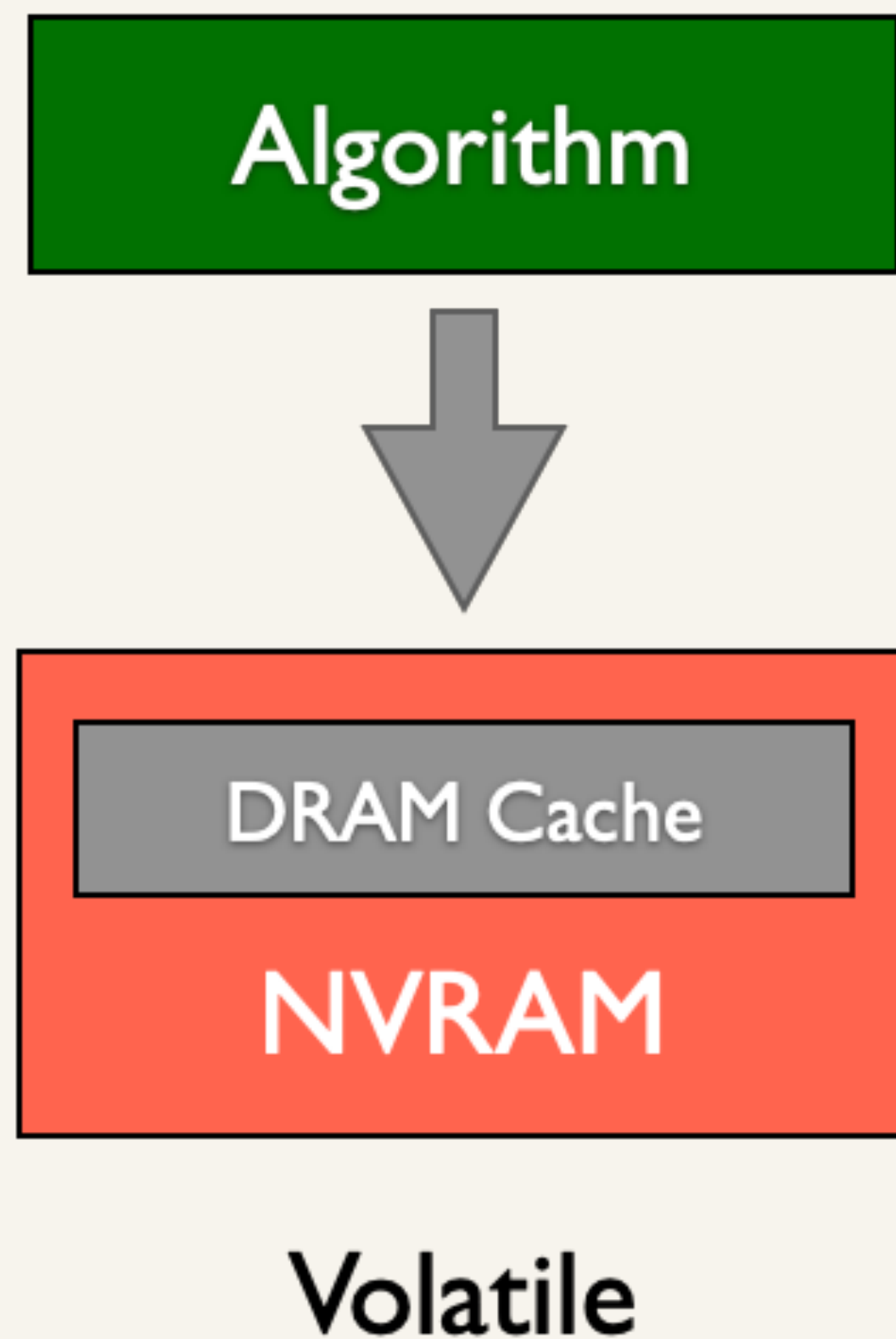


4.3 s for microbenchmark

Both graphs stored in compressed CSR format

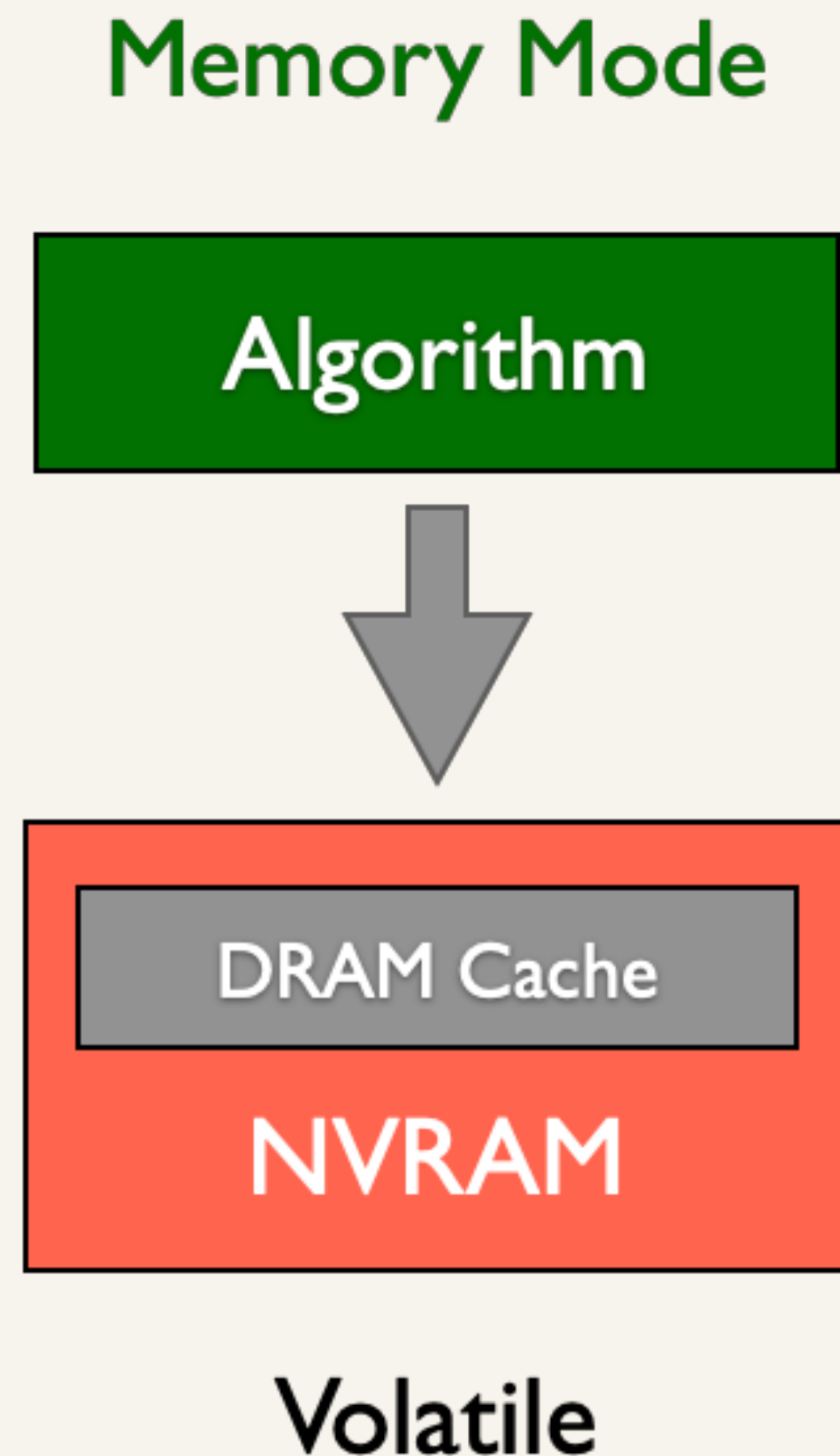
# Existing Approaches: DRAM as a Cache

## Memory Mode



- ❖ Applications do not distinguish between DRAM and NVRAM
- ❖ Existing shared-memory software does not require modification
- ❖ Workloads that are larger than DRAM can involve **costly NVRAM writes**

# Existing Approaches: DRAM as a Cache

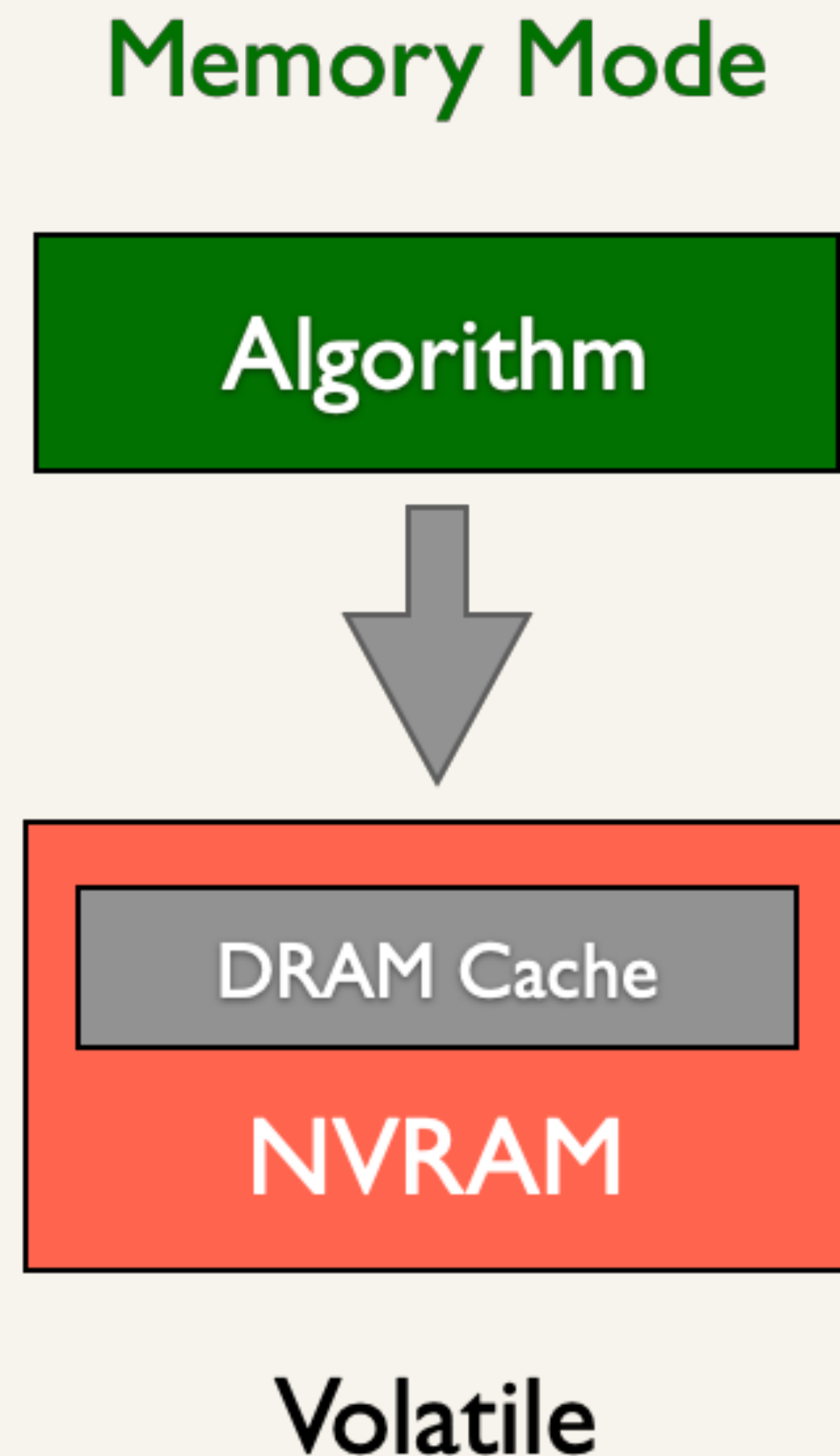


- ❖ Applications do not distinguish between DRAM and NVRAM
- ❖ Existing shared-memory software does not require modification
- ❖ Workloads that are larger than DRAM can involve **costly NVRAM writes**

## Galois (Gill et al.)

- ❖ Gill et al. study the performance of the Galois engine using MemMode
- ❖ They show promising results for scaling to larger than DRAM sizes

# Existing Approaches: DRAM as a Cache



- ❖ Applications do not distinguish between DRAM and NVRAM
- ❖ Existing shared-memory software does not require modification
- ❖ Workloads that are larger than DRAM can involve **costly NVRAM writes**

## Galois (Gill et al.)

- ❖ Gill et al. study the performance of the Galois engine using MemMode
- ❖ They show promising results for scaling to larger than DRAM sizes

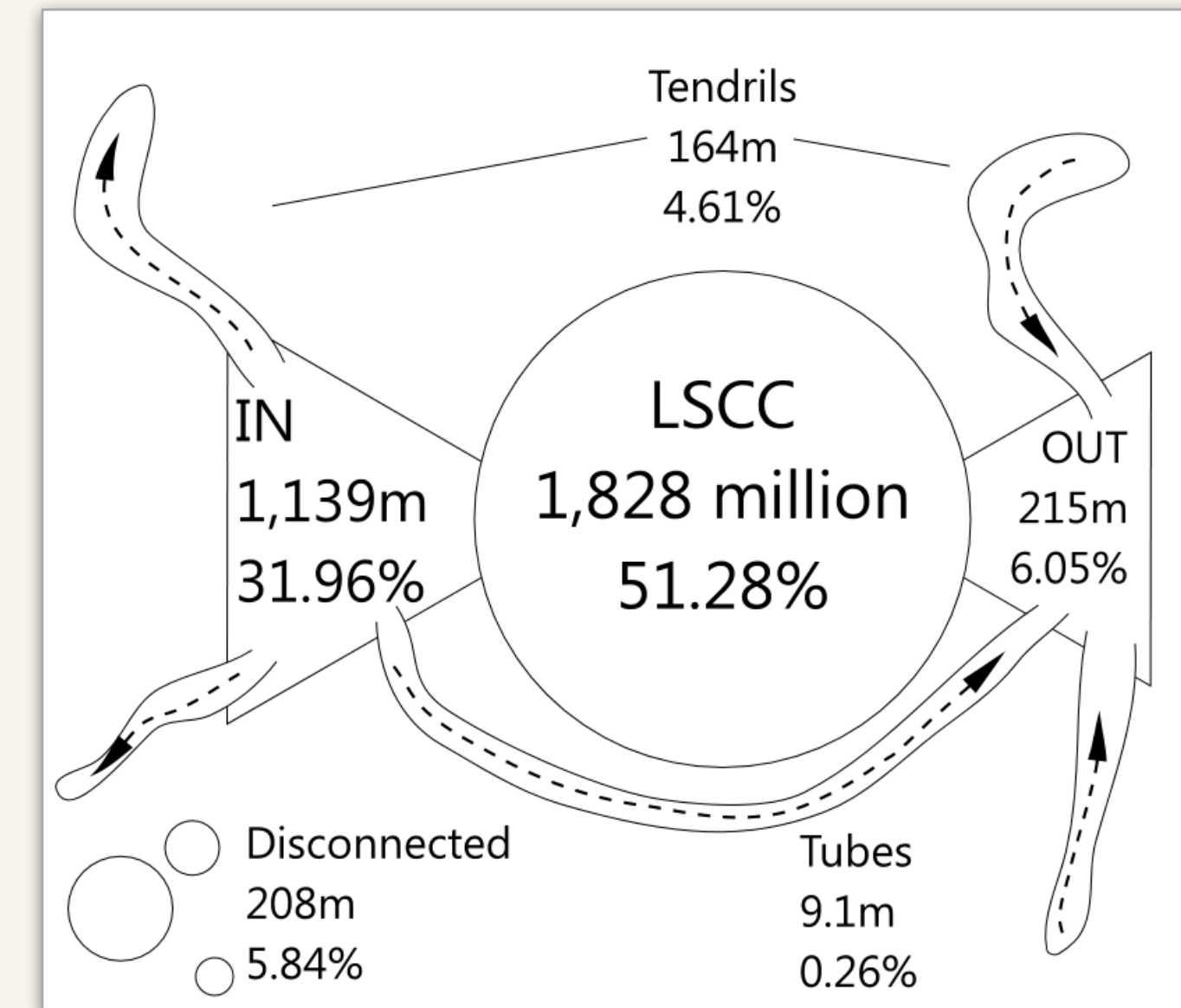
*How does our approach compare?*



# Results for Larger-than-DRAM Graphs

## WebDataCommons Graph

- ❖ Largest publicly available graph today
- ❖ 3.5B vertices connected by 128B edges (225B symmetrized)



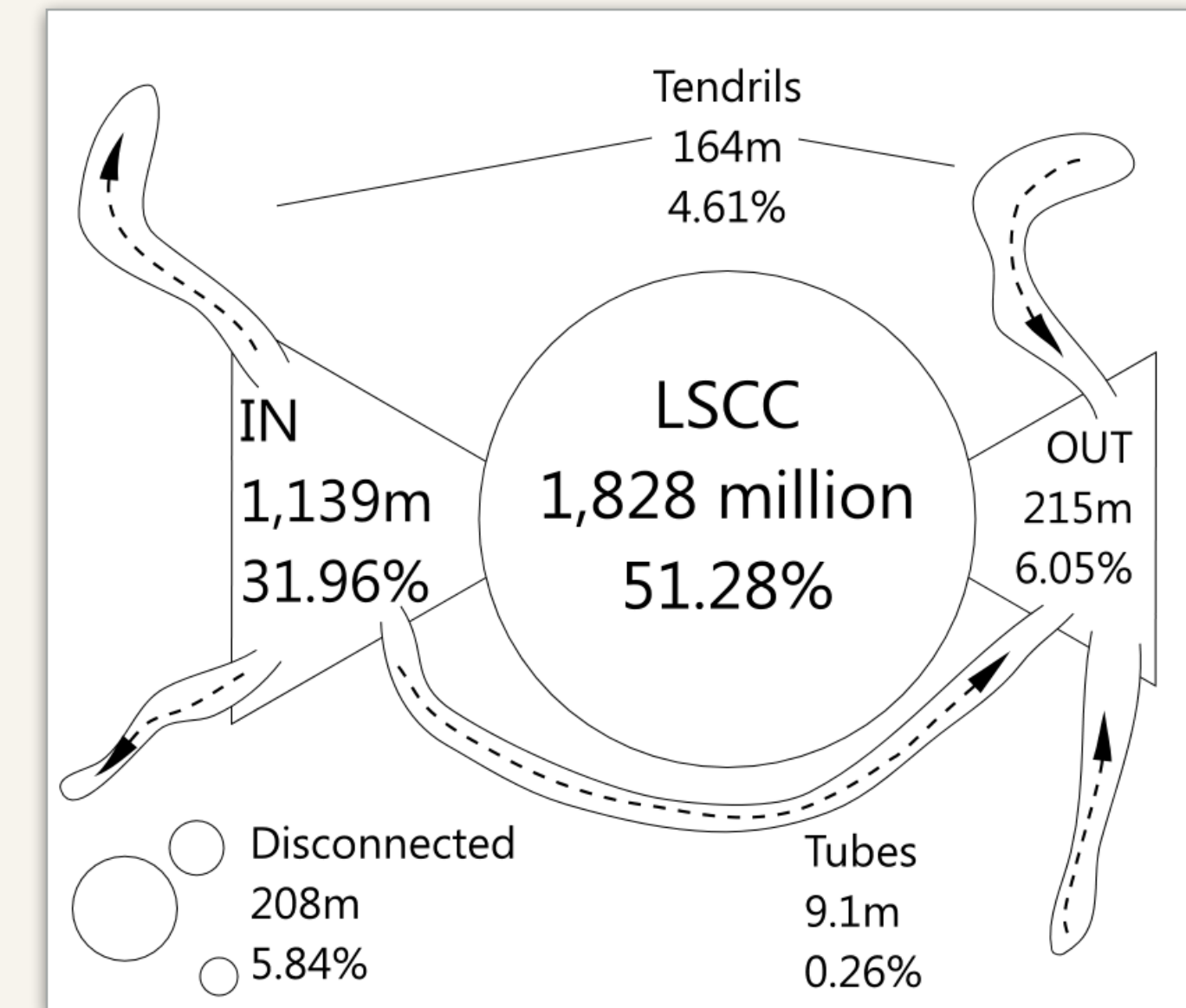
# Results for Larger-than-DRAM Graphs

## WebDataCommons Graph

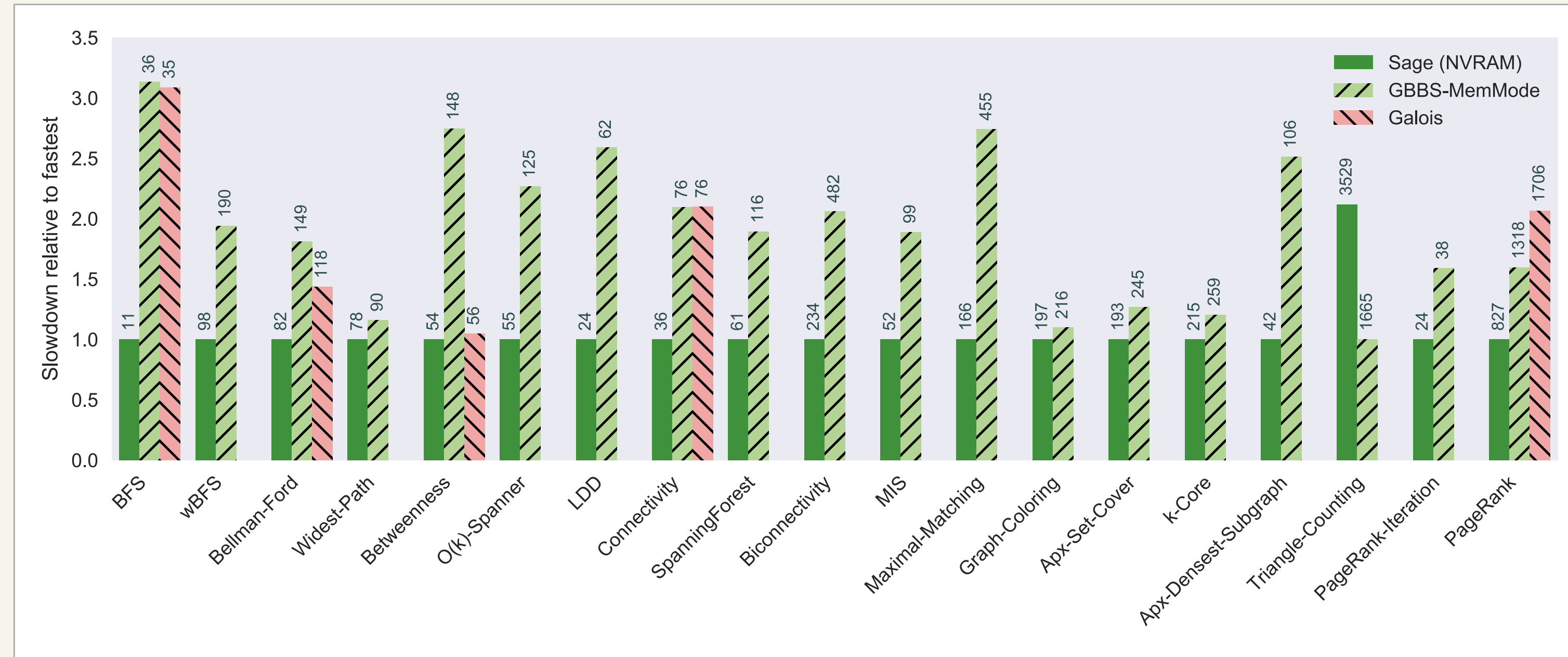
- ❖ Largest publicly available graph today
- ❖ 3.5B vertices connected by 128B edges (225B symmetrized)

## Experiment

- ❖ Compare Sage results with
  - ❖ GBBS using MemMode (existing shared-memory codes)
  - ❖ Galois using MemMode (using numbers reported by authors on the same machine)

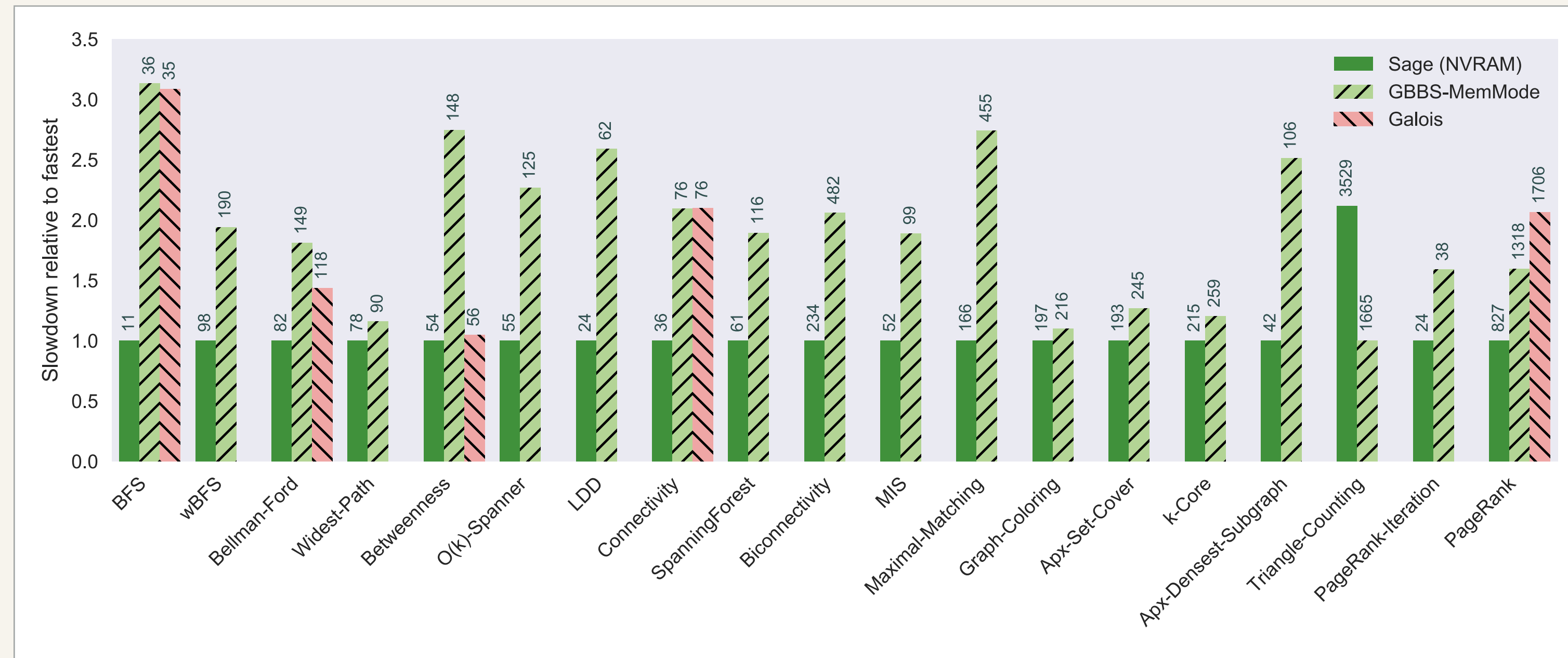


# Results for Larger-than-DRAM Graphs



Run on a 48-core machine with 2-way hyper-threading, 375 GB of DRAM and 3 TB of NVRAM

# Results for Larger-than-DRAM Graphs



Run on a 48-core machine with 2-way hyper-threading, 375 GB of DRAM and 3 TB of NVRAM

**1.94x speedup** on average over Galois (state-of-the-art existing approach to NVRAM graph processing), and **1.87x speedup** over simply running GBBS codes using MemMode

# Results for Graphs Stored in Main Memory

## ClueWeb Graph

- ❖ Large web crawl with  $\sim 1\text{B}$  vertices connected by 42B edges (74B symmetrized)
- ❖ Graph fits entirely in the main memory of our machine

# Results for Graphs Stored in Main Memory

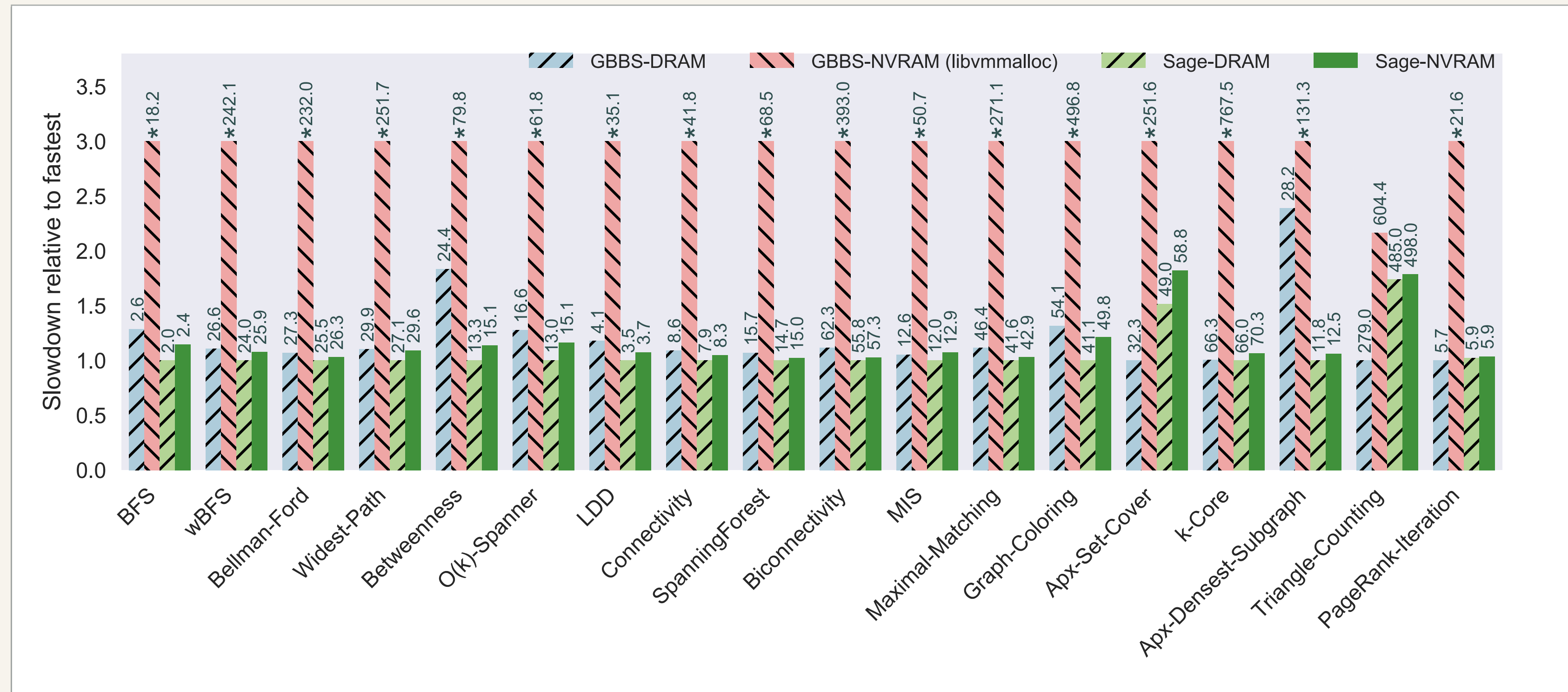
## ClueWeb Graph

- ❖ Large web crawl with  $\sim 1\text{B}$  vertices connected by 42B edges (74B symmetrized)
- ❖ Graph fits entirely in the main memory of our machine

## Experiment

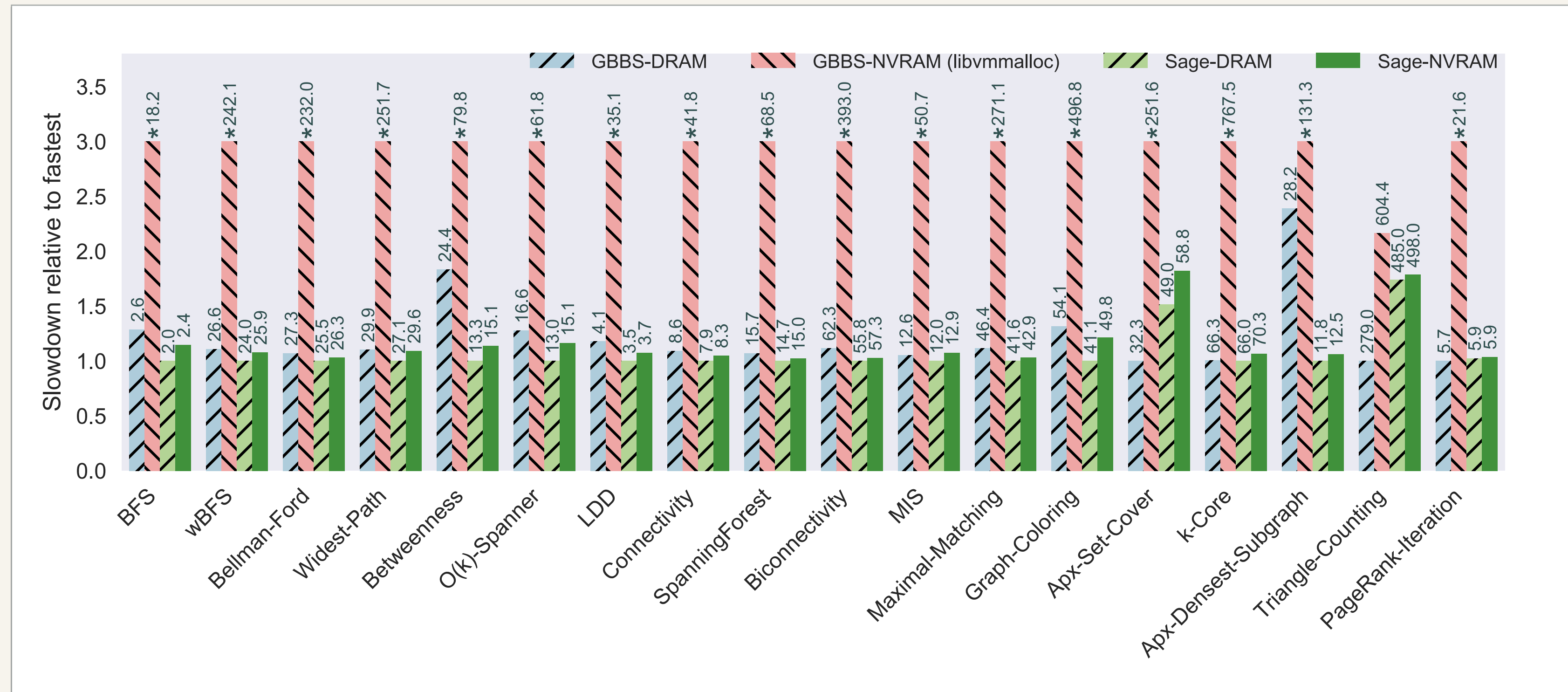
- ❖ Compare Sage (graph stored on NVRAM) with
  - ❖ Sage (graph stored in DRAM)
  - ❖ GBBS (graph stored in DRAM)
  - ❖ GBBS with libvmmalloc (graph stored on NVRAM)

# Results for Graphs Stored in Main Memory



Run on a 48-core machine with 2-way hyper-threading, 375 GB of DRAM and 3 TB of NVRAM

# Results for Graphs Stored in Main Memory



Run on a 48-core machine with 2-way hyper-threading, 375 GB of DRAM and 3 TB of NVRAM

Sage provides DRAM-competitive performance even when reading graph from NVRAM (only 5% slower on average)



# Lessons and Directions for Future Work

# Lessons and Directions for Future Work

## Avoid Cross-Socket NVRAM Traffic

- ❖ NUMA optimization which reads from the copy of the read-only graph from the same socket achieves 6x speedup over cross-socket approach

# Lessons and Directions for Future Work

## Avoid Cross-Socket NVRAM Traffic

- ❖ NUMA optimization which reads from the copy of the read-only graph from the same socket achieves 6x speedup over cross-socket approach

## Utilize App-Direct Mode

- ❖ Nearly 2x improvement for App-Direct based PSAM algorithms over two fast Memory Mode approaches

# Lessons and Directions for Future Work

## Avoid Cross-Socket NVRAM Traffic

- ❖ NUMA optimization which reads from the copy of the read-only graph from the same socket achieves 6x speedup over cross-socket approach

## Utilize App-Direct Mode

- ❖ Nearly 2x improvement for App-Direct based PSAM algorithms over two fast Memory Mode approaches

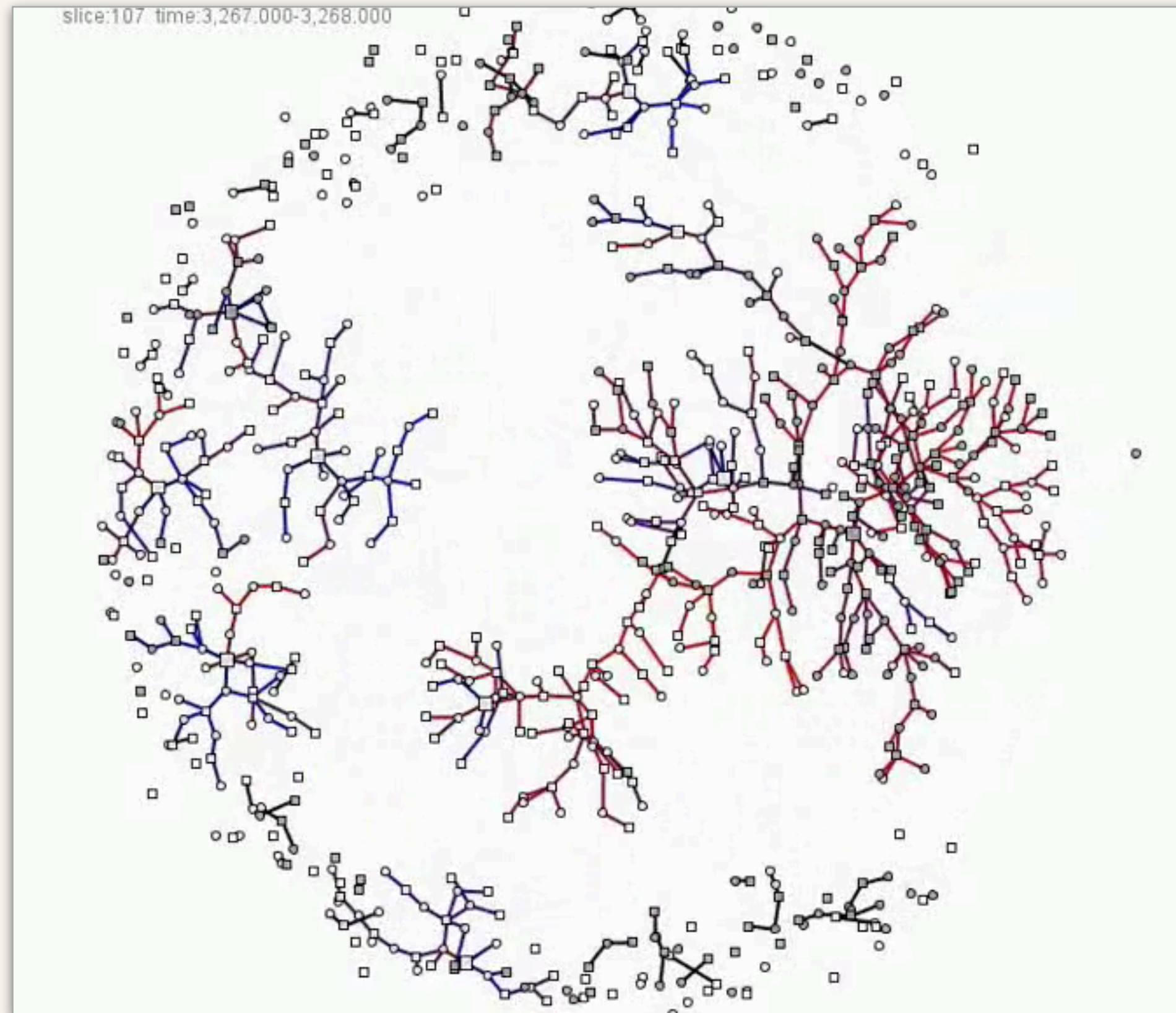
## Avoid NVRAM Writes

- ❖ PSAM implementations which only read from NVRAM are over 6x faster than our algorithms which write to NVRAM (using libvmmalloc)

# Streaming Graph Processing

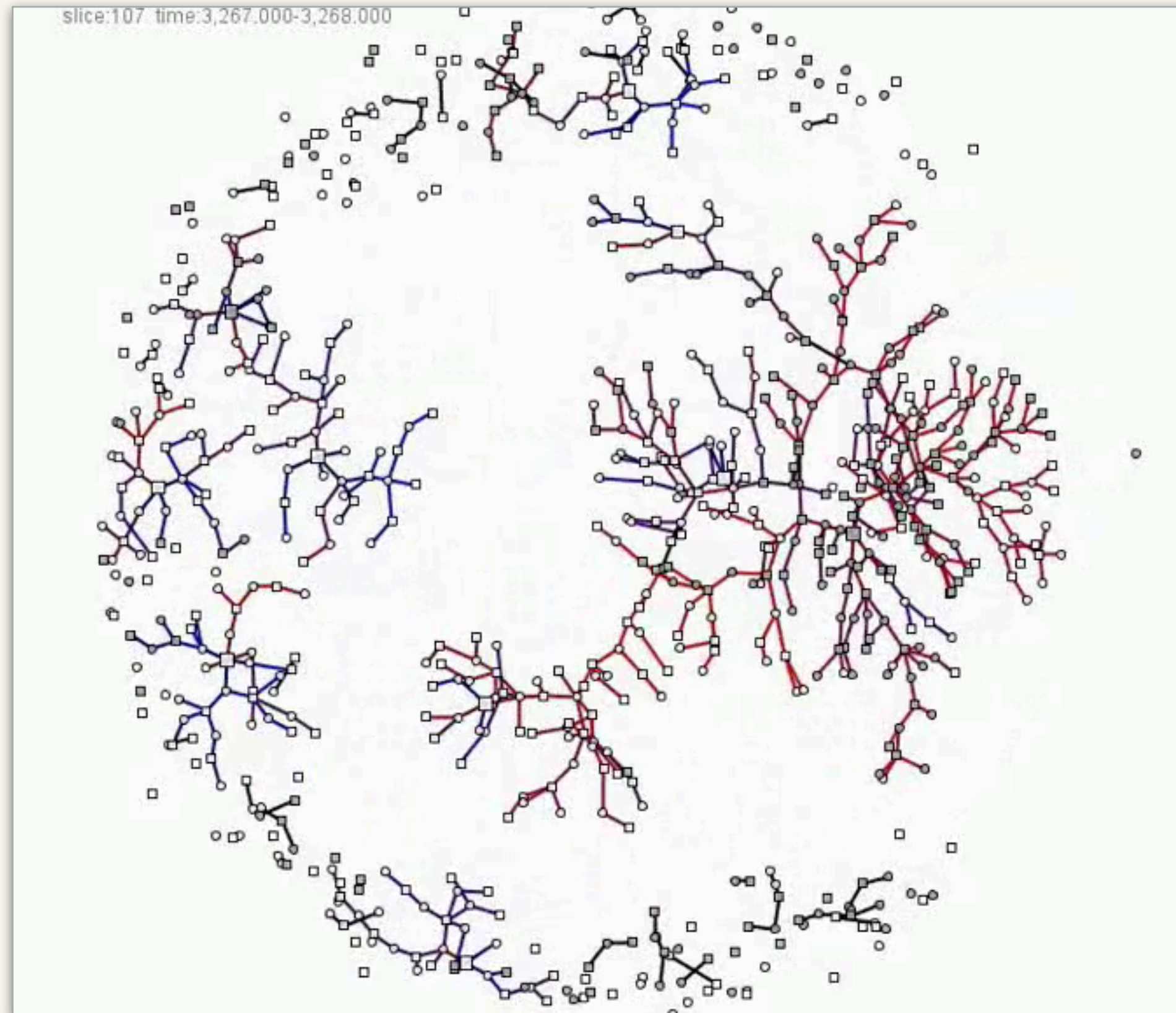
# Dynamic Graph Processing

## Measuring the spread of infections



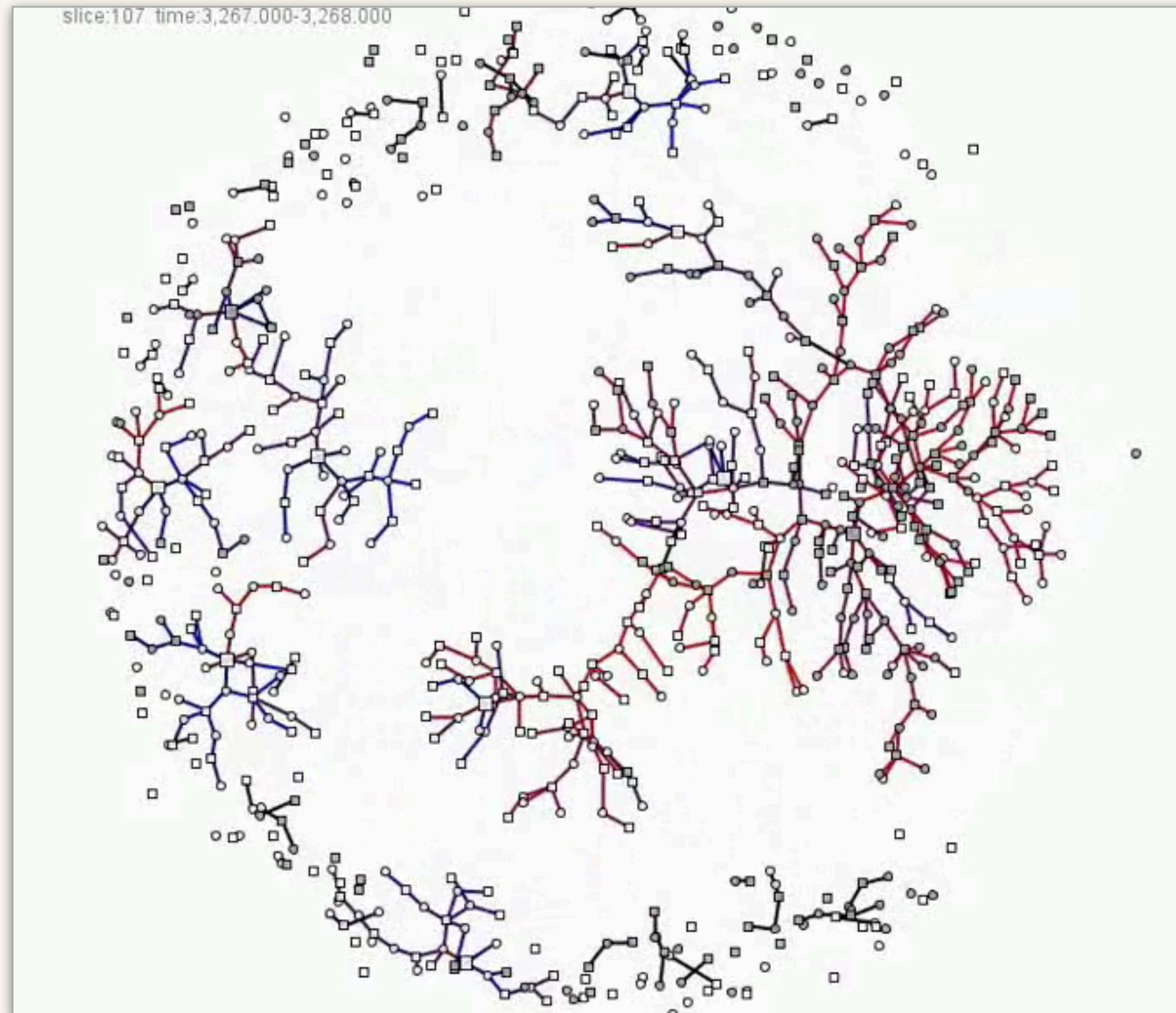
# Dynamic Graph Processing

## Measuring the spread of infections

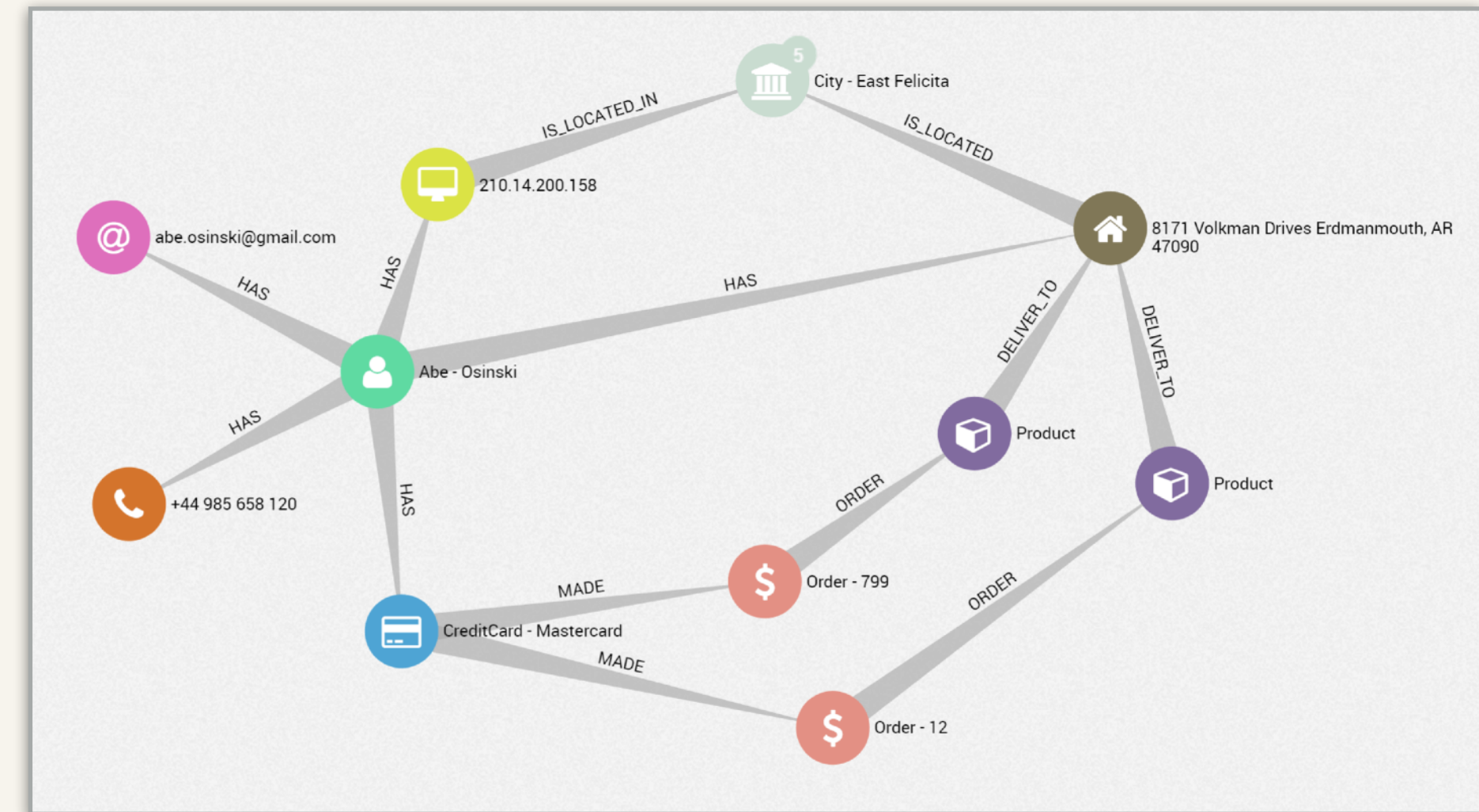


# Dynamic Graph Processing

## Measuring the spread of infections



## Preventing money laundering and fraud



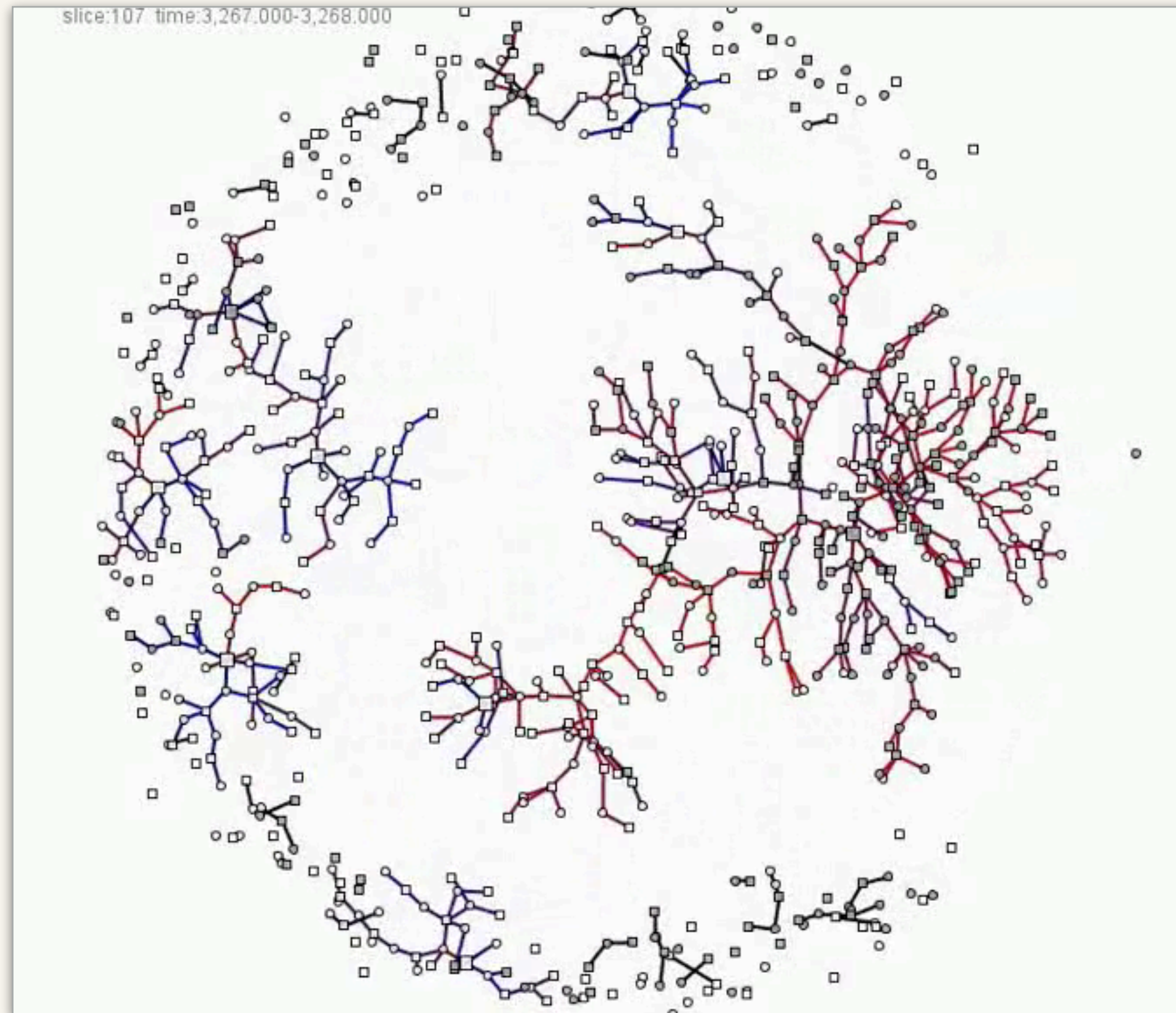
## Other Applications

- ❖ Recommendation Systems
- ❖ Geospatial Systems

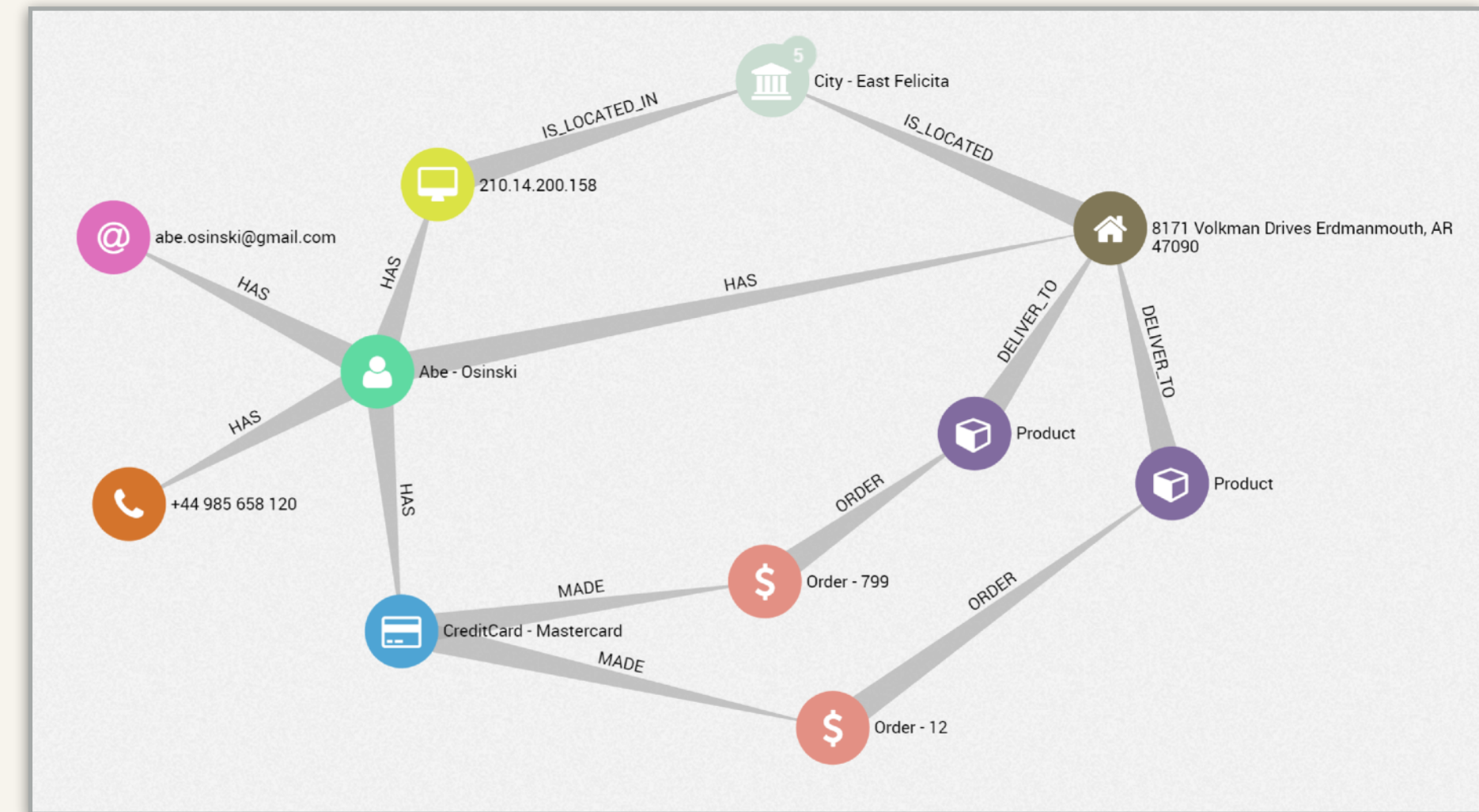


# Dynamic Graph Processing

## Measuring the spread of infections



## Preventing money laundering and fraud



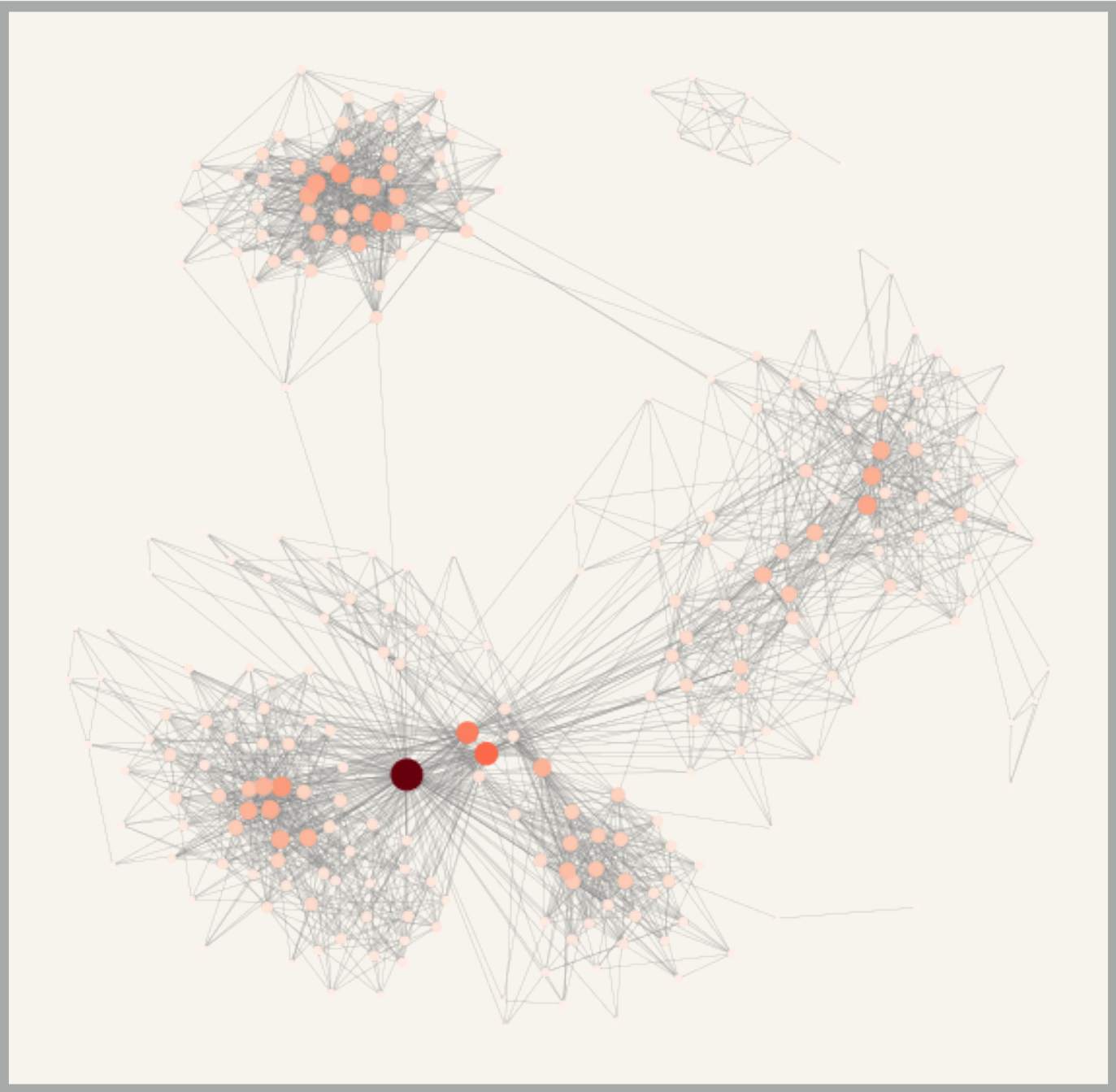
## Other Applications

- ❖ Recommendation Systems
- ❖ Geospatial Systems

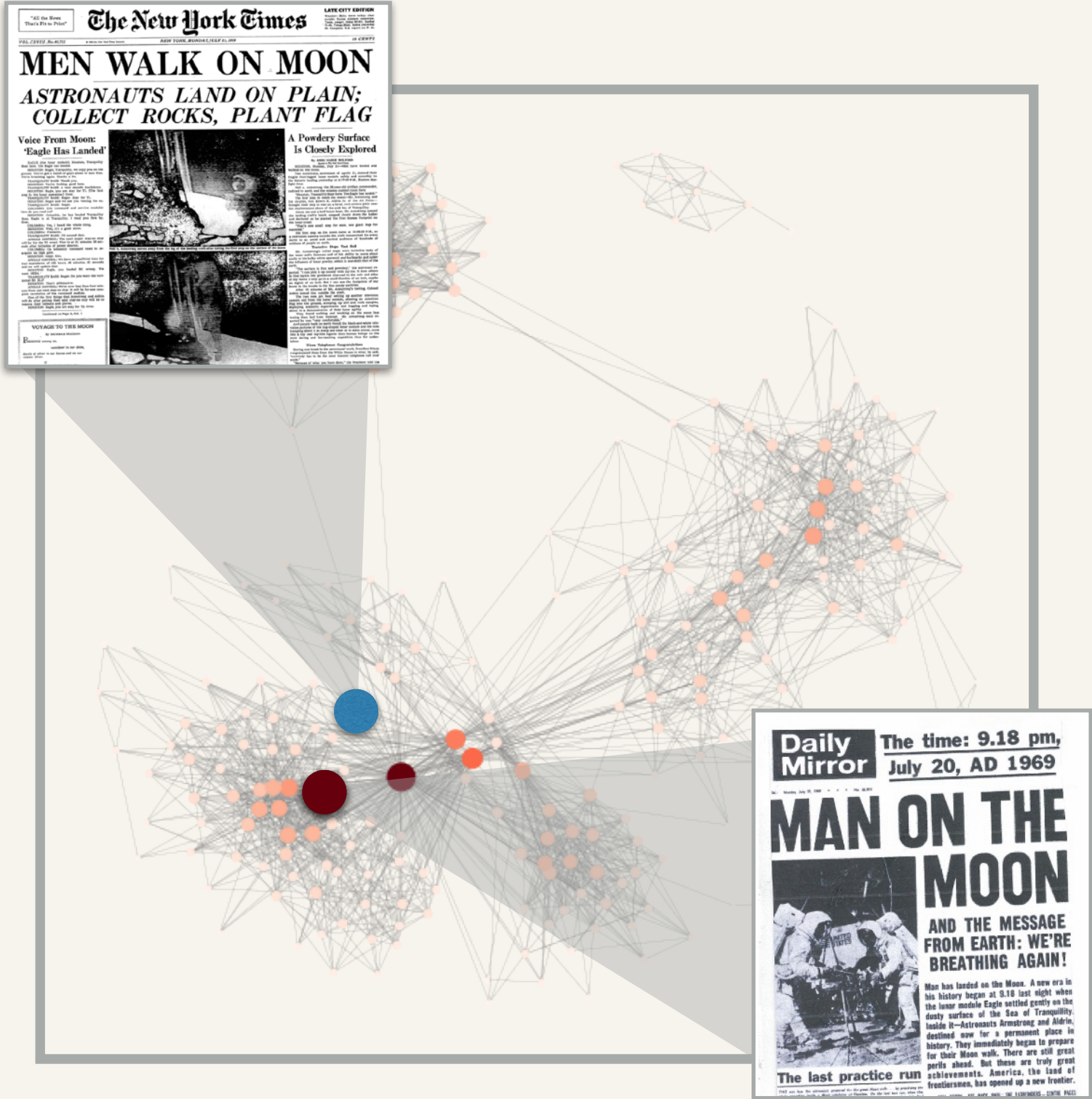
Many important applications must maintain information about evolving graphs!

# Dynamic Graph Processing: Example

# Dynamic Graph Processing: Example



# Dynamic Graph Processing: Example

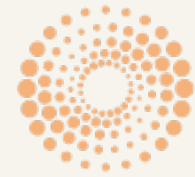


# Dynamic Graph Processing: Example

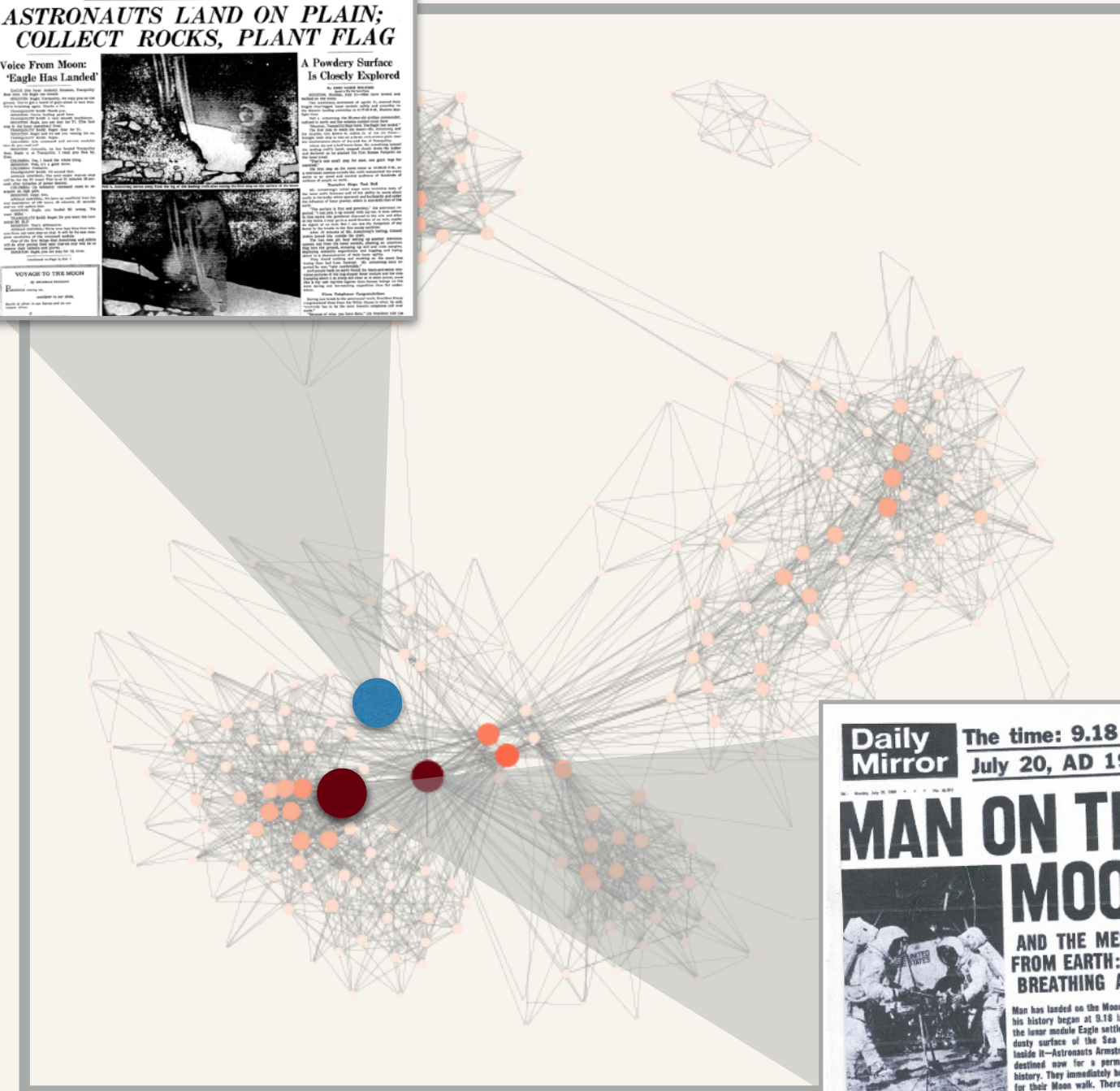
Update Stream



The Washington Post



REUTERS



# Dynamic Graph Processing: Example

Update Stream



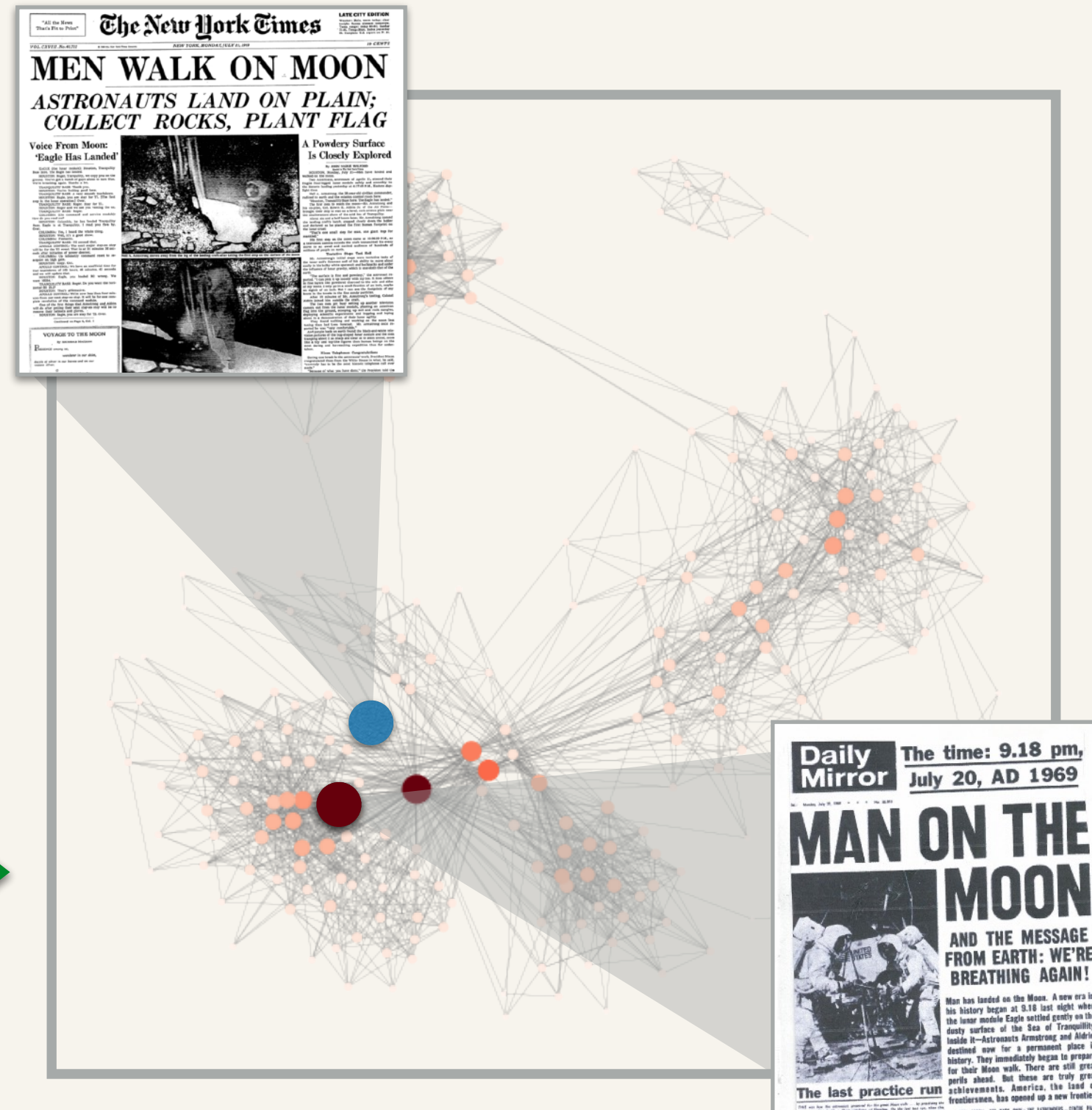
Query Stream

Fetch similar vertices

Clustering Coefficients

Fetch vertex u's neighbors

Centrality Ranking



# Dynamic Graph Processing: Example

Update Stream



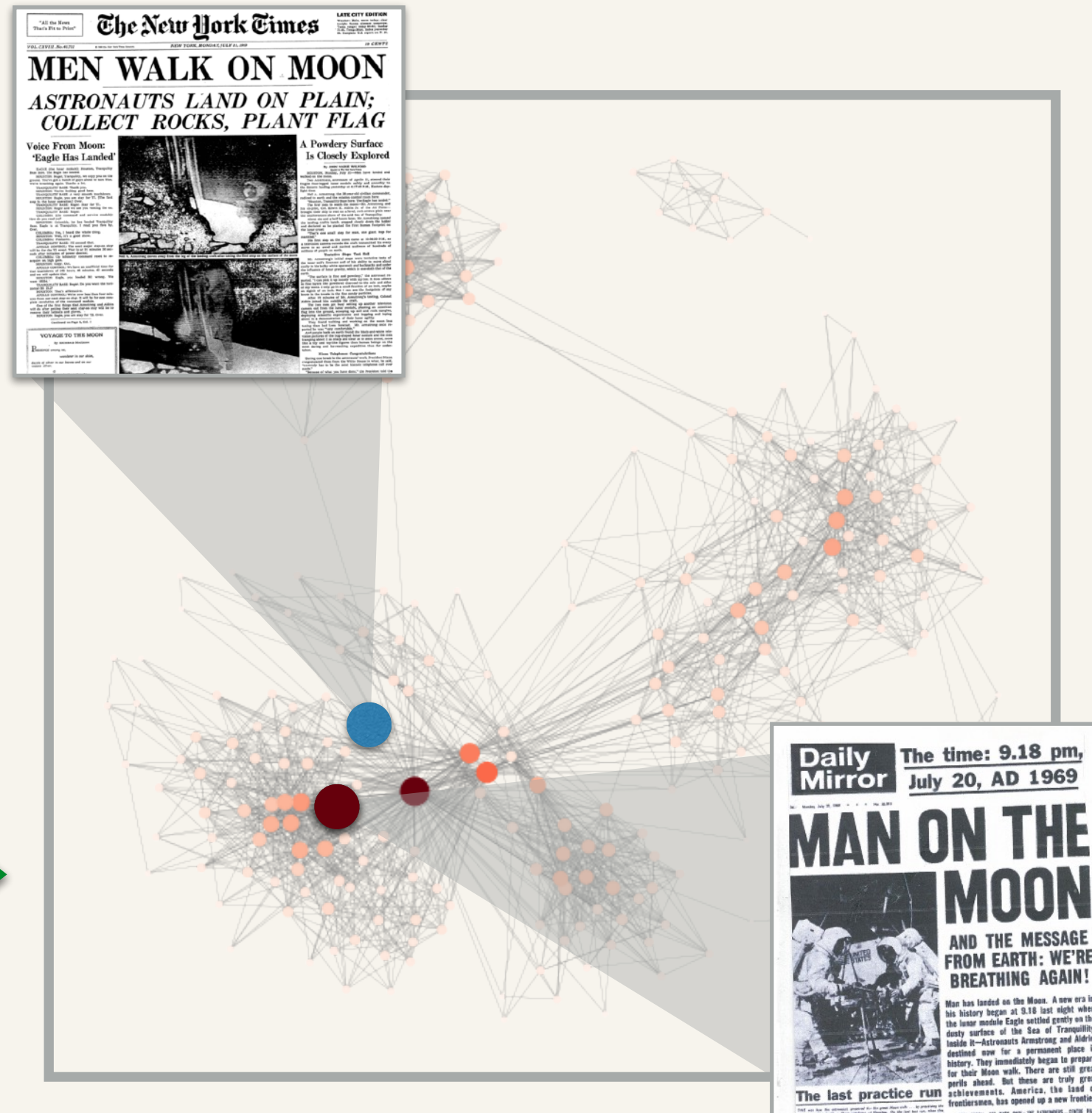
Query Stream

Fetch similar vertices

Clustering Coefficients

Fetch vertex u's neighbors

Centrality Ranking



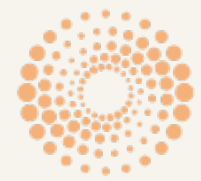
Streaming Graph Processing

Update the graph (in parallel);  
Execute arbitrary queries on snapshots.

# Dynamic Graph Processing: Example

Update Stream

The Washington Post



REUTERS

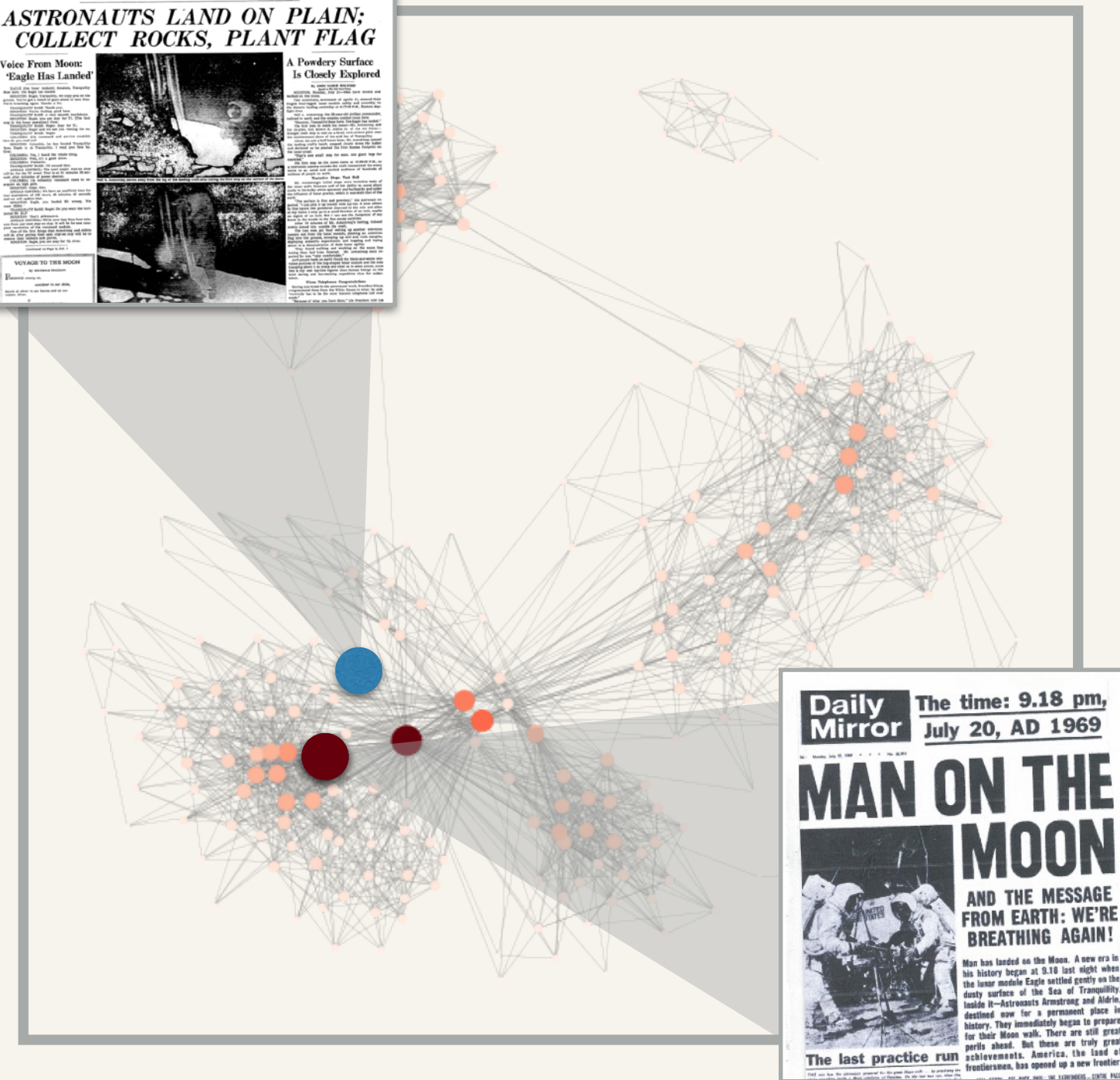
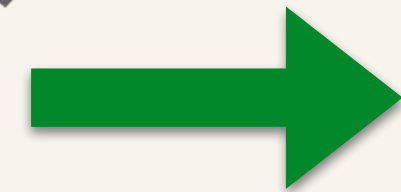
Query Stream

Fetch similar vertices

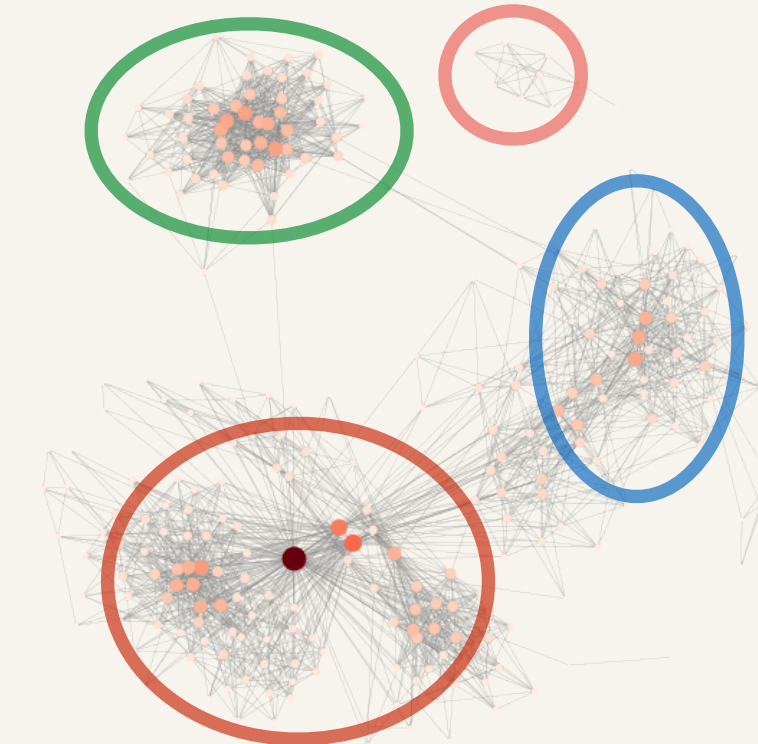
Clustering Coefficients

Fetch vertex u's neighbors

Centrality Ranking



Clustering



Triangle Counting  
Connected Components

Streaming Graph Processing

Update the graph (in parallel);  
Execute arbitrary queries on snapshots.



# Dynamic Graph Processing: Example

Update Stream



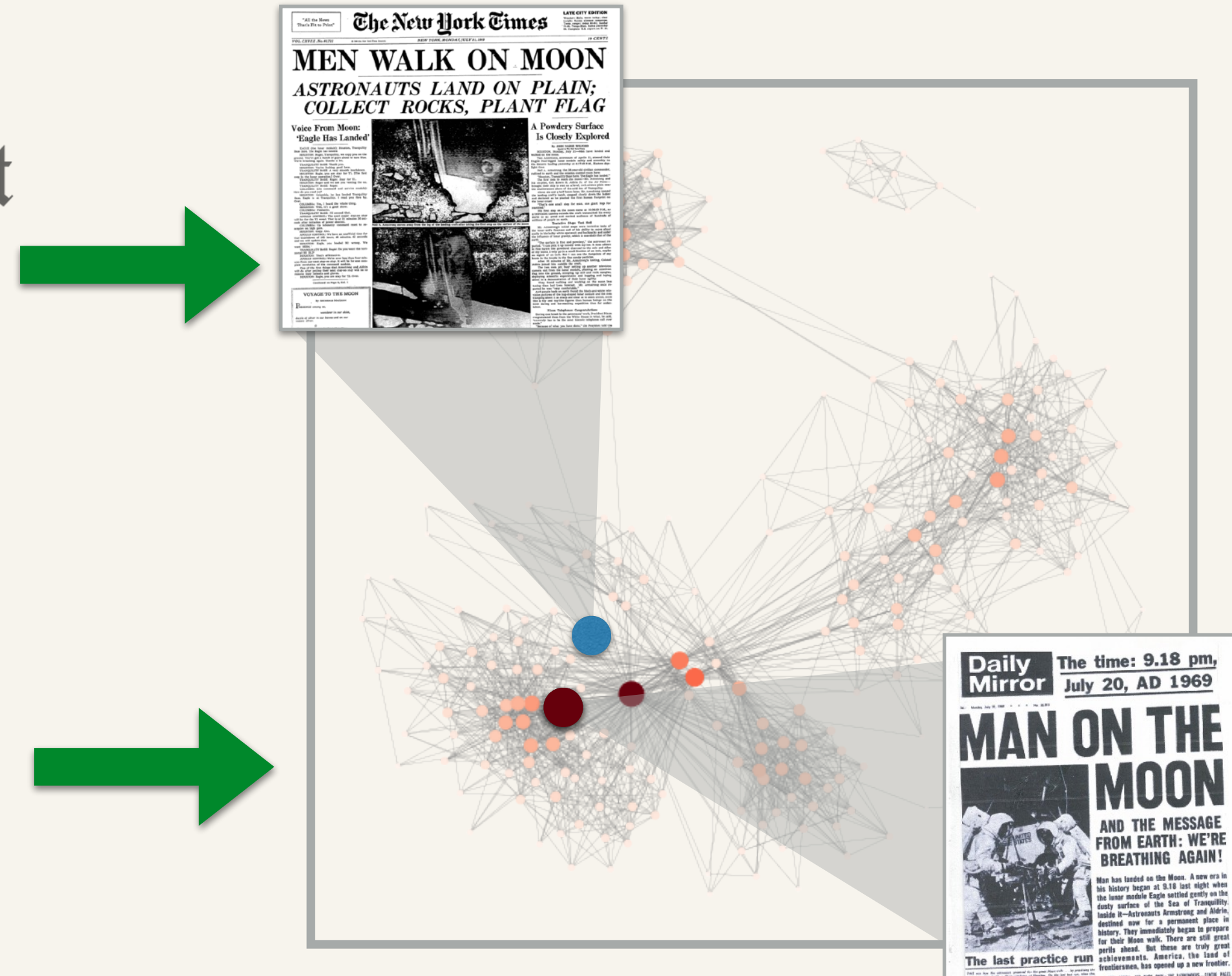
Query Stream

Fetch similar vertices

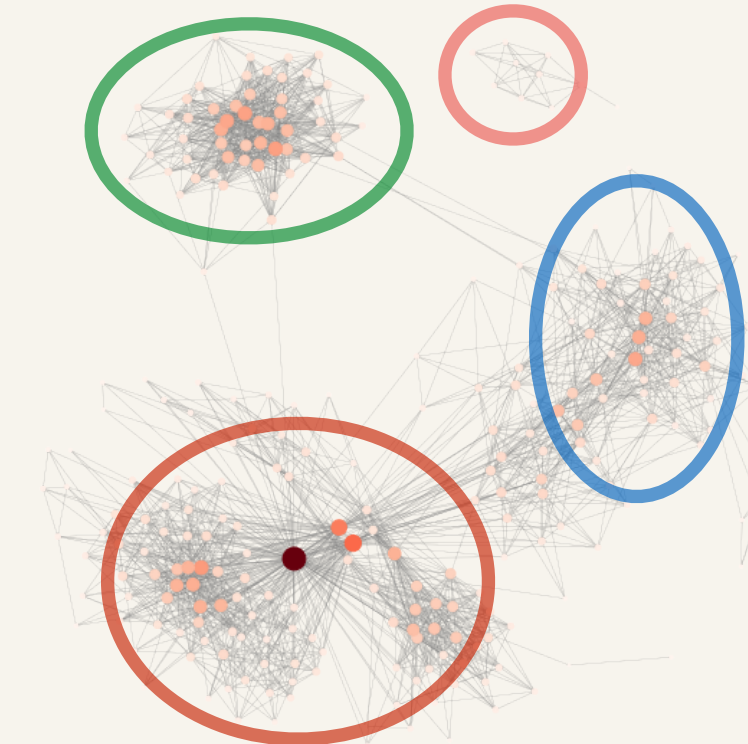
Clustering Coefficients

Fetch vertex u's neighbors

Centrality Ranking



Clustering



Triangle Counting  
Connected Components

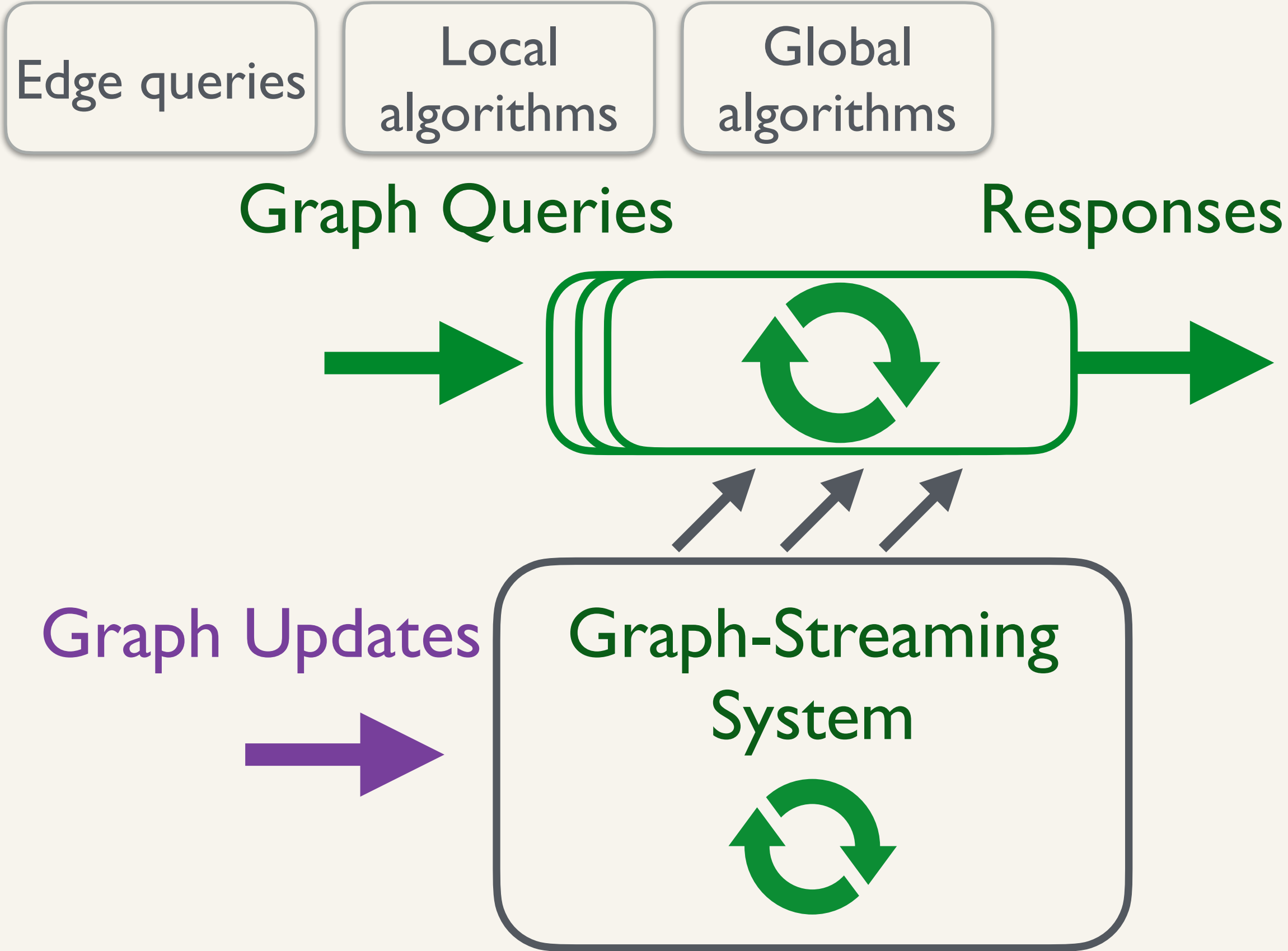
Streaming Graph Processing

Update the graph (in parallel);  
Execute arbitrary queries on snapshots.

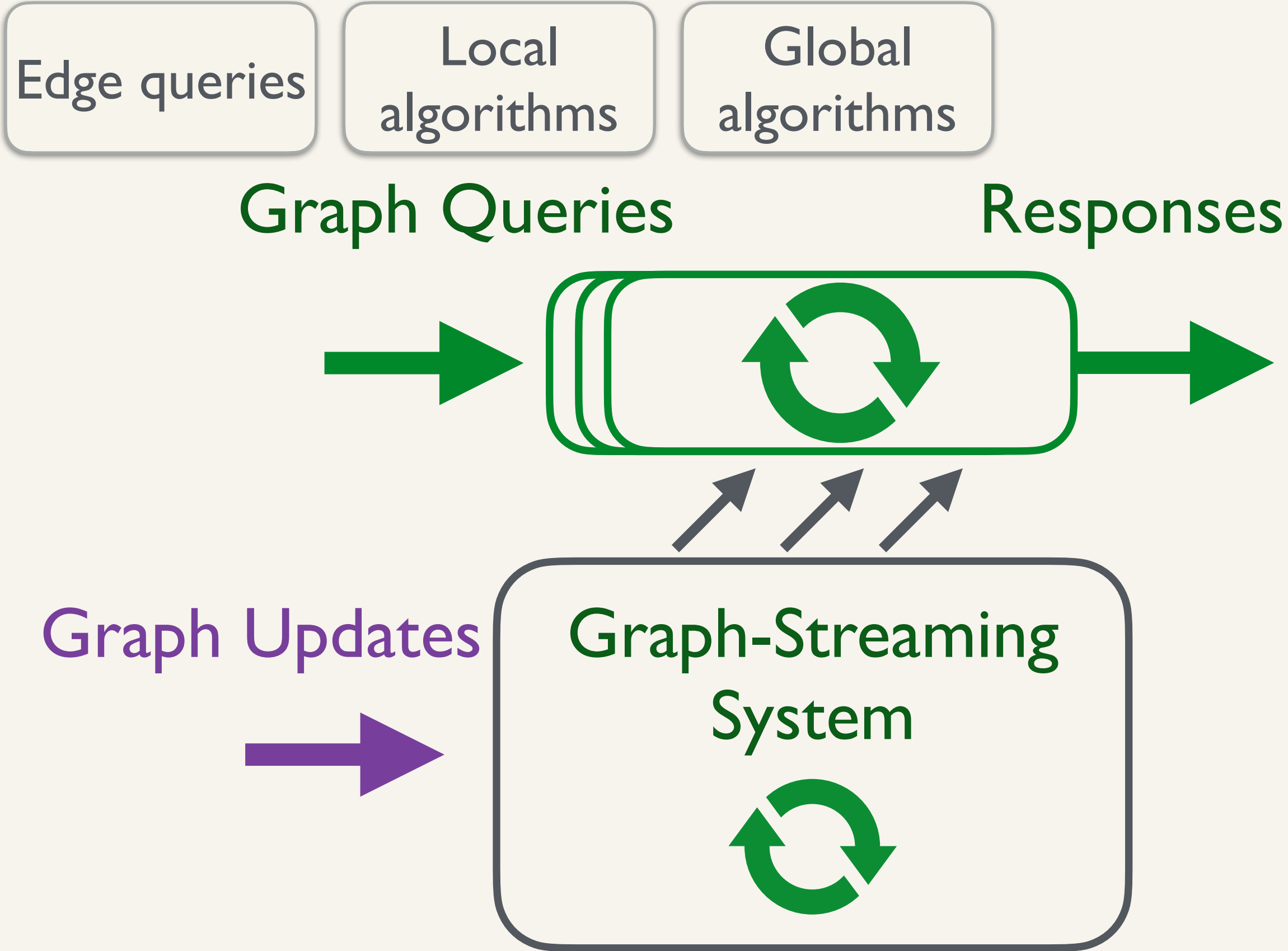
Batch-Dynamic Graph Processing

Pre-determined queries;  
Process updates faster than recomputation.

# Streaming Graph Processing

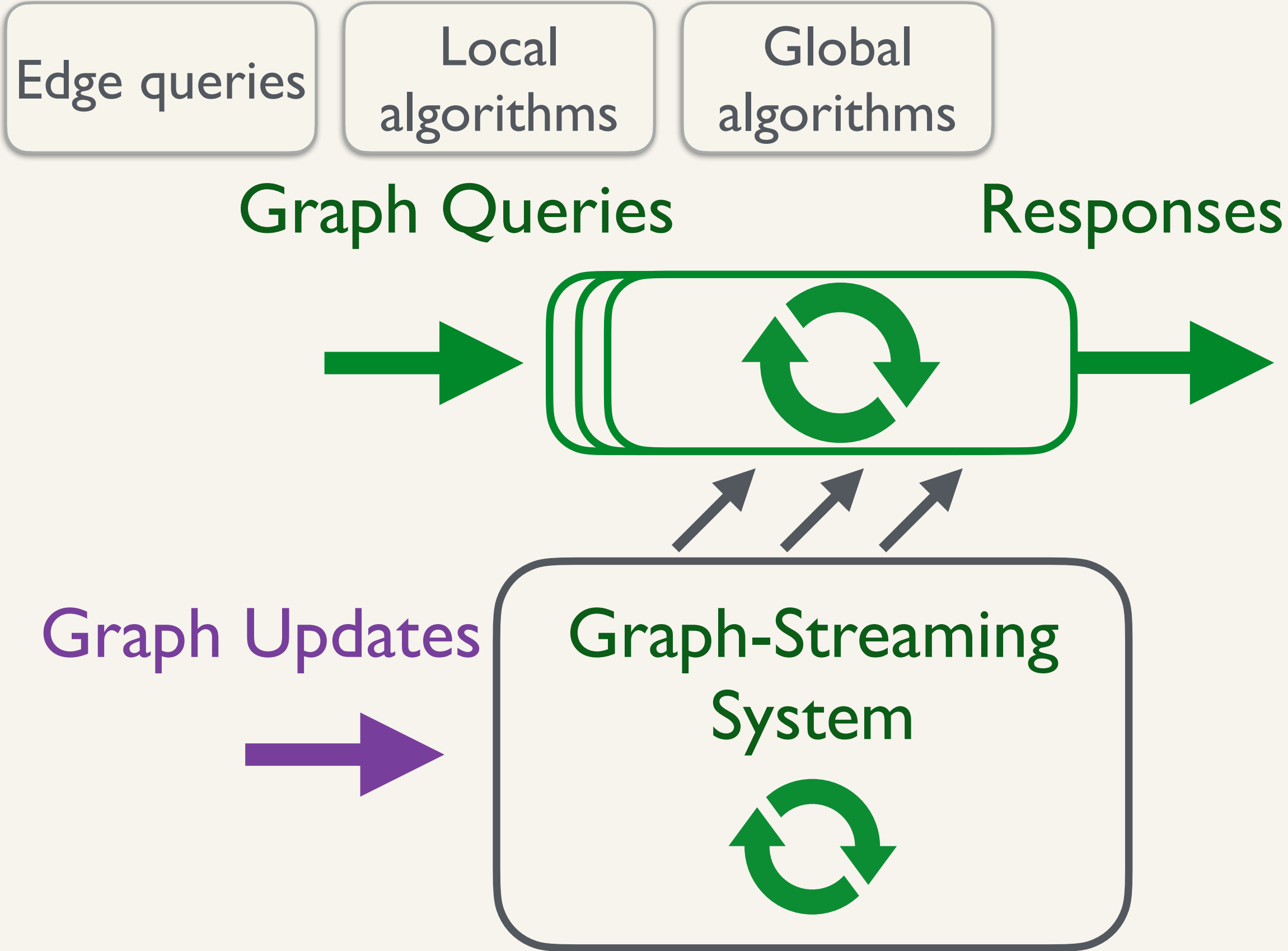


# Streaming Graph Processing



*Goal: low-latency for both updates and queries arriving concurrently to the system*

# Streaming Graph Processing



## Single-version

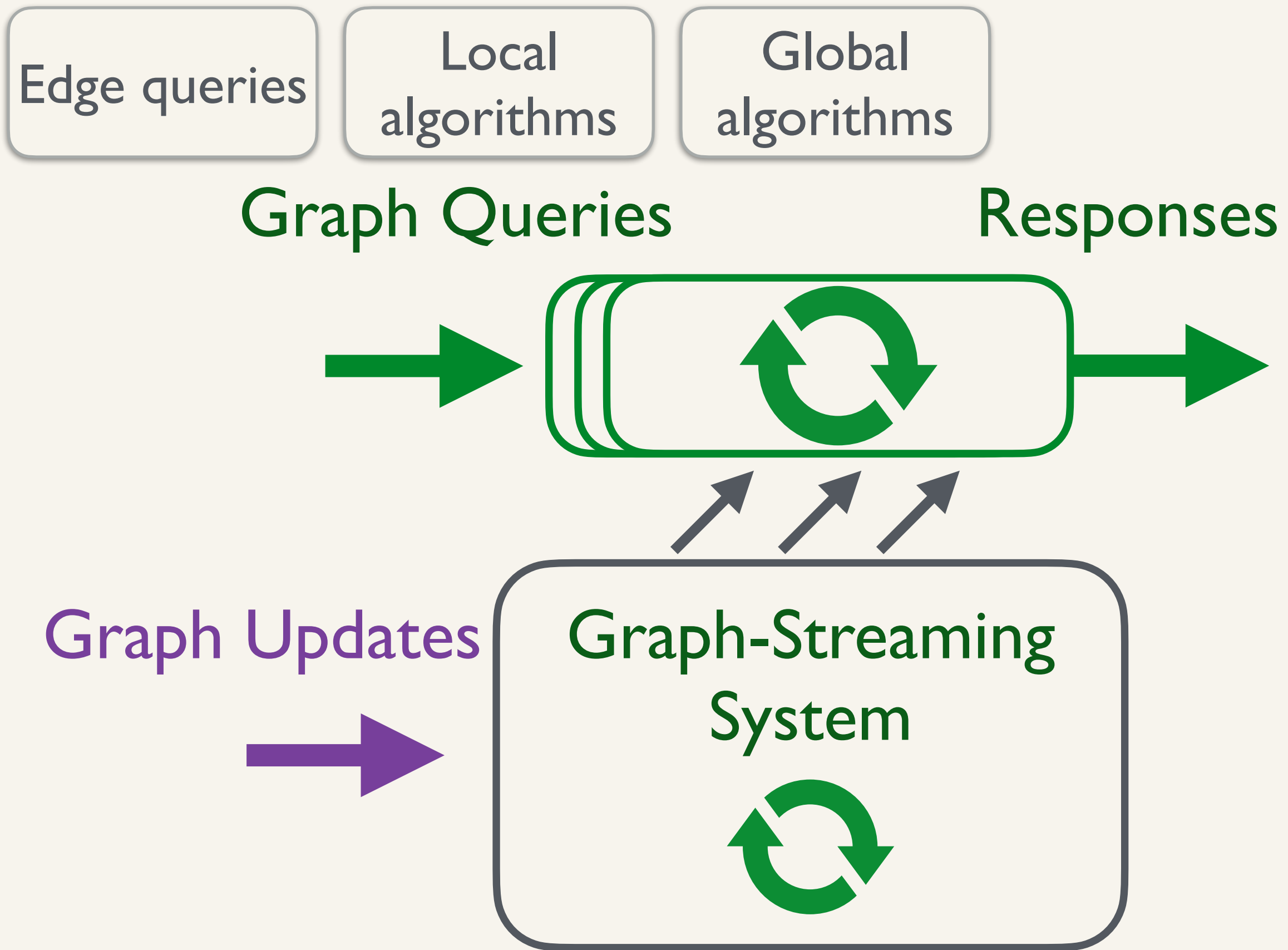
*STINGER* [EMRB'12]

*cuSTINGER* [GB'16]

*Kickstarter* [VGX'17]

*Goal: low-latency for both updates and queries arriving concurrently to the system*

# Streaming Graph Processing



## Single-version

*STINGER* [EMRB'12]

*cuSTINGER* [GB'16]

*Kickstarter* [VGX'17]

## Multi-version (snapshot-based)

*Kineograph* [CHKMW+'12]

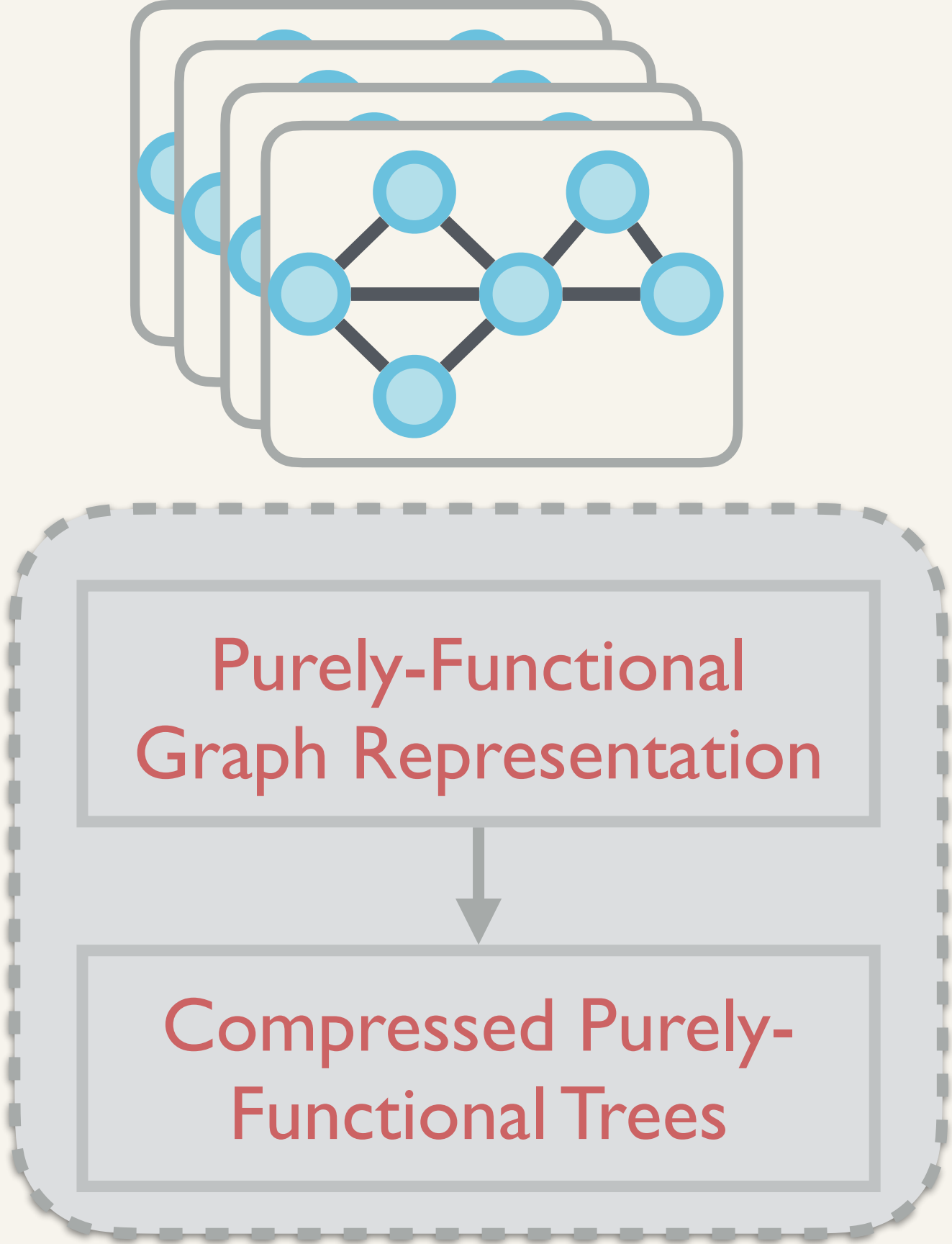
*LLAMA* [MMMS'15]

*Goal: low-latency for both updates and queries arriving concurrently to the system*

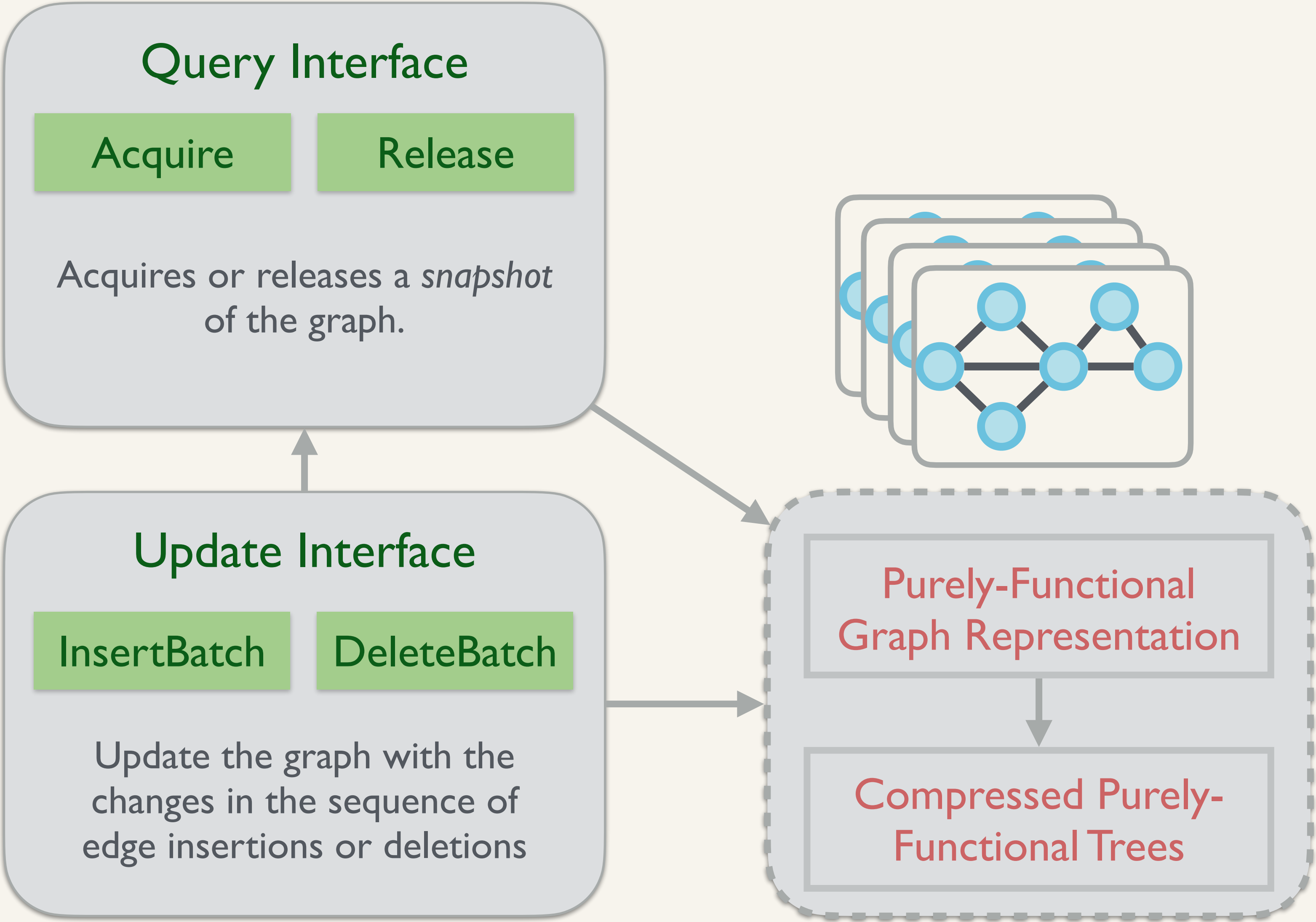
# Low-Latency Graph Streaming using Compressed Purely-Functional Trees [DBS'19]

Can we design a system that can compactly represent and concurrently update and query the largest real-world graphs?

# Aspen: A Low-Latency Graph Streaming System

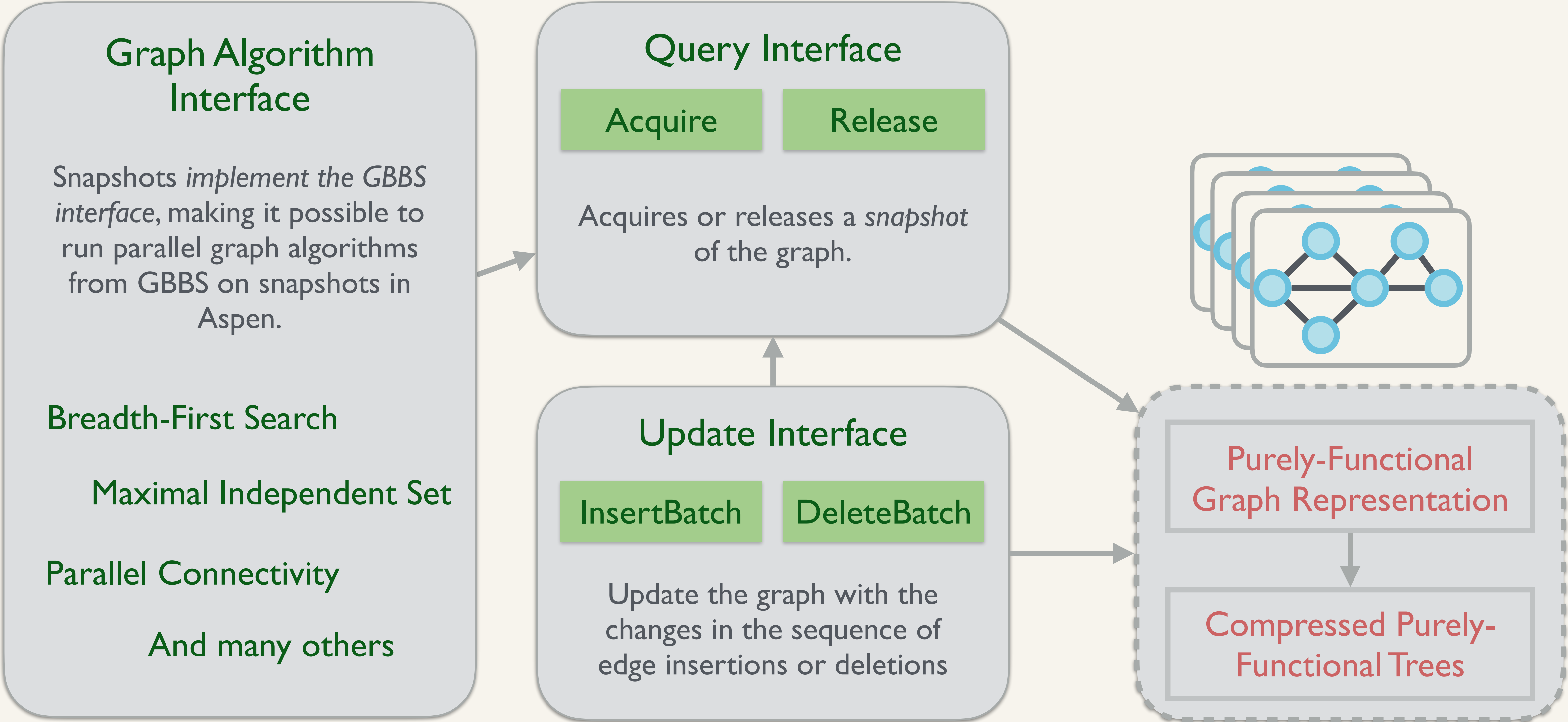


# Aspen: A Low-Latency Graph Streaming System

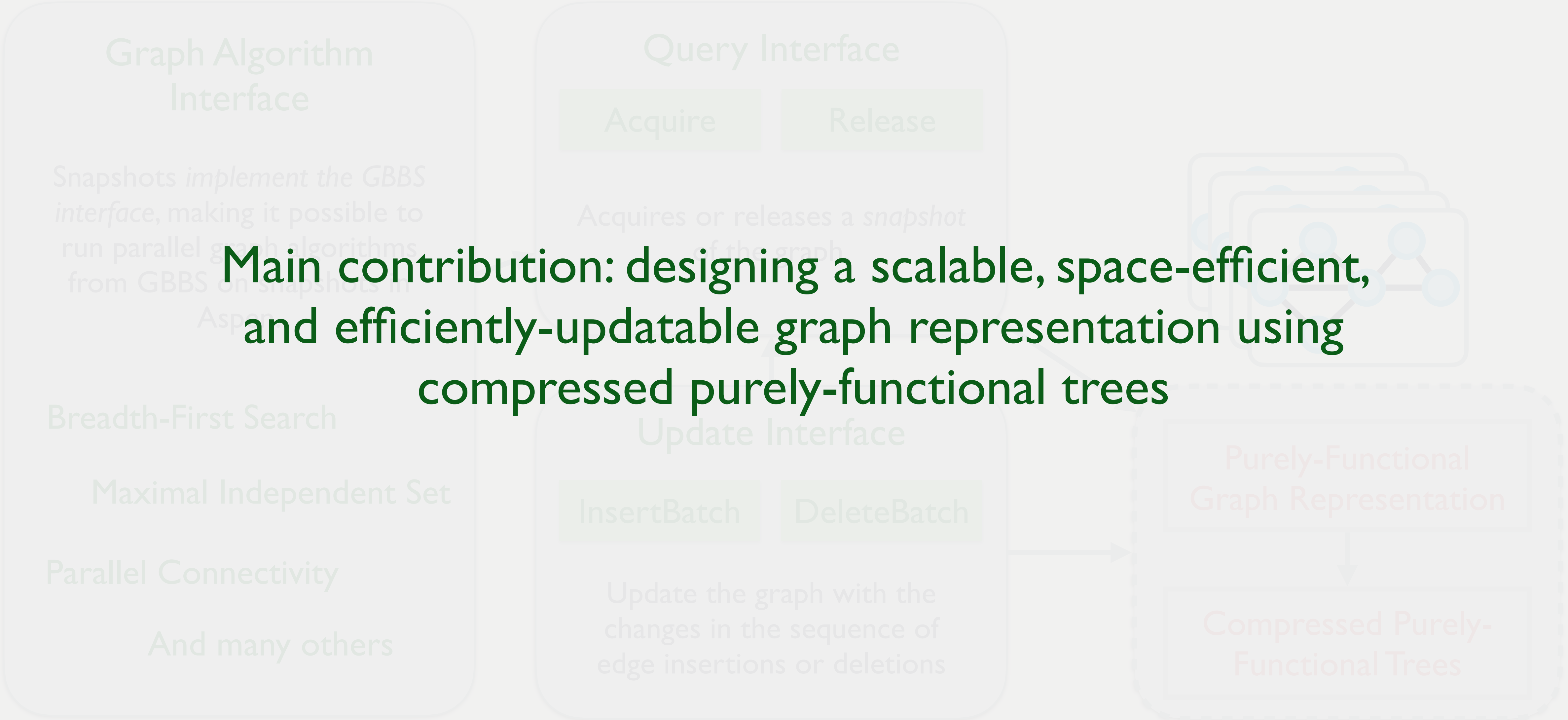




# Aspen: A Low-Latency Graph Streaming System

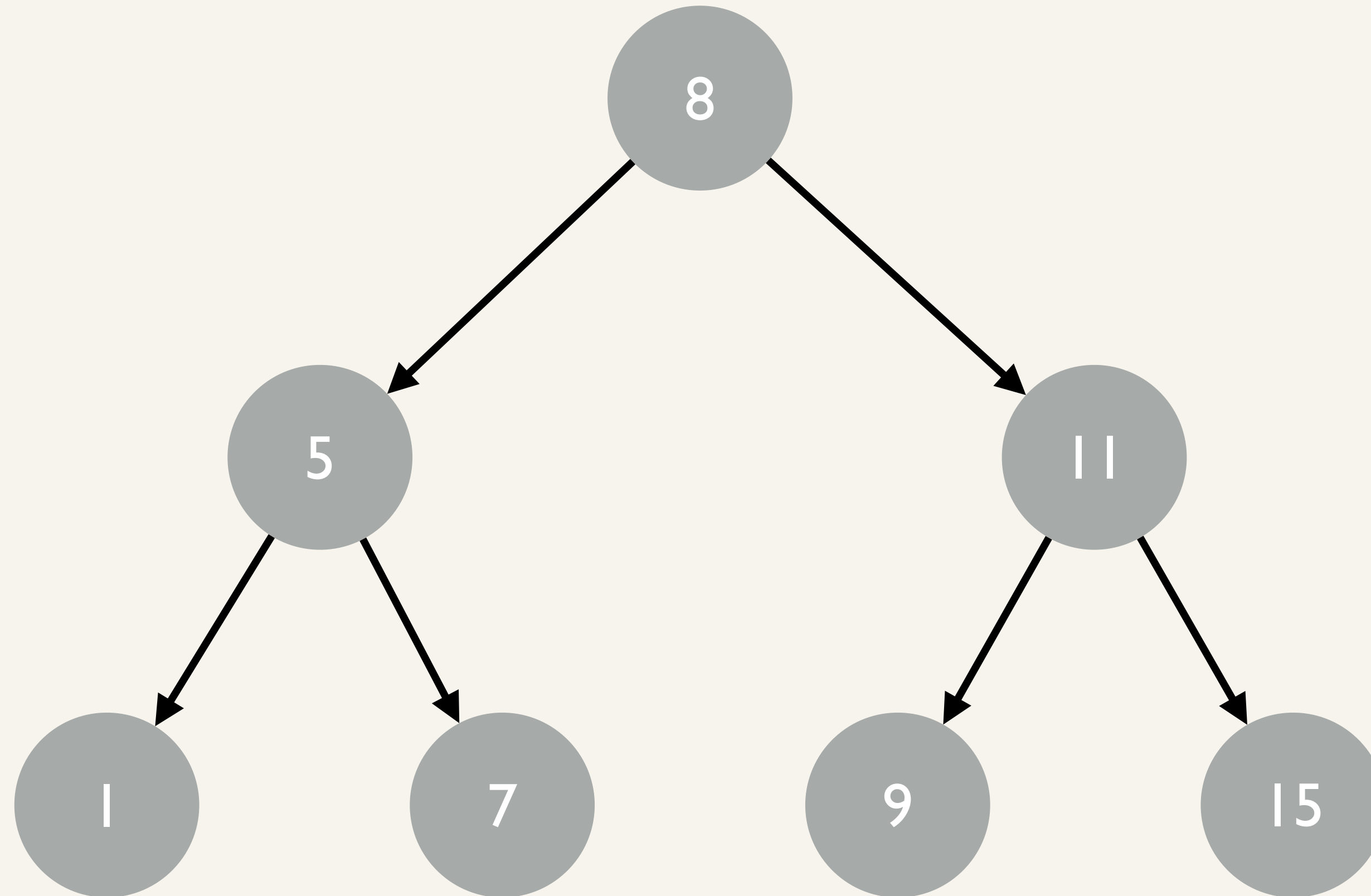


# Aspen: A Low-Latency Graph Streaming System

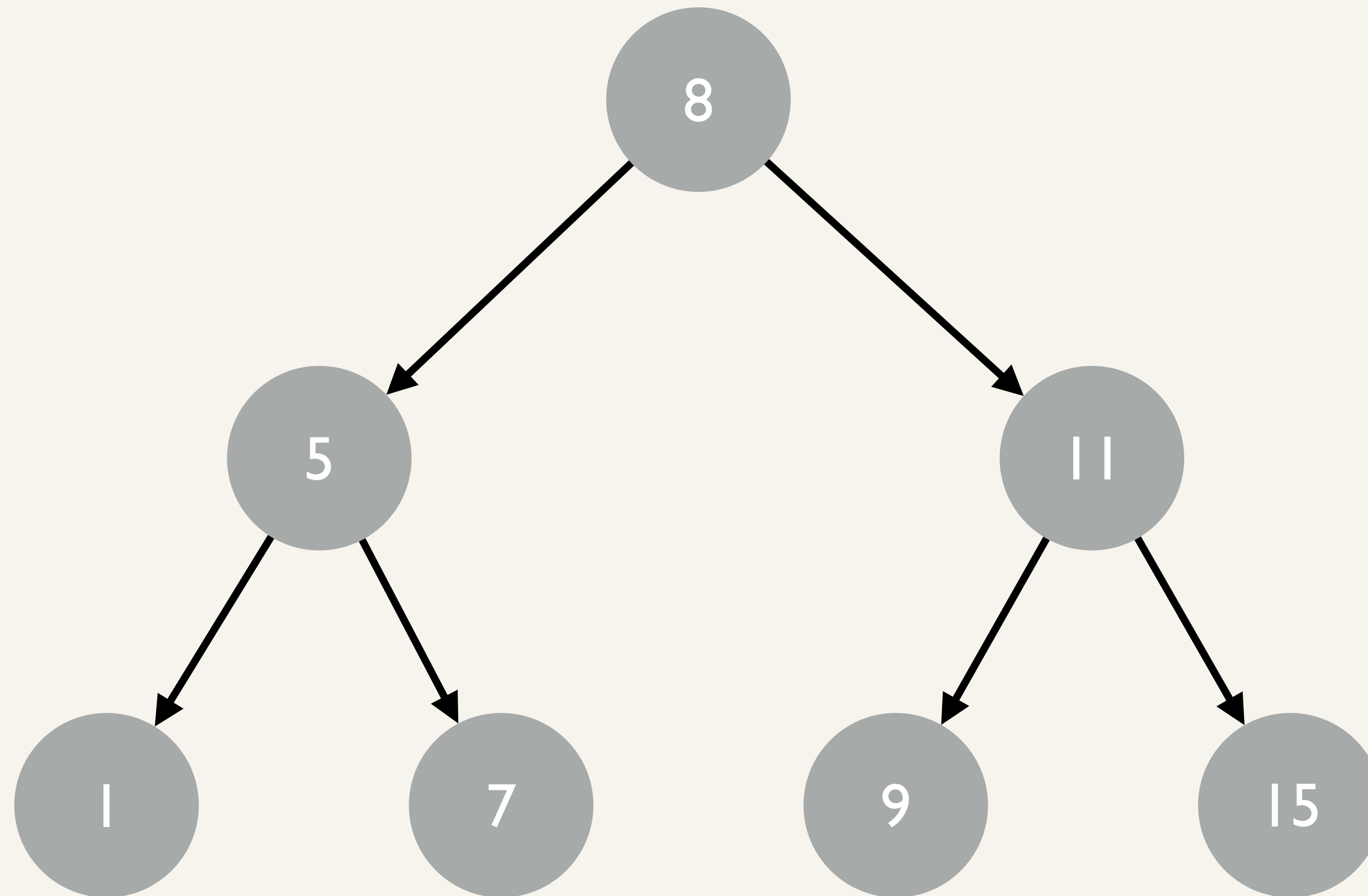


**Main contribution: designing a scalable, space-efficient, and efficiently-updatable graph representation using compressed purely-functional trees**

# Purely-Functional Trees

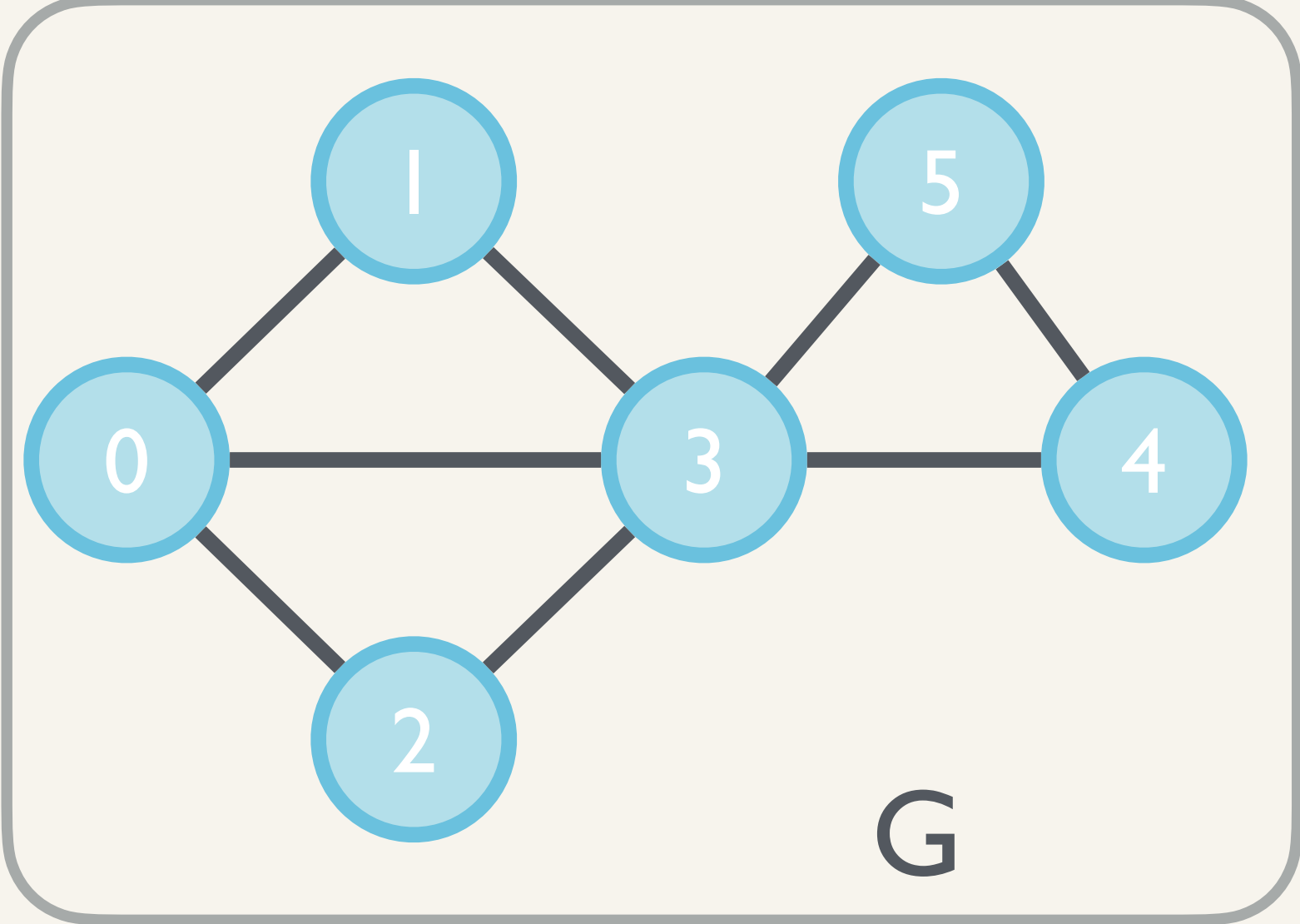


# Purely-Functional Trees

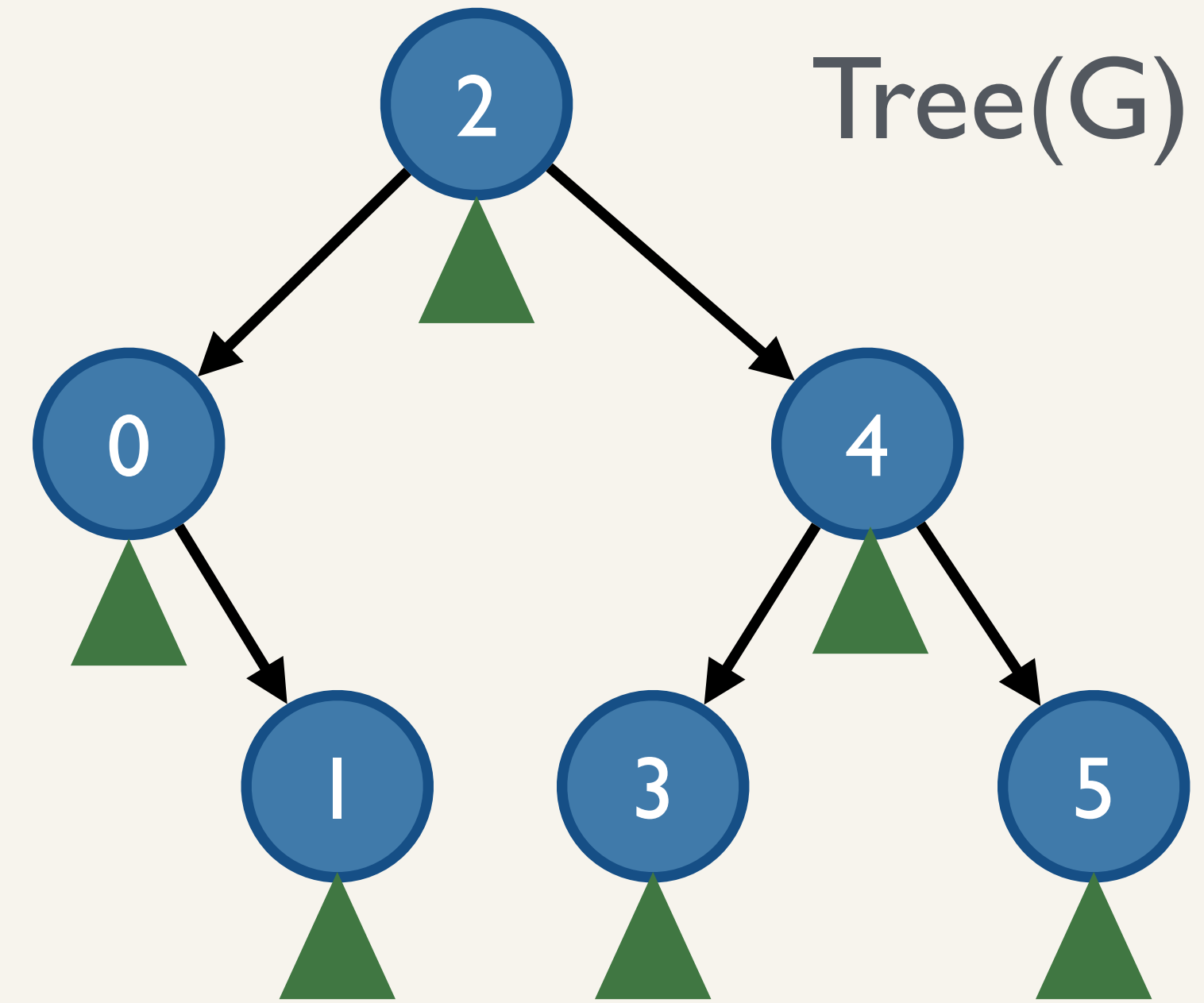
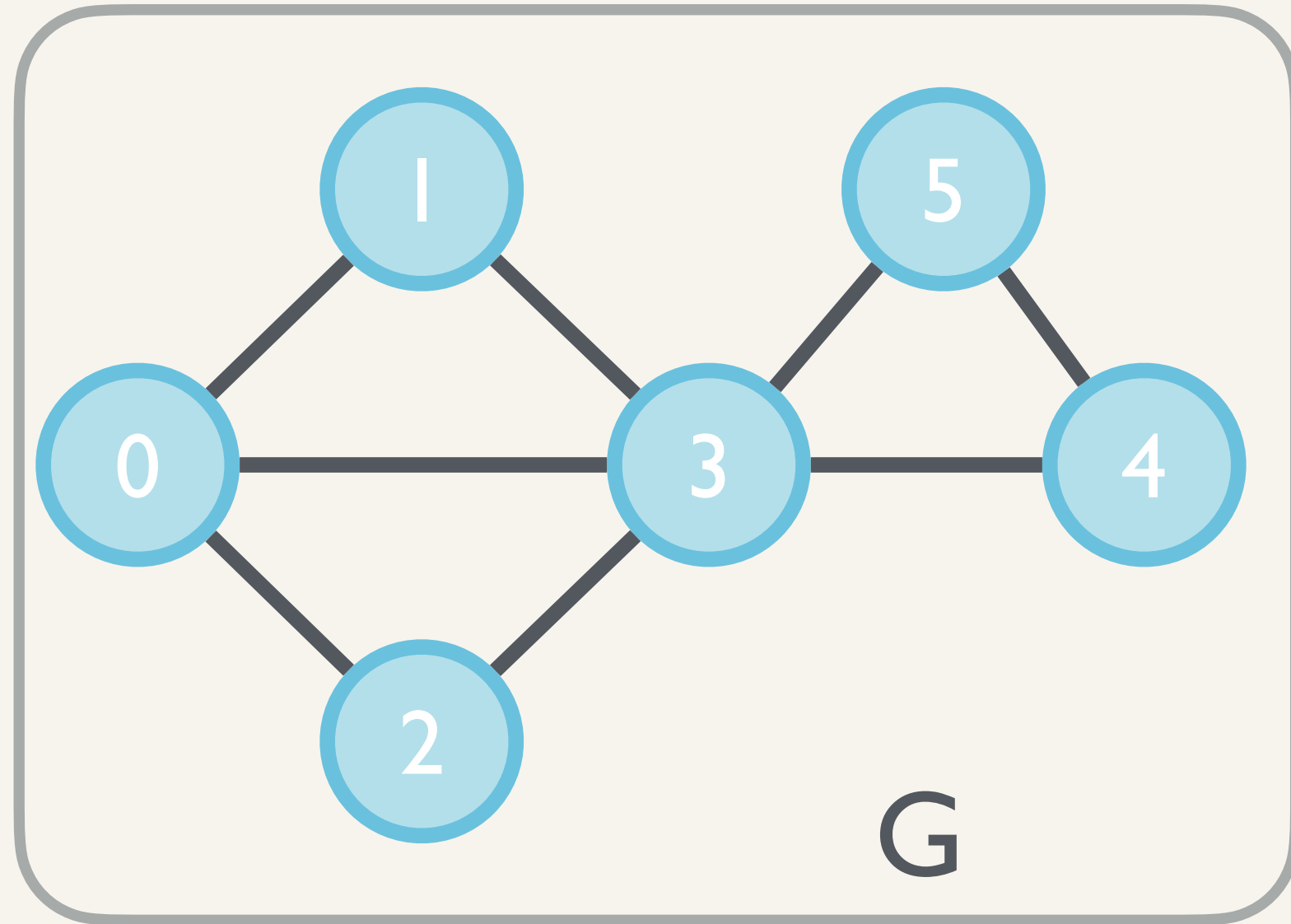


Red-black, AVL, or weight-balanced trees

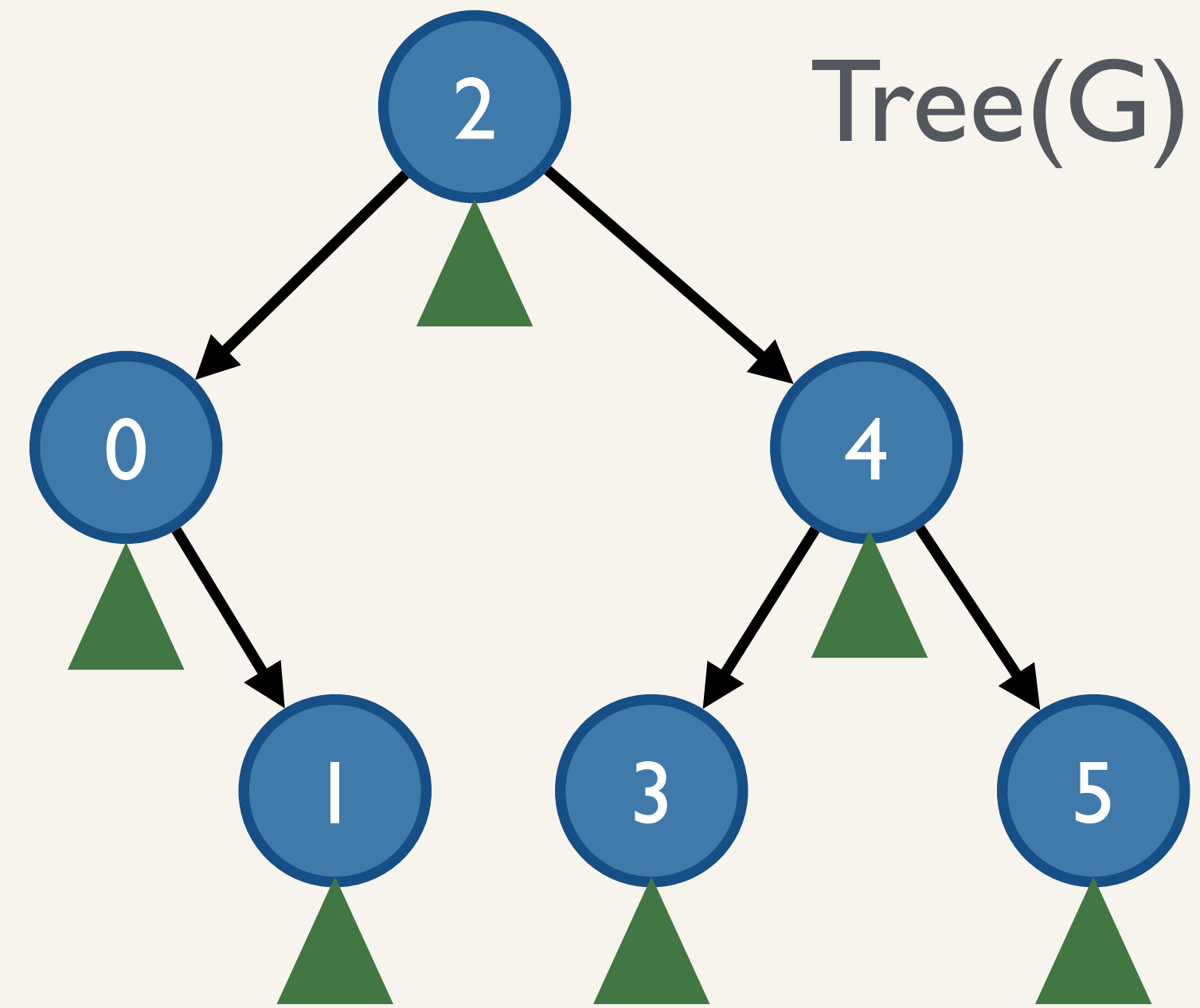
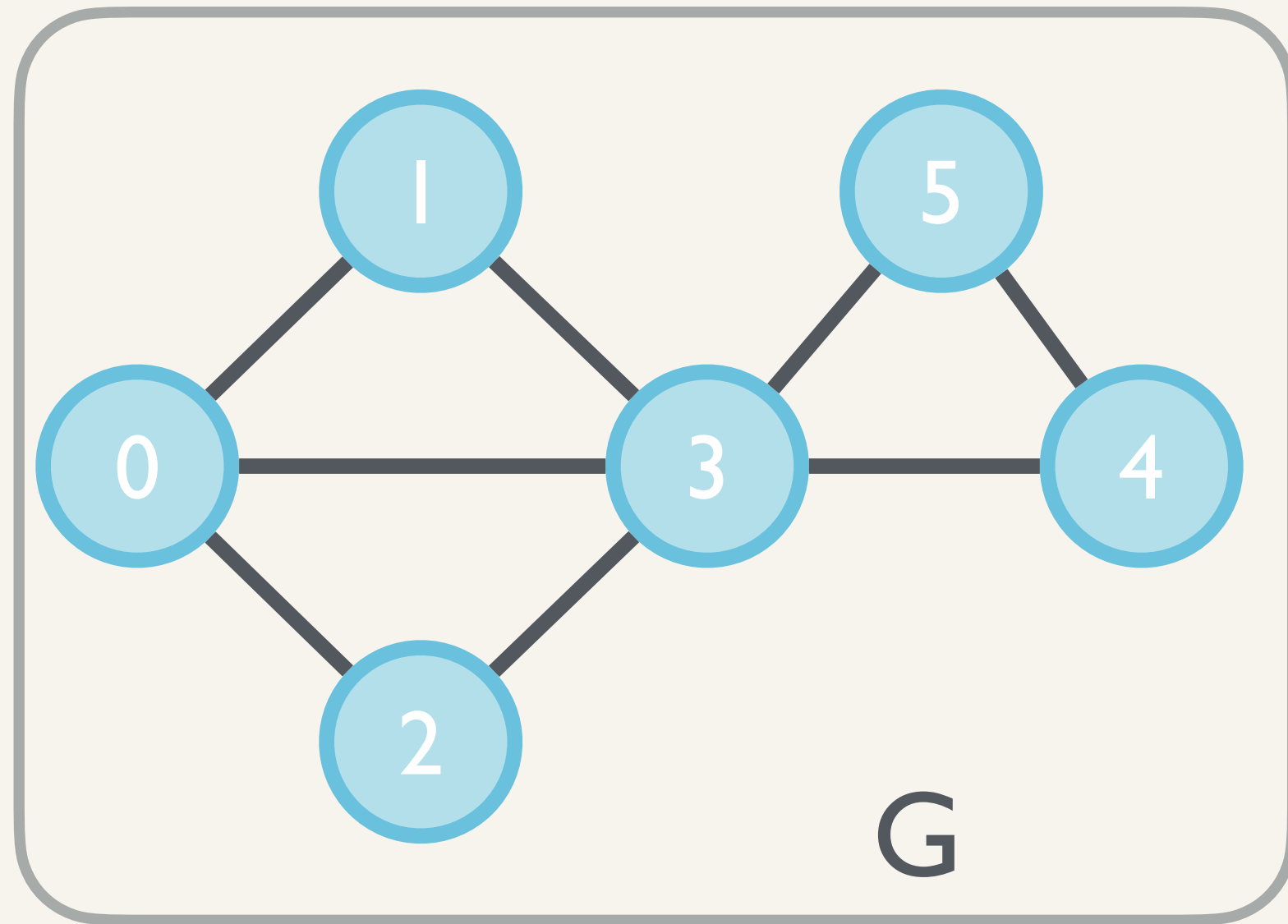
# Representing Graphs using Trees



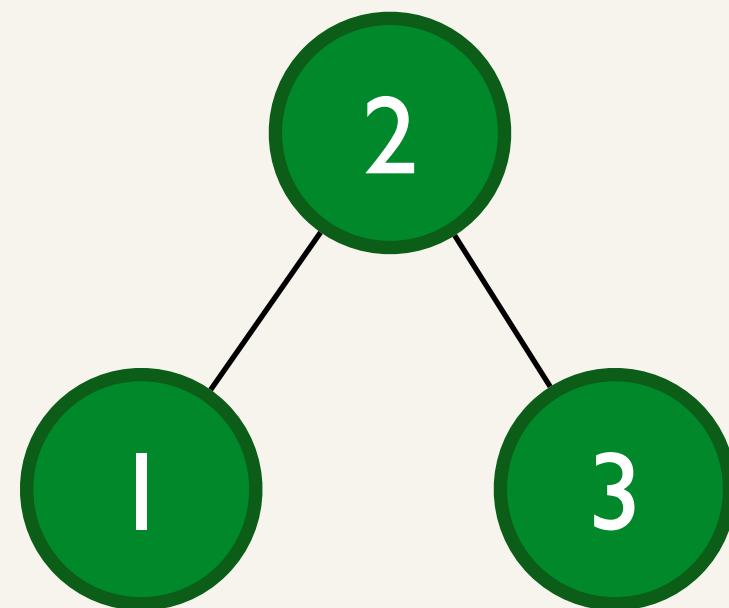
# Representing Graphs using Trees



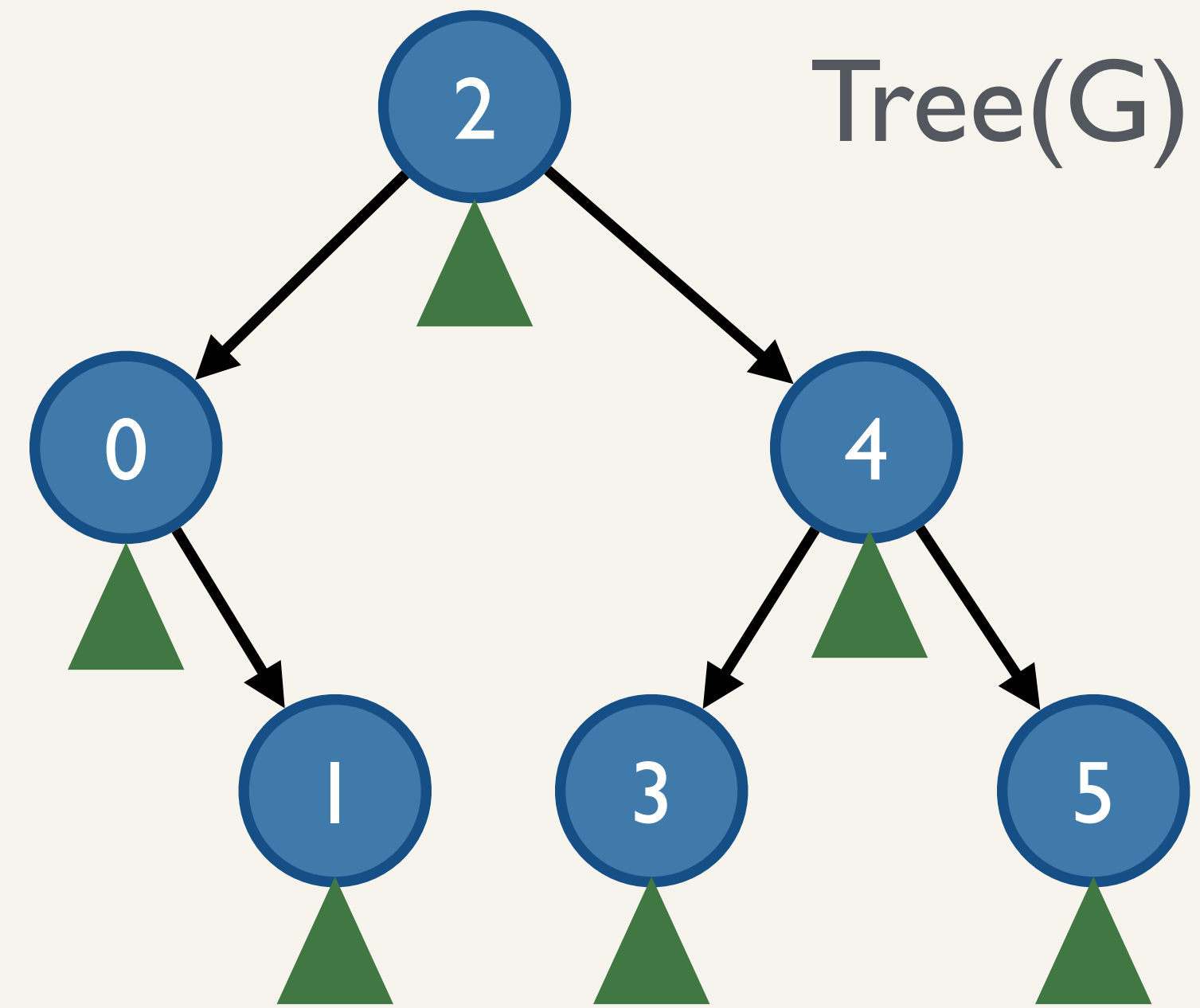
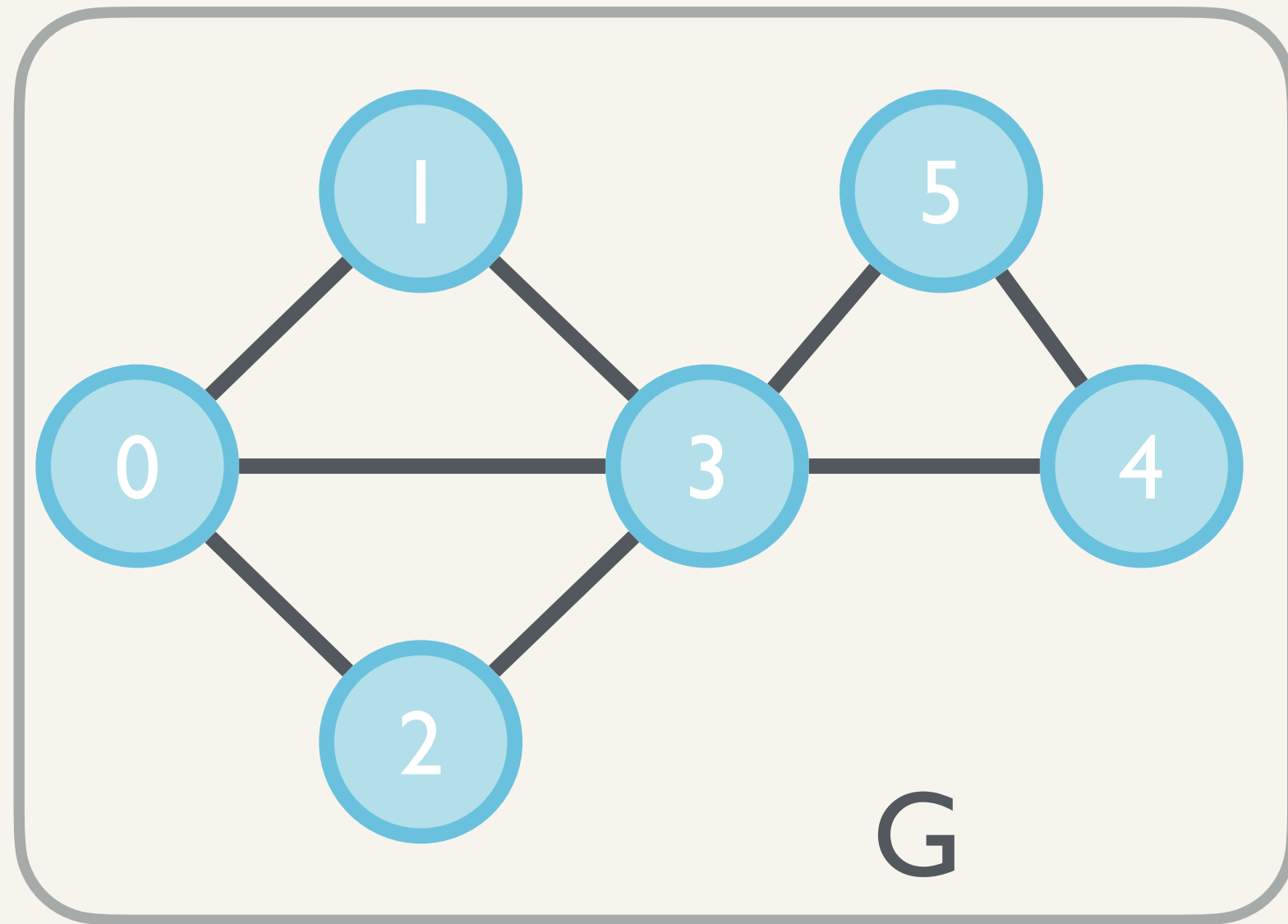
# Representing Graphs using Trees



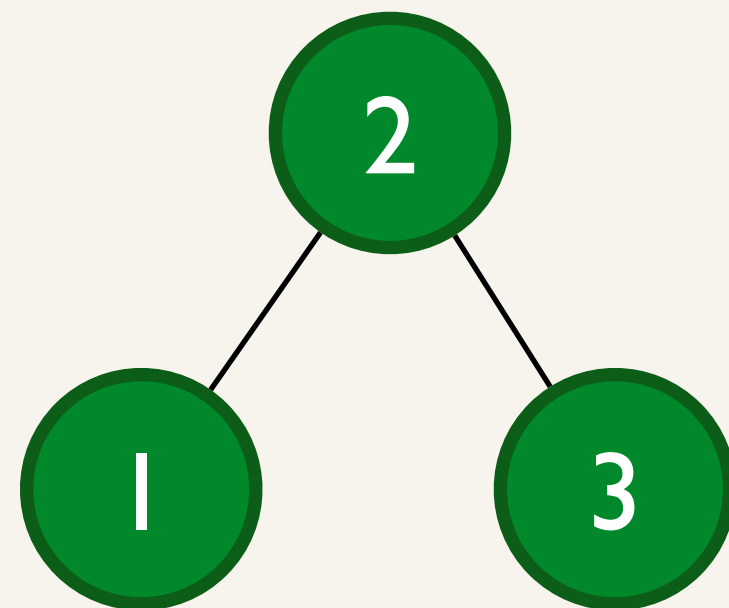
Vertex 0's  
Edge Tree



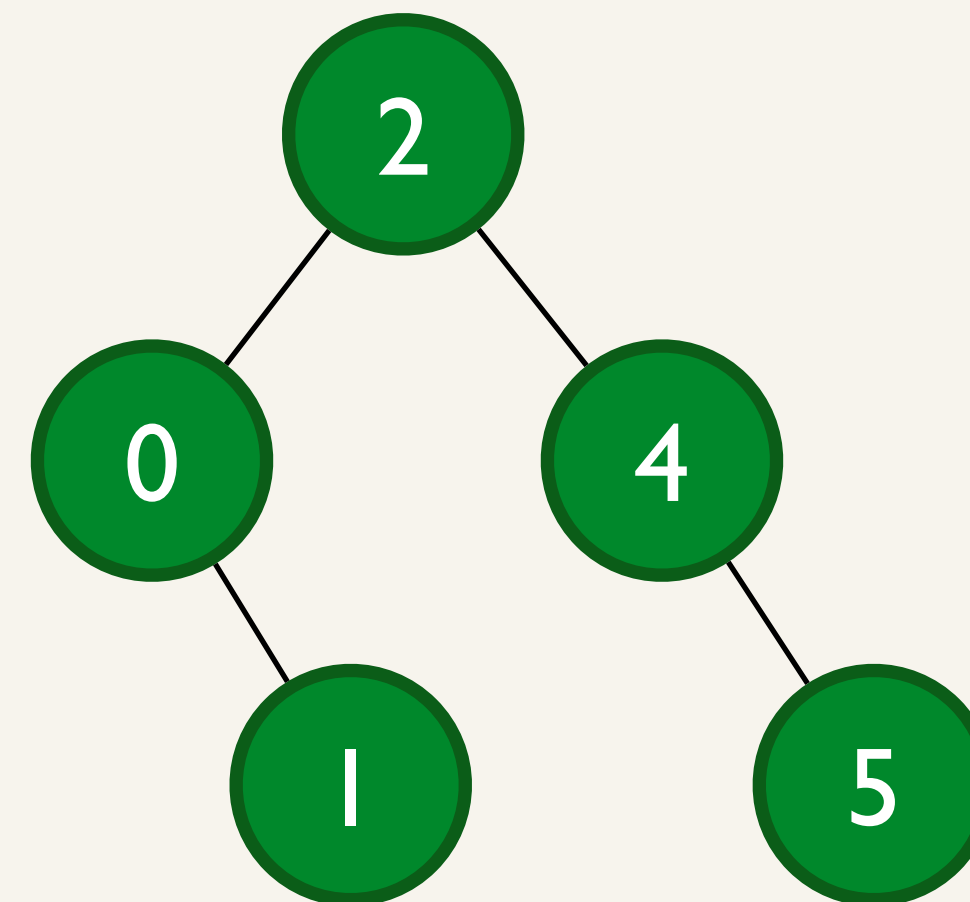
# Representing Graphs using Trees



Vertex 0's  
Edge Tree



Vertex 3's  
Edge Tree





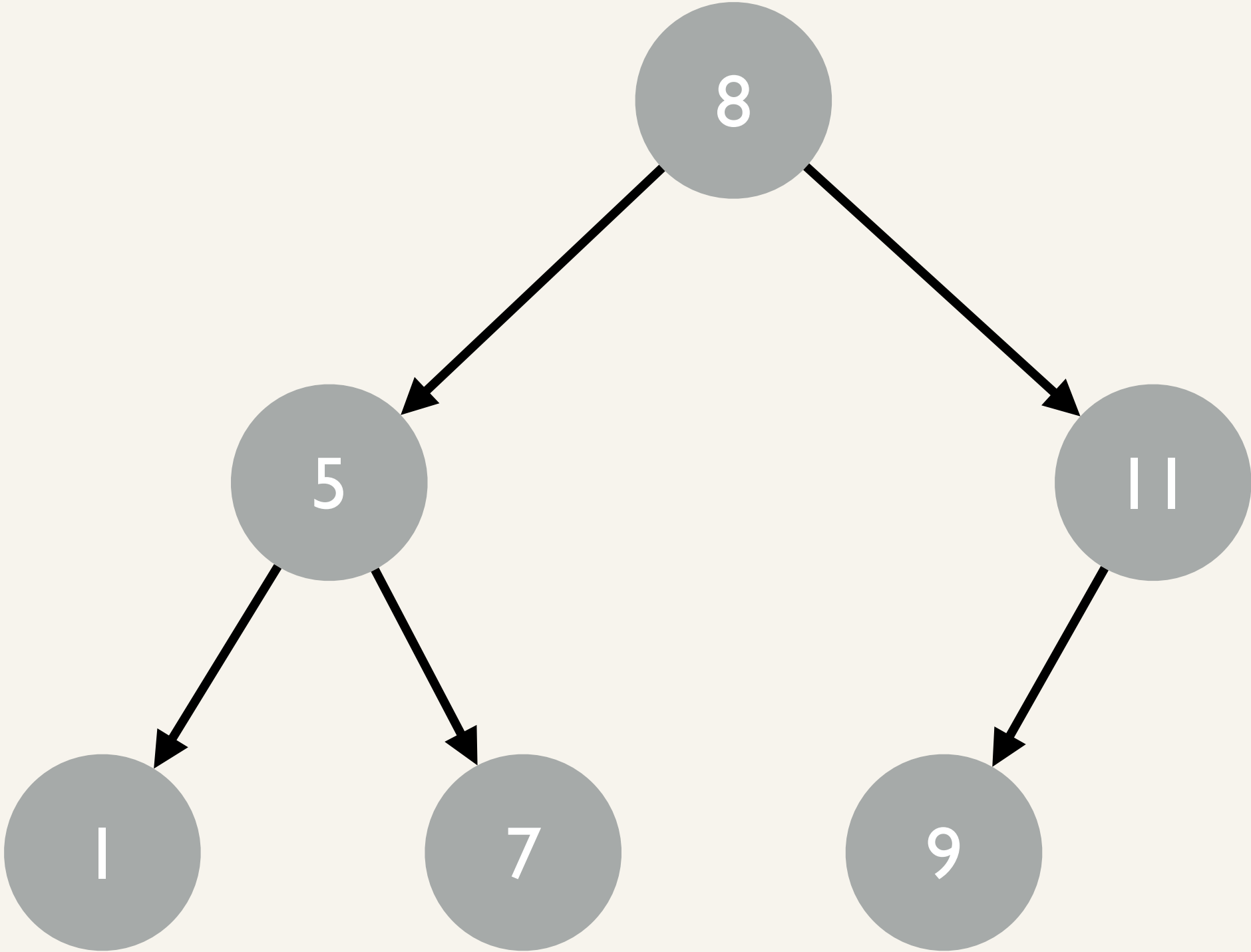
# Trees enable Simple Snapshots

# Trees enable Simple Snapshots

A snapshot is just a tree root

# Trees enable Simple Snapshots

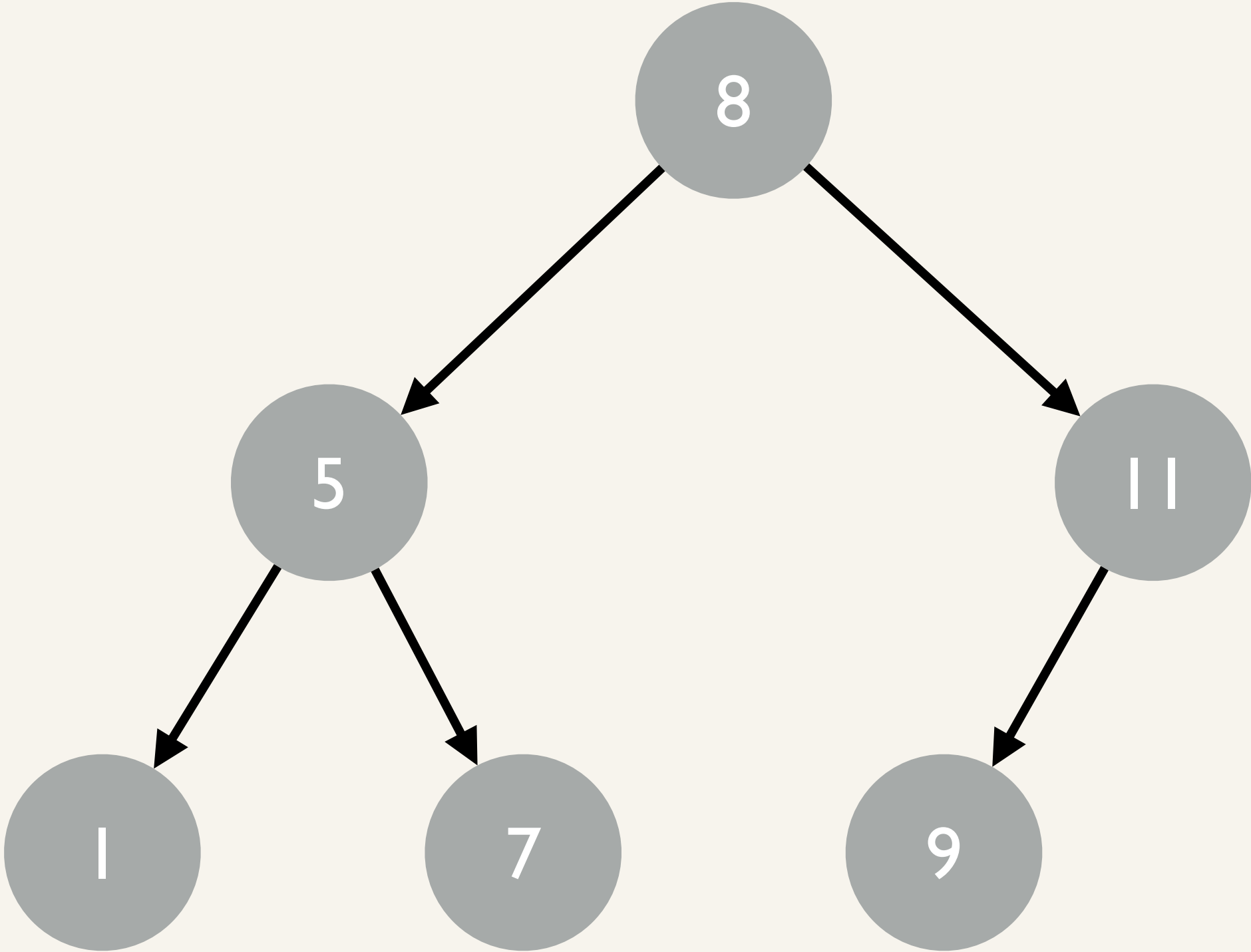
A snapshot is just a tree root



# Trees enable Simple Snapshots

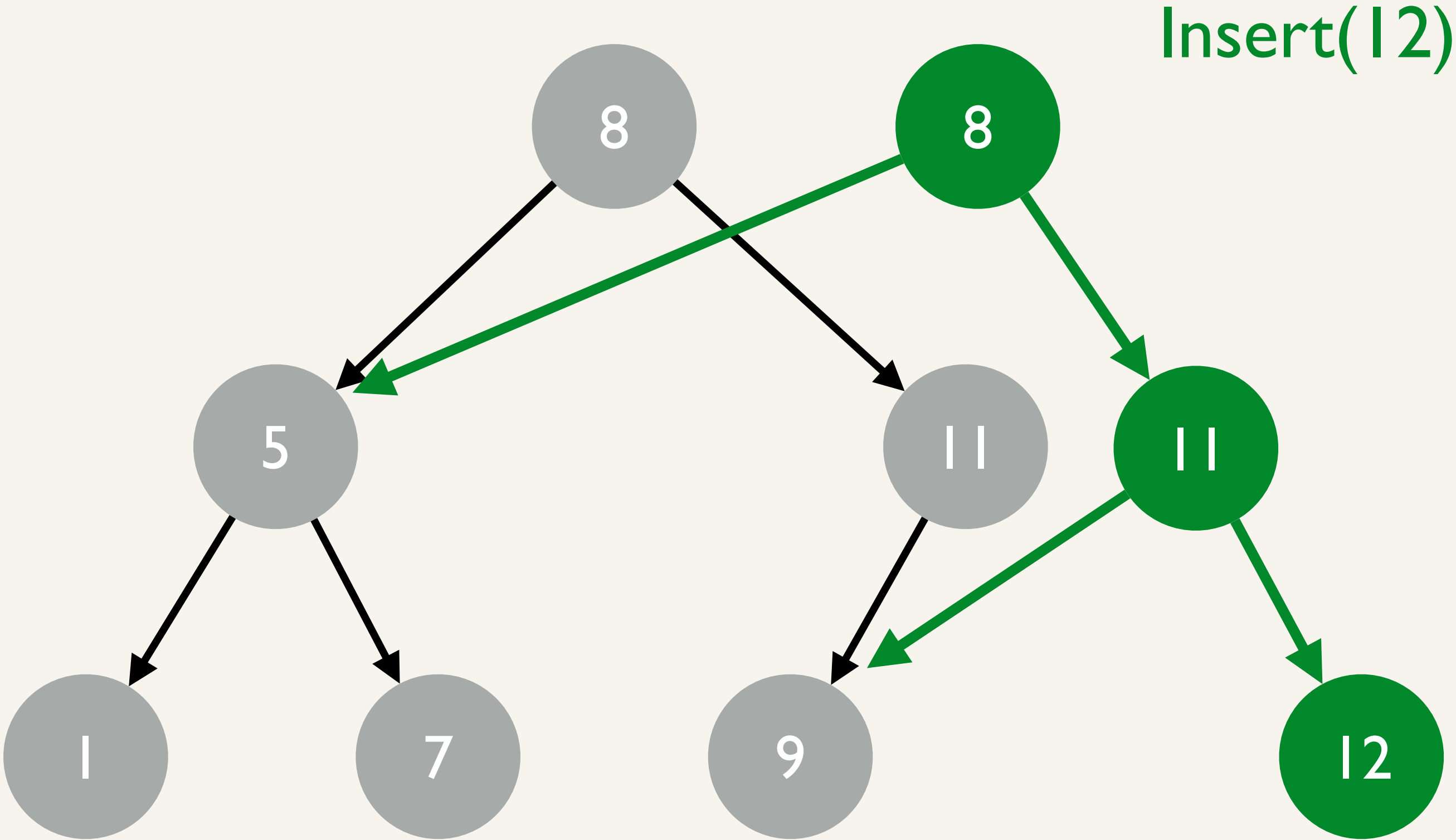
A snapshot is just a tree root

Insert(12)



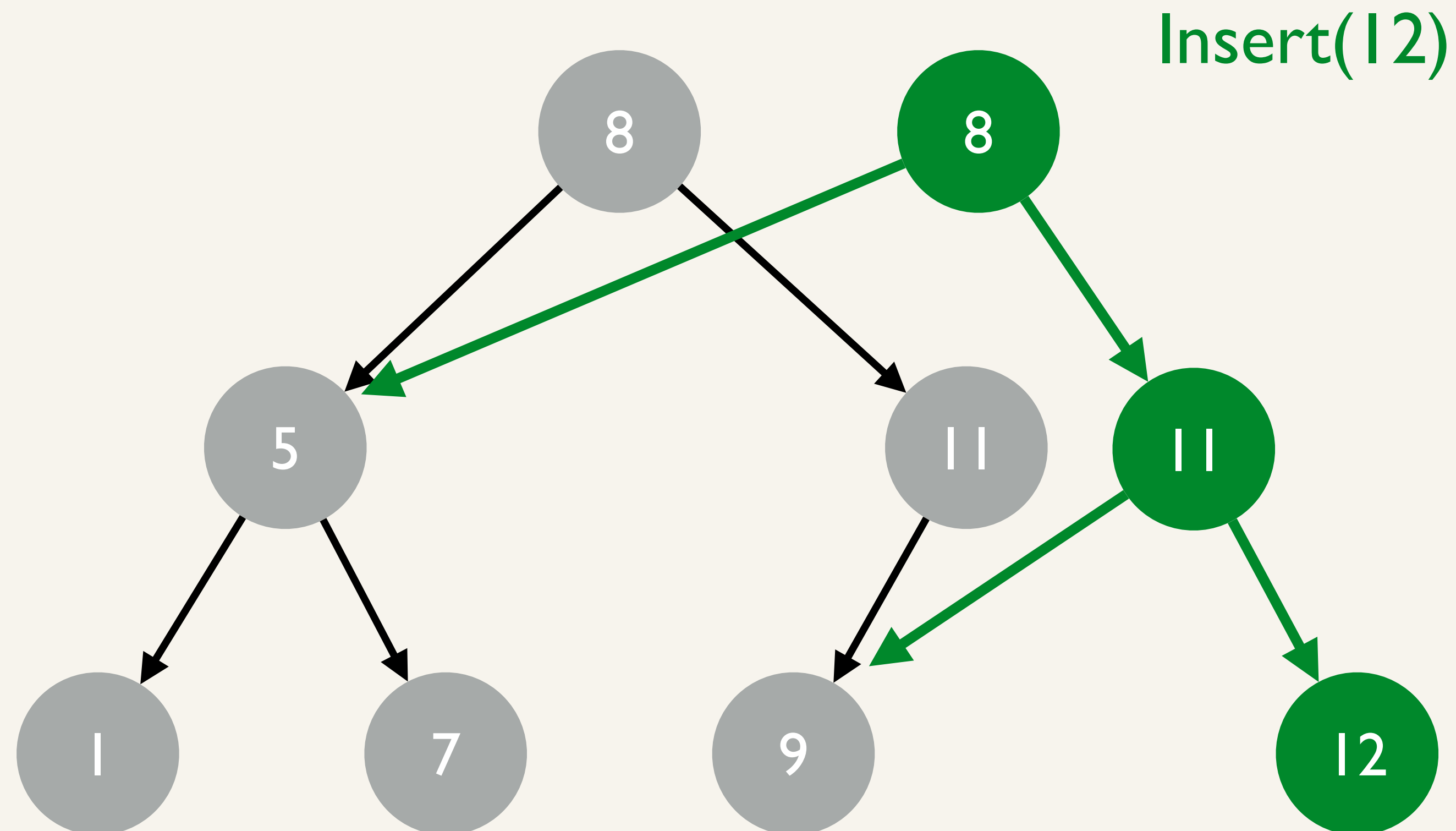
# Trees enable Simple Snapshots

A snapshot is just a tree root



# Trees enable Simple Snapshots

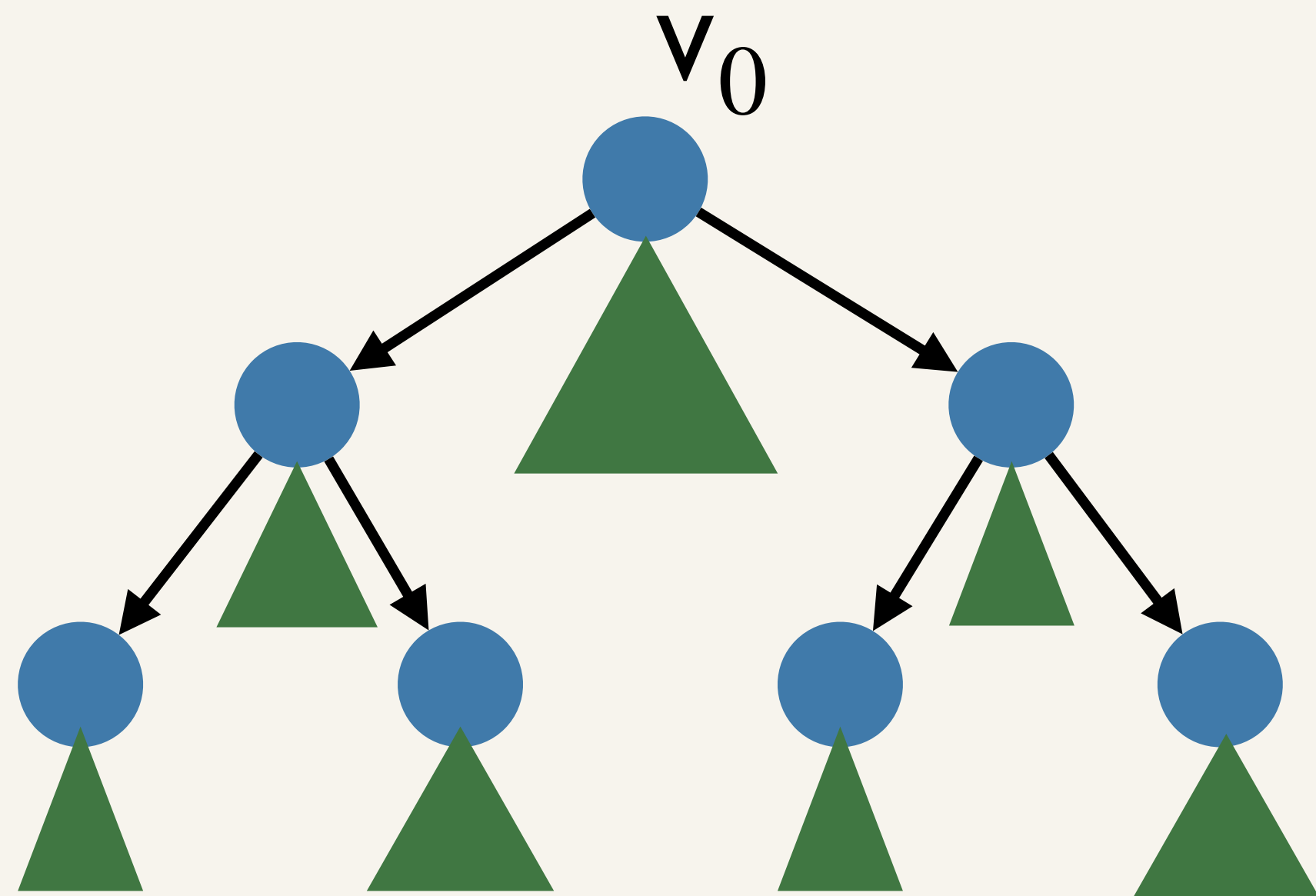
A snapshot is just a tree root



Algorithms generalize to handle *batches* of updates in low work/depth [BFS'16]

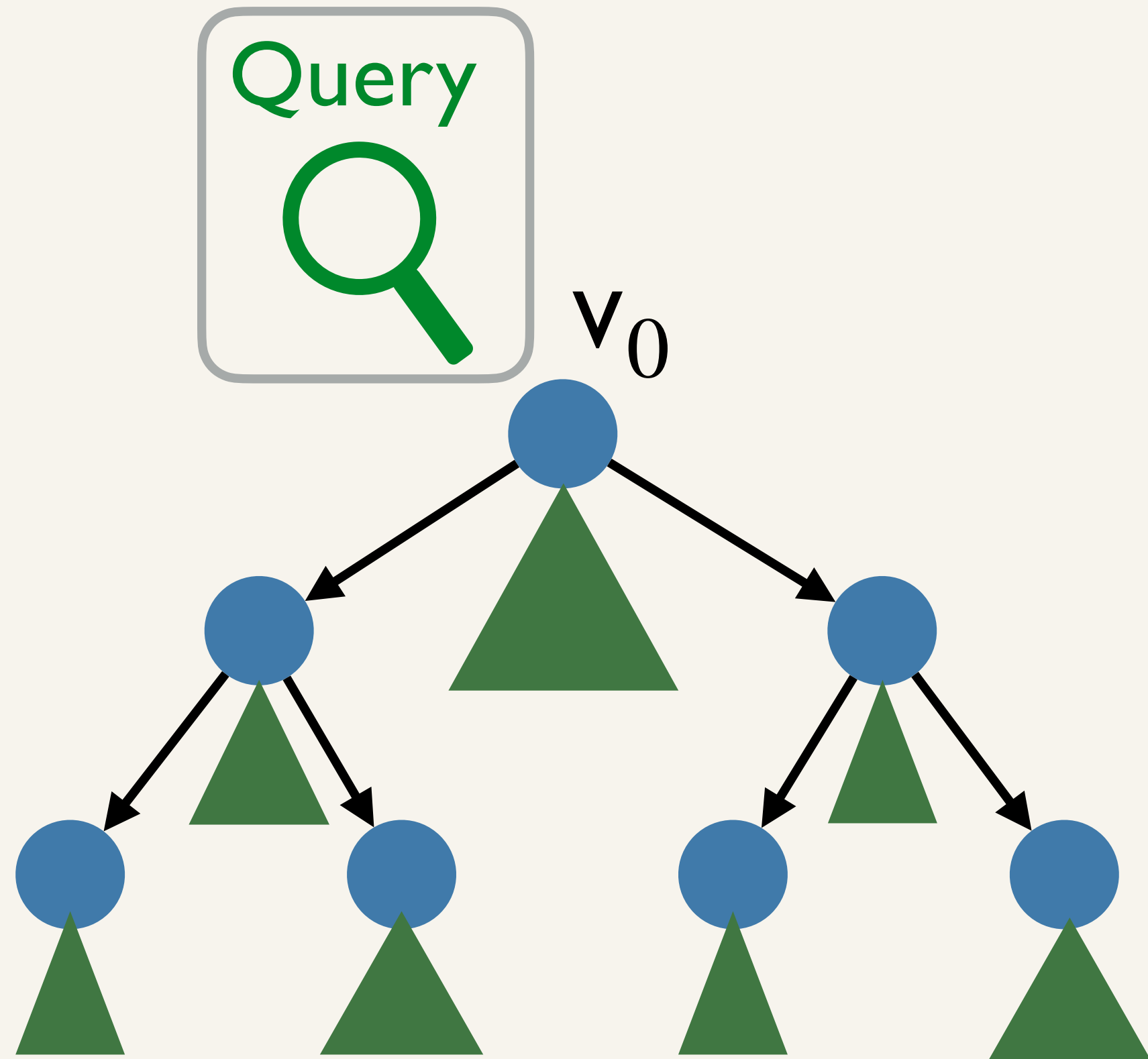
# Purely-Functional Trees are Safe for Concurrency

# Purely-Functional Trees are Safe for Concurrency

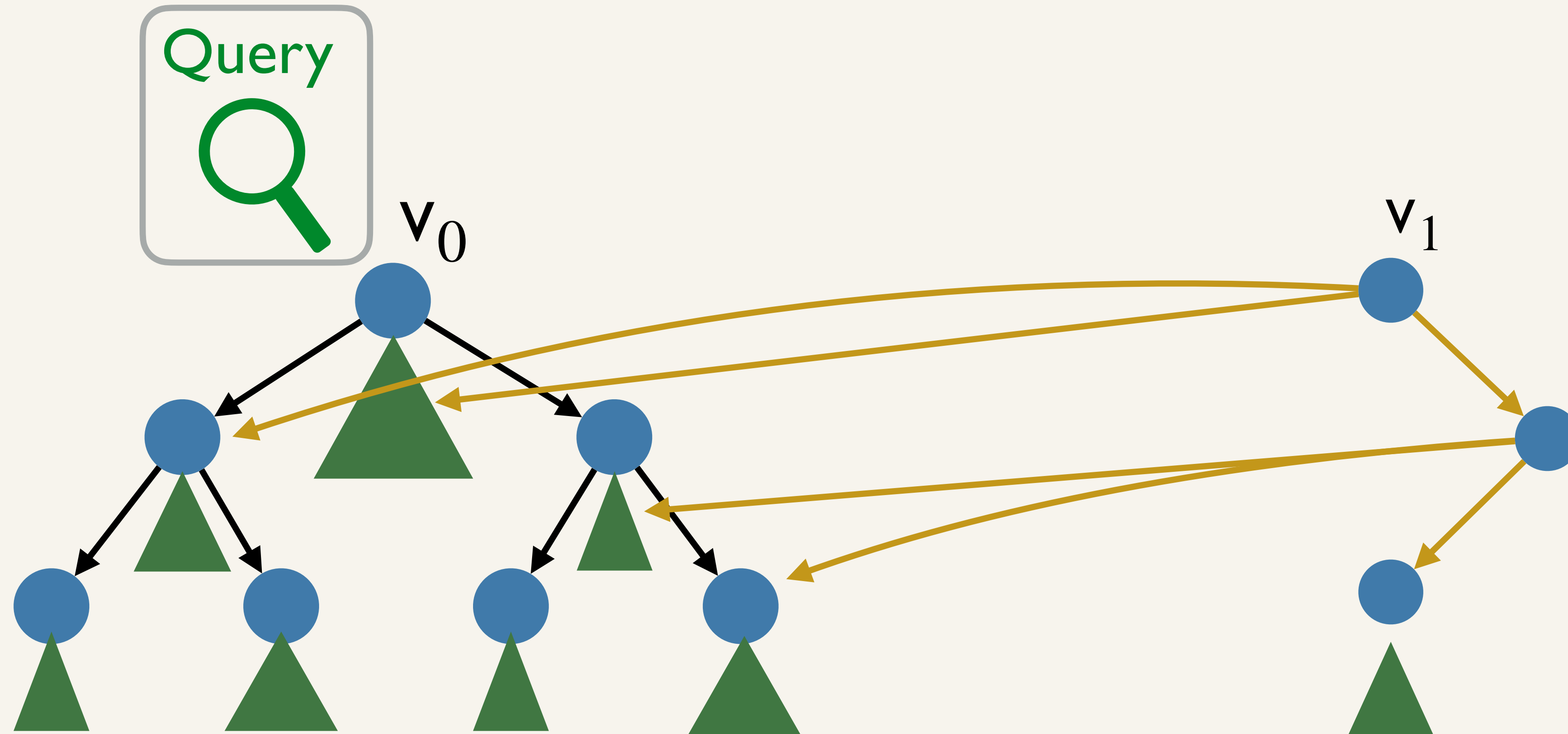




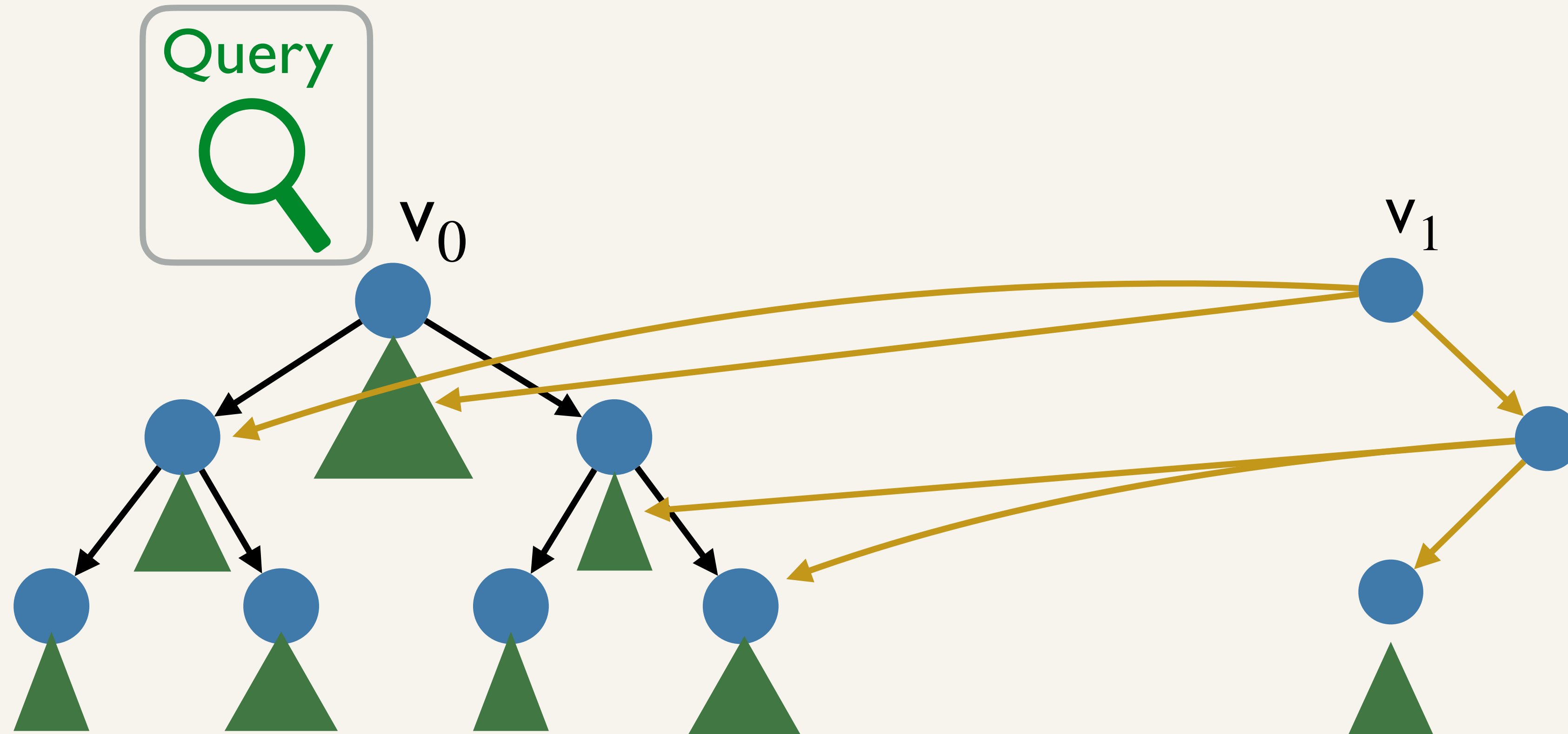
# Purely-Functional Trees are Safe for Concurrency



# Purely-Functional Trees are Safe for Concurrency



# Purely-Functional Trees are Safe for Concurrency



Queries are serialized once they  
acquire a tree root

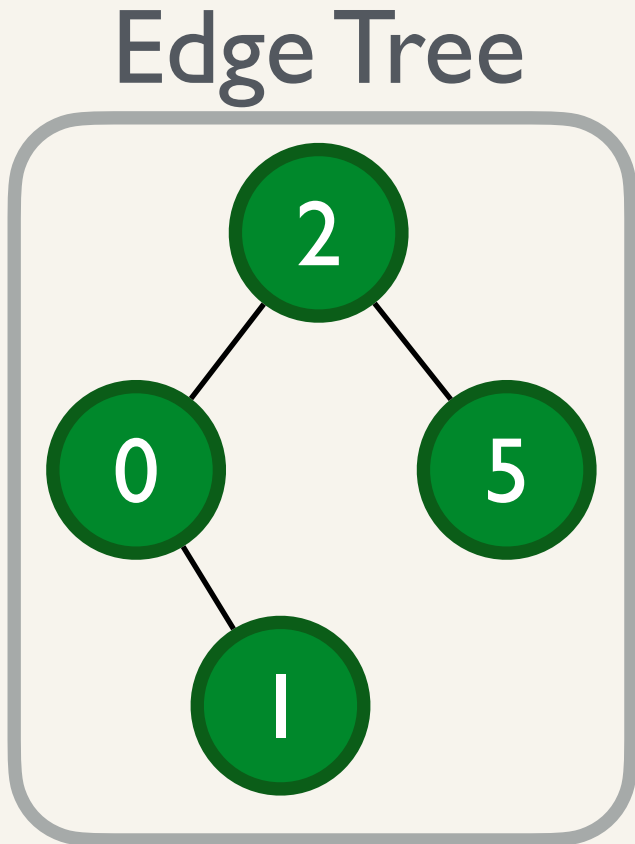
# Challenges

# Challenges

Poor Cache Usage

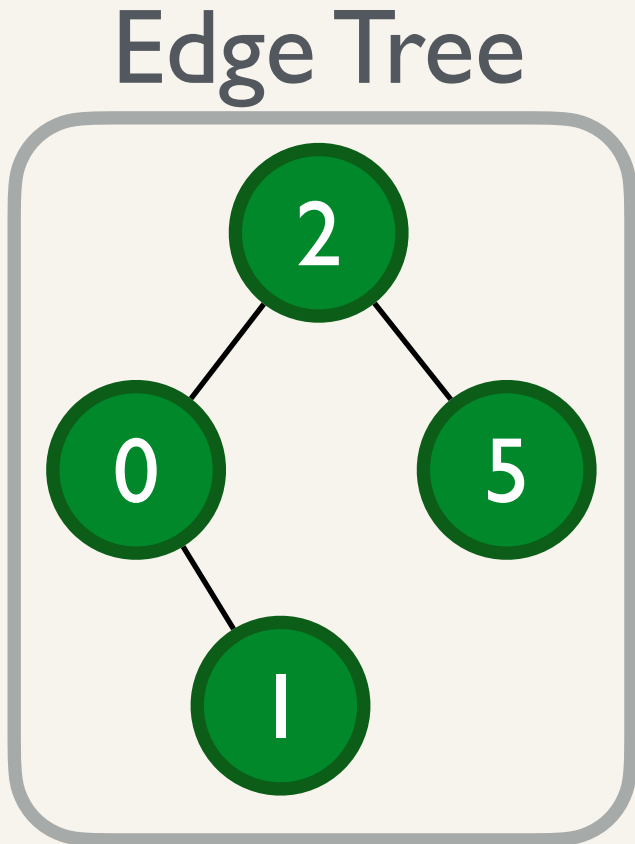
# Challenges

## Poor Cache Usage



# Challenges

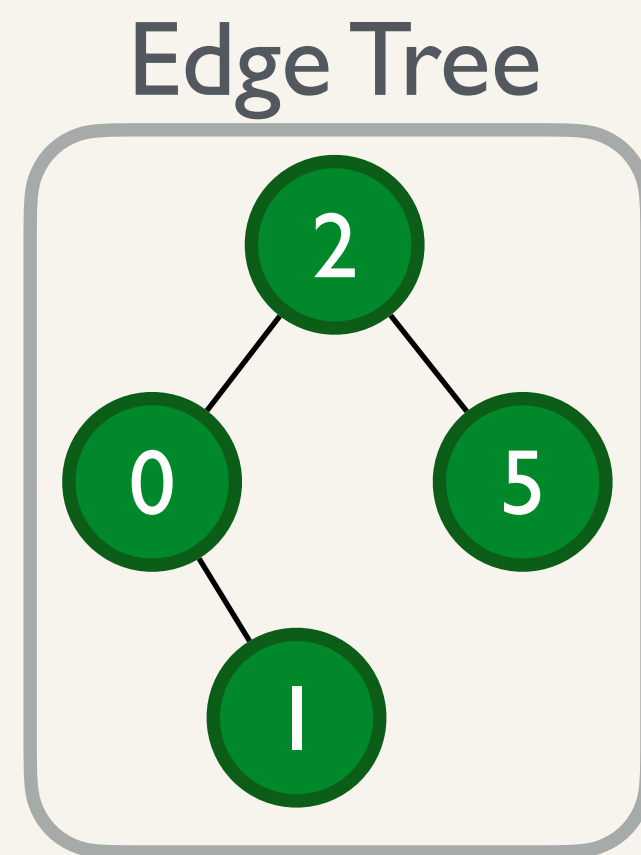
## Poor Cache Usage



## Space Inefficiency

# Challenges

## Poor Cache Usage



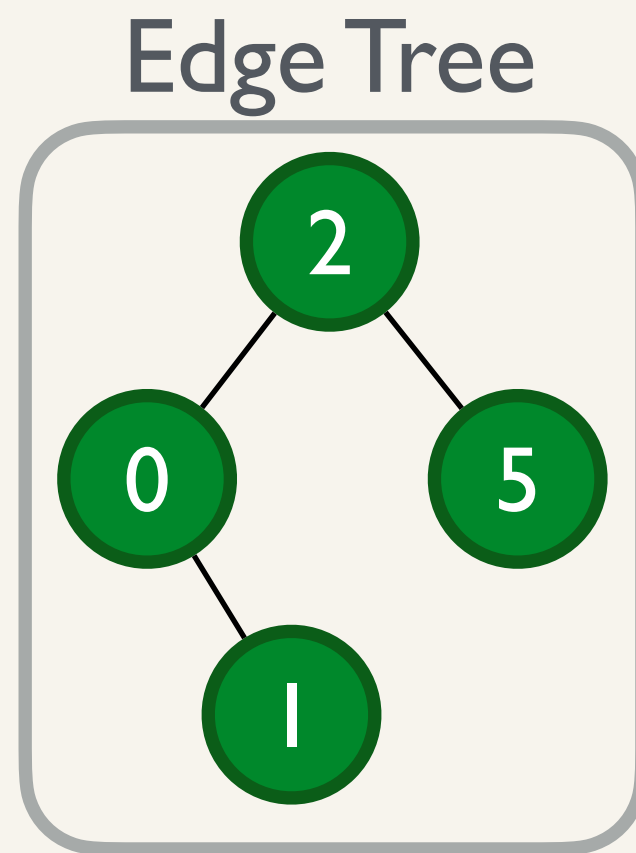
## Space Inefficiency

- ❖ Significant space overheads for tree nodes
- ❖ Lose ability to compress adjacency lists



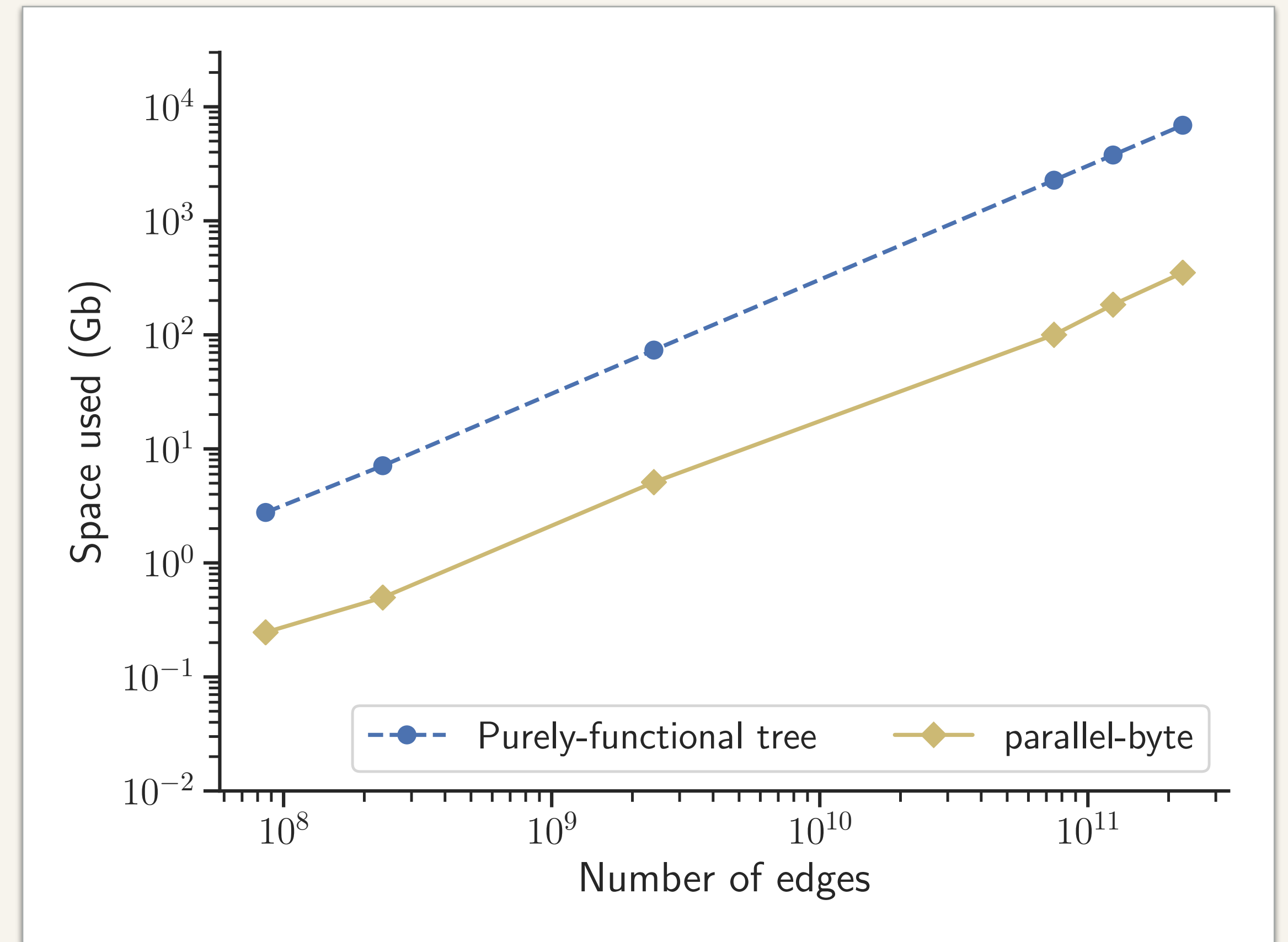
# Challenges

## Poor Cache Usage



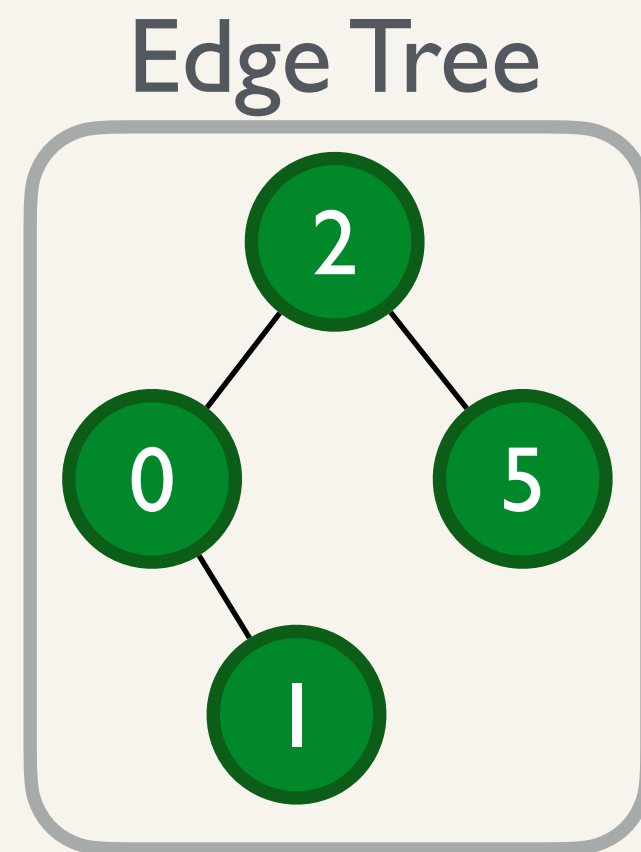
## Space Inefficiency

- ❖ Significant space overheads for tree nodes
- ❖ Lose ability to compress adjacency lists



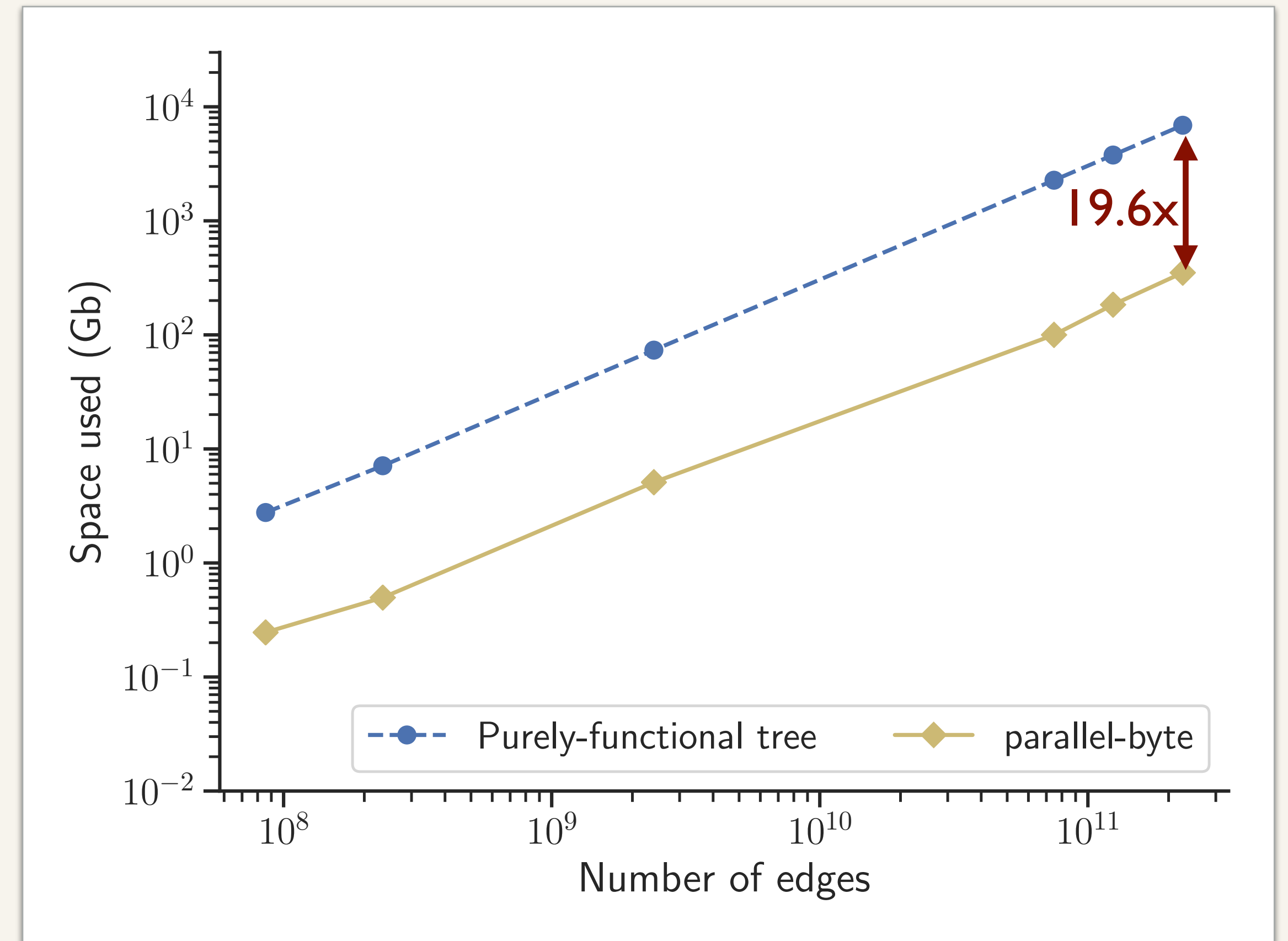
# Challenges

## Poor Cache Usage



## Space Inefficiency

- ❖ Significant space overheads for tree nodes
- ❖ Lose ability to compress adjacency lists



Aspen using naive approach requires  
7 TB of memory for WDC2012 graph

# Challenges

## Poor Cache Usage



To overcome these challenges we designed  
**C-trees: compressed purely-functional trees**

## Space Inefficiency

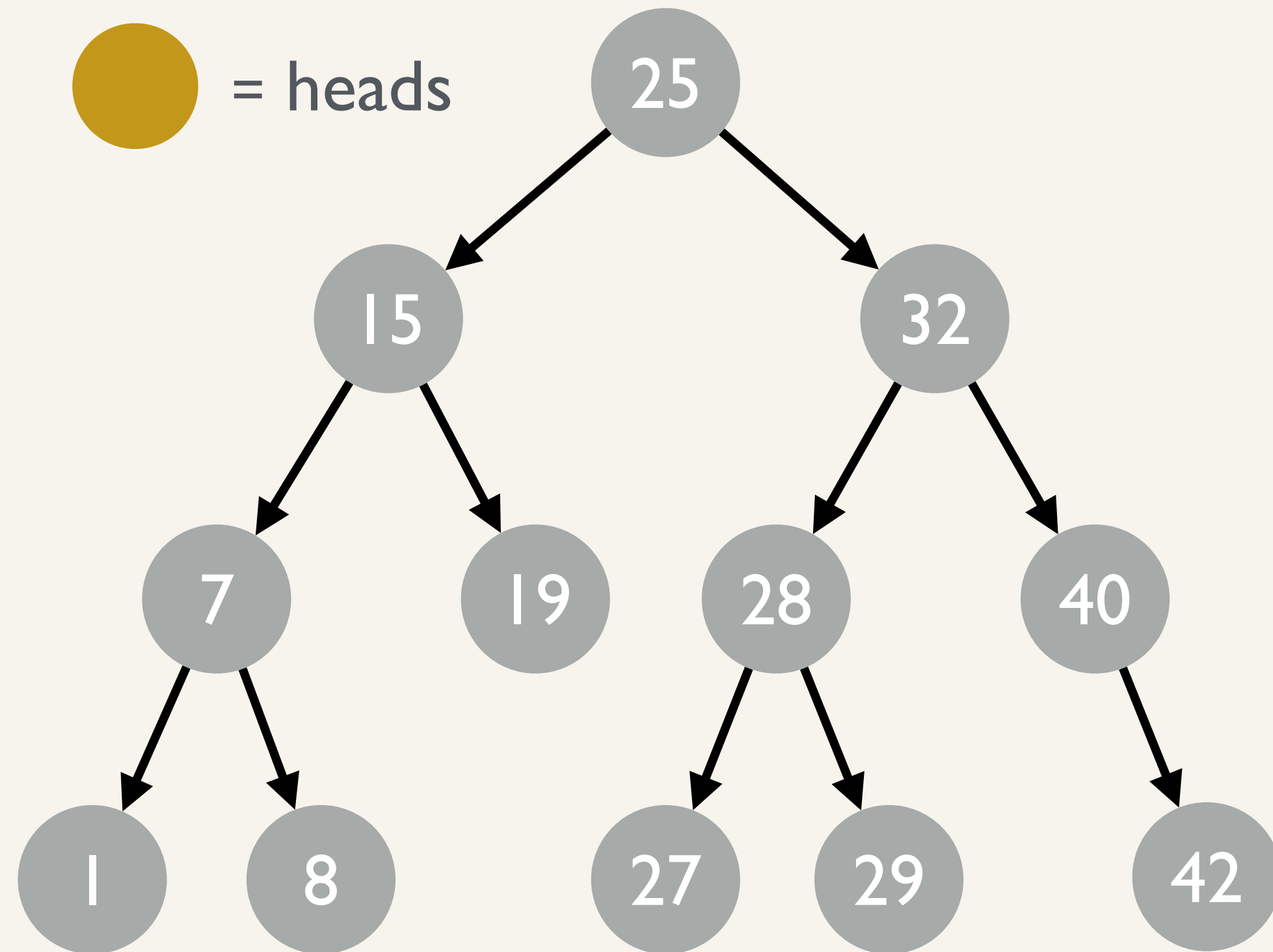
- ❖ Significant space overheads for tree nodes
- ❖ Lose ability to compress adjacency lists



Aspen using naive approach requires  
7 TB of memory for WDC2012 graph

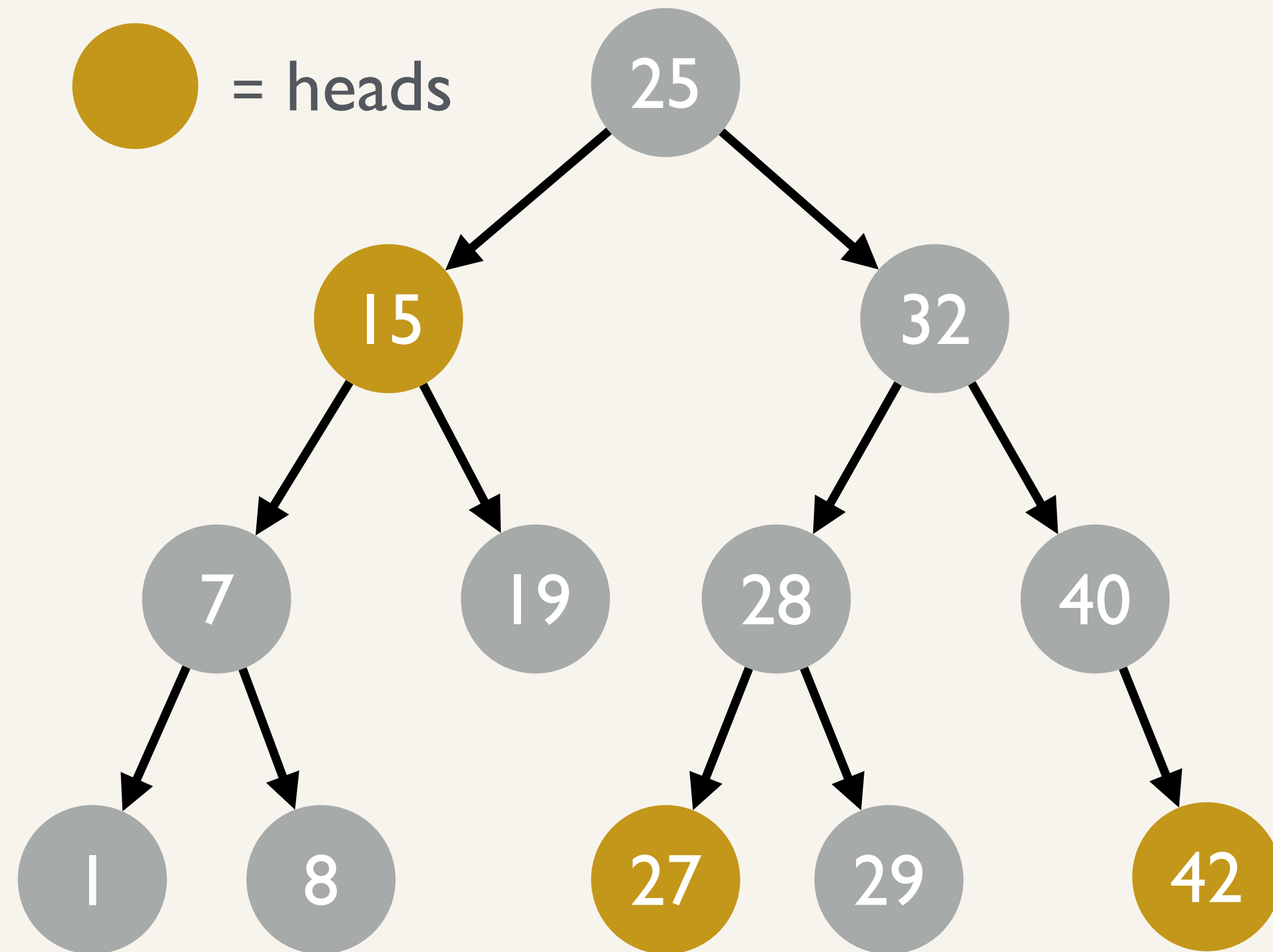
# C-trees

- ❖ Chunking parameter  $B$ . Fix a hash function,  $h$
- ❖ Select elements as *heads* with probability  $1/B$  using  $h$



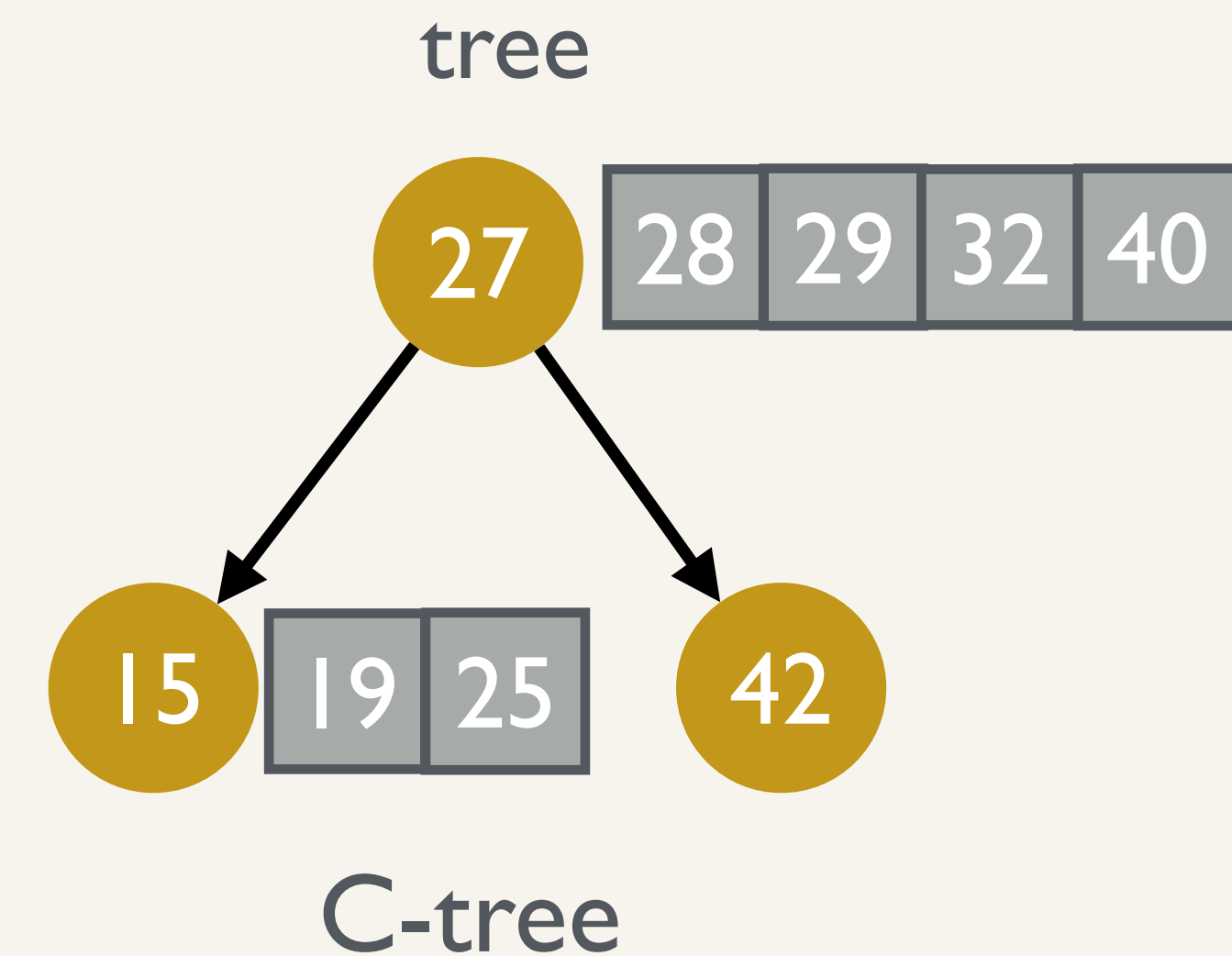
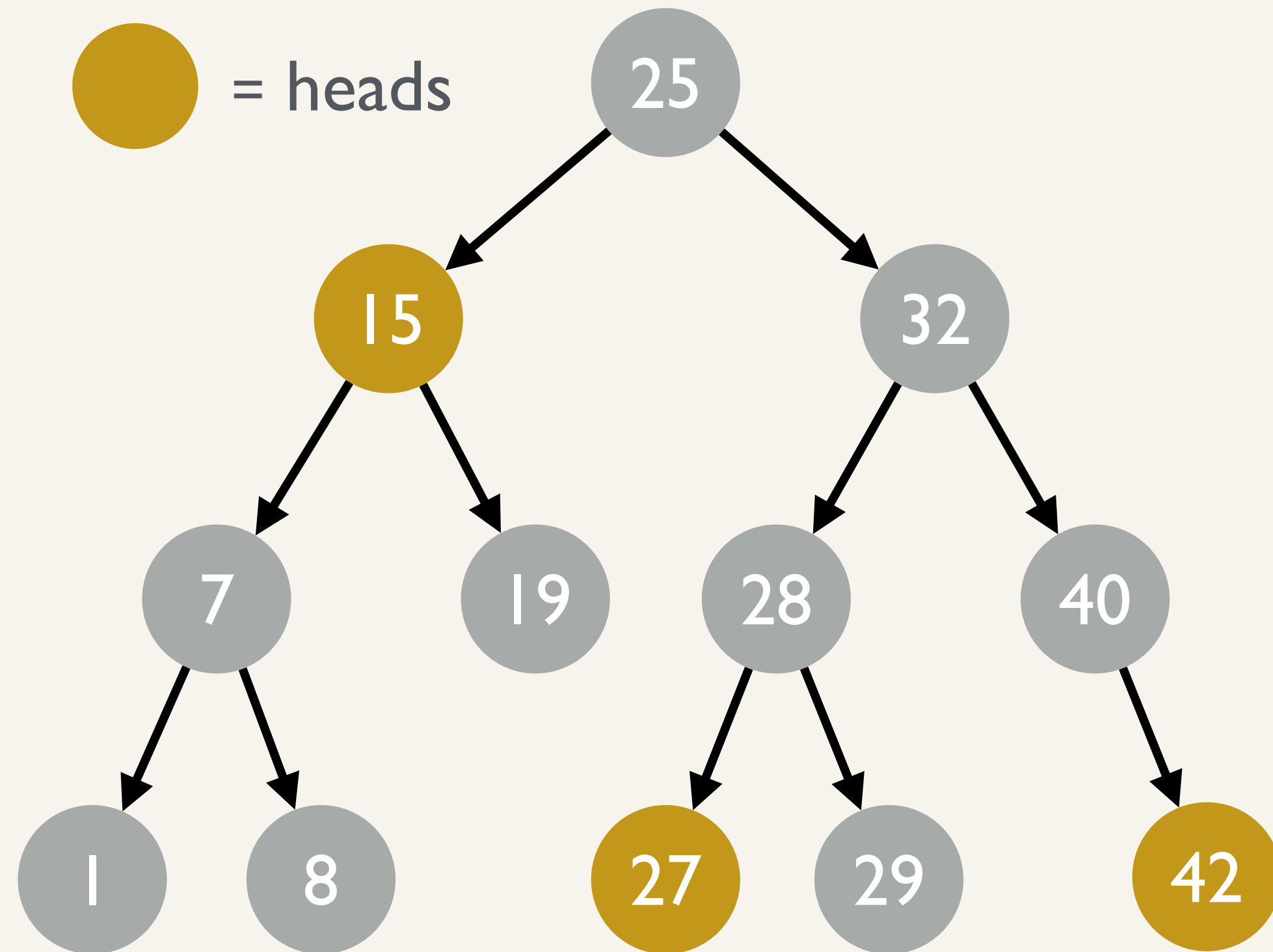
# C-trees

- ❖ Chunking parameter  $B$ . Fix a hash function,  $h$
- ❖ Select elements as *heads* with probability  $1/B$  using  $h$



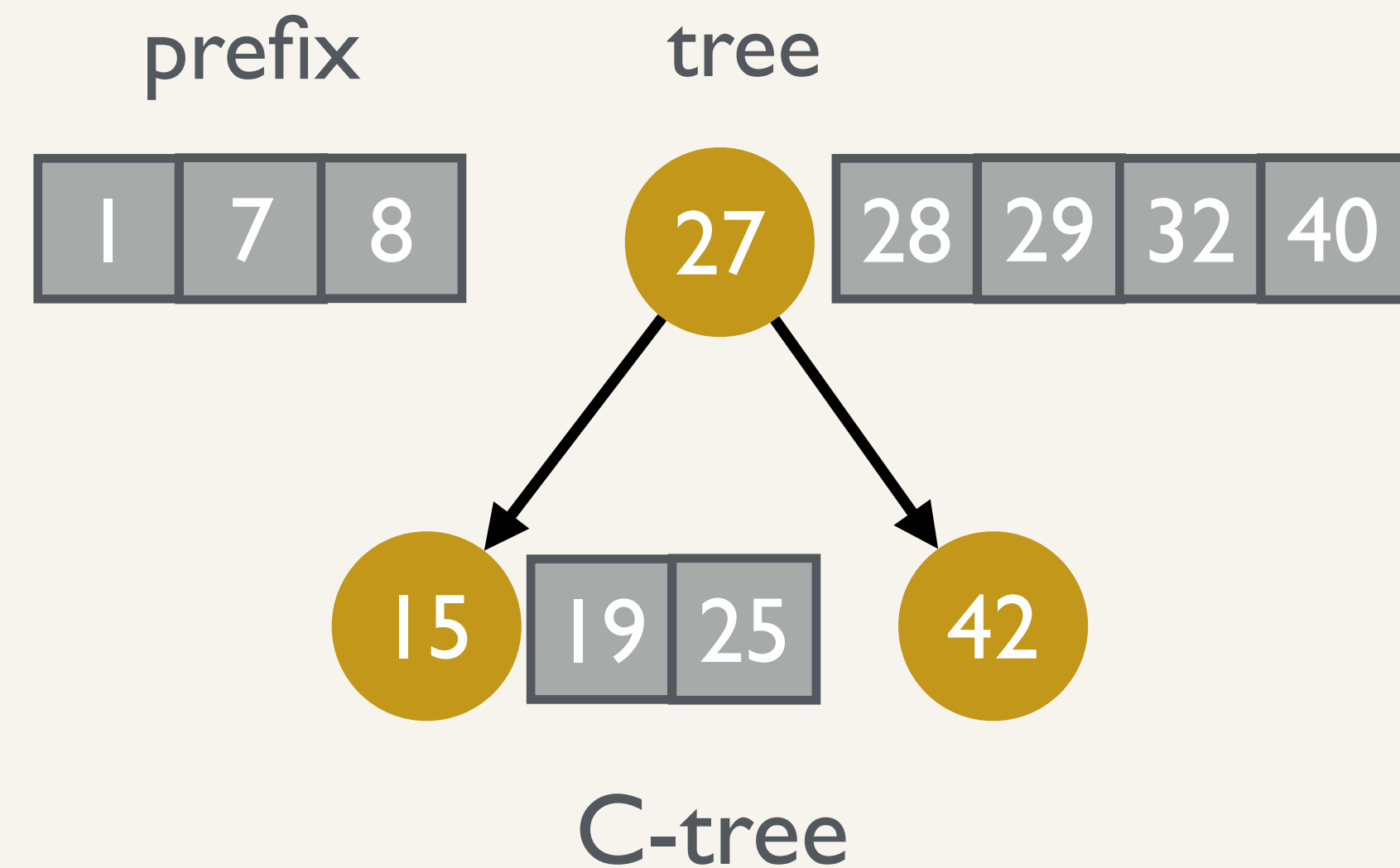
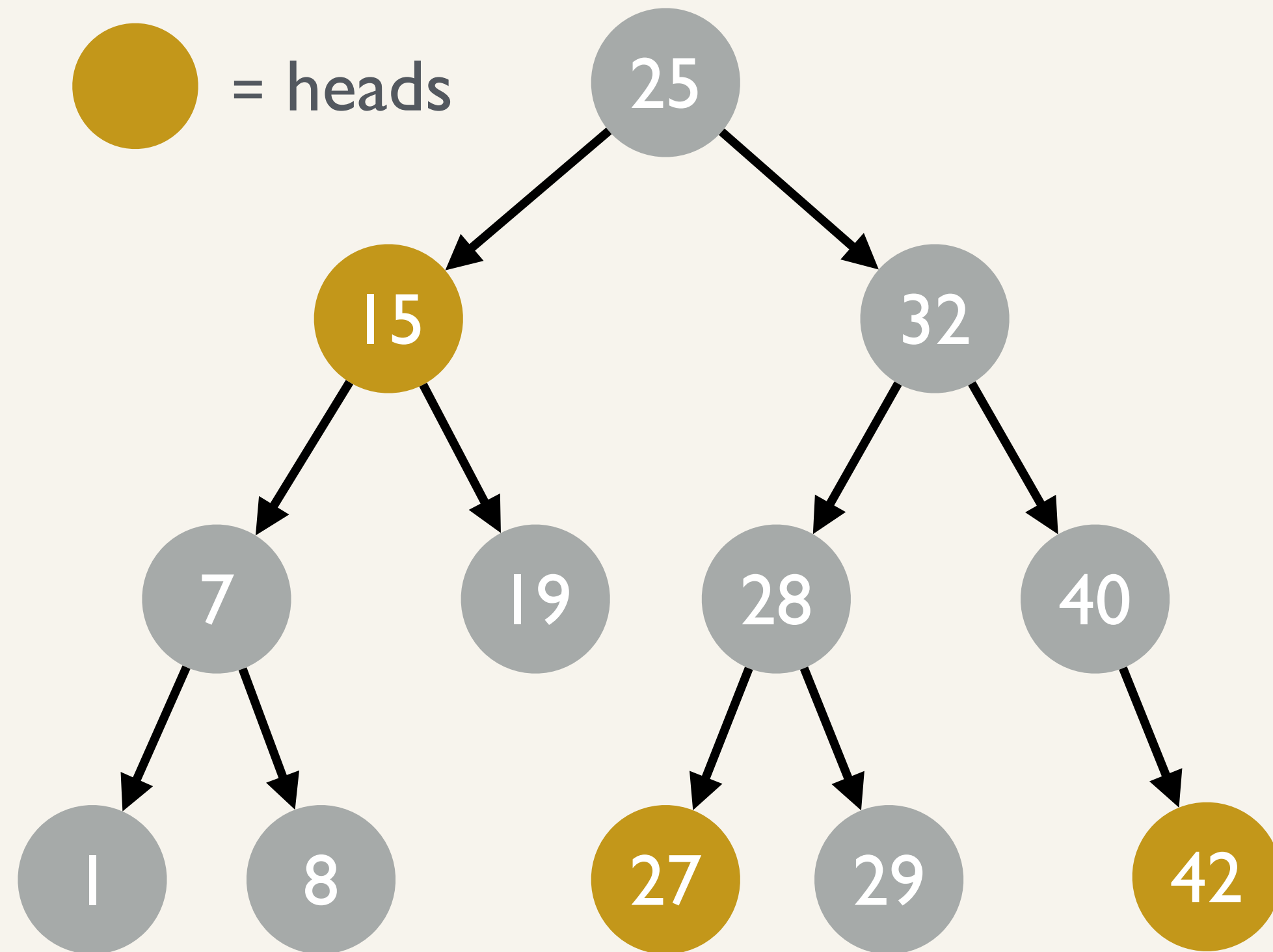
# C-trees

- ❖ Chunking parameter  $B$ . Fix a hash function,  $h$
- ❖ Select elements as *heads* with probability  $1/B$  using  $h$



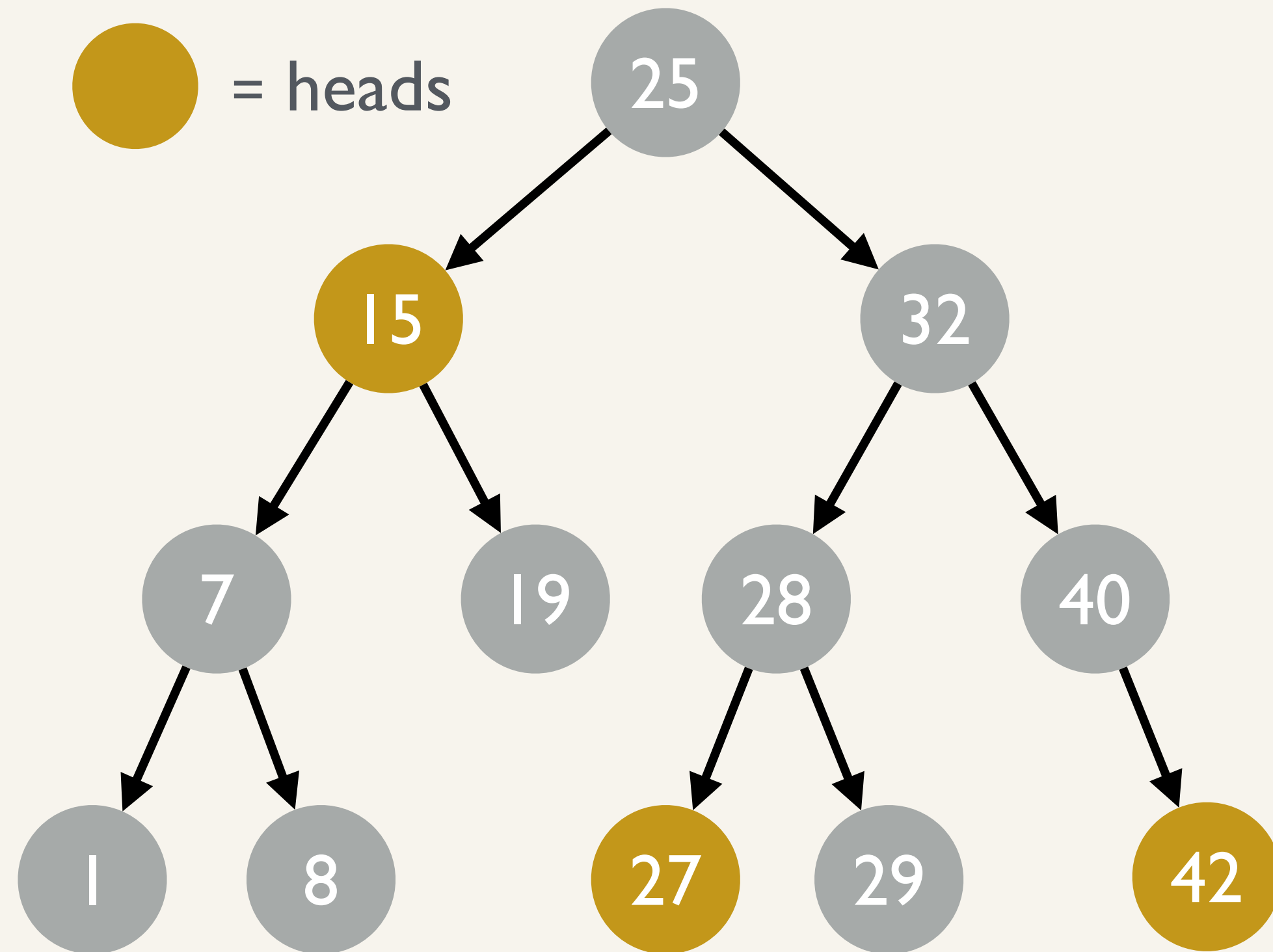
# C-trees

- ❖ Chunking parameter  $B$ . Fix a hash function,  $h$
- ❖ Select elements as *heads* with probability  $1/B$  using  $h$

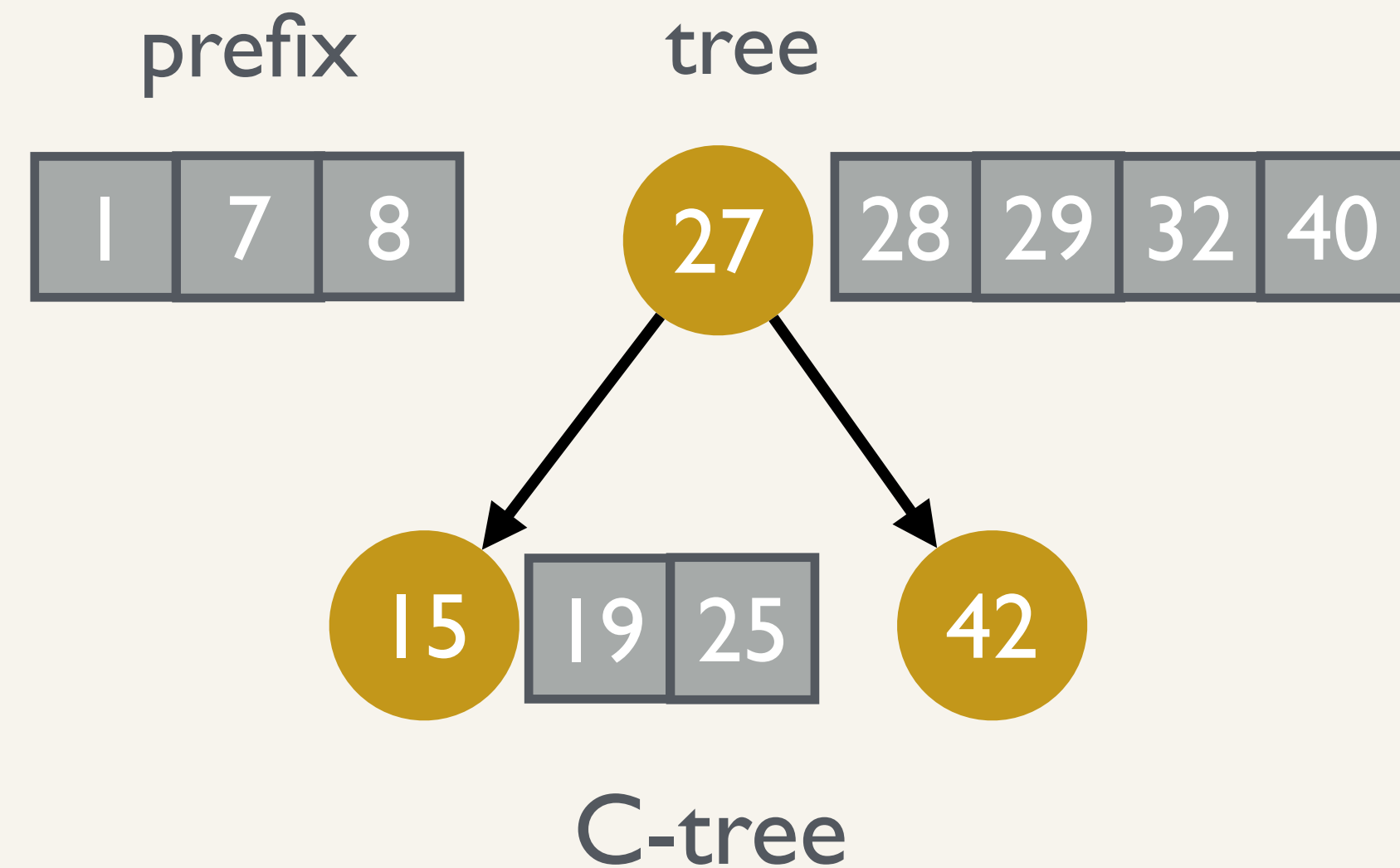


# C-trees

- ❖ Chunking parameter  $B$ . Fix a hash function,  $h$
- ❖ Select elements as *heads* with probability  $1/B$  using  $h$

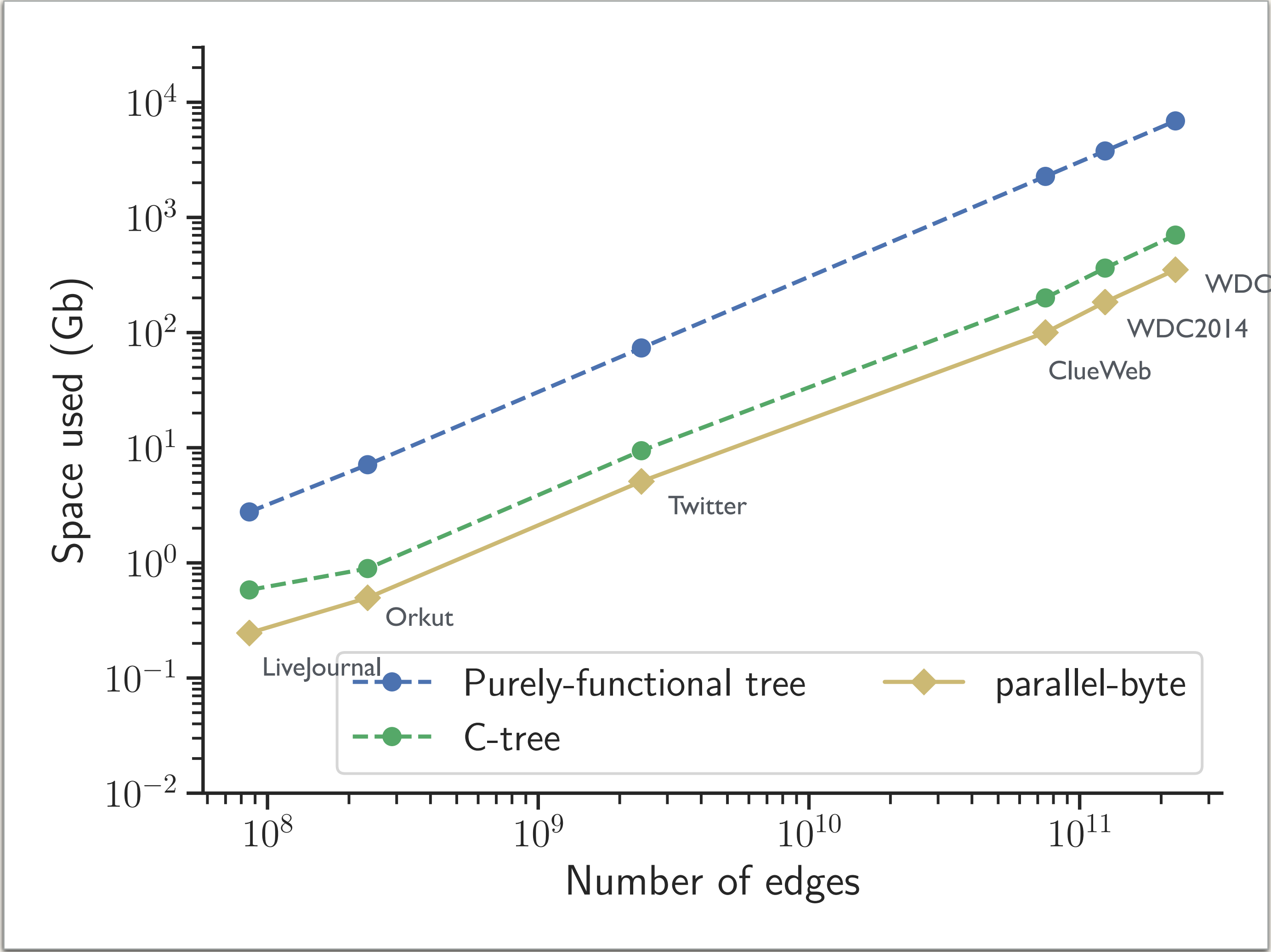


Further improve space usage for integer C-trees by difference encoding chunks

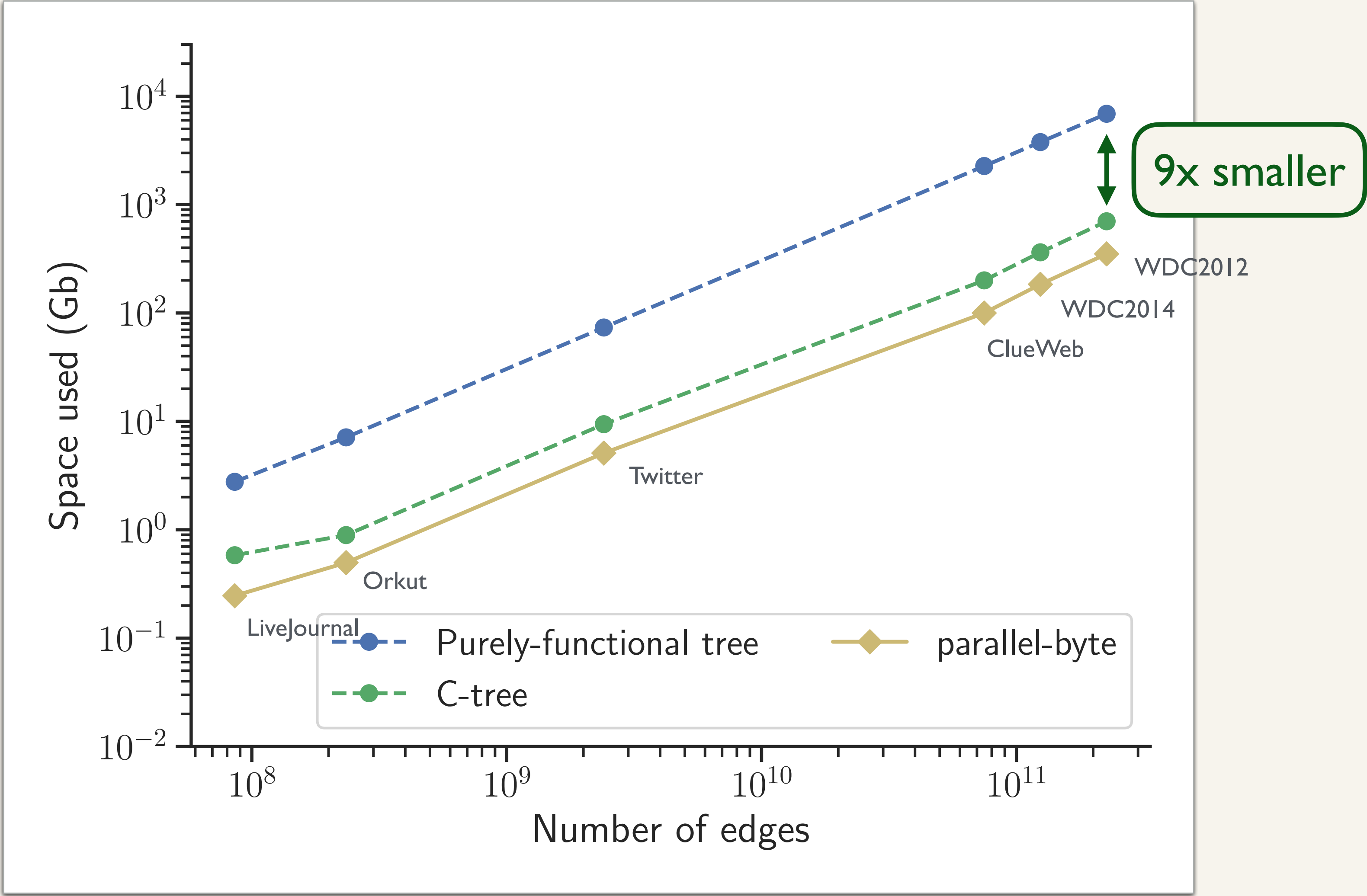




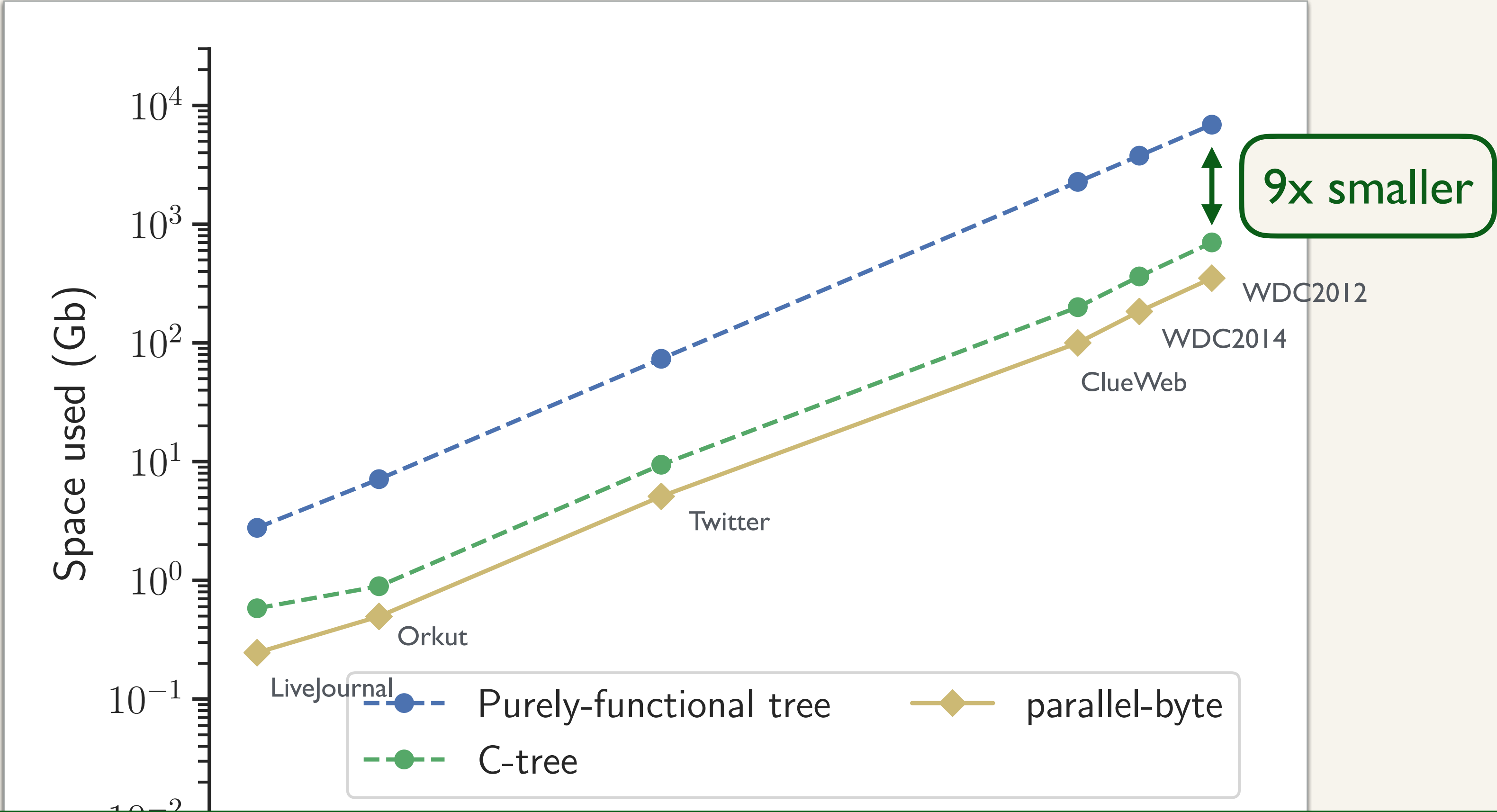
# Space Improvement in Aspen using C-trees



# Space Improvement in Aspen using C-trees



# Space Improvement in Aspen using C-trees

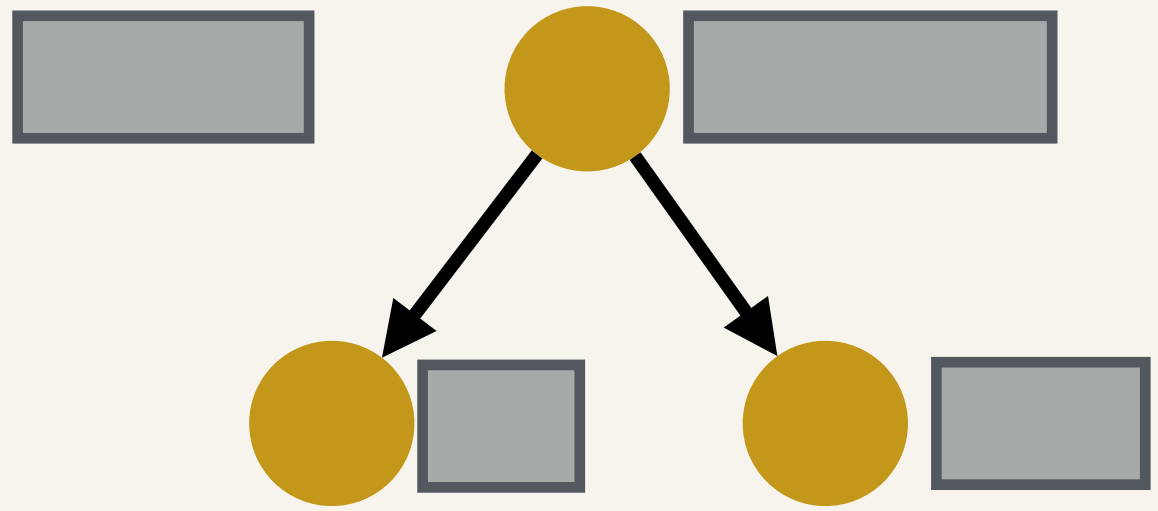
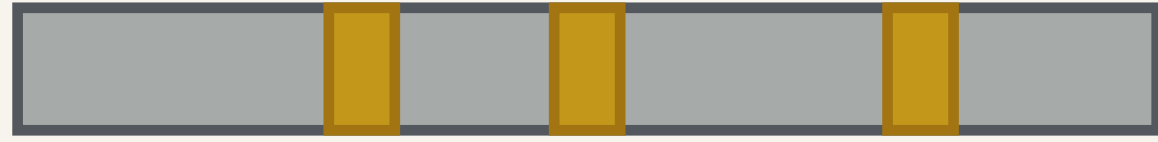


Fully-dynamic representation of the WebDataCommons hyperlink graph using 700GB of memory

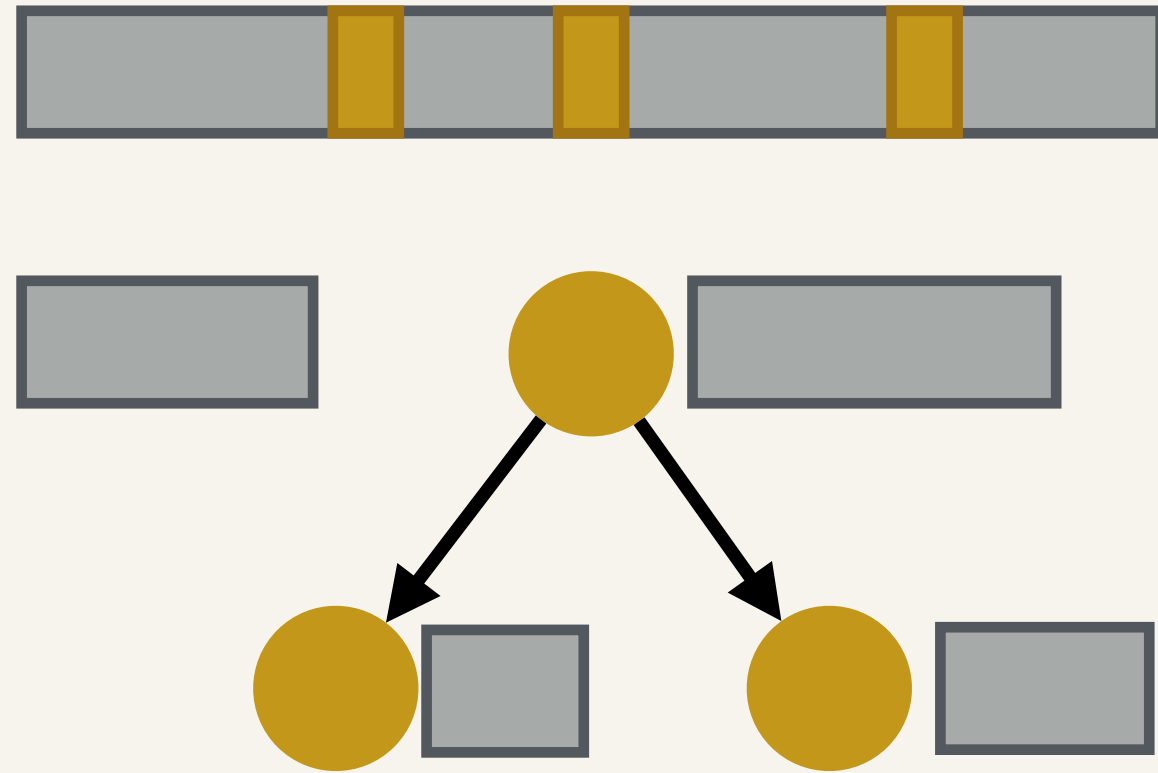
# Operations on C-trees



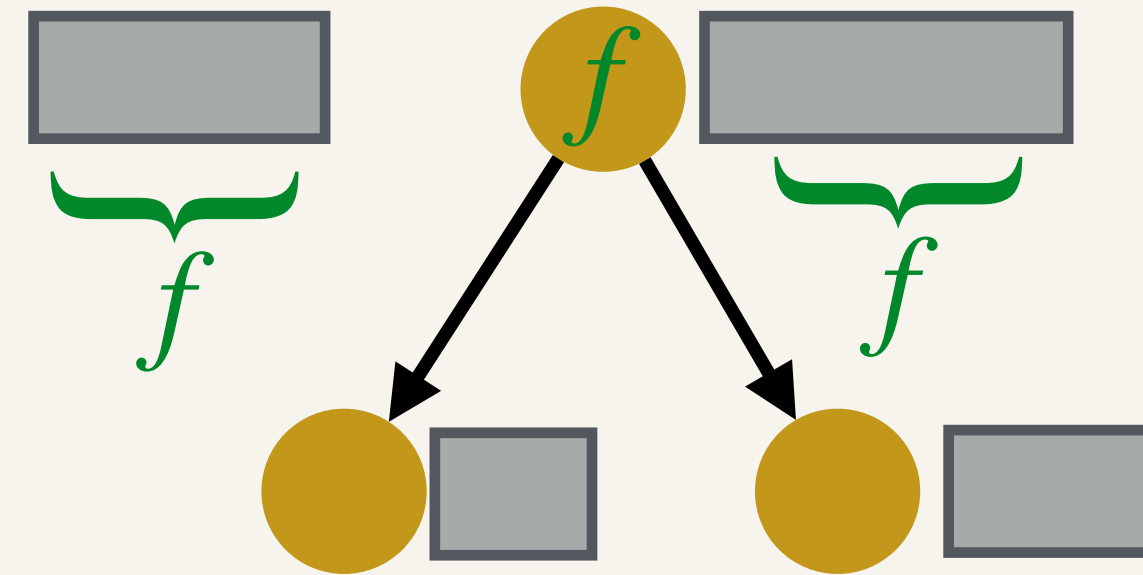
Build(Seq  $S$ )



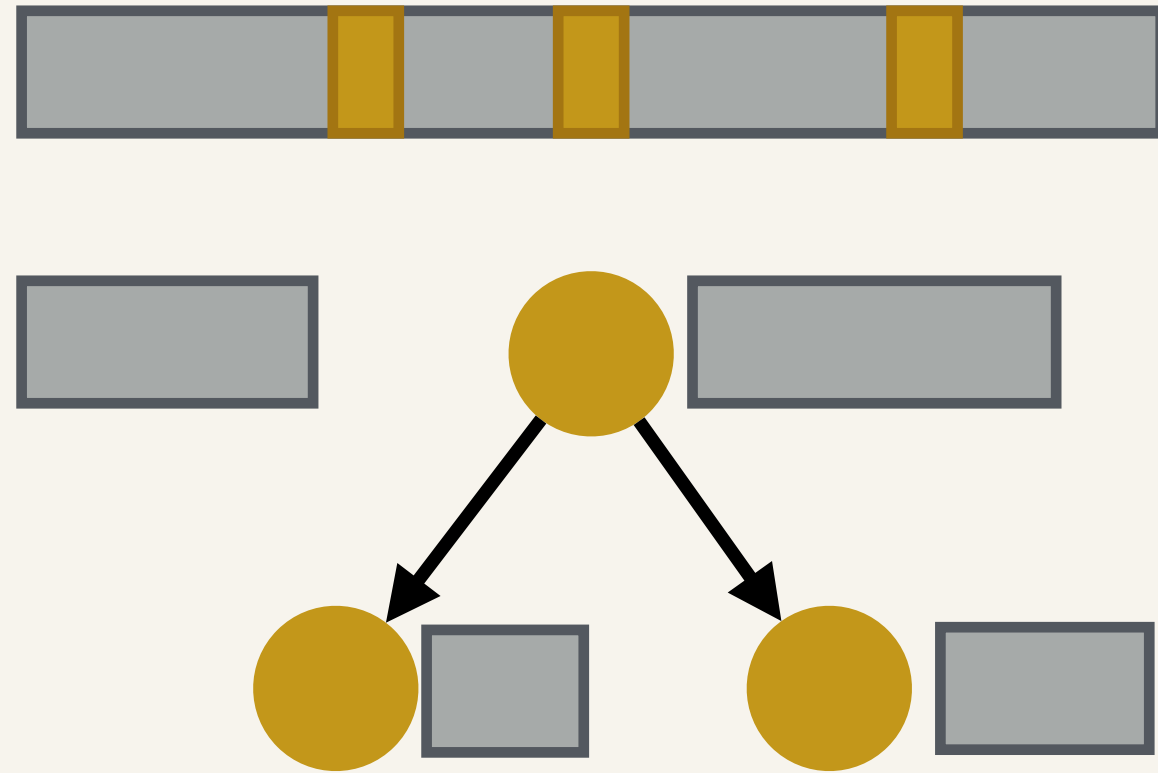
Build(Seq  $S$ )



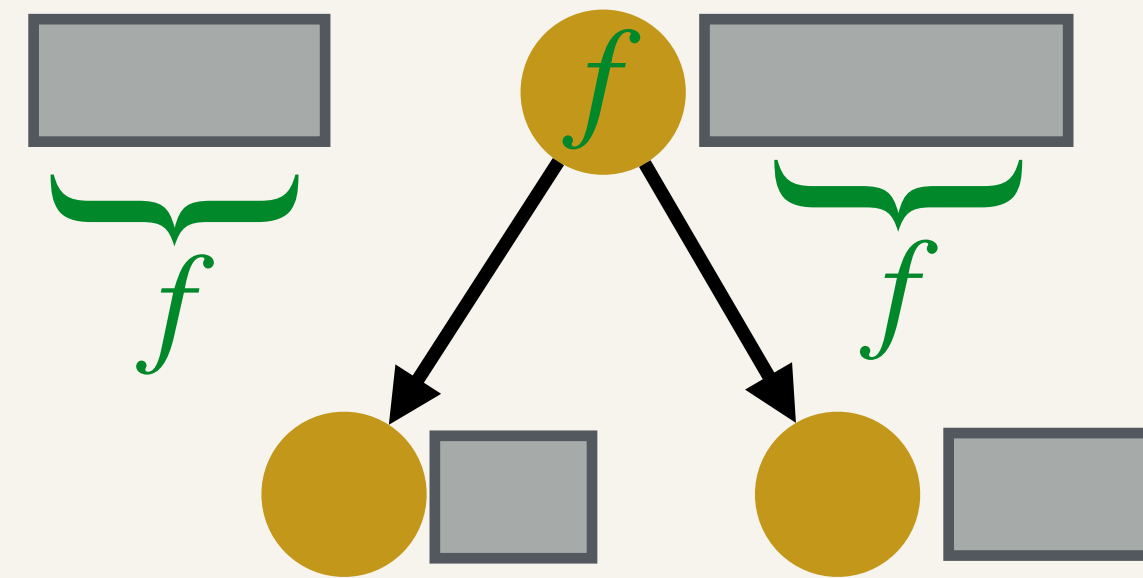
Map(Ctree  $C, f$ )



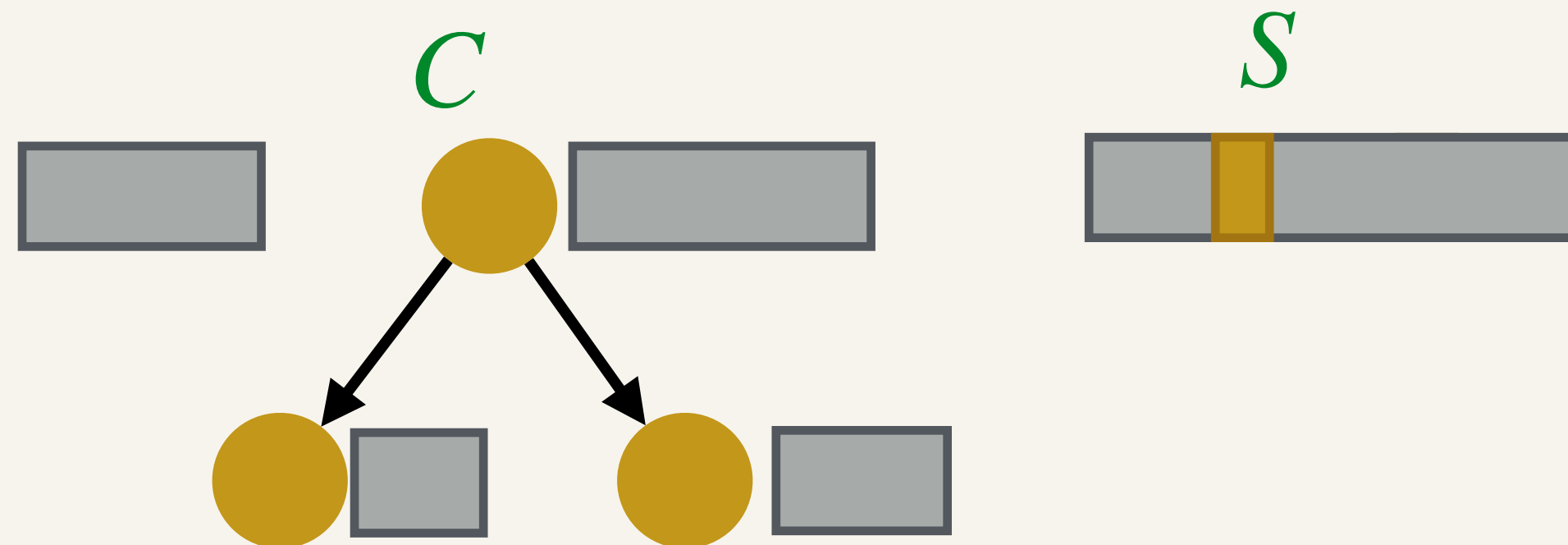
Build(Seq  $S$ )



Map(Ctree  $C, f$ )

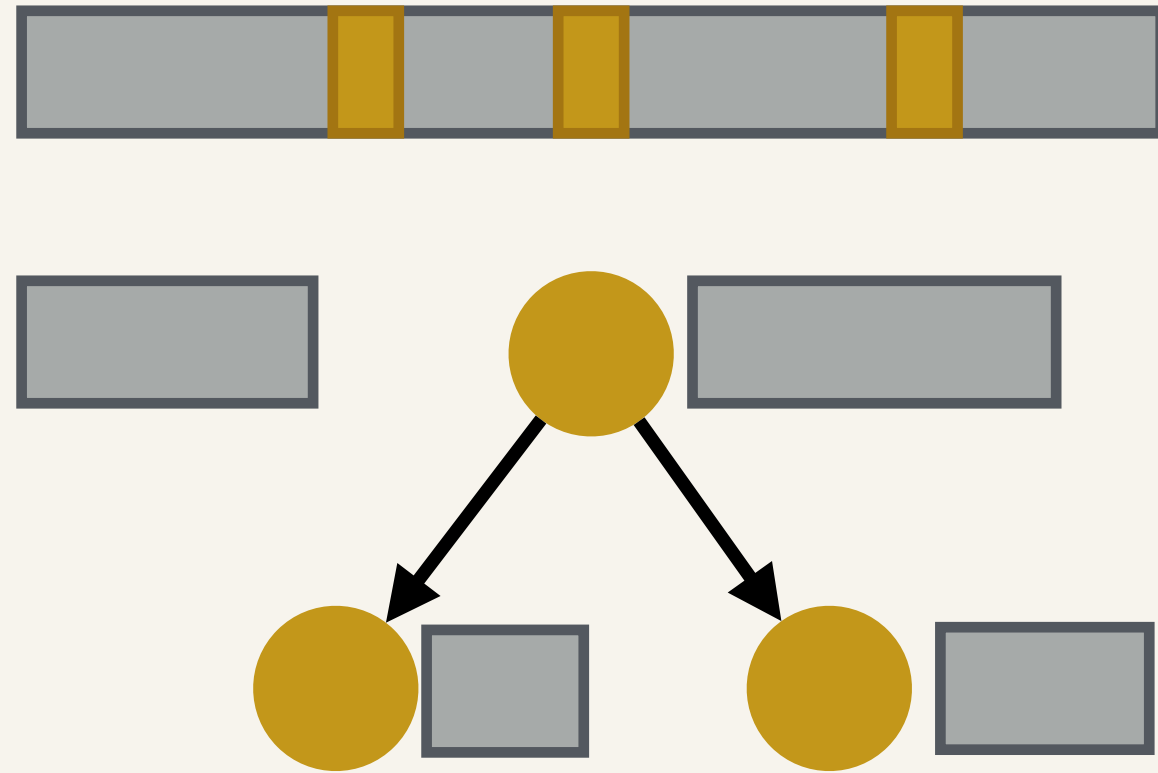


MultiInsert(Ctree  $C, \text{Seq } S$ )

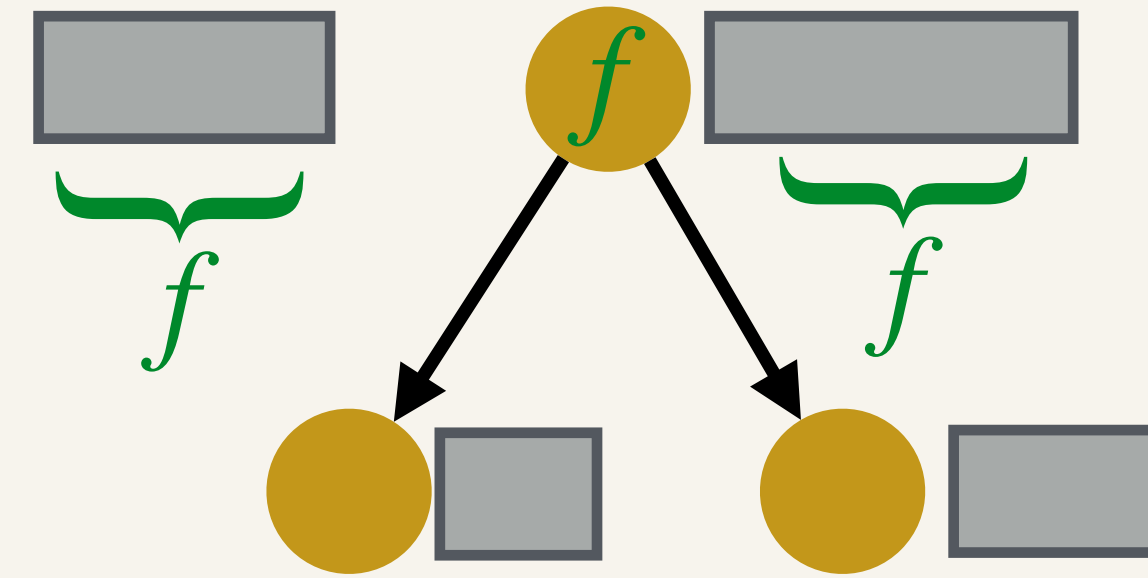




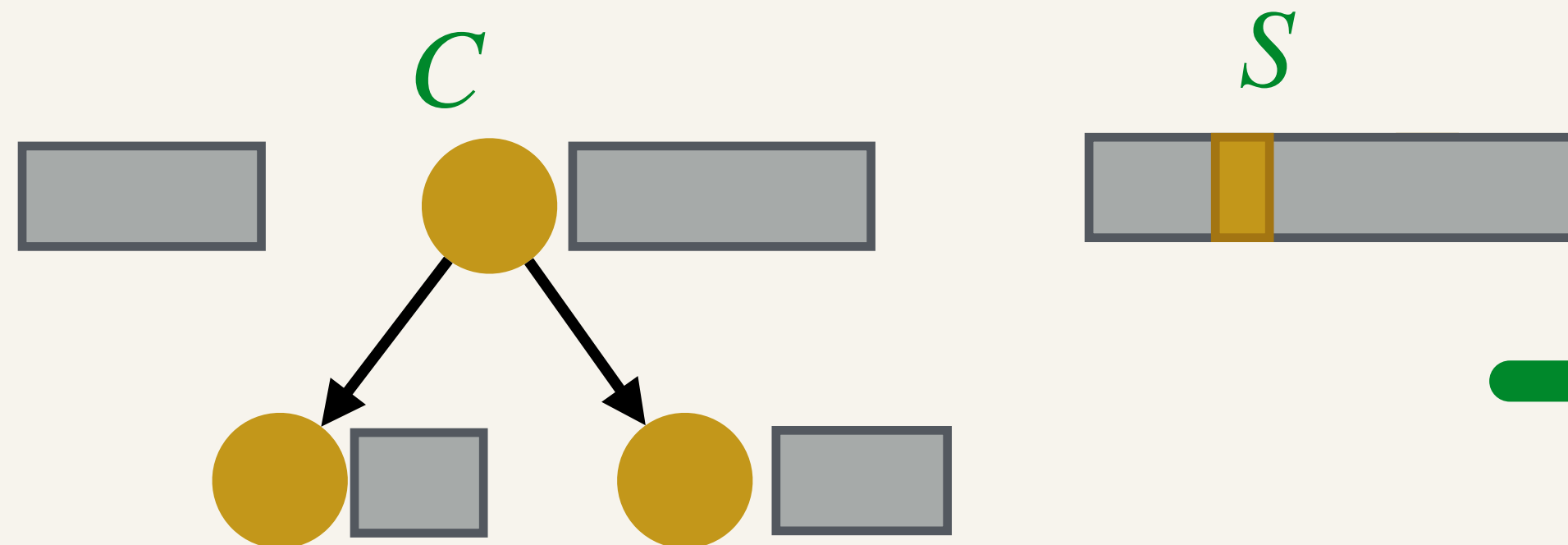
Build(Seq  $S$ )



Map(Ctree  $C, f$ )



MultiInsert(Ctree  $C$ , Seq  $S$ )



$C_S = \text{Build}(\text{Seq } S)$

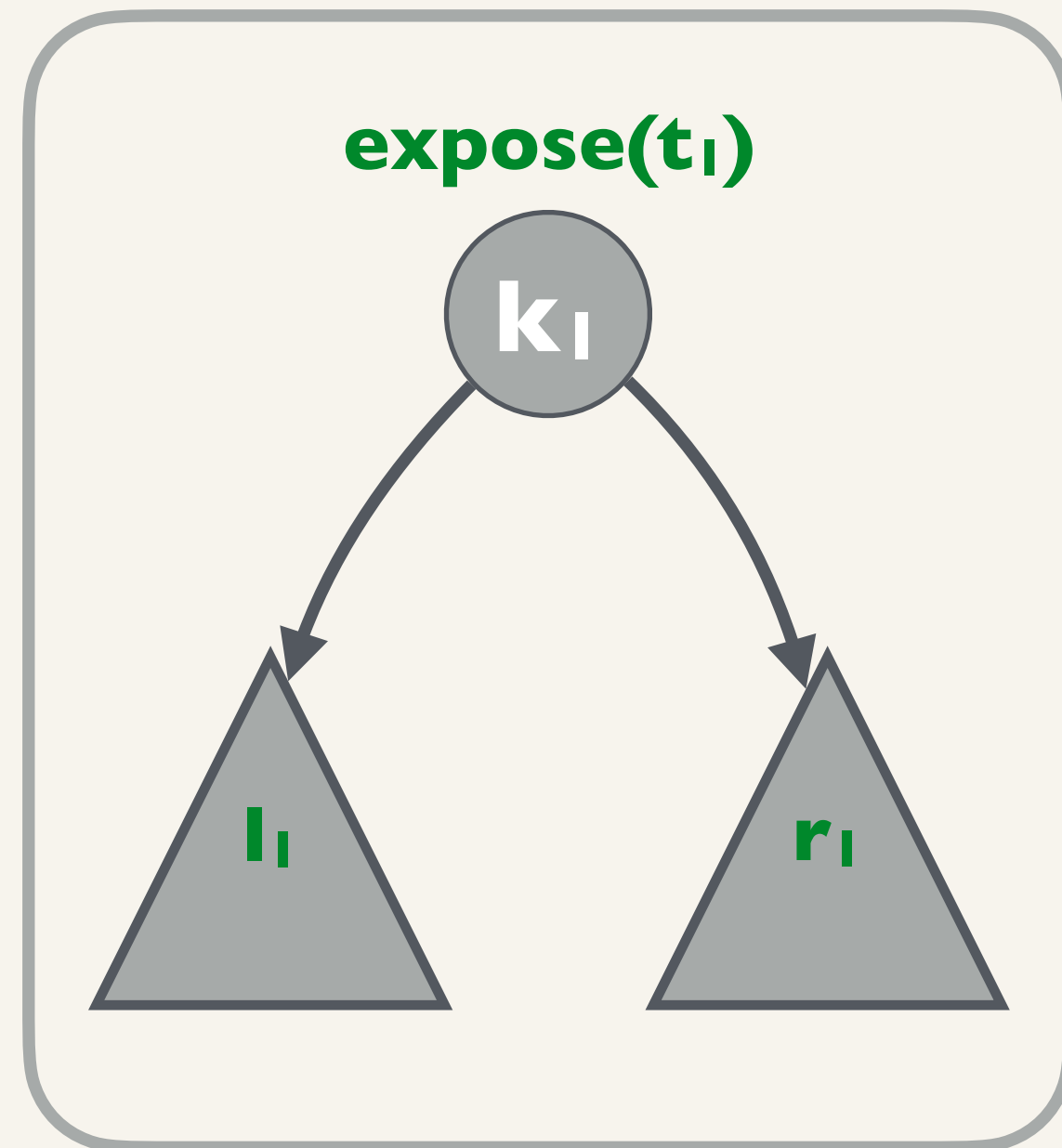
Output = Union( $C, C_S$ )

# Batch Updates on Trees

`union(t1, t2)`

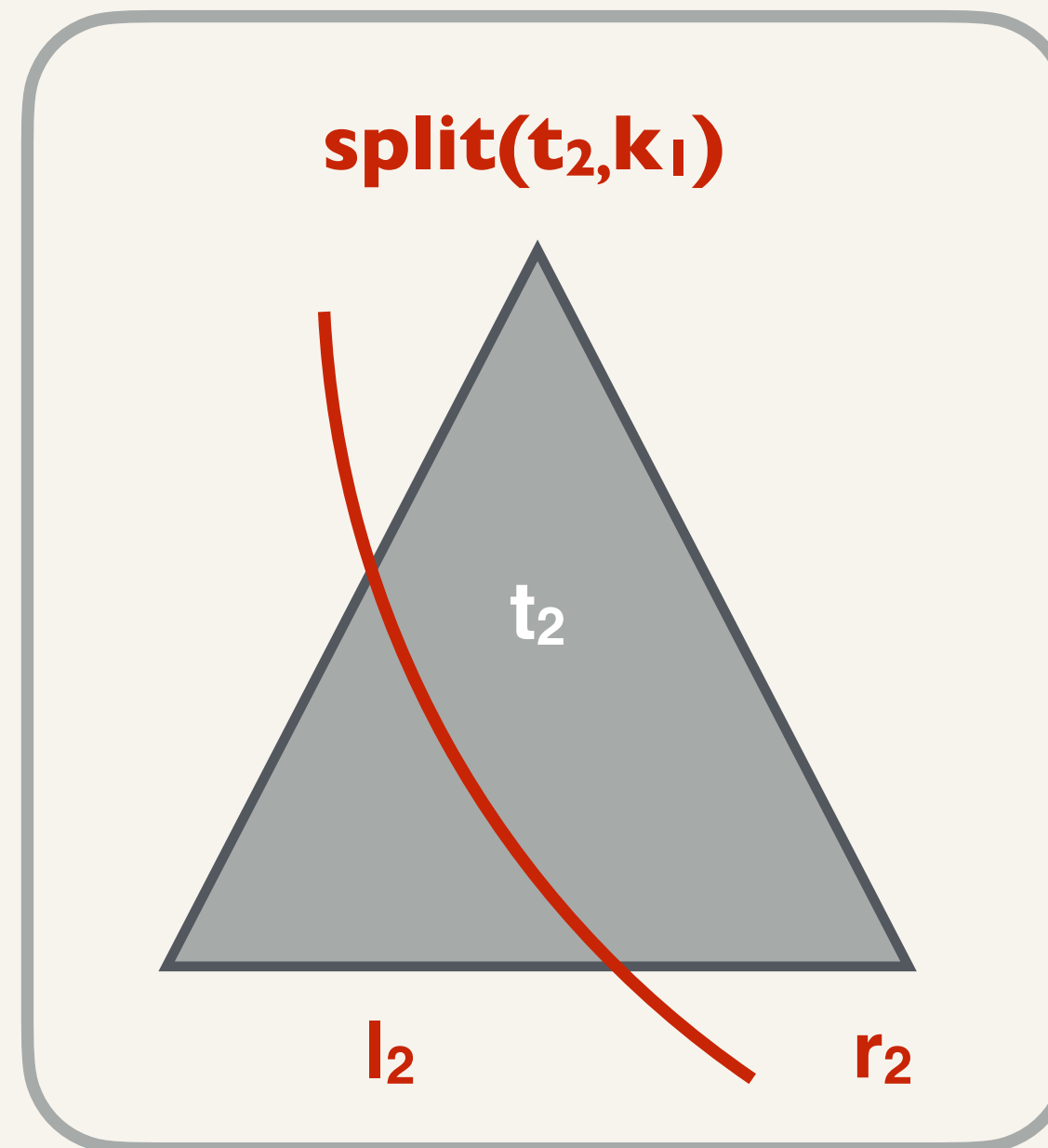
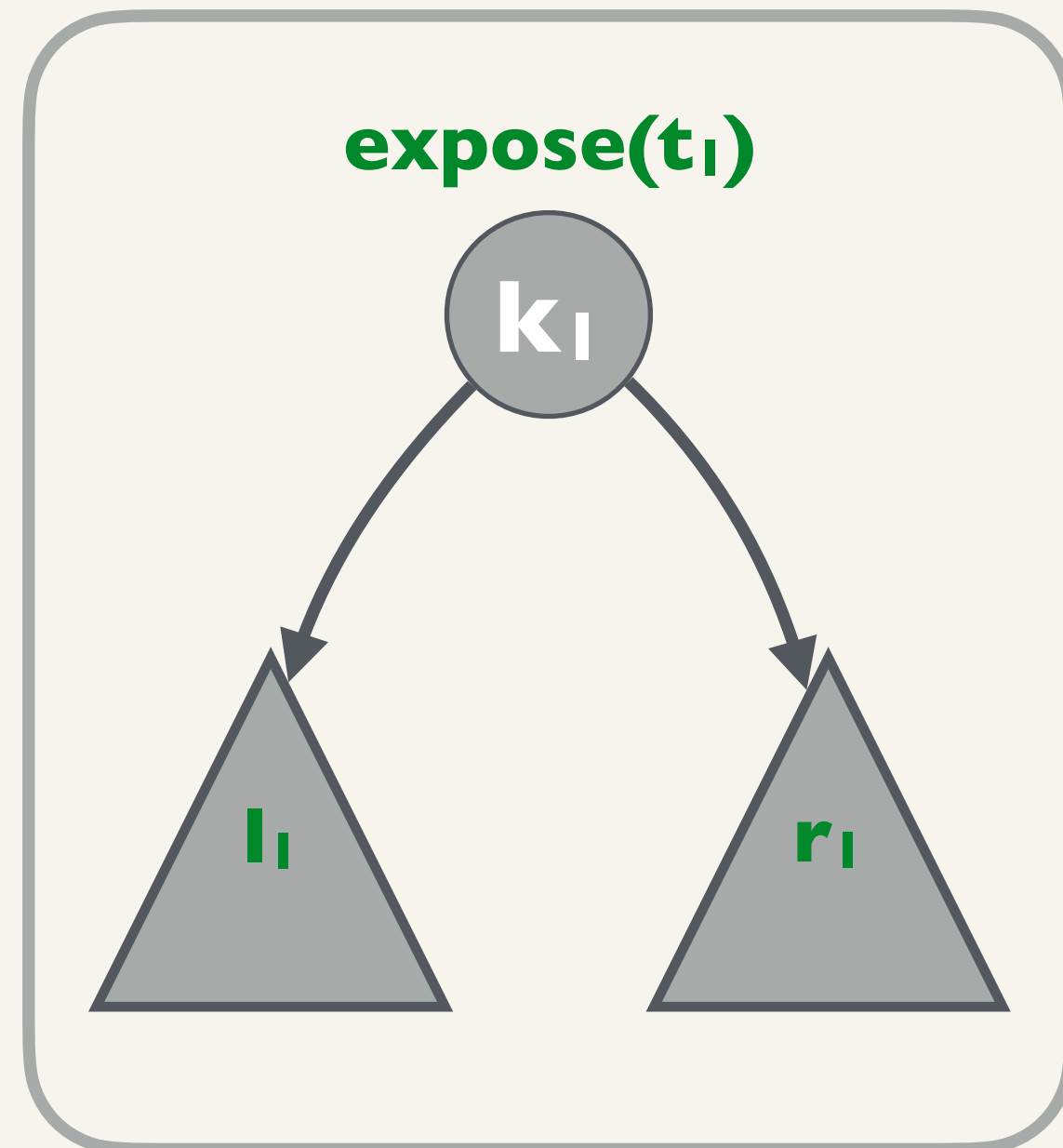
# Batch Updates on Trees

$\text{union}(t_1, t_2)$



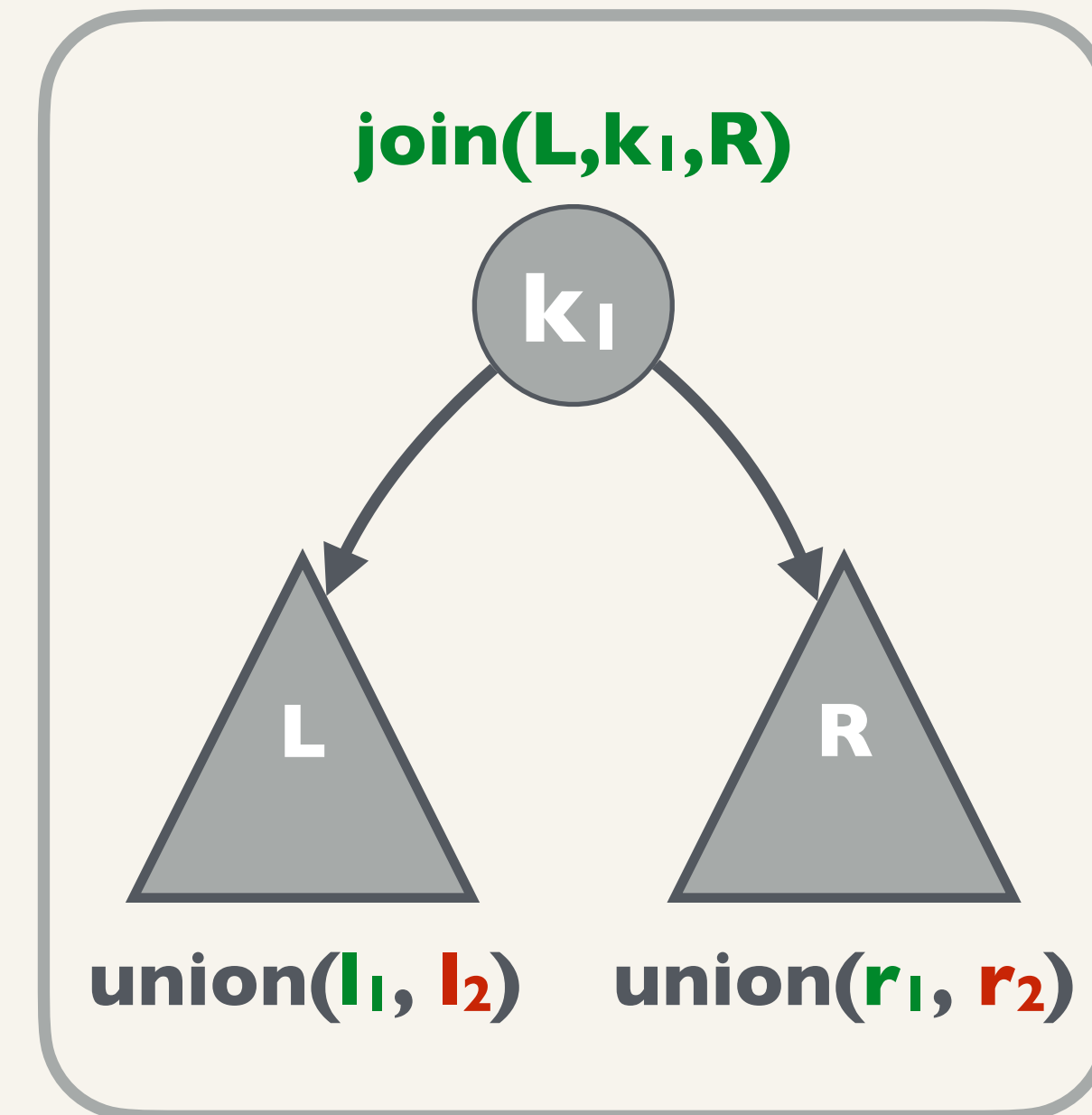
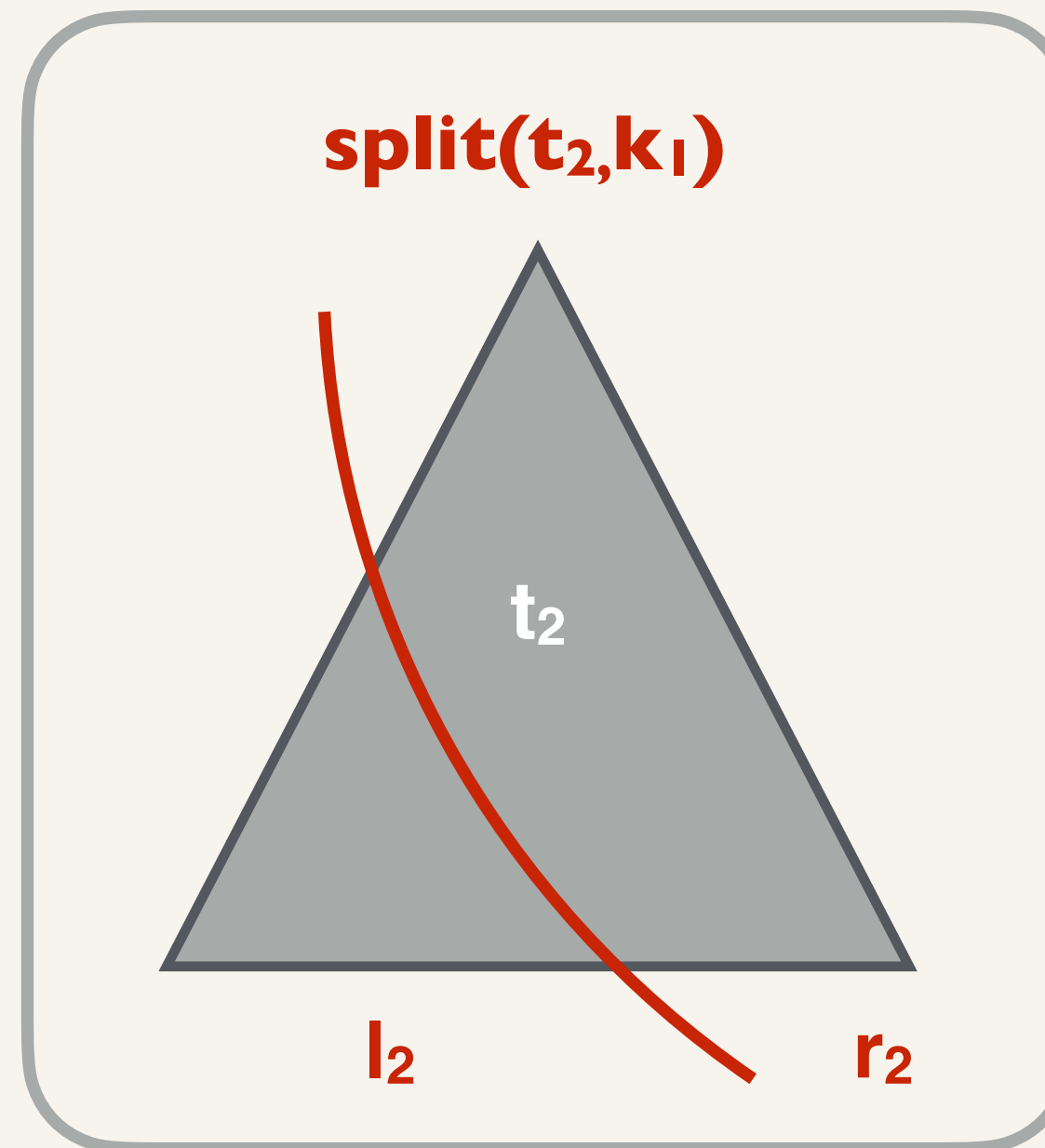
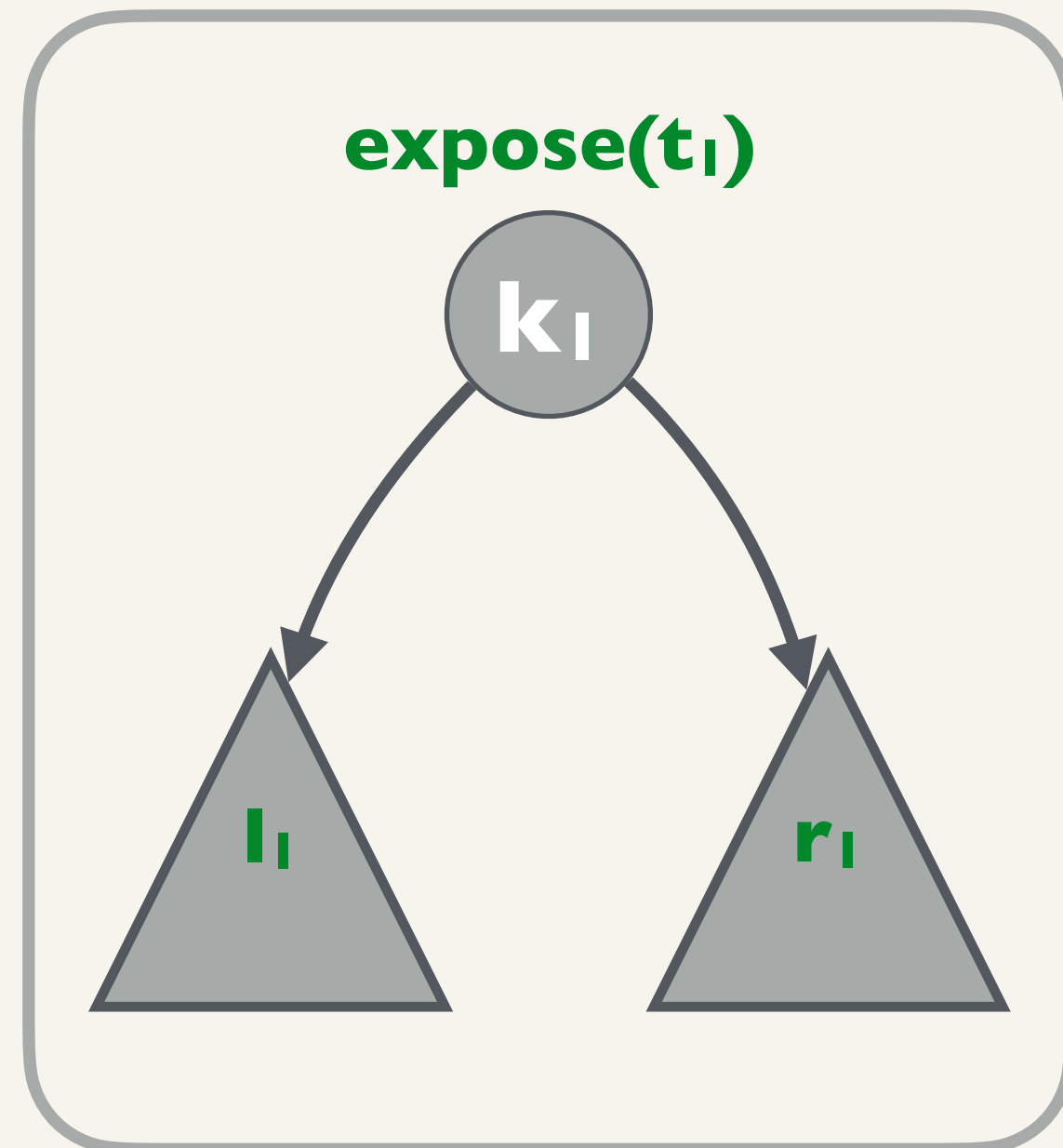
# Batch Updates on Trees

$\text{union}(t_1, t_2)$



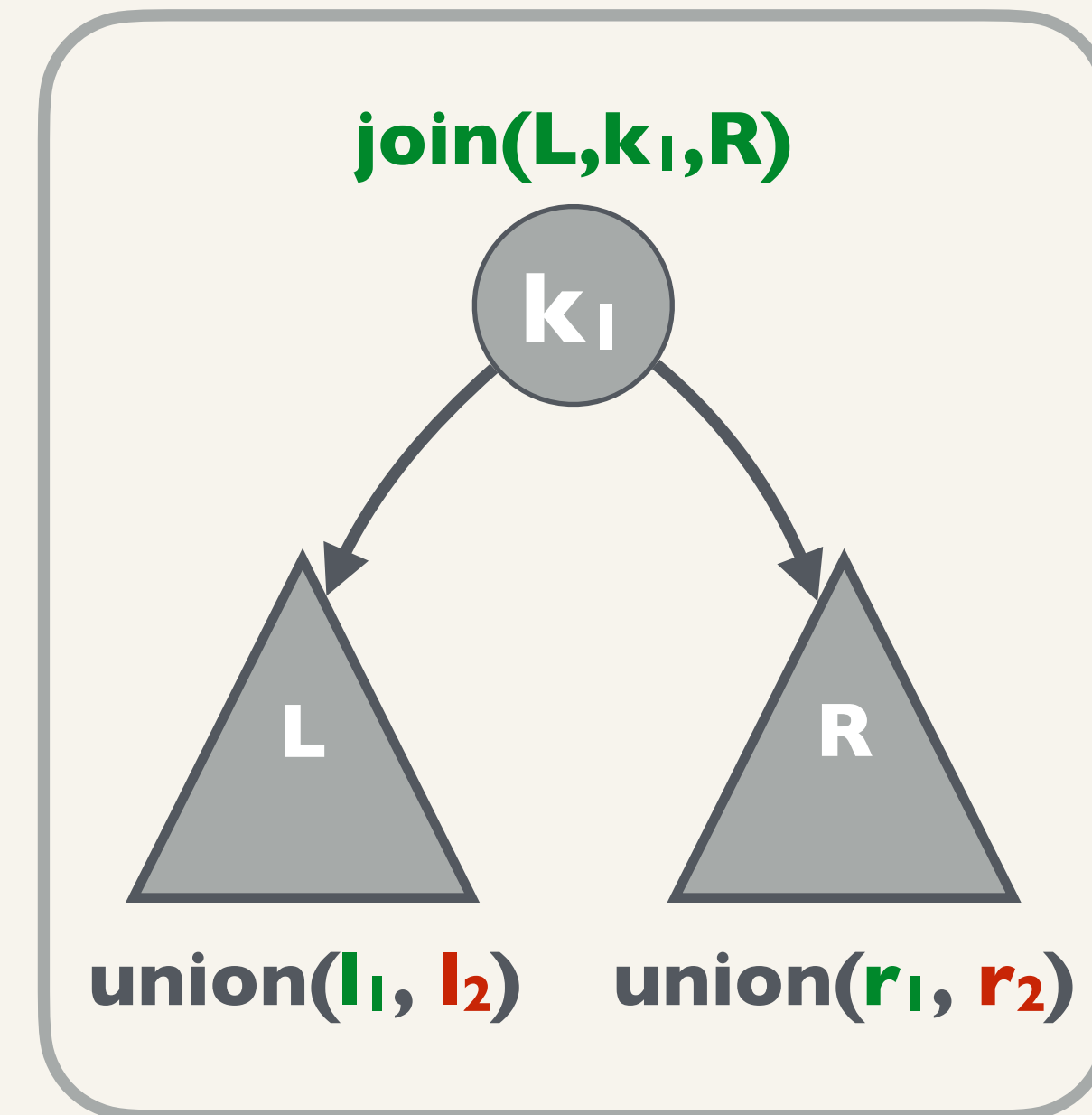
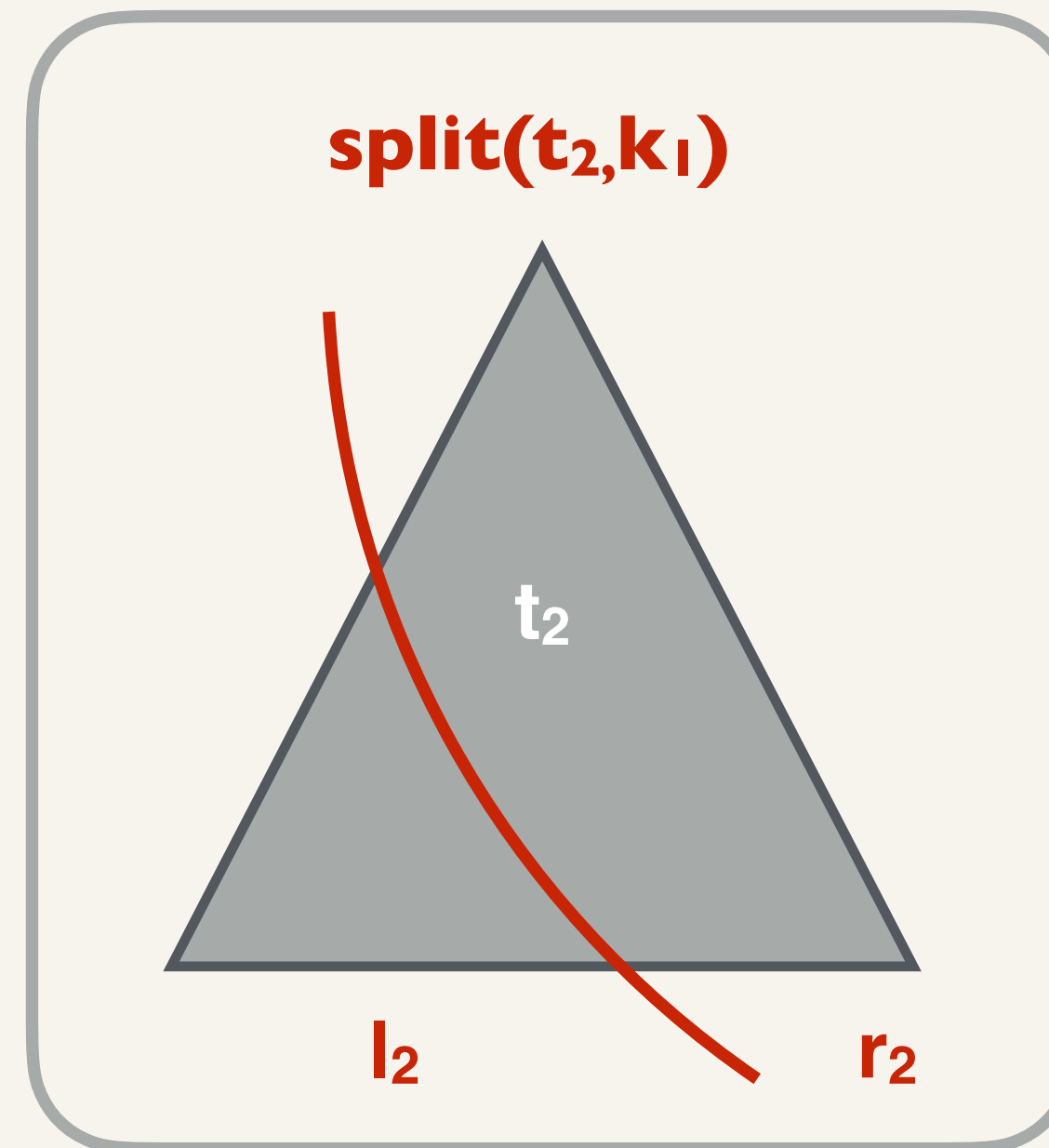
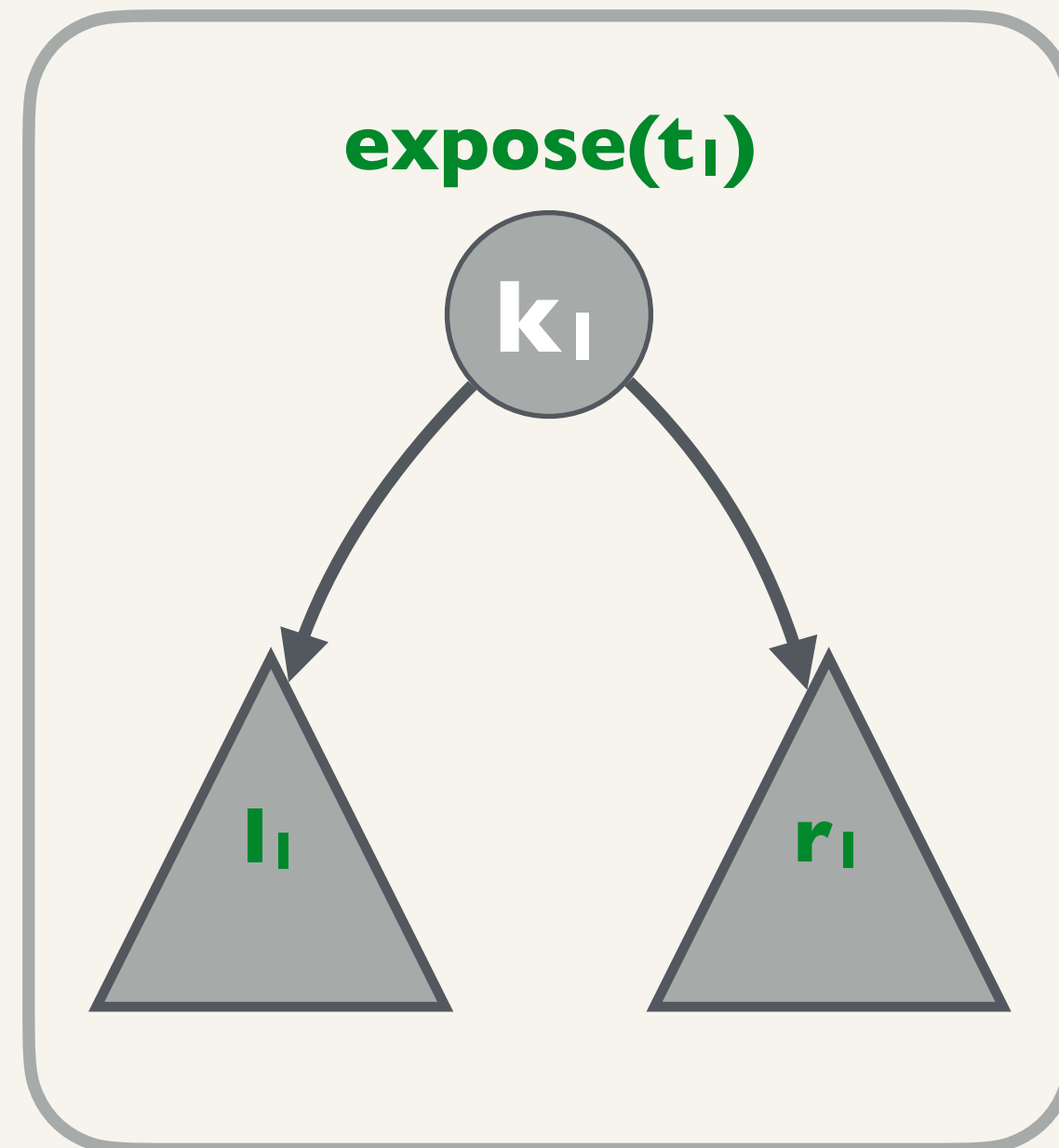
# Batch Updates on Trees

$\text{union}(t_1, t_2)$



# Batch Updates on Trees

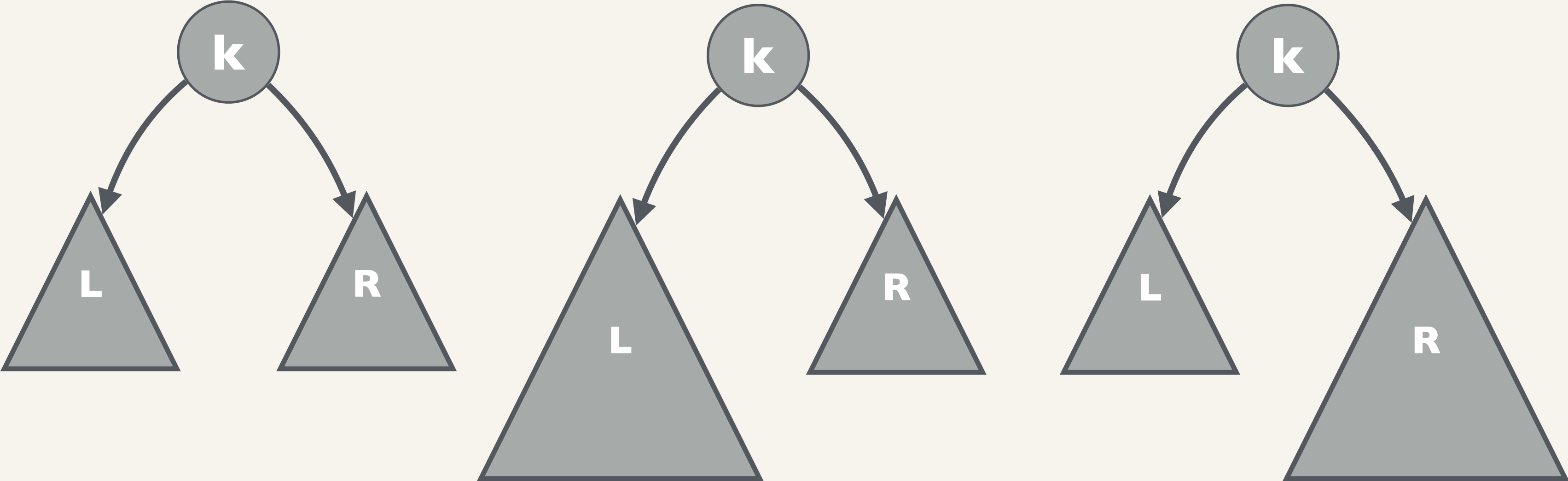
$\text{union}(t_1, t_2)$



Similar algorithms for **difference** and **intersection**

# The Join Function

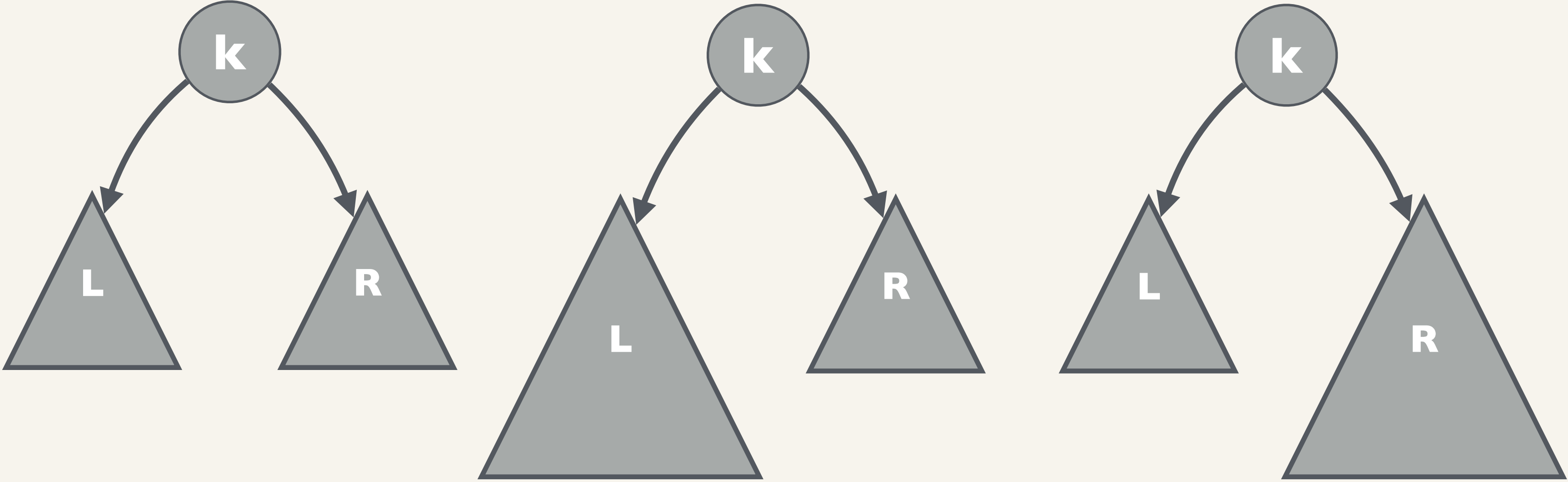
join(L, k, R)



[1] Just Join for Parallel Ordered Sets, Blelloch et al. (SPAA'16)

# The Join Function

join(L, k, R)



Join enables **balance-agnostic expression** of all other primitives[1]

[1] Just Join for Parallel Ordered Sets, Blelloch et al. (SPAA'16)

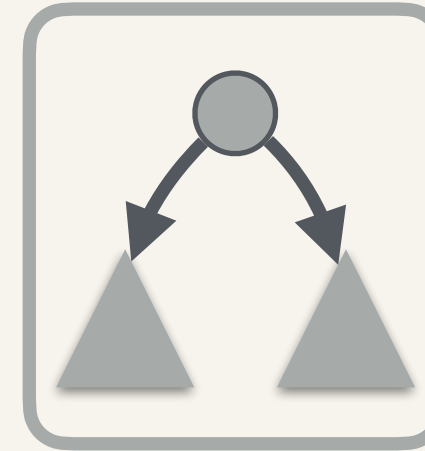


# Batch Updates on Trees

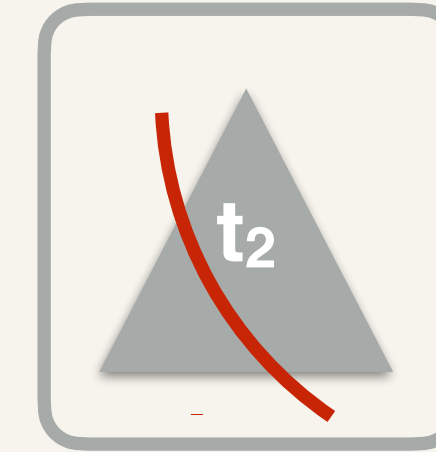
$\text{union}(t_1, t_2)$  runs in

$O\left(m \log\left(\frac{n}{m} + 1\right)\right)$  work and  $O(\log n \log m)$  depth

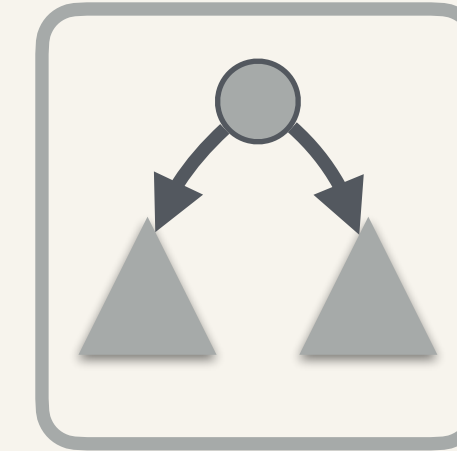
expose



split



join



# Batch Updates on Trees

$\text{union}(t_1, t_2)$  runs in

$O\left(m \log\left(\frac{n}{m} + 1\right)\right)$  work and  $O(\log n \log m)$  depth



Proof idea from [1]:

Splitting a tree costs  $O(\log |T|)$  work and depth

Overall cost = work done over all **splits**

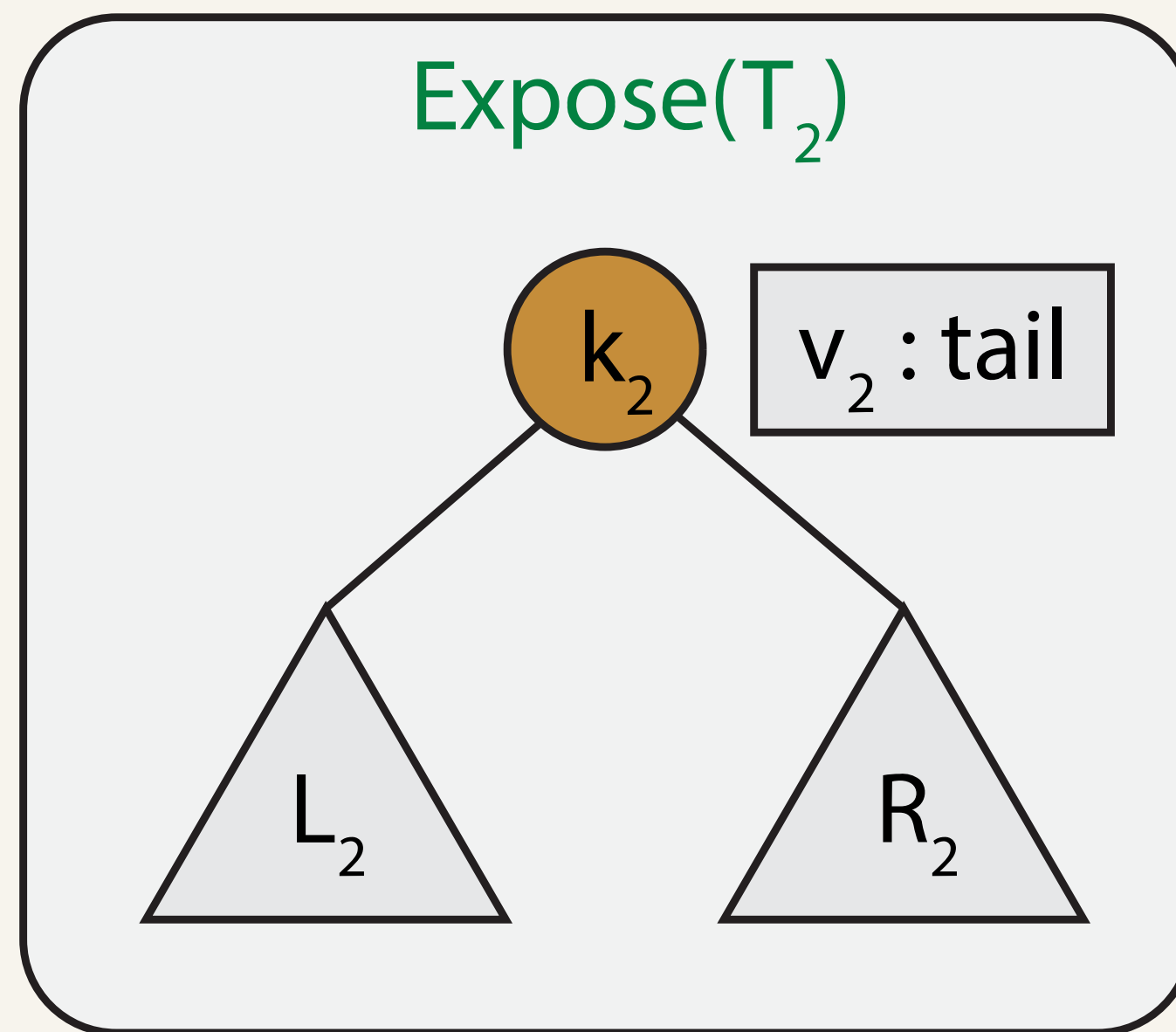
[1] Just Join for Parallel Ordered Sets, Blelloch et al. (SPAA'16)

# Batch Updates on Trees

$\text{union}(C_1=(T_1, P_1), C_2=(T_2, P_2))$

# Batch Updates on Trees

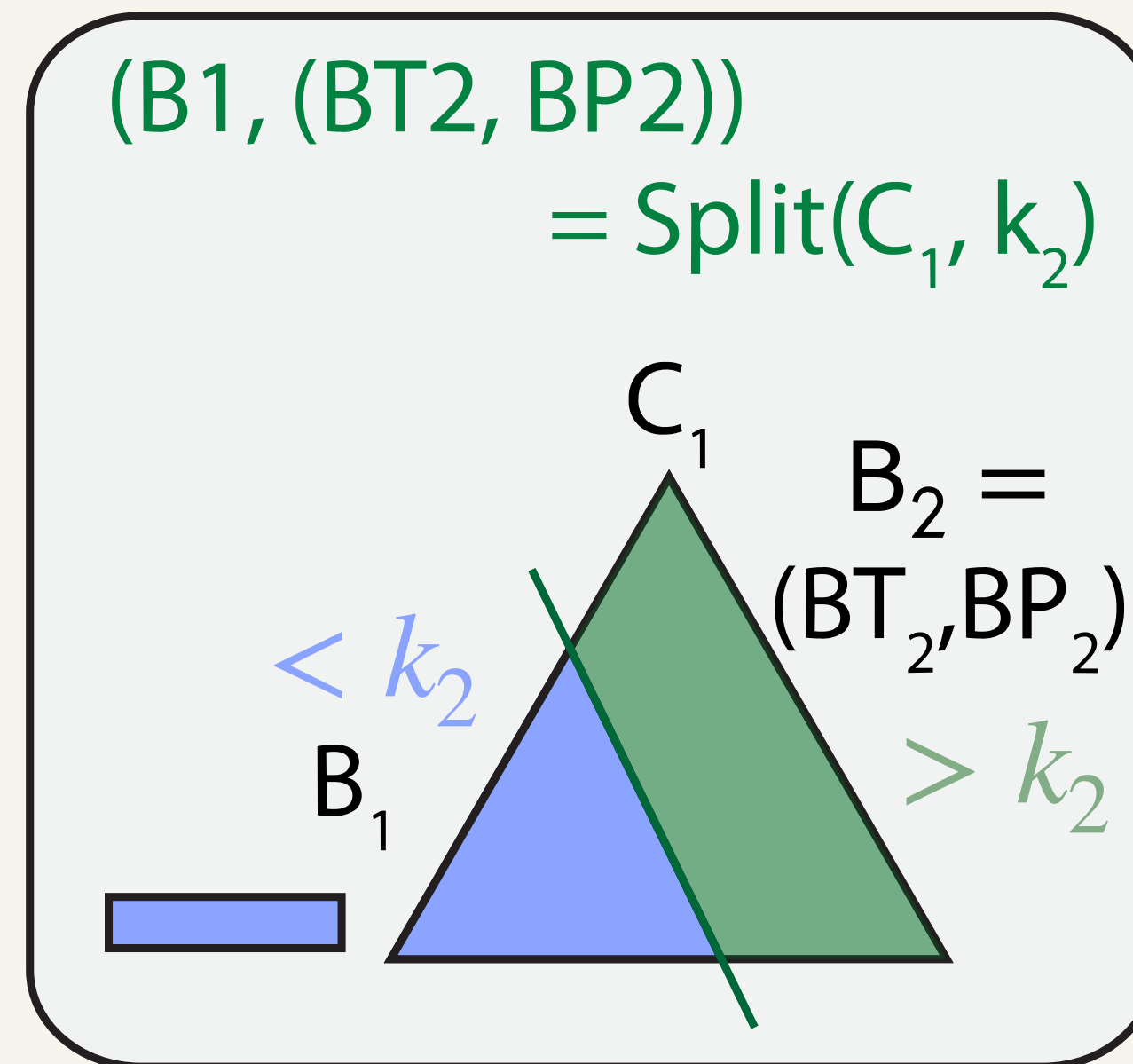
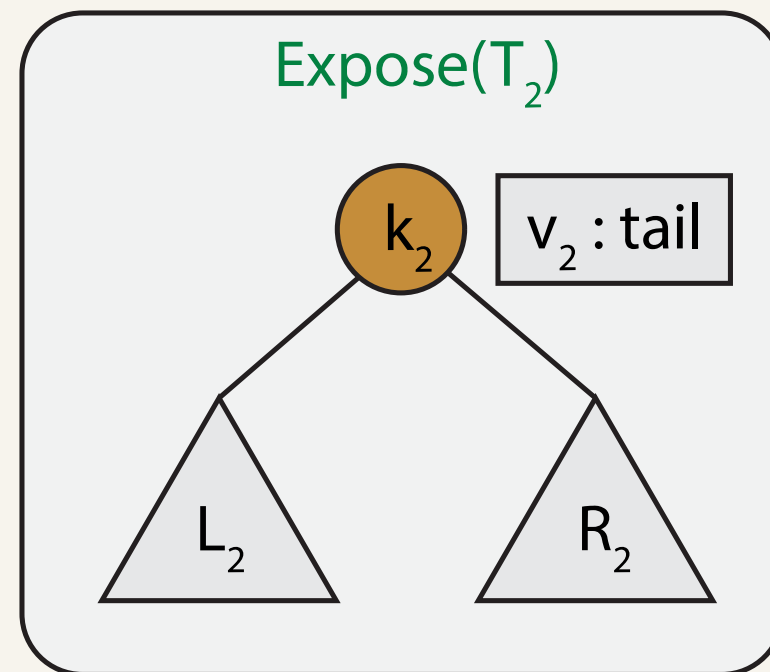
$\text{union}(C_1=(T_1, P_1), C_2=(T_2, P_2))$



Expose one of the trees

# Batch Updates on Trees

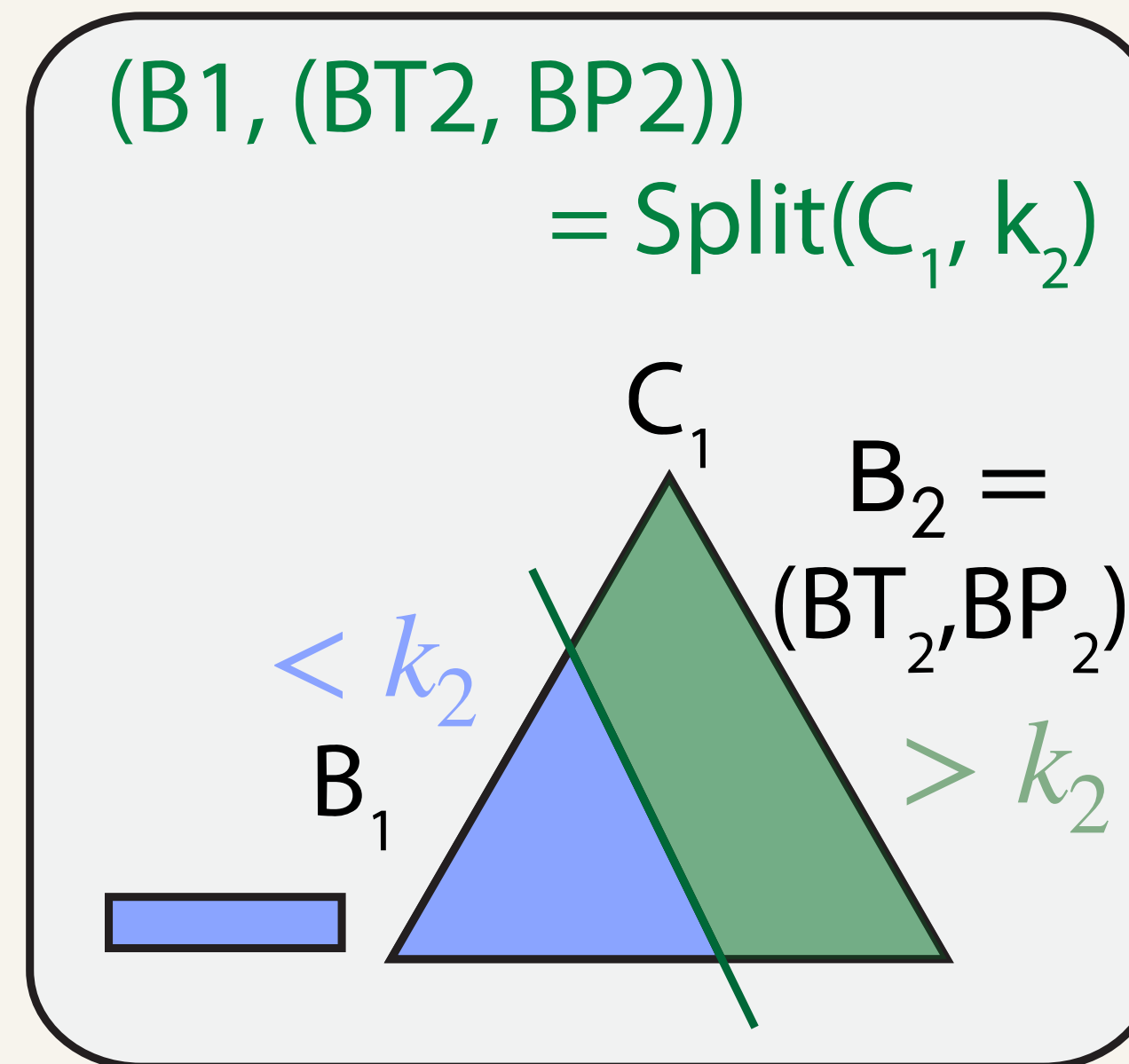
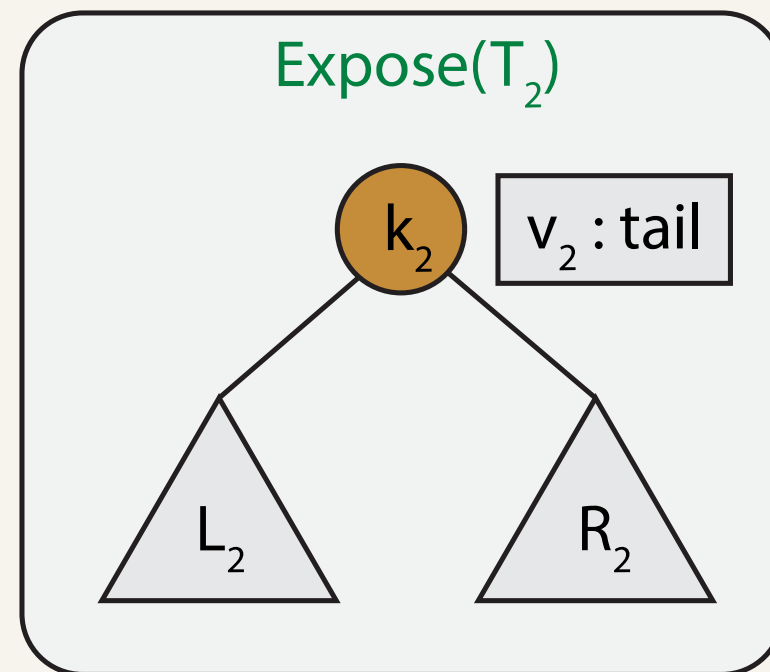
$\text{union}(C_1=(T_1, P_1), C_2=(T_2, P_2))$



Split the other  $C$ -tree with  $k_2$

# Batch Updates on Trees

$\text{union}(C_1=(T_1, P_1), C_2=(T_2, P_2))$

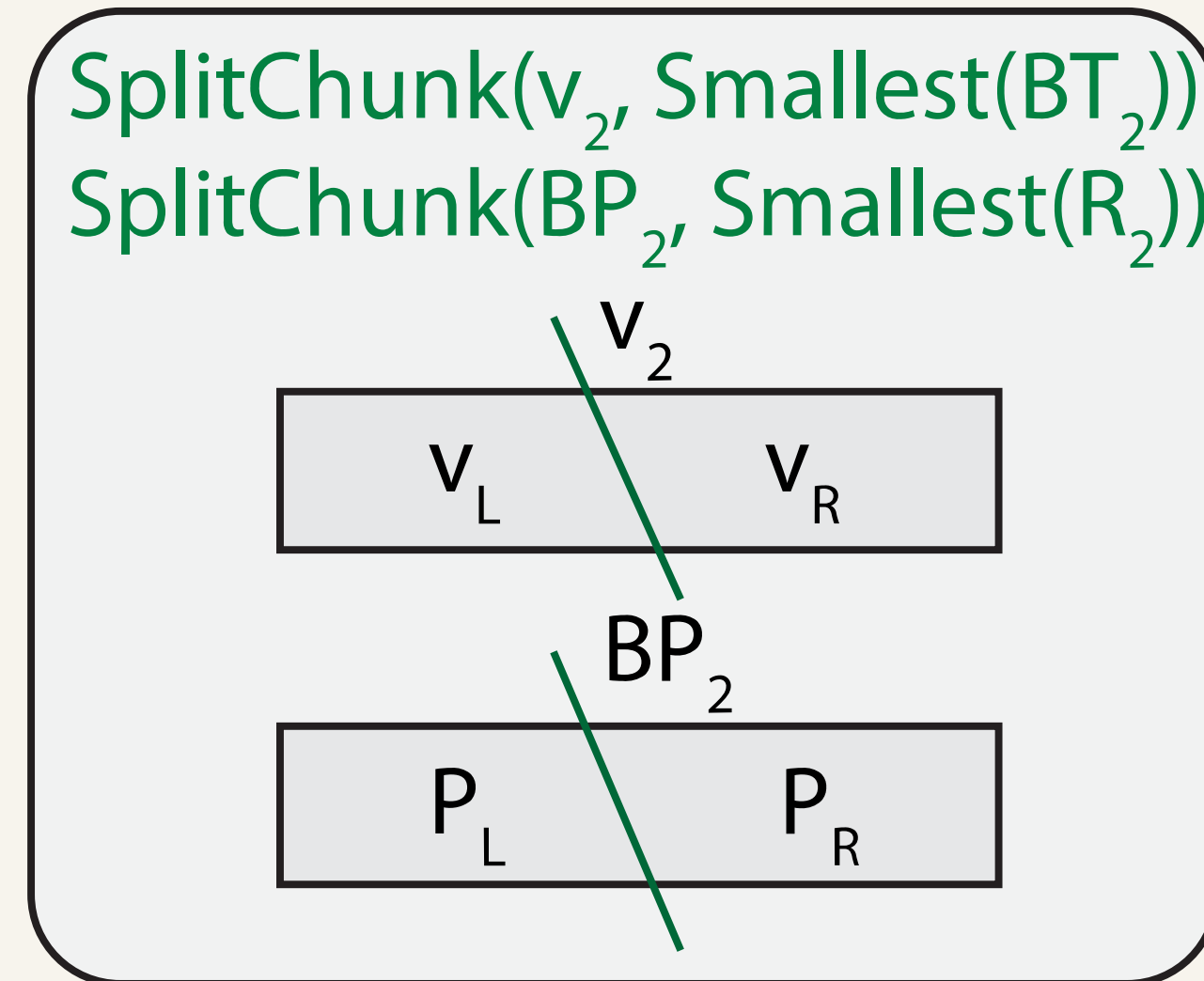
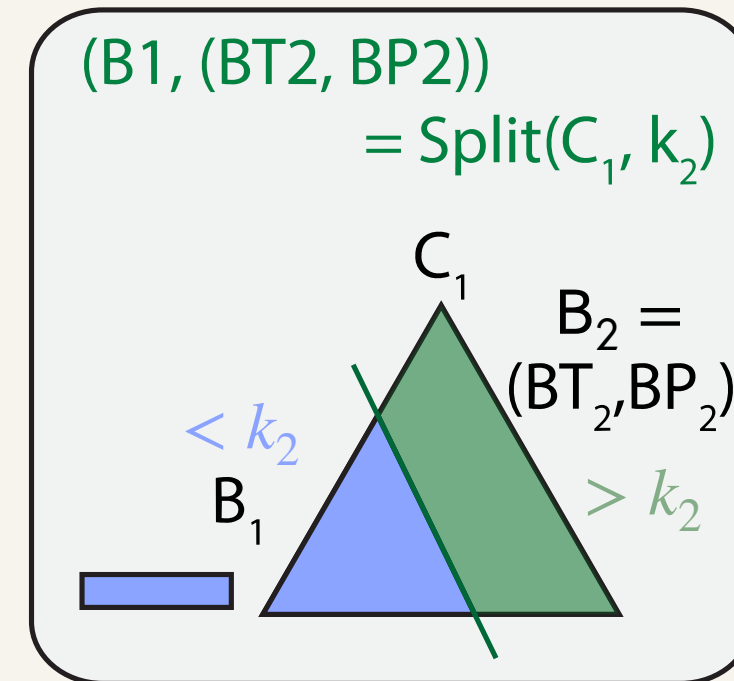
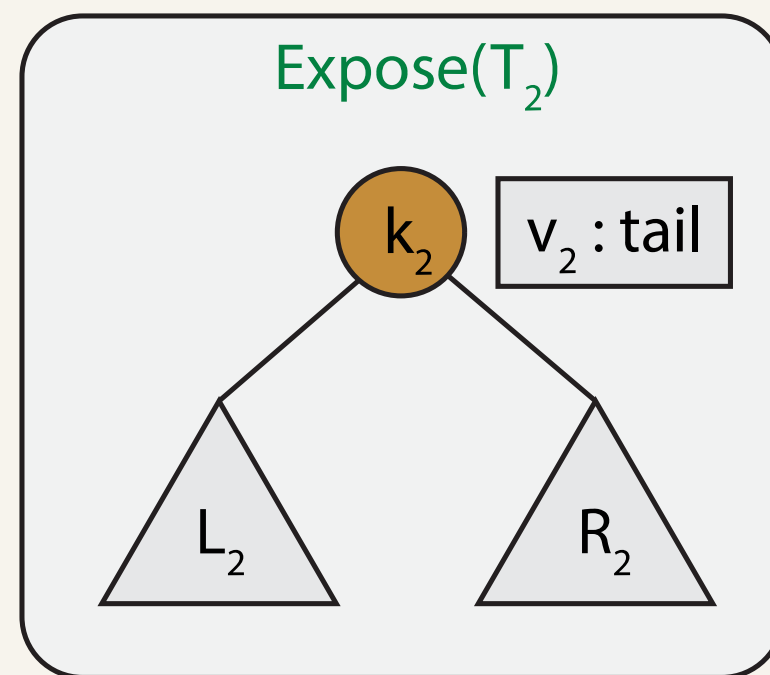


Split the other  $C$ -tree with  $k_2$

Part of  $v_2$  may belong in  $BT_2$ ,  
similarly with  $BP_2$

# Batch Updates on C-trees

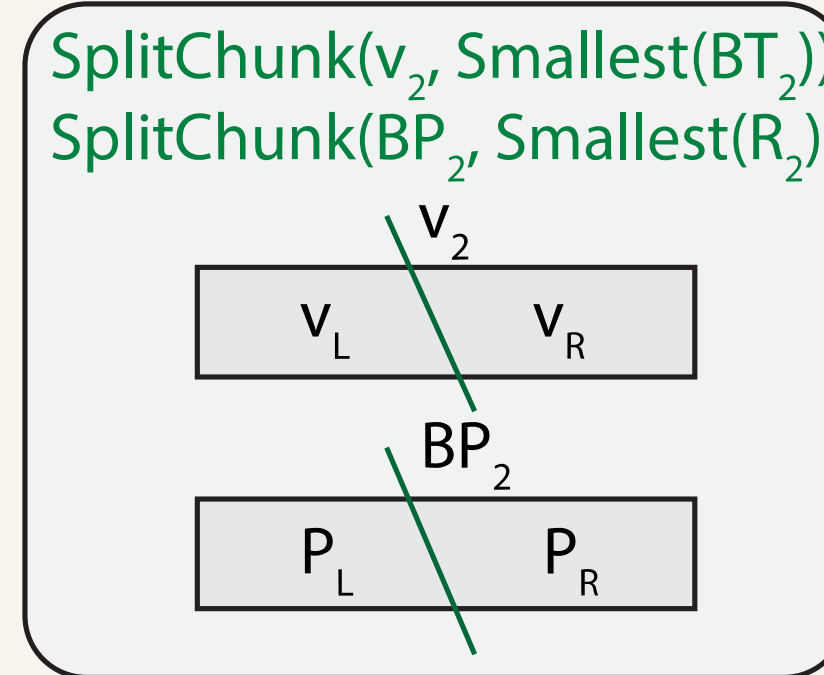
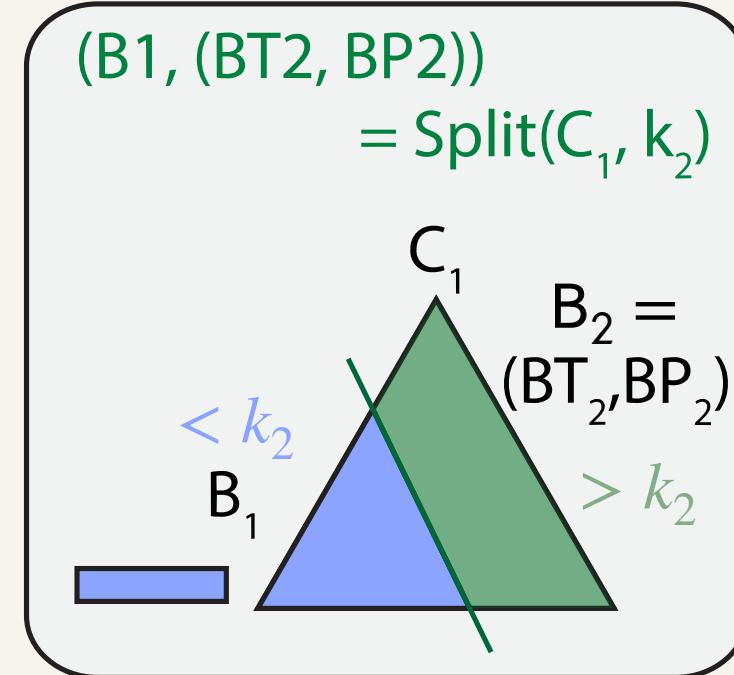
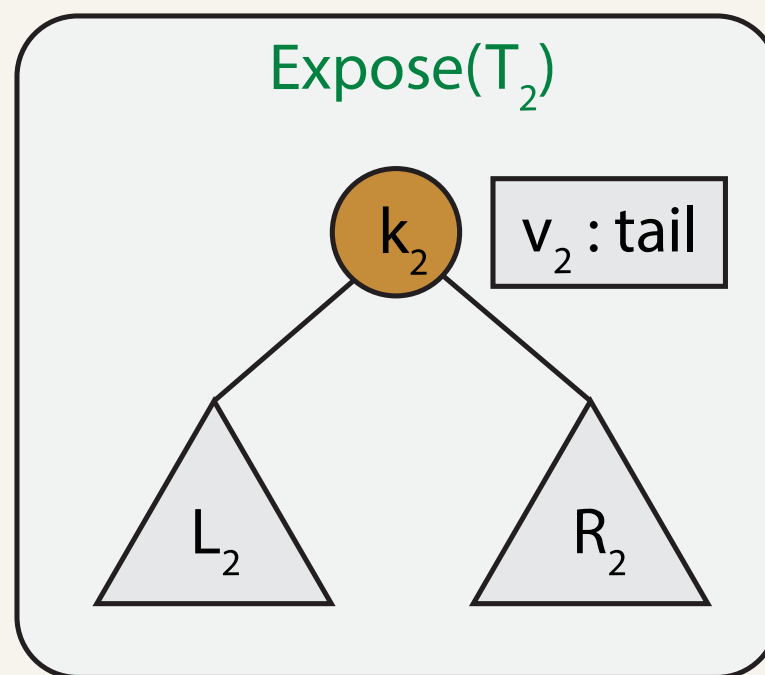
$\text{union}(C_1=(T_1, P_1), C_2=(T_2, P_2))$



Split  $v_2$  based on  $BT_2$ ,  
 $BP_2$  based on  $R_2$

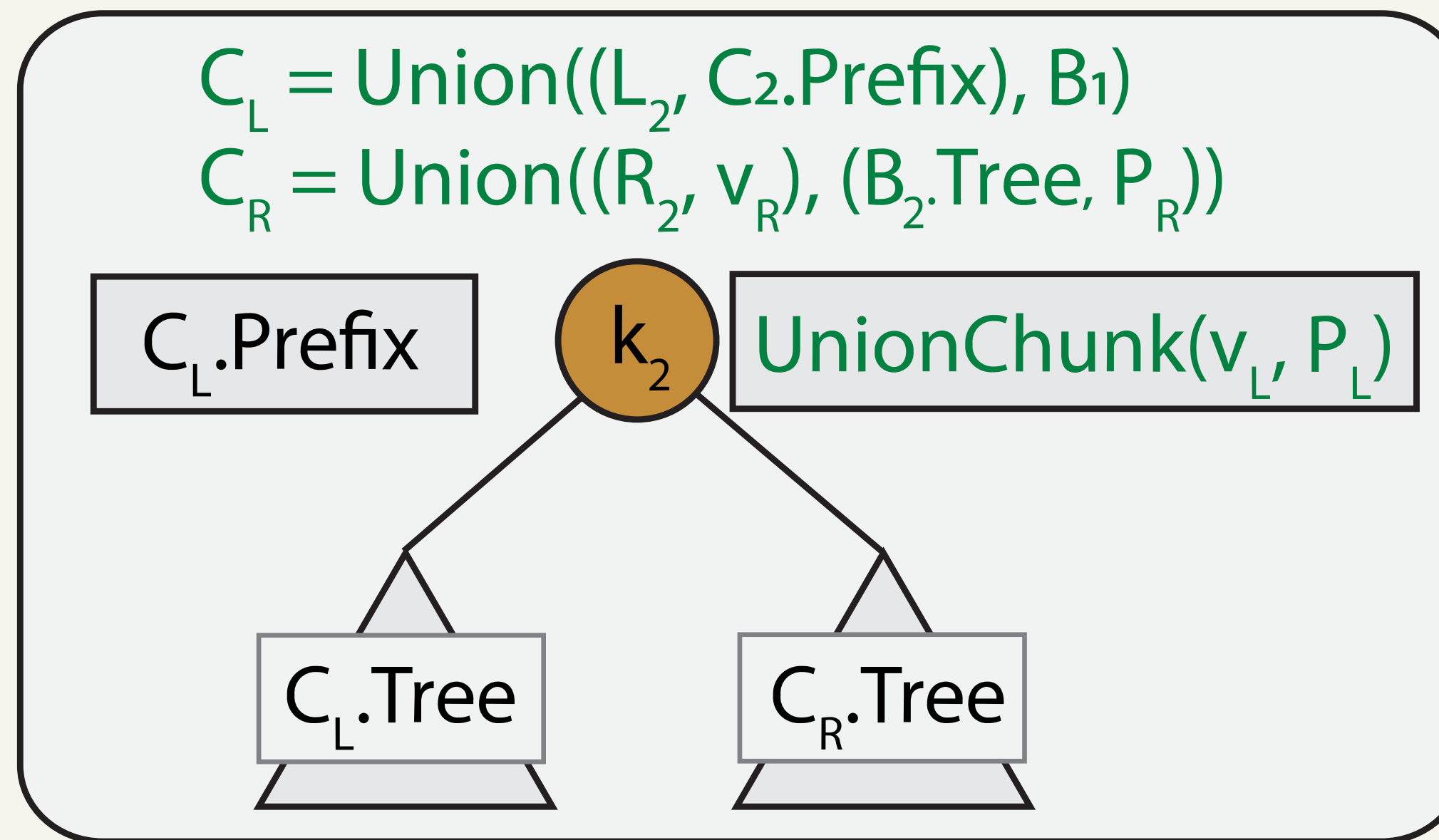
# Batch Updates on Trees

$\text{union}(C_1=(T_1, P_1), C_2=(T_2, P_2))$



Recursive union of two C-trees

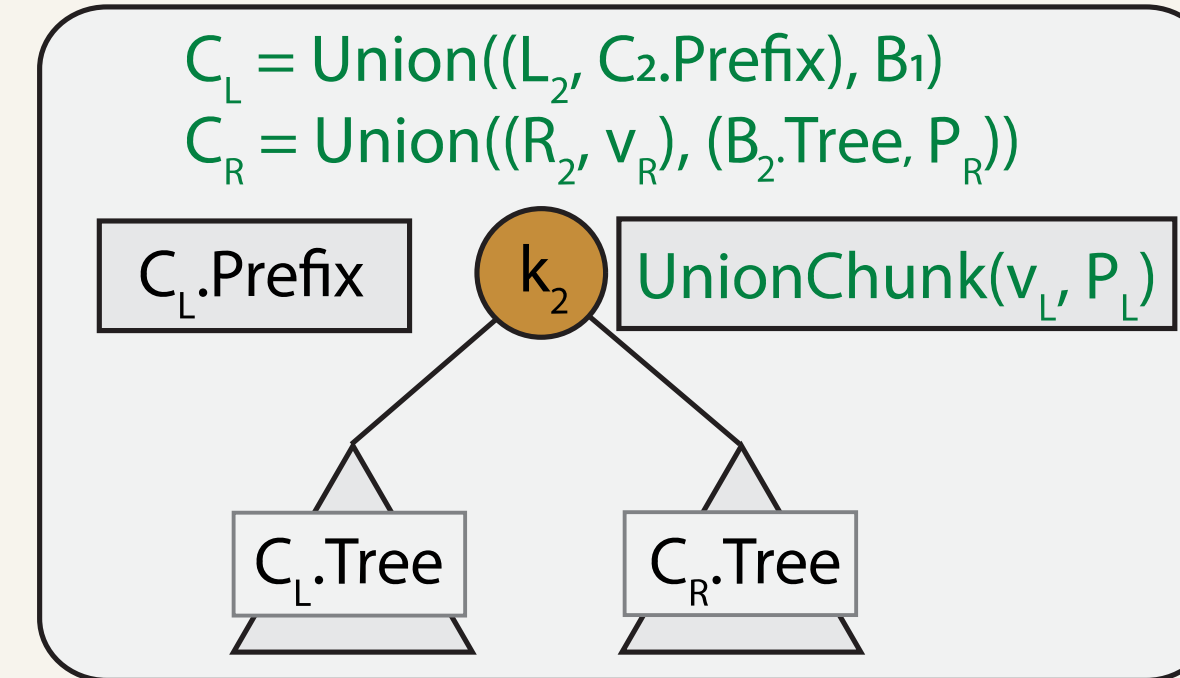
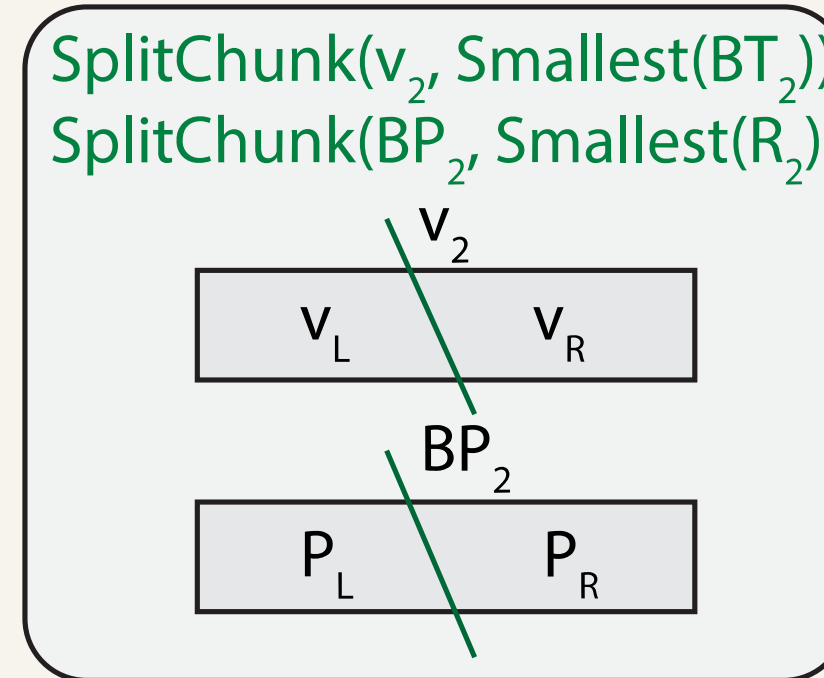
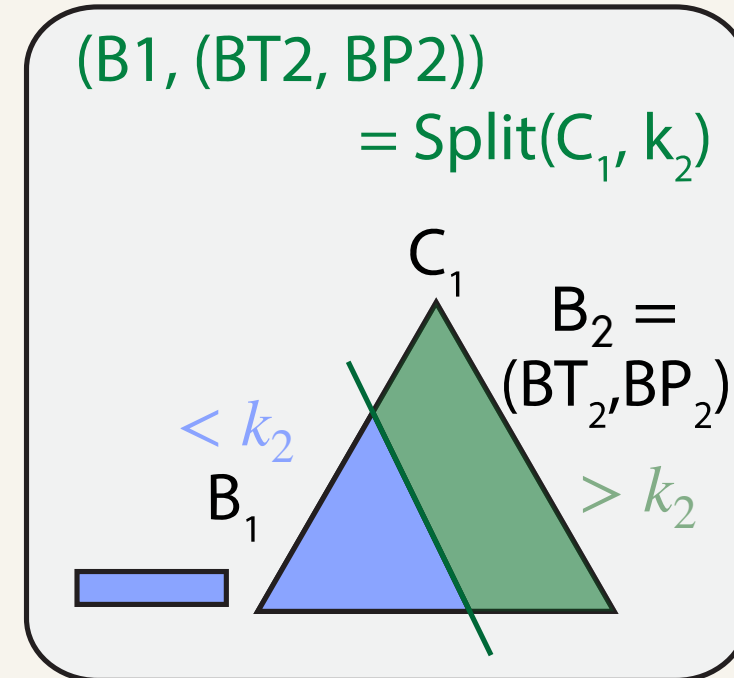
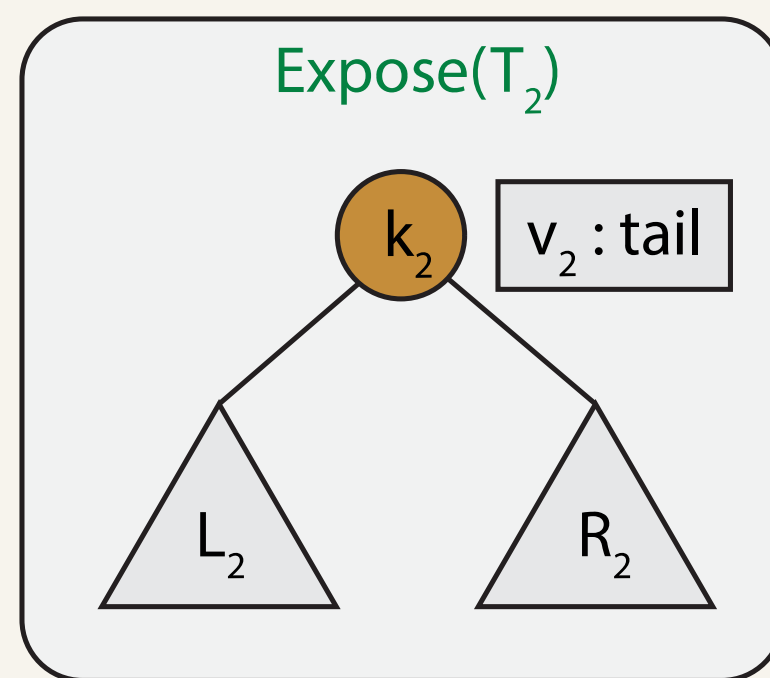
Join done on the underlying purely-functional tree





# Batch Updates on Trees

$\text{union}(C_1=(T_1, P_1), C_2=(T_2, P_2))$



$\text{union}(C_1, C_2)$  runs in

$$O\left(B^2 m \log\left(\frac{n}{m} + 1\right)\right) \text{ expected work}$$

$$O(B \log n \log m) \text{ depth whp}$$

# Experiments

# Our Machine

## Dell PowerEdge R930

- ❖ 72-cores, 2-way hyper-threaded\*
- ❖ 1TB of main memory
- ❖ Cost: about 20k USD



\* (4 x 2.4GHz 18-core E7-8867 v4 Xeon processors)

# Our Machine

## Dell PowerEdge R930

- ❖ 72-cores, 2-way hyper-threaded\*
- ❖ 1TB of main memory
- ❖ Cost: about 20k USD

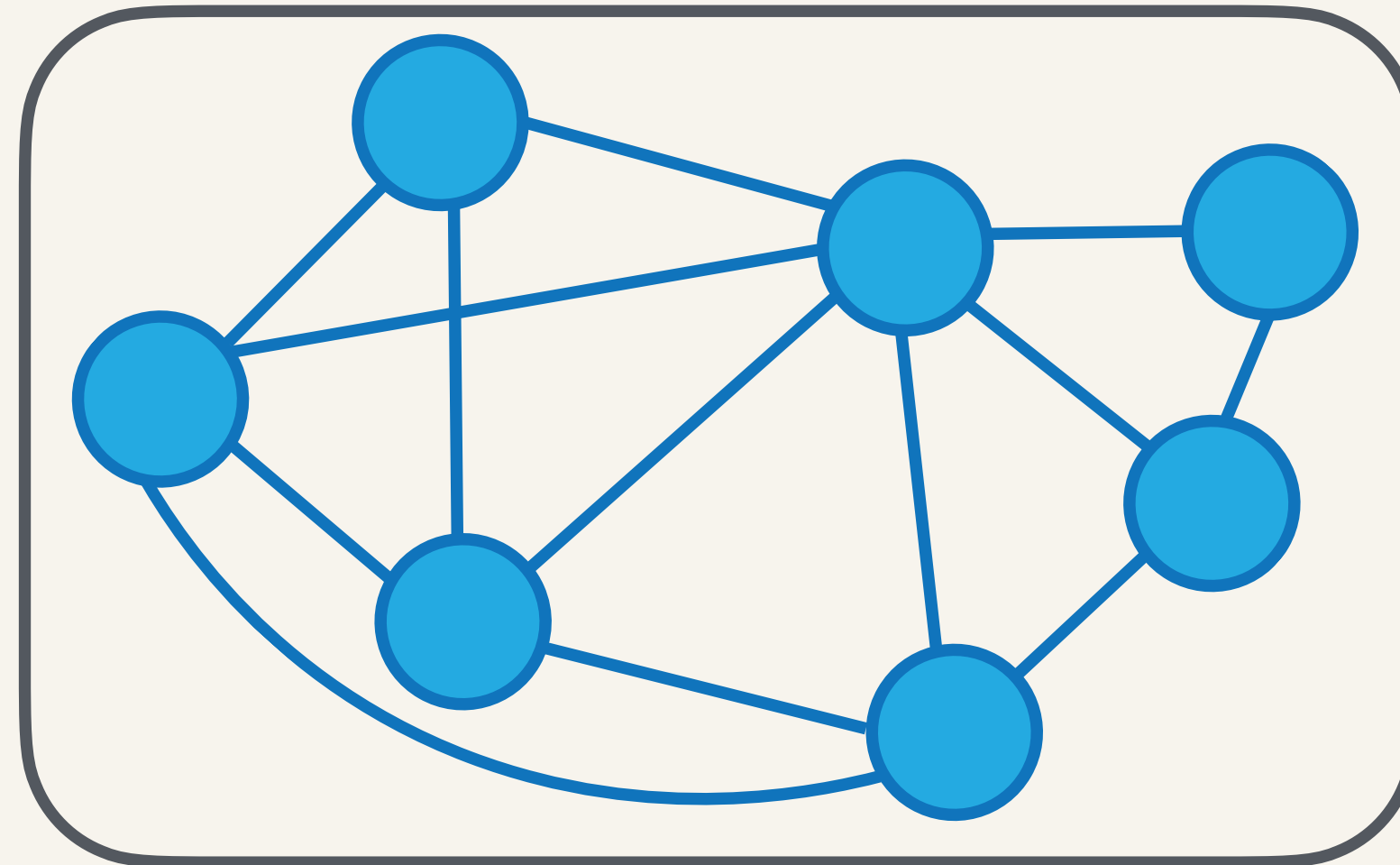


\* (4 x 2.4GHz 18-core E7-8867 v4 Xeon processors)

# Streaming Experiment

# Streaming Experiment

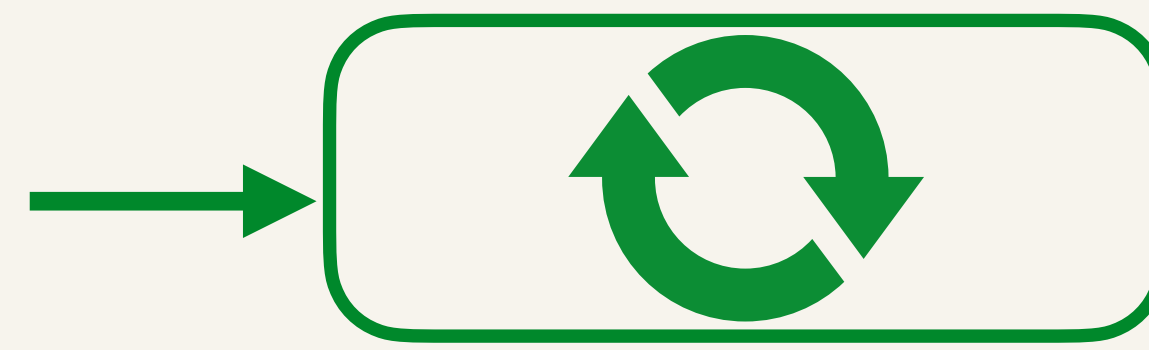
Sampled insertions  
+ deletions from  $G$



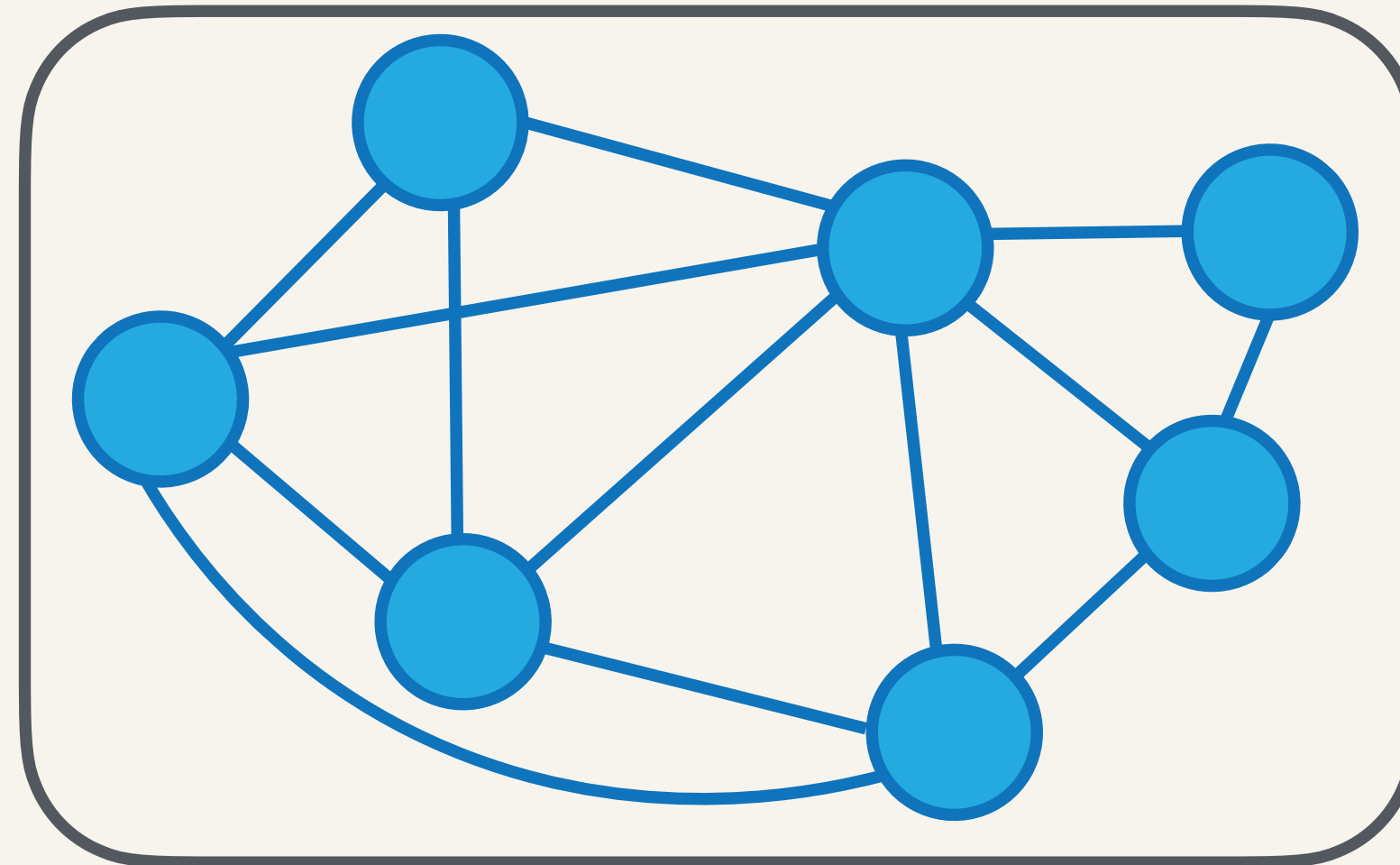
# Streaming Experiment

Stream of Parallel  
BFS queries from  
random vertices

BFS trees



Sampled insertions  
+ deletions from  $G$

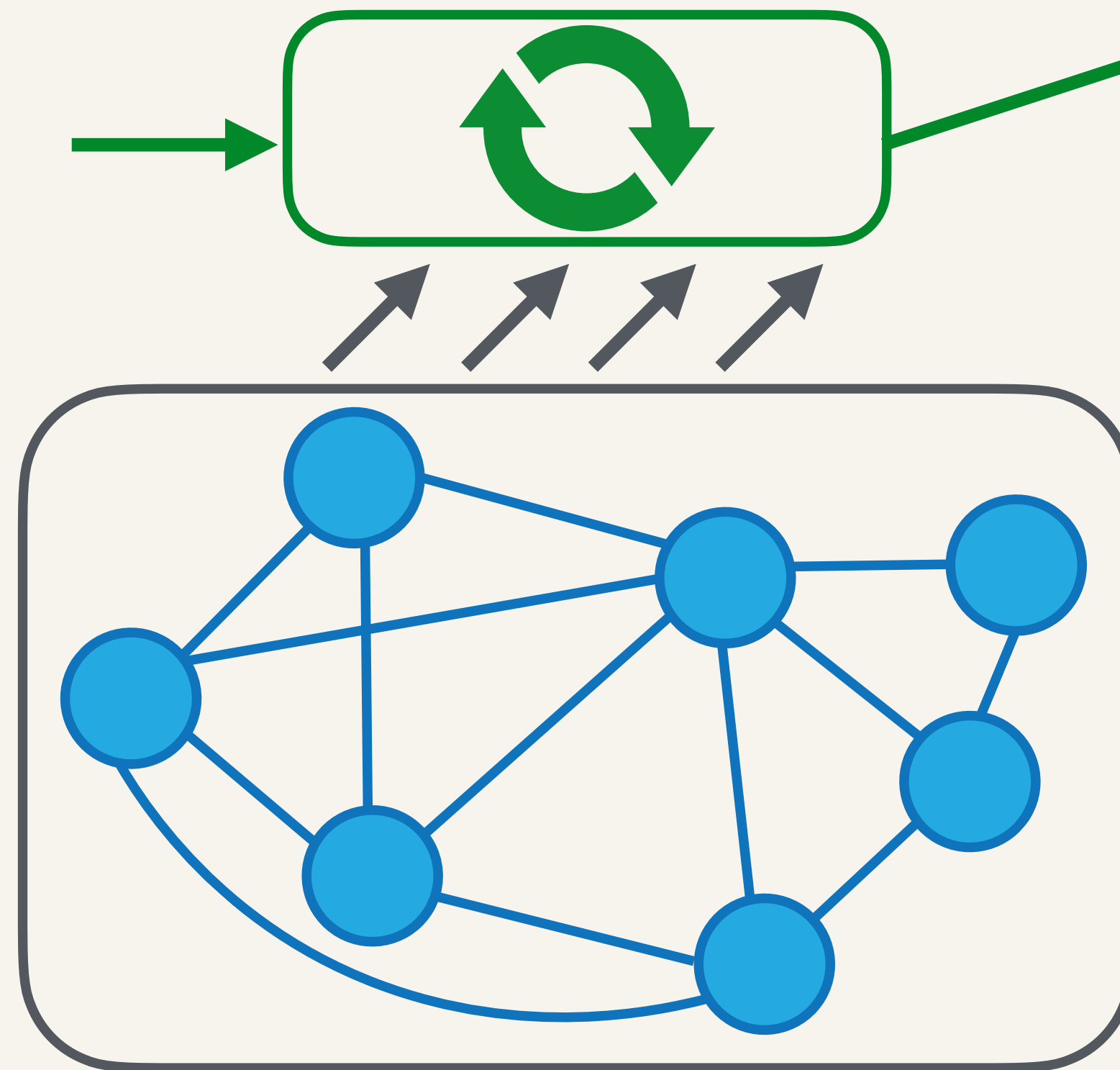


# Streaming Experiment

Stream of Parallel  
BFS queries from  
random vertices

BFS trees

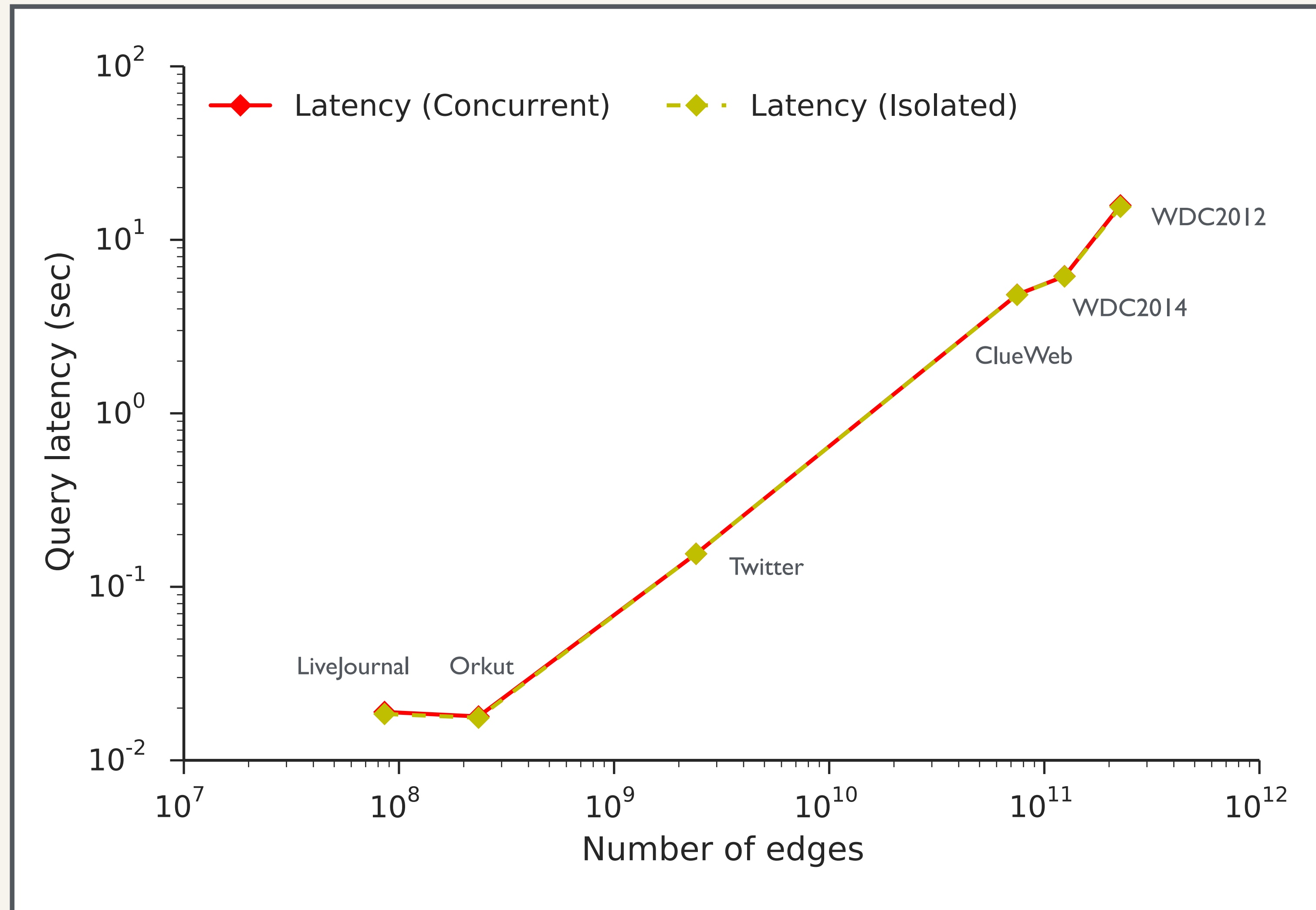
Sampled insertions  
+ deletions from  $G$



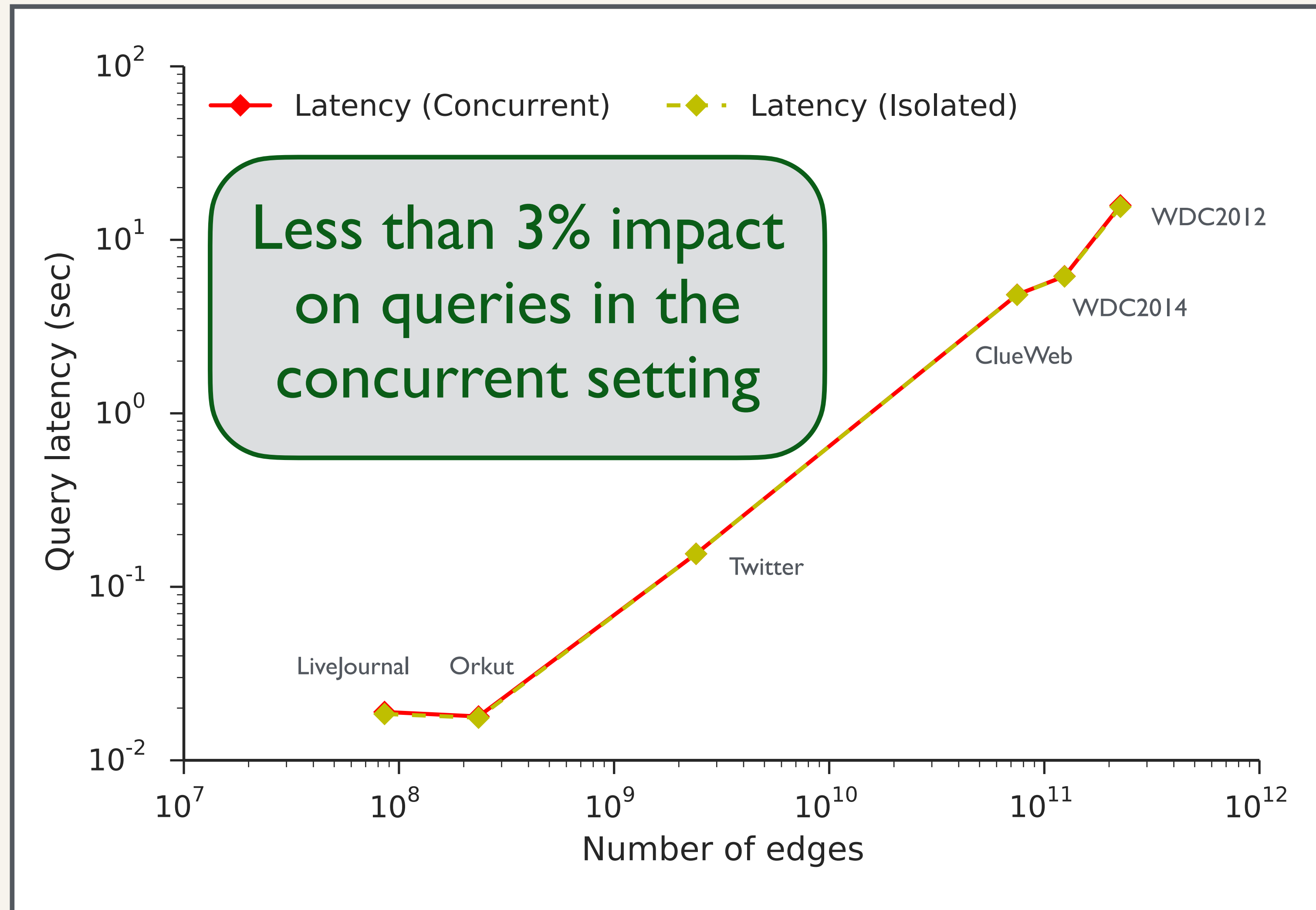
What's the impact on the  
concurrent execution on latency?



# Streaming Experiment



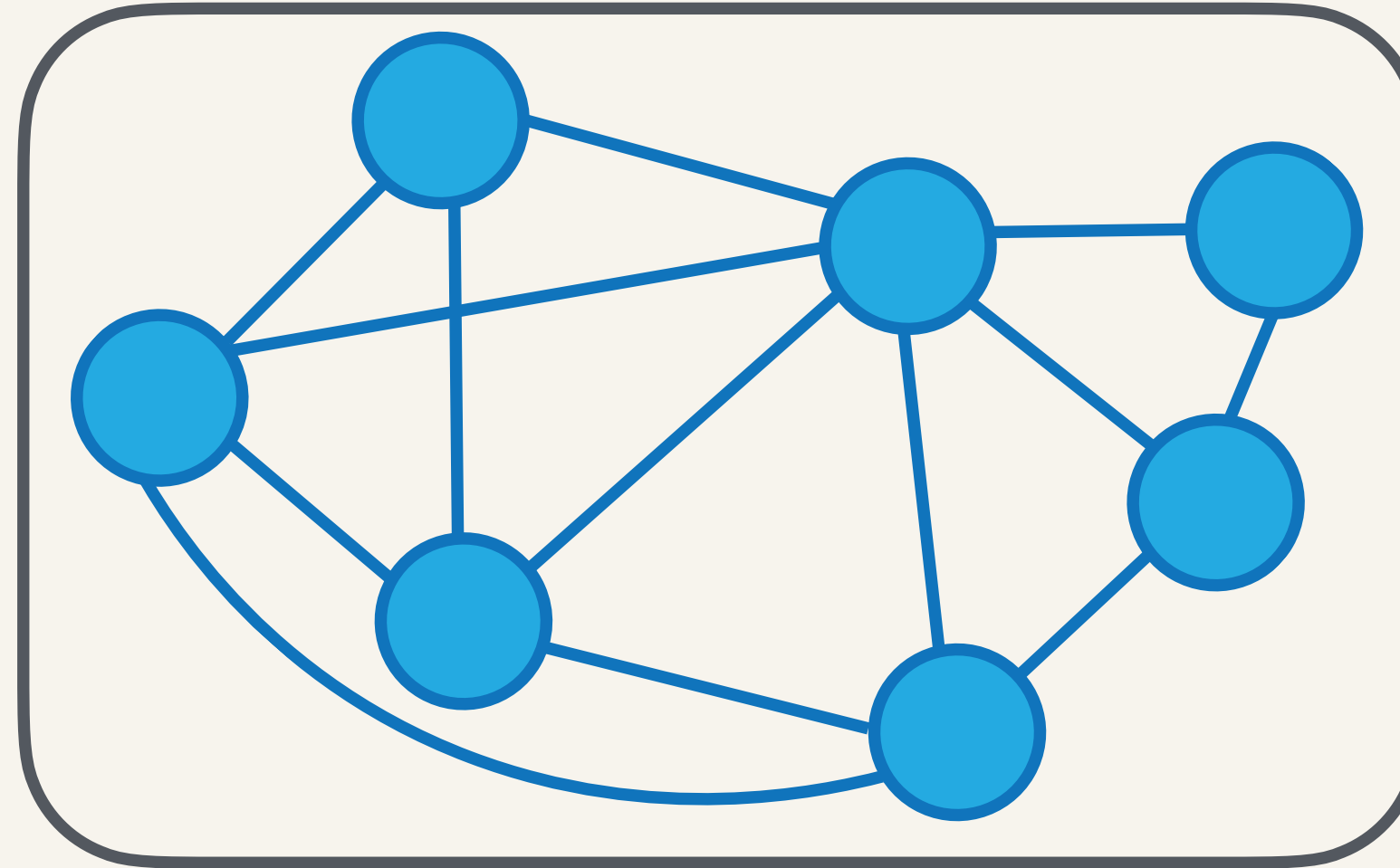
# Streaming Experiment



# Batch Update Experiment

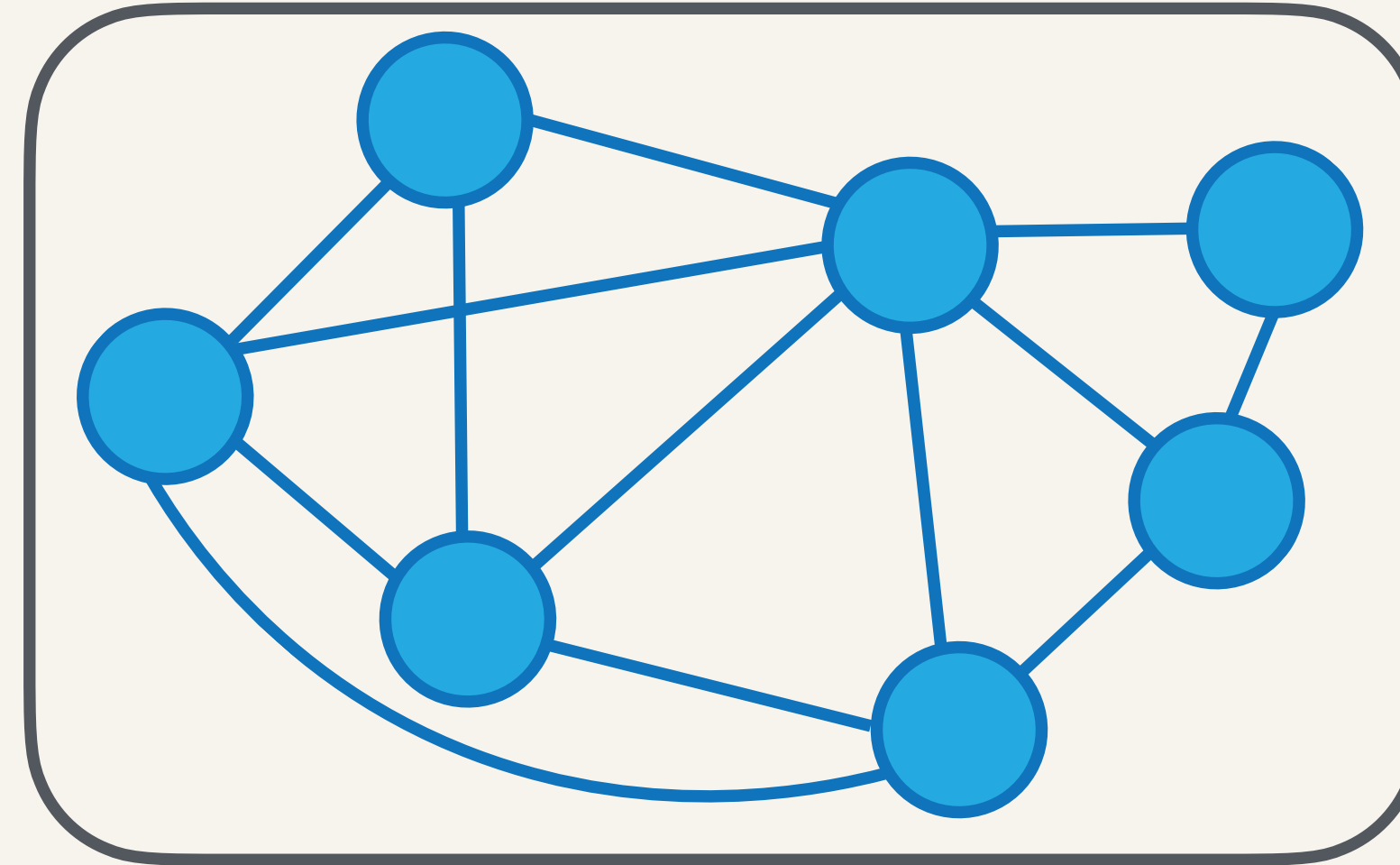
# Batch Update Experiment

Edge insertions  
drawn from RMAT



# Batch Update Experiment

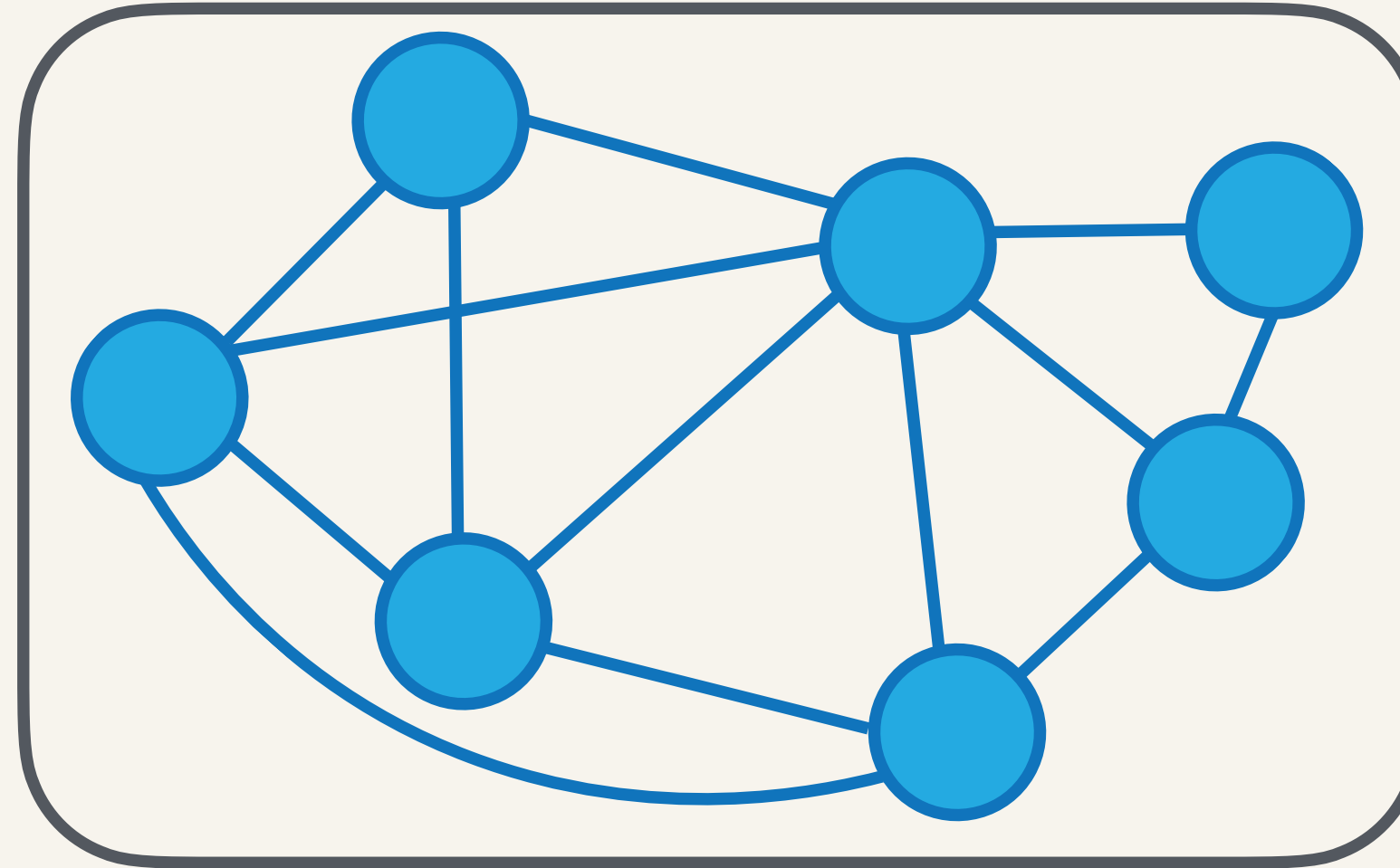
Edge insertions  
drawn from RMAT



Represent  $G$  using  
Aspen and **STINGER**

# Batch Update Experiment

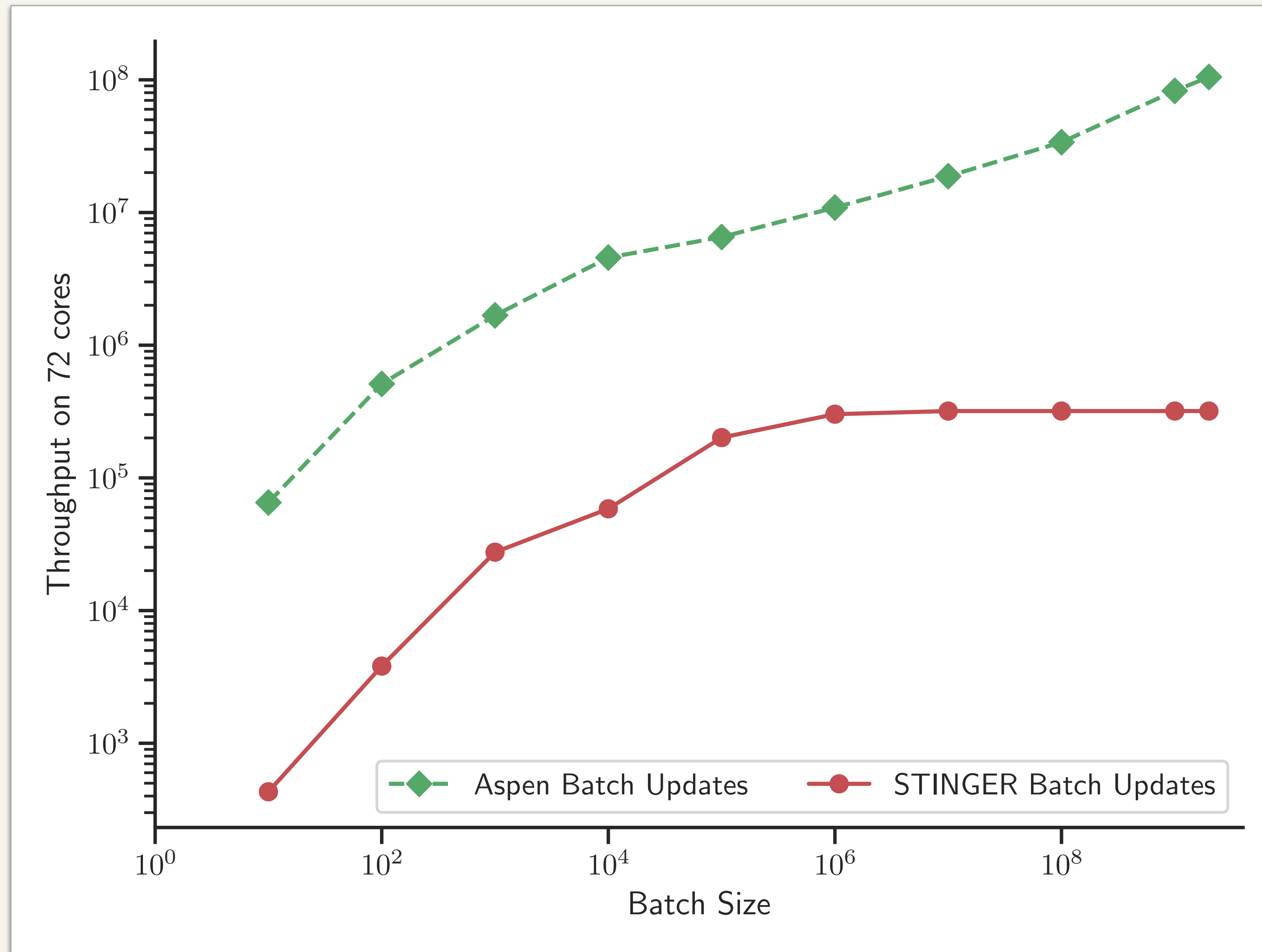
Edge insertions  
drawn from RMAT



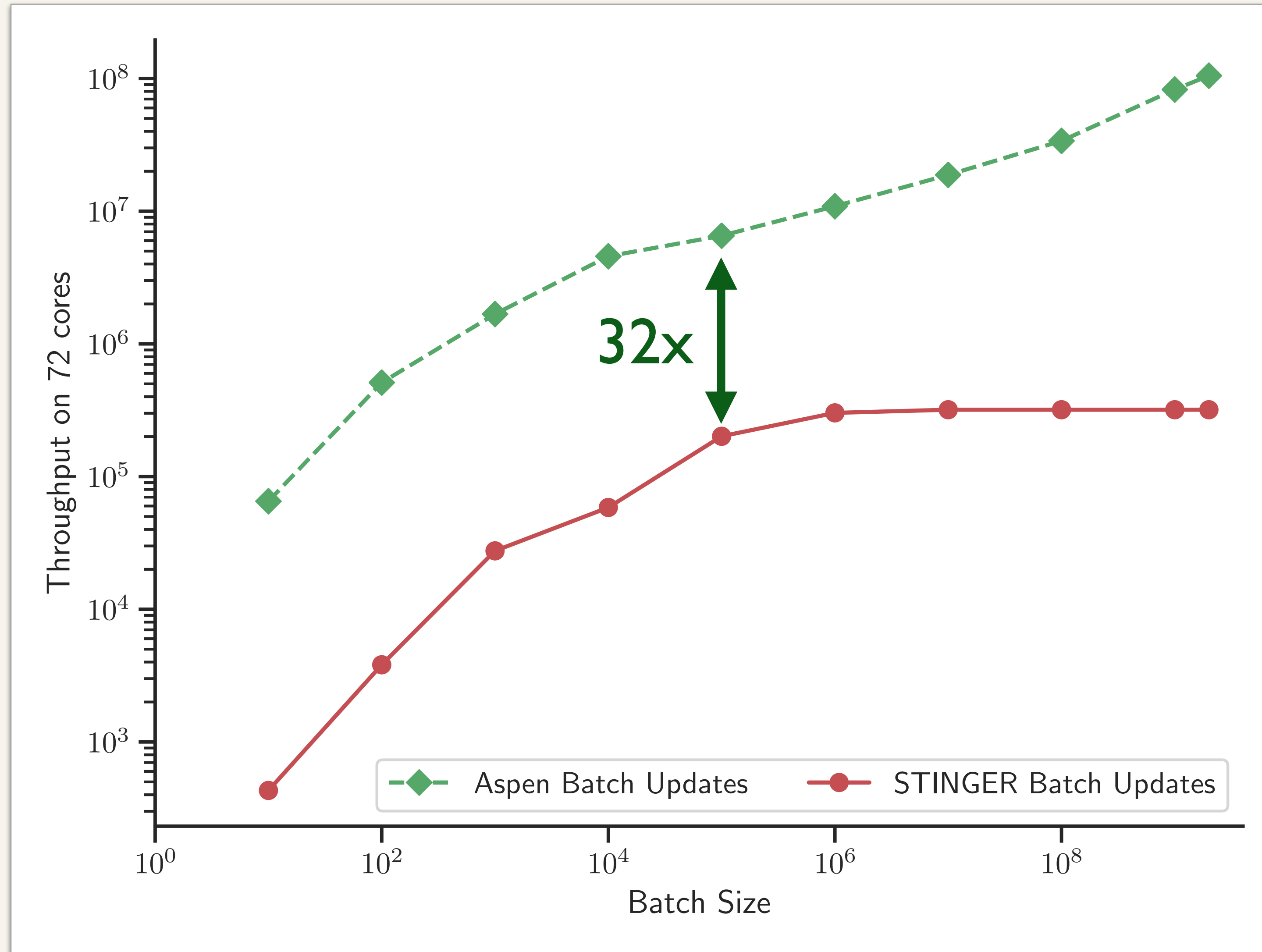
Represent  $G$  using  
Aspen and **STINGER**

How does the throughput scale  
as a function of batch size?

# Batch Update Performance

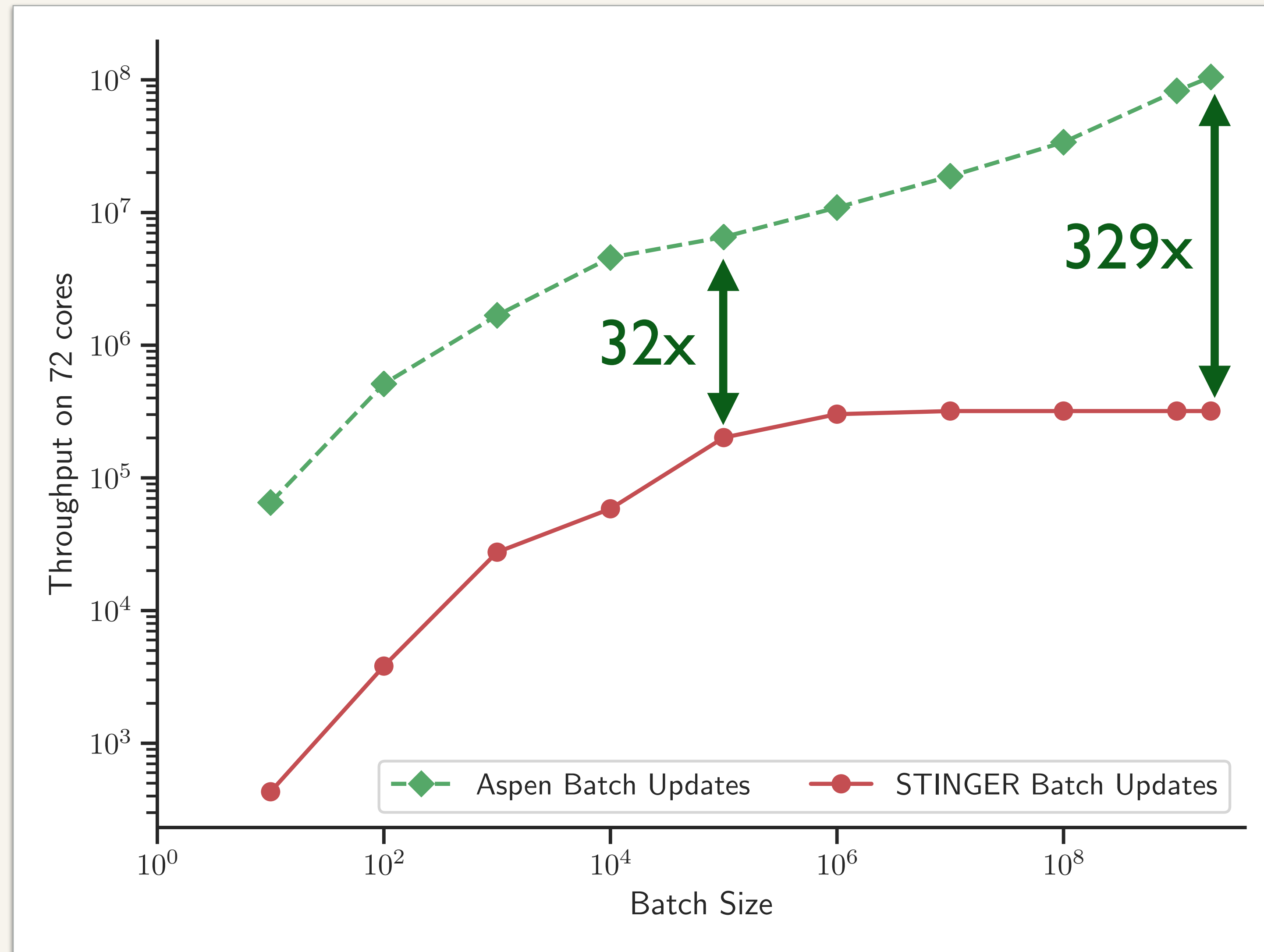


# Batch Update Performance





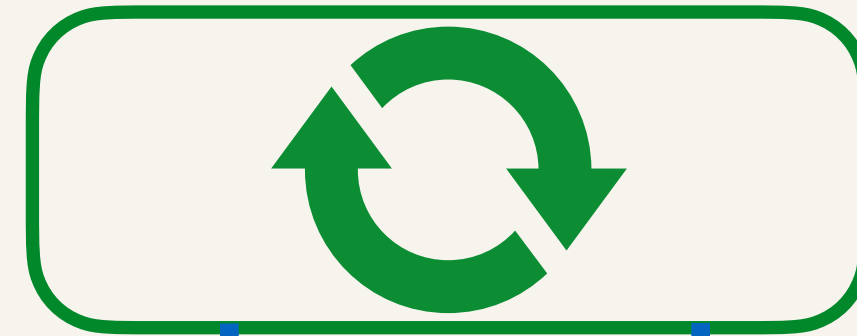
# Batch Update Performance



# Building on Aspen and C-trees

# Batch-Dynamic Graph Processing

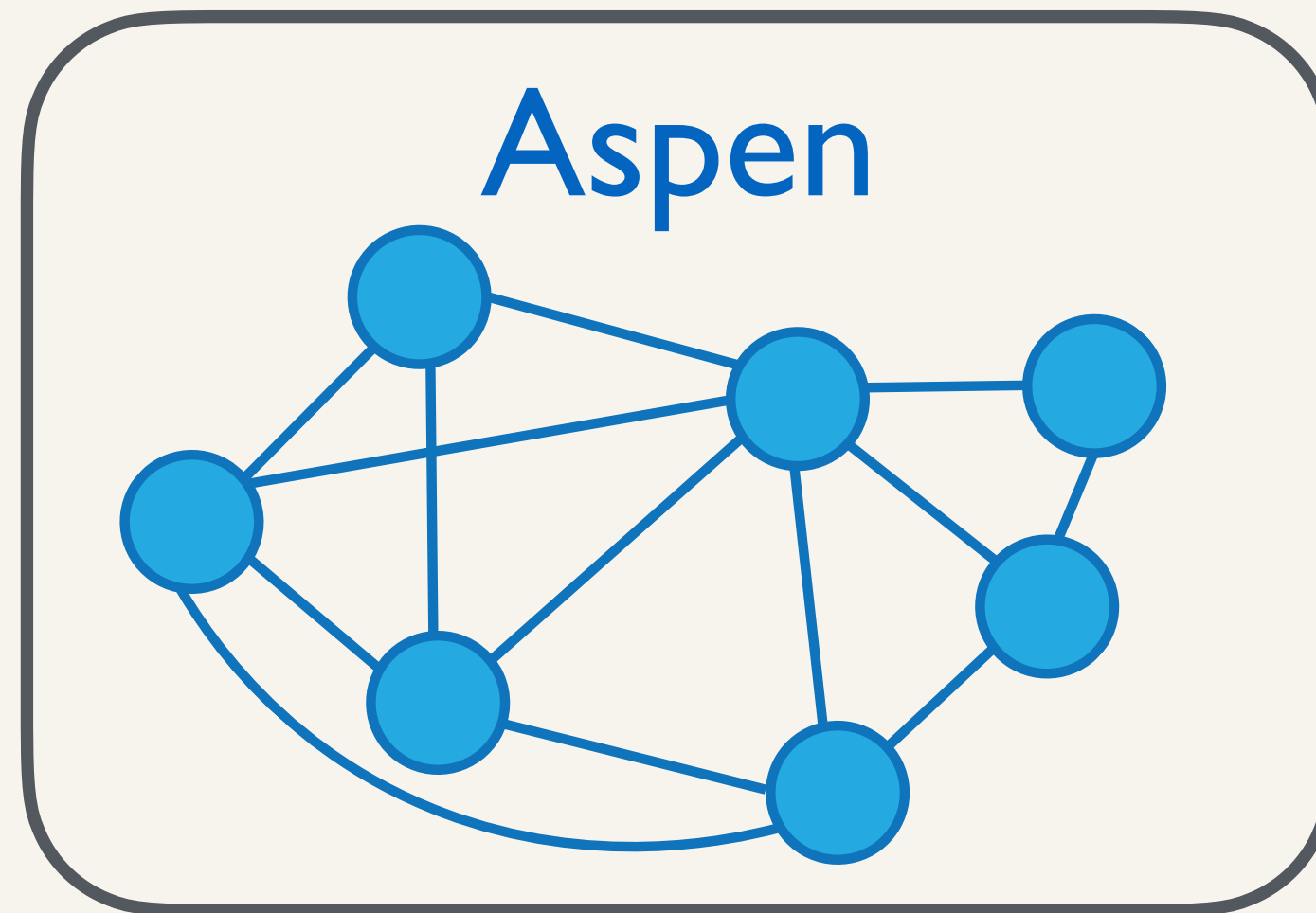
Dynamic Algorithm



$G_{\text{Prev}}$

$G_{\text{Next}}$

Updates



E.g.: Connected components, clustering coefficients, graph clusterings, etc

# Batch-Dynamic Algorithms

## Forest Conn.



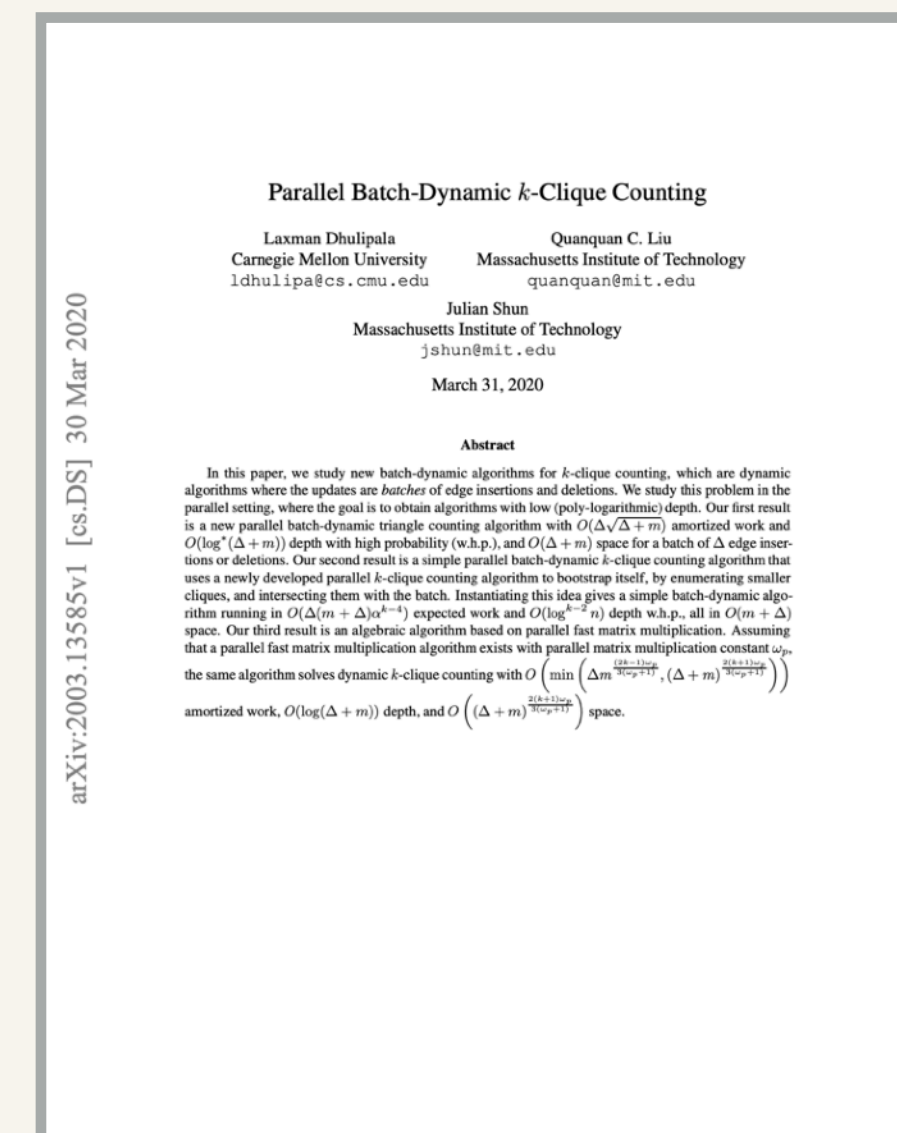
[TDB'18]

## Connectivity



[AABD'19]

## Clique-counting



[DLSY'20]

Interested in practical and memory-efficient dynamic graph algorithms

---

# Thank you!

---

## Aspen

Scalable graph data structures and interfaces for processing streaming graphs

[github.com/ldhulipala/aspen](https://github.com/ldhulipala/aspen)

- ❖ has strong theoretical bounds
- ❖ provides memory-efficient graph representations
- ❖ enables lightweight snapshots
- ❖ runs on commodity hardware
- ❖ can process the largest publicly-available graphs