

1

#### Concurrent Wait-Free Graph Snapshots using Multi-Versioning

**NETYS 2024** 

Gaurav Bhardwaj, Ayaz Ahmed, Dr. Sathya Peri Indian Institute of Technology, Hyderabad.

This project is funded by SERB Project CRG/2022/009391.



#### Content

- Background
- Motivation
- Graph Data Structure
- Graph ADT operations
- Memory Reclamation
- Results



#### Graphs

- Graphs provide a visual representation of relationships and connections in data, making complex structures easier to understand.
- Graphs are ideal for modeling various relationships such as social networks, web links, and dependencies between tasks.
- Graphs allow for flexible data representation, enabling the modeling of diverse scenarios and real-world interactions.





### **Concurrent Graphs:**

- In applications like social media analytics and network monitoring, concurrent graphs enable real-time analysis of dynamic data.
- Concurrent graphs are essential in large-scale systems, ensuring scalability and efficient processing of data in parallel.
- Concurrent graphs optimize resource utilization in multi-core processors, enhancing the performance of graph-based algorithms.





### **Correctness (Consistency)**

- The ADT operations implemented by the data structure are represented by their invocation and return steps.
- For an arbitrary concurrent execution of a set of ADT operations should satisfy the consistency framework **Linearizability**.



#### Linearizability

- Assign an atomic step as a linearization point (LP) inside the execution interval of each of the operations and show that the data structure invariants are maintained across the LPs.
- An arbitrary concurrent execution is equivalent to a valid sequential execution obtained by ordering the operations by their LPs.



#### Linearizability Example





# **Non-Blocking Progress Condition**

An execution is said to be Non-Blocking if it doesn't blocks the execution of other threads.

- **Obstruction Free:** A thread is guaranteed to finish in a finite number of steps in isolation.
- Lock-Freedom: Atleast one thread should be able to finish in the finite number of steps.
- Wait-Freedom: All threads should be able to finish in a finite number of steps.

#### **Related Work**



- Chatterjee et al.<sup>a</sup>
  - Proposed a Lock-Free dynamic unbounded unweighted graph.
  - Used Lock-Free Linked Lists to store vertices and edges.
  - Used obstruction free method to perform the partial snapshot on the graph.
- Chatterjee et al.<sup>b</sup>
  - Extended the previous work for speed up.
  - Proposed a Lock-Free dynamic unbounded weighted graph.
  - Used Lock-Free Hash Tables to store vertices and Binary Search tree to store edges.
  - Used obstruction free method to perform the partial snapshot on the graph.
- Bhardwaj et al.<sup>c</sup>
  - Extended the previous work with a wait-free snapshot operation.
  - Used wait-free snapshot algorithm to perform the complete snapshot of the graph.
  - Snapshot of the graph is used to perform the graph analytics operation.

<sup>a</sup> "A Simple and Practical Concurrent Non-blocking Unbounded Graph with Reachability Queries", Bapi Chatterjee, Sathya Peri, Muktikanta Sa, Nandini Singhal in the 20<sup>th</sup> International Conference on Distributed Computing and Networking (ICDCN), Bangalore, India, January 2019.

<sup>b</sup> "Non-blocking dynamic unbounded graphs with worst-case amortized bounds", Bapi Chatterjee, Sathya Peri, Muktikanta Sa, Manogana In: International Conference on Principles of Distributed Systems (2021)

<sup>c</sup> "Non-blocking Dynamic Unbounded Graphs with Wait-Free Snapshot", Gaurav Bhardwaj, Sathya Peri, and Pratik Shetty: 25th International Symposium on Stabilization, Safety, and Security (2023)



#### **Drawback:**

- Current Literature supports obstruction free partial snapshot or wait-free complete snapshot.
- As a result, a partial snapshot operation may continuously get interrupted
   Hence may never terminate.
- Complete Snapshot is costly as compared to the partial snapshot.
  - Graph analytics operation which requires partial snapshot still has to undergo complete snapshot.



#### **Motivation**

• To Construct a **wait-free partial** Snapshot for an Unbounded graph while allowing concurrent non-blocking (lock-free) *update* methods.



#### **Multiversioning:**

- A multiversion object maintains its previous versions, so threads can have access to the history of the object.
- Widely used in:
  - Database
  - STM
  - Concurrent Data structure
- Threads can have access to the history of the object.
- Allows to update the data structure without hindering the snapshot.



#### **Graph Data Structure:**

- Each vertex in the graph is stored as **Vnode** in a concurrent lock-free resizable hash table.
- Each **Vnode** has a key and a pointer to the **vertex version** list.
- Each node (**Vnode\_version**) in the vertex version list contains the pointer to the corresponding edge list stored in a lock-free linked list.
- Each **Vnode\_version** node contains the insert timestamp and a delete timestamp.

int global_ts	class Vnode_version {
	int k
class VNode{	int insert_ts
int k	int delete_ts
Vnode_version * vhead	Enode* elist
}	Vnode_vesion* next_ver
	}

# **Graph Data Structure:**



- Each edge is stored as **Enode** in a concurrent lock-free linked list (edge list) associated with each version node.
- Enode stores the key of the destination vertex and the pointer to the edge version list.
- Each node (Enode\_version) in the edge version list contains weight, insert timestamp and delete timestamp.
- Enode\_version points to the next version and also stores the reference to the destination Vnode\_version.

```
class ENode{
int l
ENode* enext
Enode_version* ehead
}
```

```
class Enode_version {
    int weight
    int insert_ts
    int delete_ts
    Vnode_version* dest_ver_node
    Enode_version* next_ver
```

#### **Graph Data Structure**







#### **ADT:**

- AddVertex(k)
- RemoveVertex(k)
- AddEdge(*i*, *j*,*w*)
- RemoveEdge(*i*, *j*)
- ContainsVertex(k)
- ContainsEdge(*i*, *j*)
- **BFS**(*v*)
- **SSSP**(*v*)
- Snapshot()



# Addvertex(10) : Case-1

• Vnode 10 doesn't Exists





#### Addvertex(10) : Case-1

• Vnode 10 is inserted with insert\_ts initialized to -1





### Addvertex(10) : Case-1

• Insert timestamp is set to current global timestamp





### Addvertex(15) : Case-2

• Case2: When Vnode exists but there is no active version.





### Addvertex(15) : Case-2

• A new Vnode\_version is added on the top of the list with insert timestamp as -1





# Addvertex(15) : Case-2

• Insert timestamp is set to current global timestamp



Sec 5/7 22/45



# Addvertex(18) : Case-3

• Already exists.





• Active version\_vertex exists.





• Set the delete timestamp to - 1.





• Set the delete timestamp to current global timestamp.





• Set TS if -1.





• Set the delete timestamp to current global timestamp.





• Delete Timestamp of the vertex is already set





#### Case1: Vertex 15 doesn't exists.





#### Case2: Enode(10,15) doesn't exists.





#### Case2: New Enode is created with insert timestamp as -1.



Sec 5/7 32/45



# Add Edge(10,15)

#### Case2: Set insert timestamp to current timestamp.





#### Case3: Adding a new Edge Version.



Sec 5/7 34/45



#### Case3: Adding a new Edge Version.





#### Case3: Adding a new Edge Version.



Sec 5/7 36/45



# **BFS(vertex 10) with TS 16**



Sec 5/7 37/45



# **BFS(vertex 10) with TS 16**



#### **Memory Reclamation**



- **Step-1:** Capture the smallest active snapshot timestamp among all the threads, say T.
- **Step-2:** Iterate the complete data structure and pass all the edge nodes having timestamp smaller than or equal to T-1 to DEBRA.
- **Step-3:** After removing all the reference to vertex nodes from edge nodes in step-2, re-iterate the complete data structure and remove all the vertex nodes having timestamp smaller than or equal to T-1 to DEBRA.

# IIT Hyderabad

#### **Results and Analysis :**

- Workload Distribution : The distribution is over the following ordered set of Operations (AddVertex, RemoveVertex, ContainsVertex, AddEdge, RemoveEdge, ContainsEdge, and Critical Operation(Snapshot/BFS/Betweenness Centrality).
  - Read Heavy Workload : 3%, 2%, 44%, 3%, 2%, 44% , 2%
  - Update Heavy Workload: 13%, 12%, 24%, 13%, 12%, 24% , 2%
- **Algorithms :** We compare our implementation to Obstruction-free and wait-free implementations of same operations namely
  - Obst-Free: "A Simple and Practical Concurrent Non-blocking Unbounded Graph with Reachability Queries", Bapi Chatterjee, Sathya Peri, Muktikanta Sa, Nandini Singhal in the 20<sup>th</sup> International Conference on Distributed Computing and Networking (ICDCN), Bangalore, India, January 2019.
  - **PANIGRAHAM:** "Non-blocking dynamic unbounded graphs with worst-case amortized bounds", Bapi Chatterjee, Sathya Peri, Muktikanta Sa, Manogana In: International Conference on Principles of Distributed Systems (2021)
  - **Graph\_Iterator:** "Non-blocking Dynamic Unbounded Graphs with Wait-Free Snapshot", Gaurav Bhardwaj, Sathya Peri, and Pratik Shetty: 25th International Symposium on Stabilization, Safety, and Security (2023)
  - LIGRA



### **Results: Point Operation**



**Update Heavy Workload** 

**Read Heavy Workload** 

Sec 7/7 41/45



#### **Results: Snapshot**



**Update Heavy Workload** 

**Read Heavy Workload** 

Sec 7/7 42/45

#### **Results: BFS**





**Update Heavy Workload** 

**Read Heavy Workload** 

Sec 7/7 43/45





**Update Heavy Workload** 

**Read Heavy Workload** 



#### **Conclusion:**

- First algorithm which supports partial wait-free snapshot for the graph.
- Significant improvement in the performance when compared to its counterparts.



