Lock-free Augmented Trees

Panagiota Fatourou, U. of Crete and FORTH, Greece Eric Ruppert, York University, Canada

> DISC 2024 Madrid, Spain



Panagiota Fatourou and Eric Ruppert Lo

Lock-free Augmented Trees

A D N A D N A D N A D

Main Result

Technique to augment lock-free search trees to support more operations.

- Simple to implement using single-word CAS
- General: can handle any augmentation
- Efficient: queries as fast as in sequential system, minimal overhead for updates
- Wait-free: additional work for augmentation is wait-free
- Snapshots of tree easily support complex queries

4 **A** N A **A** N A **A**



Panagiota Fatourou and Eric Ruppert Lock-free Augmented Trees



Leaf-oriented BST. Set = $\{A, C, E, G, L\}$.

・ロト ・部 ト ・ヨト ・ヨト



Augmentation: Each node stores number of leaves in subtree.

YORK

Panagiota Fatourou and Eric Ruppert Lock-free Augmented Trees



Augmentation: Each node stores number of leaves in subtree.

YORK





Augmentation: Each node stores number of leaves in subtree.

YORK





Augmentation: Each node stores number of leaves in subtree.

YORK





Keys have weights; each node stores sum of subtree's weights.

YORK

(日)

Panagiota Fatourou and Eric Ruppert Lock-free Augmented Trees



Keys have weights; each node stores sum of subtree's weights.

YORK

Panagiota Fatourou and Eric Ruppert Lock-free Augmented Trees

Many other ways to augment a BST.

- For database of employees, number of women in subtree. How many women's salaries are more than €100,000?
- Store min key, max key, and smallest gap in subtree. Find two closest keys in the set.

Key Property of Augmentations

Values of a node's new field(s) can be computed from information in node and its children.



A D N A D N A D N A D

Many other ways to augment a BST.

- For database of employees, number of women in subtree. How many women's salaries are more than €100,000?
- Store min key, max key, and smallest gap in subtree. Find two closest keys in the set.

Key Property of Augmentations

Values of a node's new field(s) can be computed from information in node and its children.



A D N A D N A D N A D

Many other ways to augment a BST.

- For database of employees, number of women in subtree. How many women's salaries are more than €100,000?
- Store min key, max key, and smallest gap in subtree. Find two closest keys in the set.

Key Property of Augmentations

Values of a node's new field(s) can be computed from information in node and its children.



(日)

Applications of Augmented Trees

Augmented BSTs are basis of many other data structures.

- Interval tree
- Tango tree
- Measure tree
- Priority search tree
- Segment tree
- Link/cut tree

• . . .

<text>

Lots of applications in computational geometry, databases, graph algorithms,

Main Result

Technique to augment lock-free search trees.

Example: Lock-free Binary Search Tree

- Can handle any augmentation.
- Adds only O(height) steps to insert, delete.
- Supports simple snapshots.
- Wait-free queries run sequential code.
- Based on BST from PODC 2014.



Main Result

Technique to augment lock-free search trees.

Example: Lock-free Binary Search Tree

- Can handle any augmentation.
- Adds only O(height) steps to insert, delete.
- Supports simple snapshots.
- Wait-free queries run sequential code.
- Based on BST from PODC 2014.



Updating an Augmented Tree



Panagiota Fatourou and Eric Ruppert

Updating an Augmentee Tree



Panagiota Fatourou and Eric Ruppert Loc

Updating an Augmented Tree



Panagiota Fatourou and Eric Ruppert

Updating an Augmented Tree



Panagiota Fatourou and Eric Ruppert

Goal

Create lock-free augmented tree.

Challenges

- An update changes many nodes along a path All changes must appear atomic
- Queries traverse a path while concurrent updates change it
- Contention: all updates need to modify root's size



A B > A B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

Goal

Create lock-free augmented tree.

Challenges

- An update changes many nodes along a path All changes must appear atomic
- Queries traverse a path while concurrent updates change it
- Contention: all updates need to modify root's size





Node stores pointer to current version of augmented field.

- Old versions can still be used by queries in progress.
- Pointers between versions provide consistent view.





Node stores pointer to current version of augmented field.

- Old versions can still be used by queries in progress.
- Pointers between versions provide consistent view.





Node stores pointer to current version of augmented field.

- Old versions can still be used by queries in progress.
- Pointers between versions provide consistent view.



Node stores pointer to current version of augmented field.

- Old versions can still be used by queries in progress.
- Pointers between versions provide consistent view.



Accessing root's version provides snapshot of version tree.

- Versions also store keys to direct searches.
- Supports any sequential query operation.
- Old versions are unreachable when no longer needed.





Accessing root's version provides snapshot of version tree.

- Versions also store keys to direct searches.
- Supports any sequential query operation.
- Old versions are unreachable when no longer needed.





Accessing root's version provides snapshot of version tree.

- Versions also store keys to direct searches.
- Supports any sequential query operation.
- Old versions are unreachable when no longer needed.



Propagate for each node x on path to root do (at most) twice // Refresh x's version create new version v v.left \leftarrow x.left.version v.right \leftarrow x.right.version compute contents of v CAS x.version to v



(日)







Image: A matrix





Image: A matrix

Propagatefor each node x on path to root
do (at most) twice
// Refresh x's version
create new version v
v.left \leftarrow x.left.version
compute contents of v
CAS x.version to v



Image: Image:

Propagatefor each node x on path to root
do (at most) twice
// Refresh x's version
create new version v
v.left \leftarrow x.left.version
v.right \leftarrow x.right.version•compute contents of v
CAS x.version to v



Image: Image:



Propagate for each node x on path to root do (at most) twice // Refresh x's version create new version v $v.left \leftarrow x.left.version$ $v.right \leftarrow x.right.version$ compute contents of v• CAS x.version to v



Image: Image:

Propagatefor each node x on path to rootdo (at most) twice// Refresh x's versioncreate new version vv.left \leftarrow x.left.versionv.right \leftarrow x.right.versioncompute contents of vCAS x.version to v



A B A B A
 B A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

Contents of version never change once it is attached to tree. Propagation is wait-free.





Panagiota Fatourou and Eric Ruppert





Panagiota Fatourou and Eric Ruppert

Lock-free Augmented Trees





Panagiota Fatourou and Eric Ruppert





Panagiota Fatourou and Eric Ruppert

Refresh on each node *x* uses CAS to update *x*'s version. What if the CAS fails?

Try again.

What if the CAS fails again?

Stop; someone else's refresh has propagated your change to *x*.

Cooperation and Contention

- Updates are propagated cooperatively
- One change can propagate many operations together

• All update operations perform CAS on root BUT not all have to succeed

Double Refresh

Refresh on each node *x* uses CAS to update *x*'s version. What if the CAS fails?

Try again.

What if the CAS fails again?

Stop; someone else's refresh has propagated your change to *x*.

Cooperation and Contention

- Updates are propagated cooperatively
- One change can propagate many operations together

• All update operations perform CAS on root BUT not all have to succeed Refresh on each node x uses CAS to update x's version.

What if the CAS fails?

Try again.

What if the CAS fails again?

Stop; someone else's refresh has propagated your change to *x*.

Cooperation and Contention

- Updates are propagated cooperatively
- One change can propagate many operations together
- All update operations perform CAS on root BUT not all have to succeed

A D N A D N A D N A D

Refresh on each node x uses CAS to update x's version.

What if the CAS fails?

Try again.

What if the CAS fails again?

Stop; someone else's refresh has propagated your change to x.

Cooperation and Contention

- Updates are propagated cooperatively
- One change can propagate many operations together

A D N A D N A D N A D

• All update operations perform CAS on root BUT not all have to succeed Refresh on each node *x* uses CAS to update *x*'s version.

What if the CAS fails?

Try again.

What if the CAS fails again?

Stop; someone else's refresh has propagated your change to x.

Cooperation and Contention

- Updates are propagated cooperatively
- One change can propagate many operations together
- All update operations perform CAS on root BUT not all have to succeed

Run original algorithm to perform update Refresh each ancestor (at most) twice

Adds O(height) to step complexity of updates.

Query operation

Read *Root.version* to get snapshot of version tree Run standard *sequential* algorithm on that snapshot

Step complexity same as *sequential* query time.



Run original algorithm to perform update Refresh each ancestor (at most) twice

Adds O(height) to step complexity of updates.

Query operation

Read *Root.version* to get snapshot of version tree Run standard *sequential* algorithm on that snapshot

Step complexity same as *sequential* query time.



Run original algorithm to perform update Refresh each ancestor (at most) twice

Adds O(height) to step complexity of updates.

Query operation

Read *Root.version* to get snapshot of version tree Run standard *sequential* algorithm on that snapshot

Step complexity same as *sequential* query time.



• • • • • • • • • • • • •

Run original algorithm to perform update Refresh each ancestor (at most) twice

Adds O(height) to step complexity of updates.

Query operation

Read *Root.version* to get snapshot of version tree Run standard *sequential* algorithm on that snapshot

Step complexity same as *sequential* query time.



A B A B A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 B
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

Challenges

- Updates on original tree, queries on version tree
- The two trees might look very different

(Strong) Linearizability

- Define *arrival point* of each update at each node on its leaf-to-root path.
- Invariant: tree rooted at x.version reflects all operations that have arrived at x (done in order of their arrival points).
- Linearization point = arrival point at the root.
- Show that propagate ensures all operations arrive at root.



Challenges

- Updates on original tree, queries on version tree
- The two trees might look very different

(Strong) Linearizability

- Define *arrival point* of each update at each node on its leaf-to-root path.
- Invariant: tree rooted at x.version reflects all operations that have arrived at x (done in order of their arrival points).
- Linearization point = arrival point at the root.
- Show that propagate ensures all operations arrive at root.



Challenges

- Updates on original tree, queries on version tree
- The two trees might look very different

(Strong) Linearizability

- Define arrival point of each update at each node on its leaf-to-root path.
- Invariant: tree rooted at x.version reflects all operations that have arrived at x (done in order of their arrival points).
- Linearization point = arrival point at the root.
- Show that propagate ensures all operations arrive at root.



Challenges

- Updates on original tree, queries on version tree
- The two trees might look very different

(Strong) Linearizability

- Define *arrival point* of each update at each node on its leaf-to-root path.
- Invariant: tree rooted at x.version reflects all operations that have arrived at x (done in order of their arrival points).
- Linearization point = arrival point at the root.
- Show that propagate ensures all operations arrive at root.

YOR

Improving Query Time

We can improve query time for unbalanced trees.

Simply plug in a different augmentation: *x.version* is a red-black tree of all keys in *x*'s subtree



Join can be done in $O(\log n)$ time. Update takes additional $O(height \cdot \log n)$ steps. Queries take $O(\log n)$ steps, even if tree height is $\Theta(n)$.



Related Work

• Lock-free BST augmented with size

[Kokorin, Alistarh, Aksenov IPDPS 2024]

- Each operation must join a queue at each node and help all those ahead.
- Not generalizable to other augmentations.
- $\Omega((\# processes) \cdot height)$ steps per operation.
- Lock-based tree augmentation [Sela, Petrank DISC 2024]
- Much work on taking snapshots of shared data structures
 - They are more complicated, and have slower queries
 - Those based on multiversioning have complex GC
- Double refresh has been used in other ways [Afek, Dauber, Touitou 1995]

• Lock-free BST augmented with size

[Kokorin, Alistarh, Aksenov IPDPS 2024]

- Each operation must join a queue at each node and help all those ahead.
- Not generalizable to other augmentations.
- $\Omega((\# processes) \cdot height)$ steps per operation.
- Lock-based tree augmentation [Sela, Petrank DISC 2024]
- Much work on taking snapshots of shared data structures
 - They are more complicated, and have slower queries
 - Those based on multiversioning have complex GC
- Double refresh has been used in other ways [Afek, Dauber, Touitou 1995]

• Lock-free BST augmented with size

[Kokorin, Alistarh, Aksenov IPDPS 2024]

- Each operation must join a queue at each node and help all those ahead.
- Not generalizable to other augmentations.
- $\Omega((\# processes) \cdot height)$ steps per operation.
- Lock-based tree augmentation [Sela, Petrank DISC 2024]
- Much work on taking snapshots of shared data structures
 - They are more complicated, and have slower queries
 - Those based on multiversioning have complex GC
- Double refresh has been used in other ways [Afek, Dauber, Touitou 1995]

• Lock-free BST augmented with size

[Kokorin, Alistarh, Aksenov IPDPS 2024]

- Each operation must join a queue at each node and help all those ahead.
- Not generalizable to other augmentations.
- $\Omega((\# processes) \cdot height)$ steps per operation.
- Lock-based tree augmentation [Sela, Petrank DISC 2024]
- Much work on taking snapshots of shared data structures
 - They are more complicated, and have slower queries
 - Those based on multiversioning have complex GC
- Double refresh has been used in other ways [Afek, Dauber, Touitou 1995]



Our scheme for augmenting concurrent trees

- is simple to implement
- works for any augmentation
- adds O(height) to step complexity of updates
- preserves lock- or wait-freedom of updates
- has wait-free, fast queries
- supports simple snapshots



< 17 ▶