

More Scalable Ordered Set for ETS Using Adaptation

Kostis Sagonas
Kjell Winblad

Department of Information Technology
Uppsala University, Sweden

5 Sept 2014





What is ETS?

Context

Motivation

CA Trees

Performance

Future Work

Conclusion

- Erlang Term Storage
- Key-value store
- In-memory database
- Shared memory



ETS Table Types and their Options

Context

Motivation

CA Trees

Performance

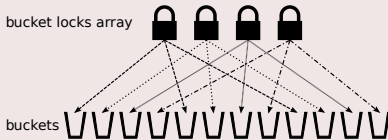
Future Work

Conclusion

- Creates a new shared **hash based** table with **write concurrency** enabled

```
TID1 = ets:new(tab, [set, public,  
                    {write_concurrency, true}])
```

- **write concurrency** enables fine grained locking



- Other hash based table types are **bag** and **duplicate bag**



ETS Table Types and their Options

Context

Motivation

CA Trees

Performance

Future Work

Conclusion

- Creates a new **search tree based** table with **read concurrency** enabled

```
TID2 = ets:new(tab, [ordered_set, public,  
                    {read_concurrency, true}])
```

- Current ordered set does not support fine grained locking
- read concurrency enables frequent-read-optimized readers-writer locks
 - One cache line per scheduler
 - Read only operations do not interfere each other



Use of ETS

Context

Motivation

CA Trees

Performance

Future Work

Conclusion

- Insert an Erlang tuple into a table

```
ets:insert(TID, {42, "A value"}),
```



Use of ETS

Context

Motivation

CA Trees

Performance

Future Work

Conclusion

- Insert an Erlang tuple into a table

```
ets:insert(TID, {42, "A value"}),
```

- Lookup a value

```
ets:lookup(TID, 42),
```



Use of ETS

Context

Motivation

CA Trees

Performance

Future Work

Conclusion

- Insert an Erlang tuple into a table

```
ets:insert(TID, {42, "A value"}),
```

- Lookup a value

```
ets:lookup(TID, 42),
```

- Find first element

```
FirstElem = ets:first(TID),
```

- Find next element

```
NextElem = ets:next(TID, FirstElem).
```



Use of ETS

Context

Motivation

CA Trees

Performance

Future Work

Conclusion

- Insert an Erlang tuple into a table

```
ets:insert(TID, {42, "A value"}),
```

- Lookup a value

```
ets:lookup(TID, 42),
```

- Find first element

```
FirstElem = ets:first(TID),
```

- Find next element

```
NextElem = ets:next(TID, FirstElem).
```

- Many other functions like **delete**, **match**, **foldl**, etc.



Current ETS Scalability Benchmark

Benchmark

- ets_bench from bencheri
- key range $[1, 2^{21}]$, $2^{21} \approx 2 * 10^6$
- Three phases:
 - 1 Insert phase: inserts 2^{20} random keys, $2^{20} \approx 10^6$
 - 2 Update and read phase, parameter for percentage update
 - 3 Delete phase: deletes 10^6 random keys

Machine

- **Four** Intel(R) Xeon(R) CPU E5-4650 CPUs (2.70GHz), **eight cores each**
 - total of 32 physical cores, each **with hyperthreading**
- The machine has 128GB of RAM and is running Debian Linux 3.10.17-amd64 and **Erlang/OTP release 17.0**



100% Updates

Context

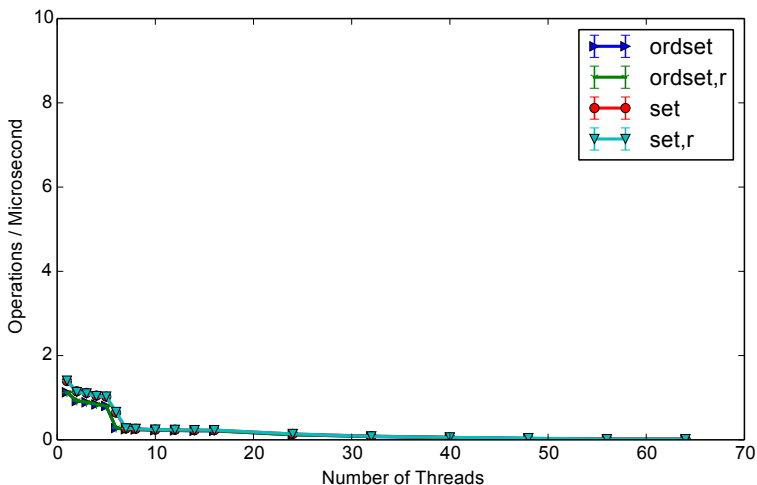
Motivation

CA Trees

Performance

Future Work

Conclusion





100% Updates

Context

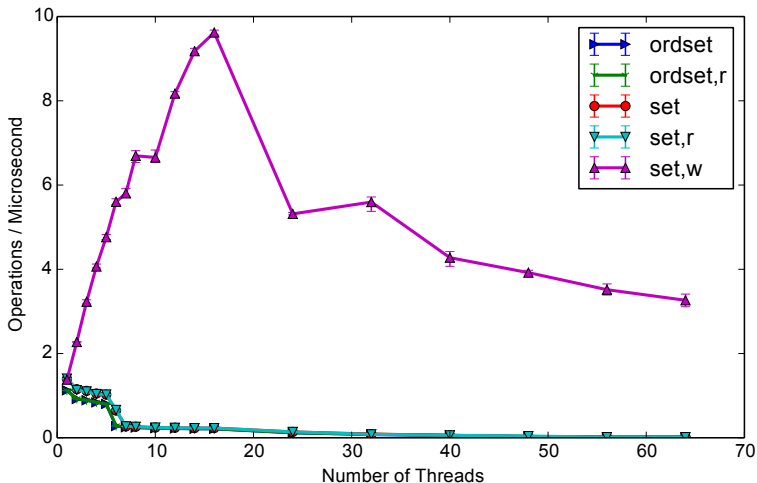
Motivation

CA Trees

Performance

Future Work

Conclusion





100% Updates

Context

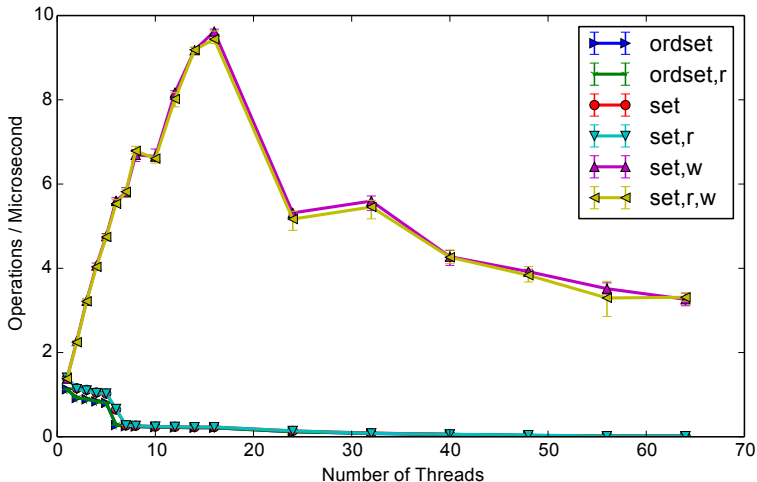
Motivation

CA Trees

Performance

Future Work

Conclusion





100% Lookups

Context

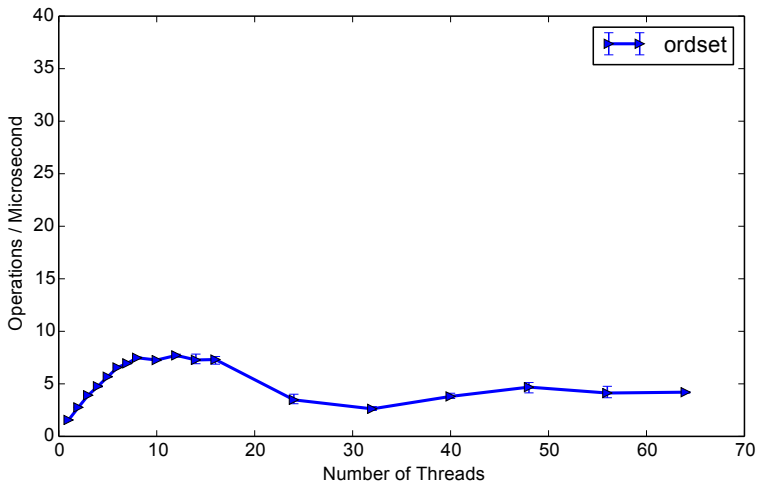
Motivation

CA Trees

Performance

Future Work

Conclusion





100% Lookups

Context

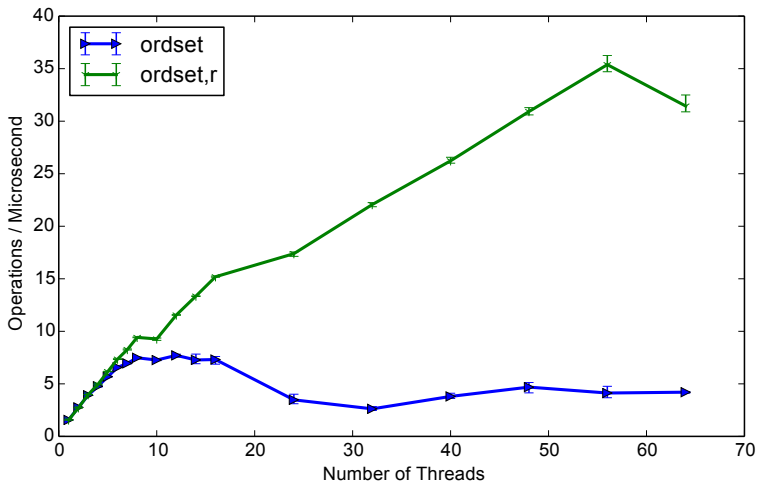
Motivation

CA Trees

Performance

Future Work

Conclusion





100% Lookups

Context

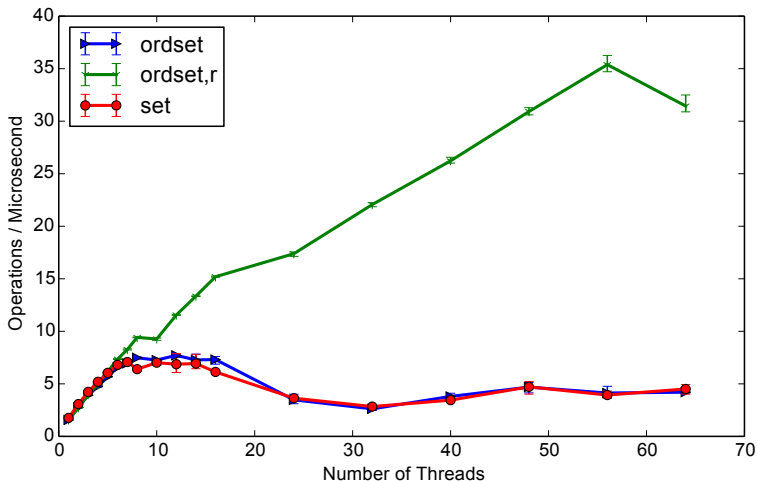
Motivation

CA Trees

Performance

Future Work

Conclusion





100% Lookups

Context

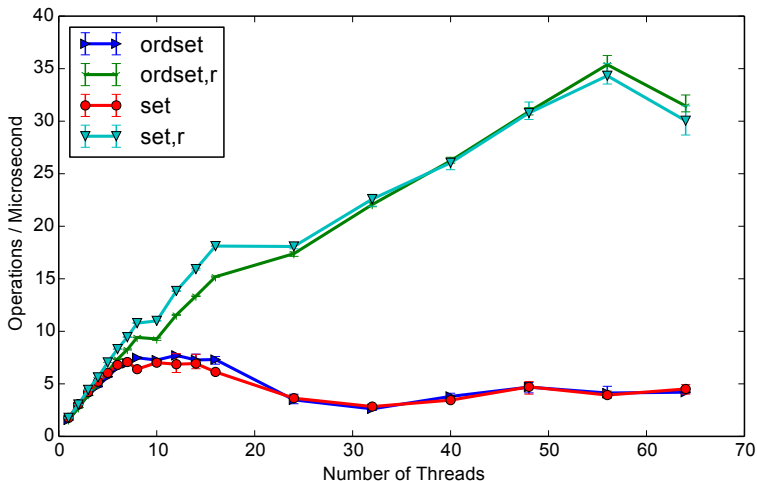
Motivation

CA Trees

Performance

Future Work

Conclusion





100% Lookups

Context

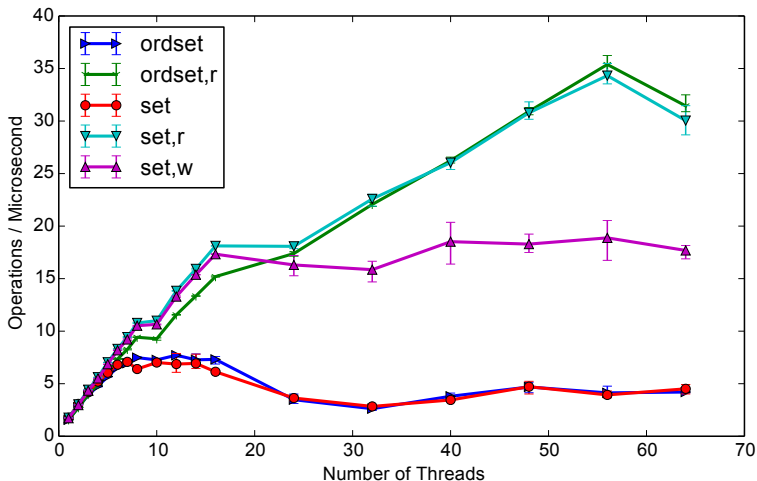
Motivation

CA Trees

Performance

Future Work

Conclusion





100% Lookups

Context

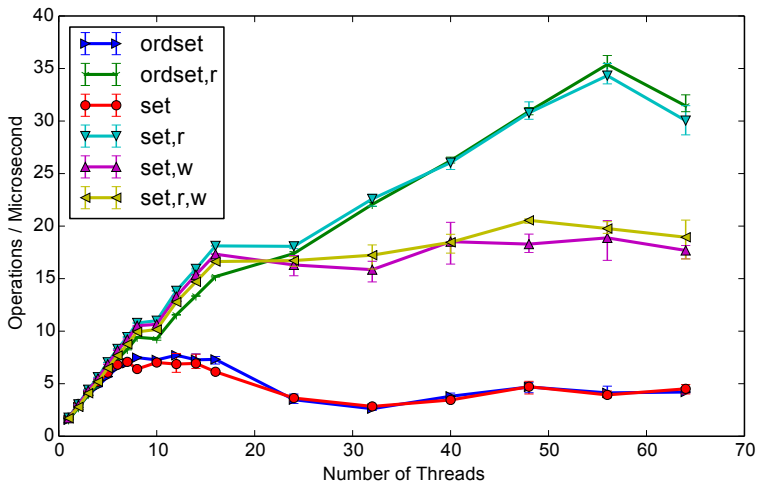
Motivation

CA Trees

Performance

Future Work

Conclusion





50% Lookups, 50% Updates

Context

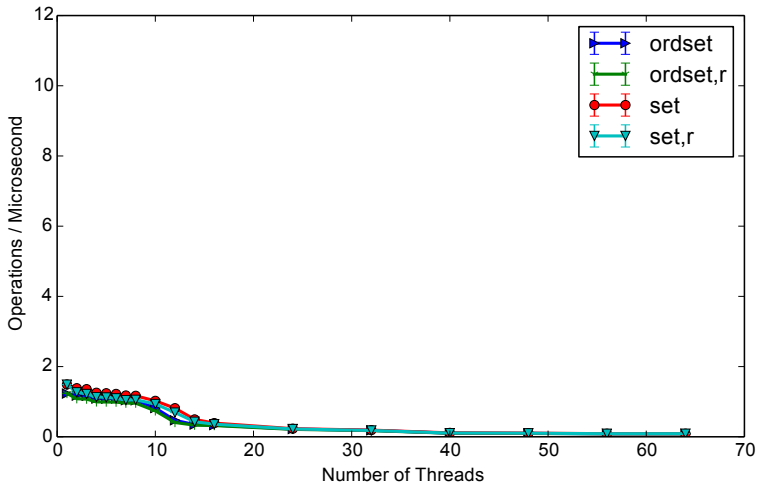
Motivation

CA Trees

Performance

Future Work

Conclusion





50% Lookups, 50% Updates

Context

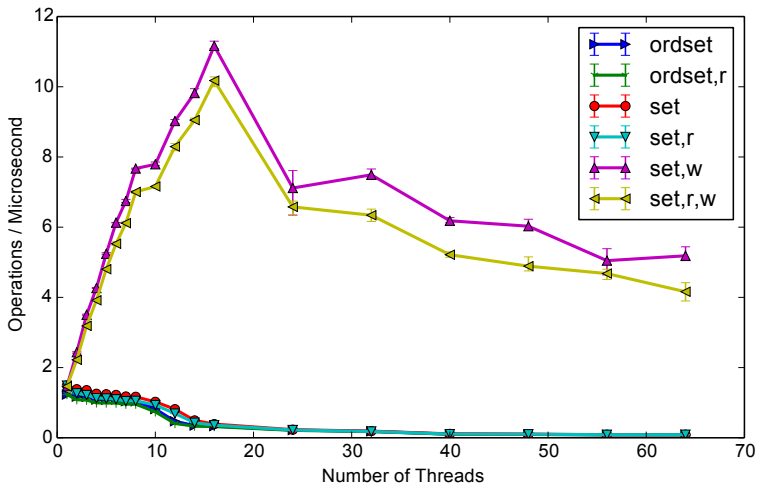
Motivation

CA Trees

Performance

Future Work

Conclusion





99% Lookups, 1% Updates

Context

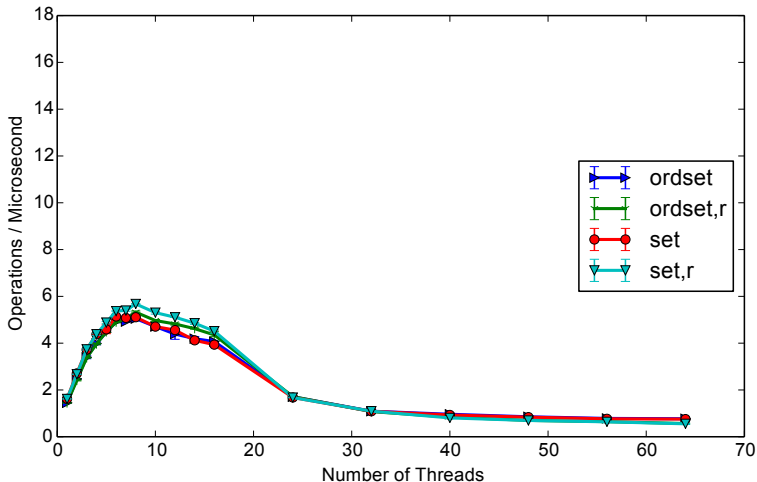
Motivation

CA Trees

Performance

Future Work

Conclusion





99% Lookups, 1% Updates

Context

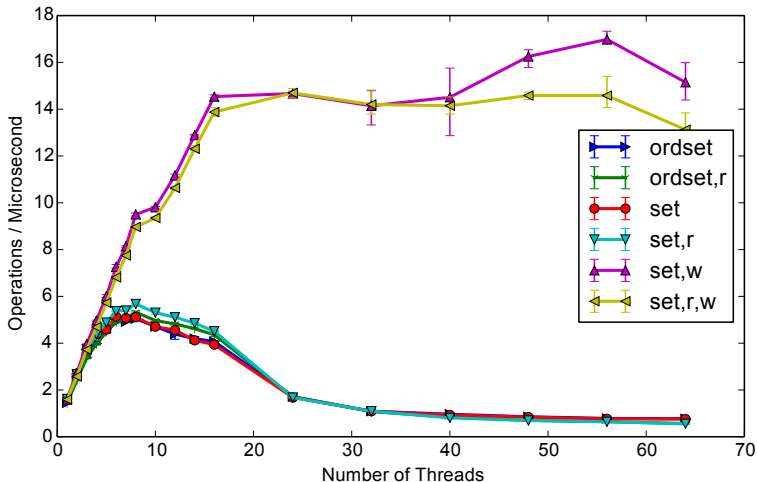
Motivation

CA Trees

Performance

Future Work

Conclusion





Summary of Current ETS Scalability

Context

Motivation

CA Trees

Performance

Future Work

Conclusion

- **Something needs to be done with ordered set to make it scale when there are parallel writes!**
 - Huge slow down even with 99% reads



We Want

Context

Motivation

CA Trees

Performance

Future Work

Conclusion

Wish list for ordered set

- Good scalability
- Reuse code from the current ETS implementation
- Low overhead in sequential case
 - Current algorithms for concurrent ordered sets sacrifice sequential performance and memory consumption for scalability



Contention Adapting Binary Search Trees (CA trees)

Context

Motivation

CA Trees

Performance

Future Work

Conclusion

Key Ideas

- Start with sequential binary search tree protected by a lock
- Collect statistics from the lock
- Adapt the tree according to the statistics



CA Tree Components

Context

Motivation

CA Trees

Performance

Future Work

Conclusion

Statistics Collecting Lock

- The lock has an associated counter C
- $C += 250$ **if** needed to wait to acquire the lock
- $C -= 1$ **if not** needed to wait to acquire the lock
- Adapt when C reach thresholds
 - E.g. -1000 and 1000



CA Tree Components

Sequential Ordered Set Data Structure

- Requires support for split and join
 - The **split** operation splits a tree into two so that all keys in one tree are smaller than the keys in the other
 - The **join** operation merges two trees given that all keys in one tree are smaller than the keys in the other
 - $\mathcal{O}(\log(N))$ implementations for **AVL trees**, **Red-Black trees**, **Treaps**, etc.



CA Tree Structure

Context

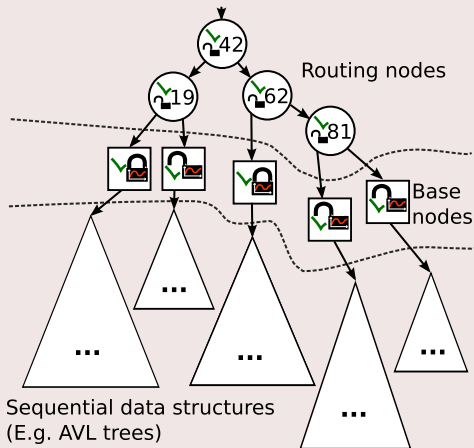
Motivation

CA Trees

Performance

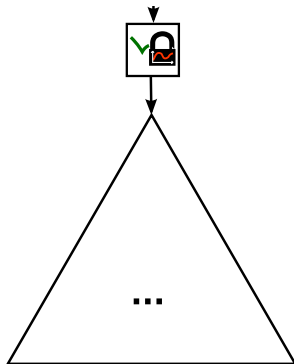
Future Work

Conclusion





CA Tree Animation



Context

Motivation

CA Trees

Performance

Future Work

Conclusion



CA Tree Animation

Context

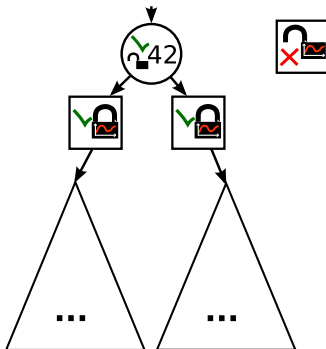
Motivation

CA Trees

Performance

Future Work

Conclusion





CA Tree Animation

Context

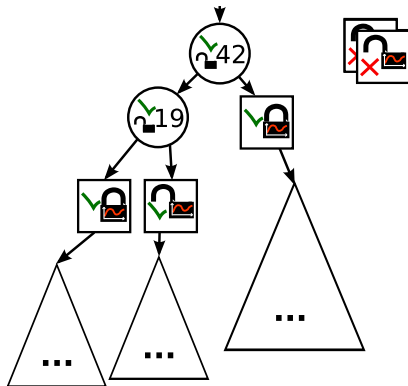
Motivation

CA Trees

Performance

Future Work

Conclusion





CA Tree Animation

Context

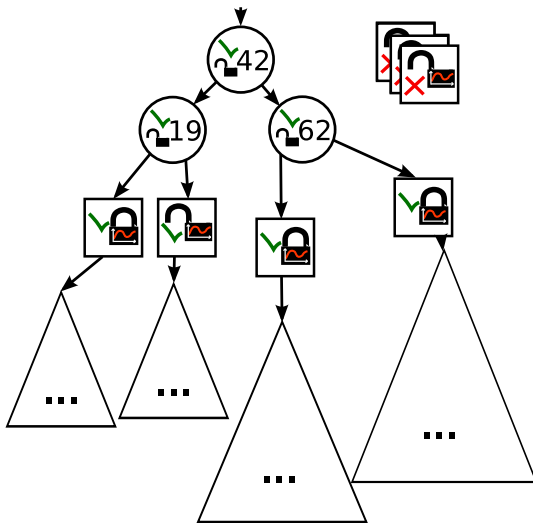
Motivation

CA Trees

Performance

Future Work

Conclusion





CA Tree Animation

Context

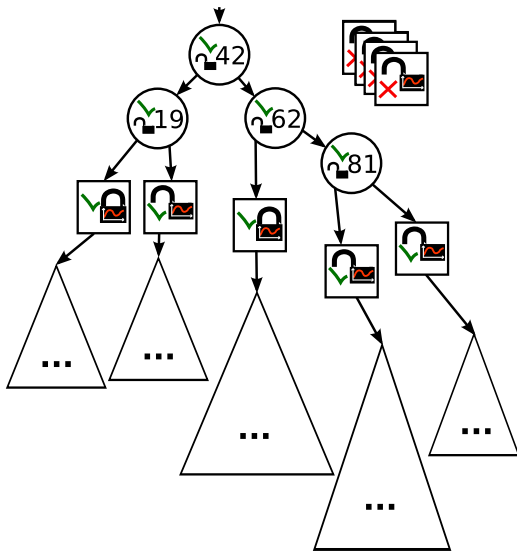
Motivation

CA Trees

Performance

Future Work

Conclusion





CA Tree Animation

Context

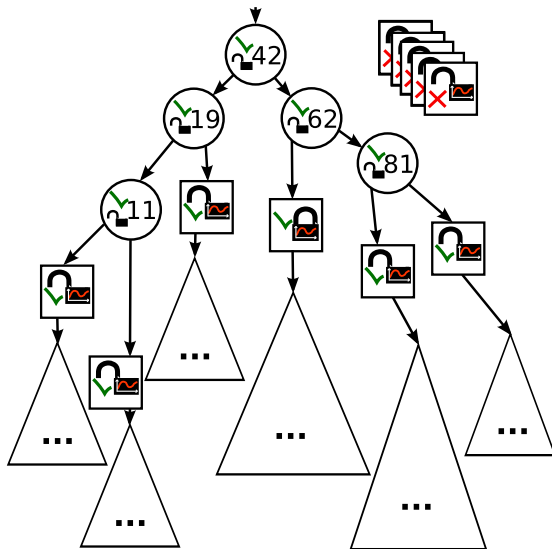
Motivation

CA Trees

Performance

Future Work

Conclusion





CA Tree Animation

Context

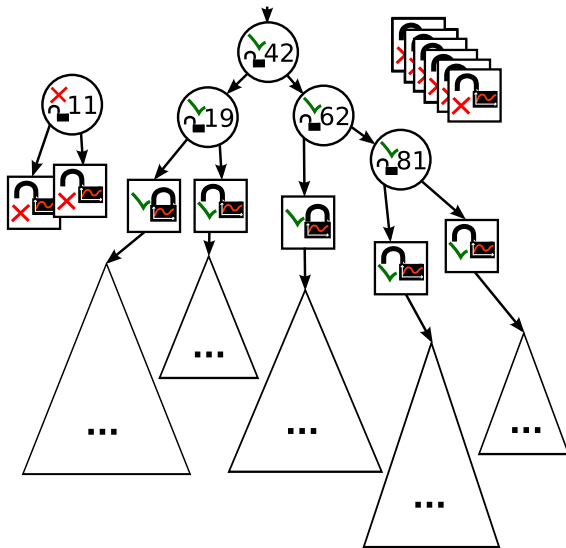
Motivation

CA Trees

Performance

Future Work

Conclusion





CA Tree Animation

Context

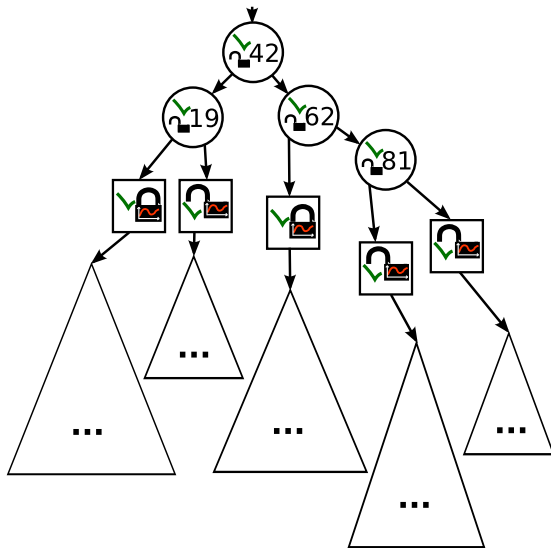
Motivation

CA Trees

Performance

Future Work

Conclusion





Integration into ETS

Context

Motivation

CA Trees

Performance

Future Work

Conclusion

- Reuse code from current `ordered_set`



Integration into ETS

Context

Motivation

CA Trees

Performance

Future Work

Conclusion

- Reuse code from current `ordered_set`
- Routing layer needs special memory management
 - Currently quiescent state based reclamation
 - Better to reuse memory reclamation system for lock free data structures that is already integrated into the Erlang runtime system



Performance of CA Tree Optimized ETS

Context

Motivation

CA Trees

Performance

Future Work

Conclusion

The same benchmark again



100% Updates

Context

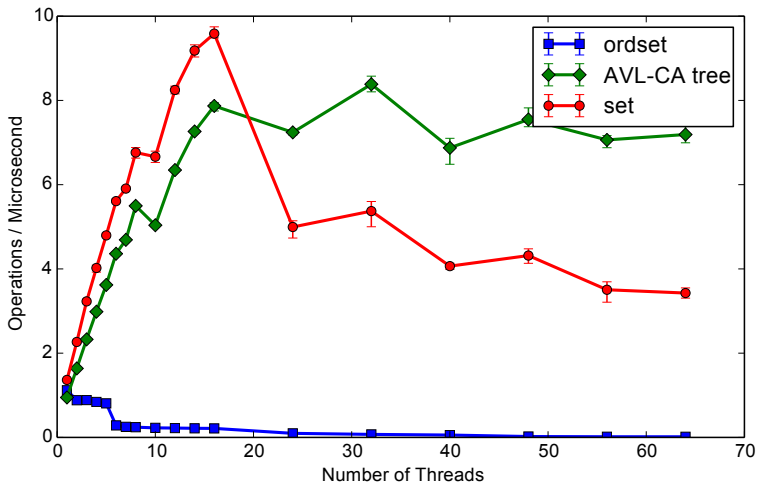
Motivation

CA Trees

Performance

Future Work

Conclusion





100% Lookups

Context

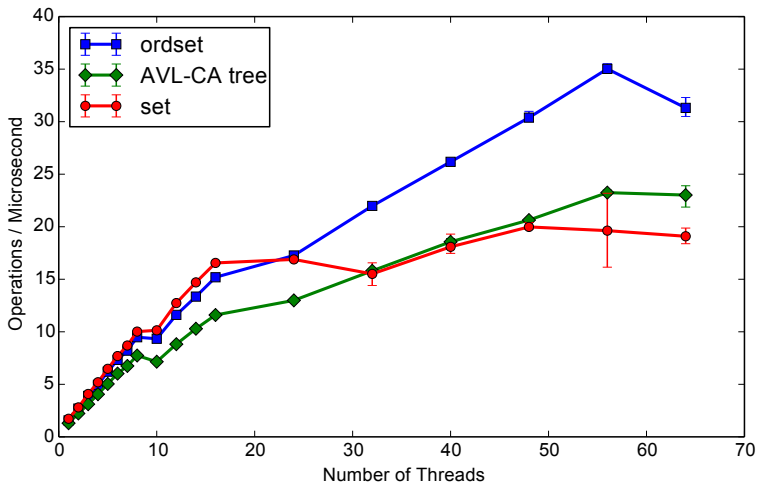
Motivation

CA Trees

Performance

Future Work

Conclusion





50% Update 50% Lookup Scalability

Context

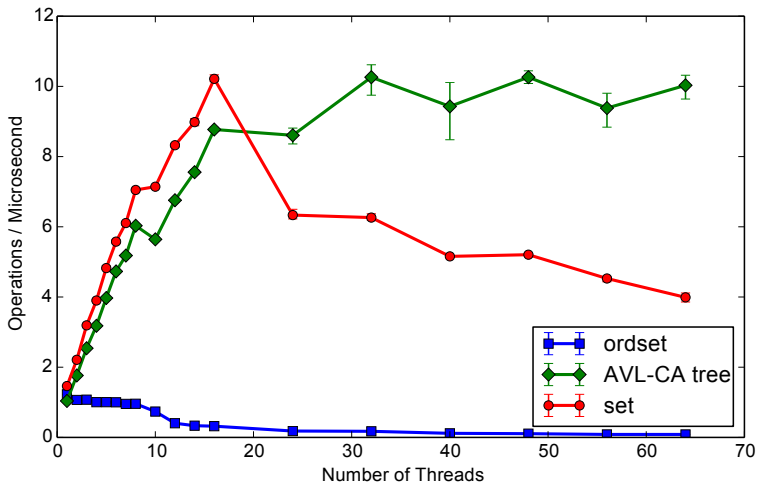
Motivation

CA Trees

Performance

Future Work

Conclusion





20% Update 80% Lookup Scalability

Context

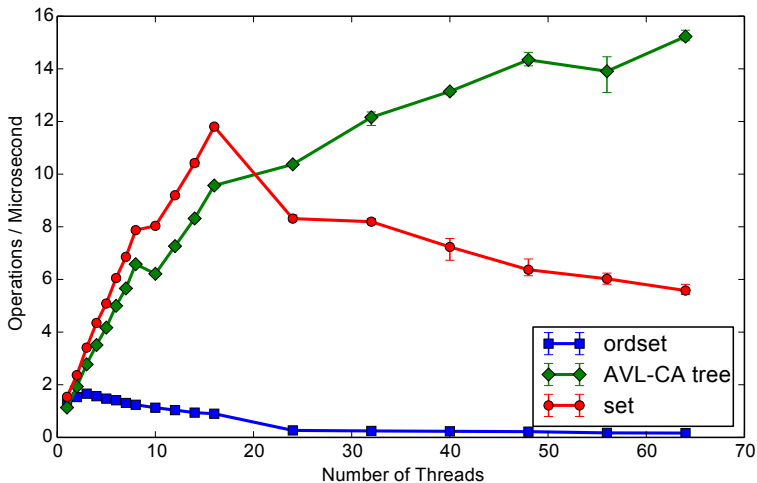
Motivation

CA Trees

Performance

Future Work

Conclusion





1% Update 99% Lookup Scalability

Context

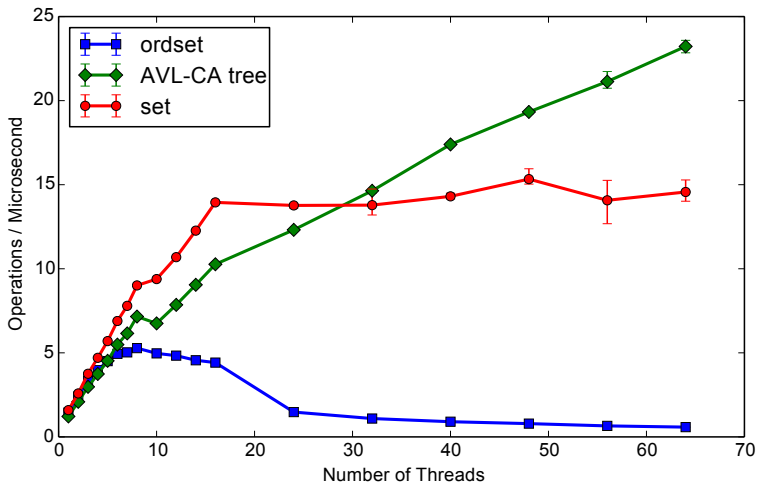
Motivation

CA Trees

Performance

Future Work

Conclusion





Summary of Scalability Improvements

Context

Motivation

CA Trees

Performance

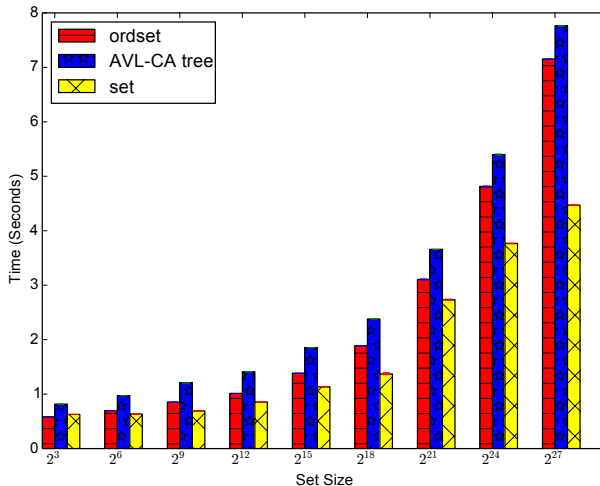
Future Work

Conclusion

- The CA trees does not suffer from large slow down
- Scales reasonably well on one chip
- Update heavy scenarios scale far from perfect on NUMA
 - Centralized statistics counter in memory management
 - set has even more problems on NUMA



Sequential Performance, 80% reads





Future Work

Context

Motivation

CA Trees

Performance

Future Work

Conclusion

To integrate into ETS

- Implement the whole ETS interface
 - Code can be reused from current implementation
- Decide how to integrate it into ETS
 - Always activate on public tables
 - ▶ Read only case might suffer
 - Only activate when `write_concurrency` is specified



More Information

Context

Motivation

CA Trees

Performance

Future Work

Conclusion

Work Already Done

- Compare to other concurrent ordered set data structures
- Investigate optimization for read heavy scenarios
 - RW-locks, Sequence Lock, Hardware Lock Elision (HLE)
- Discuss algorithm in detail

More in a Technical Report

[http://www.it.uu.se/research/group/languages/
software/ca_tree](http://www.it.uu.se/research/group/languages/software/ca_tree)



Concluding Remarks

Context

Motivation

CA Trees

Performance

Future Work

Conclusion

- Performance of concurrent writes on ordered set can be substantially improved in ETS
- CA tree based table scales even better than the current hash based implementation on a NUMA system
- good scalability with low sequential overhead



Questions?