

INTERRUPTING SNAPSHOTS AND THE JAVA SIZE() METHOD

Moran Tzafrir
Tel-Aviv University

Joint work with Yehuda Afek and Nir Shavit

Concurrent Data-Structures

2

- Multicore machines are the mainstream computing platform.
- Inter thread communication is the bottleneck to application performance
- Keeping concurrent data-structures scalable is key to applications' scalability.

Concurrent Data-Structure Libraries

3

- Operations inherited from the sequential API
- Typical global-operations:

ADT: LINKED-LIST

- Insert()**
- Remove()**
- Find()**
- Size()**
- IsEmpty()**
- Clear()**

ADT: QUEUE

- Enq()**
- Deq()**
- Peek()**
- Size()**
- IsEmpty()**
- Clear()**

ADT: RB-TREE

- Insert()**
- Remove()**
- Find()**
- Size()**
- IsEmpty()**
- Clear()**

- Global-operations, show little or no scalability.

The Size() Operation

4

- *Size()* linearizable operation added to any existing data-structure.
- Interface:
 - *Update()* - inc/dec with add/delete of elements.
 - *Size()* - get the number of elements in the data-structure.
- *Update()* more frequent than *Size()* by order of magnitude.

Known Size() Solutions

5

	Scalable	Update Performance	Size Performance
<i>Global-Counter + CAS</i>	X	X	√
<i>Local-Counters + LOCKs</i>	X	f	f
<i>Weaken the spec (JDK)</i>	√	√	√

Interrupting-Snapshots

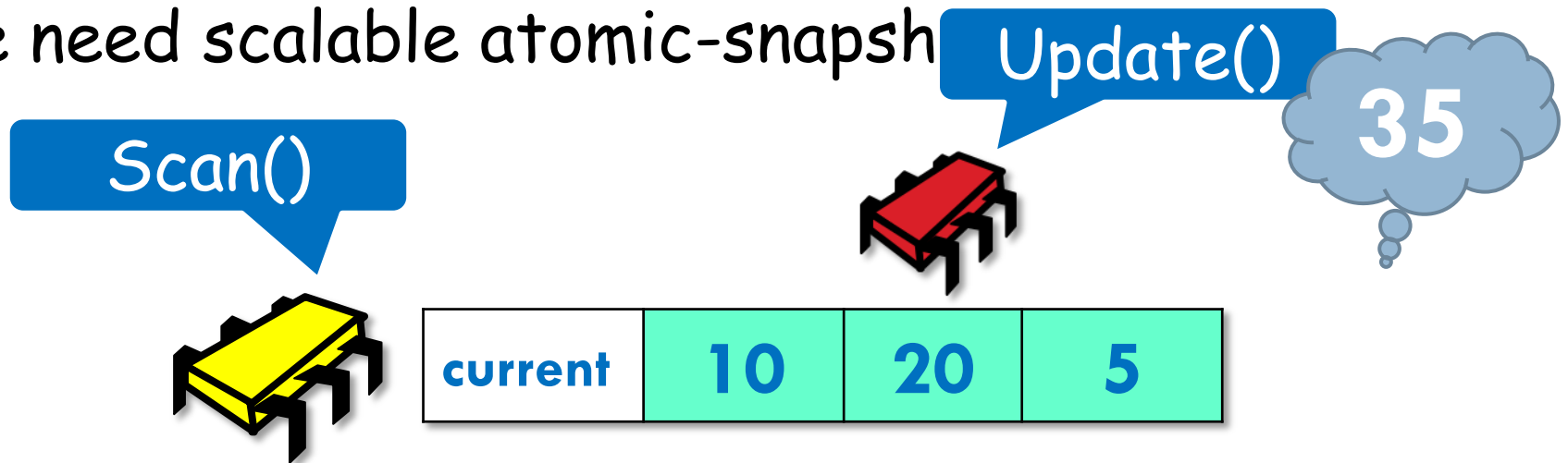
6

- New *Size()* algorithm based on atomic-snapshots.
- Unlike existing solutions:
 - Linearizable** (and not an approximation)
 - Scalable**
 - Wait-free**
- Potentially first example of using atomic-snapshots in industrial code.

Local Counters + Snapshots

7

- **Local Counters:** one per thread.
- *Size()* = $\sum_{i=1,n}$ (snapshot of local-counters)
- We need scalable atomic-snapsh



Single Scanner Algorithm

[Kirousis, Spirakis, T: 95]

(curr.seq == thread seq) →
update current value

95]

Read scan-seq (#6)



scan -seq	#6
-----------	----

current	10	20	5
curr.seq	#6	#6	#6
previous	8	4	0

Update() - If (curr.seq == thread seq), update current value. Else copy current to previous with new seq before updating current.

Single Scanner Algorithm

10

Update(22)
(Read scan-seq #6)



scan -seq	#6
-----------	----

current	10	22	5
curr.seq	#6	#6	#6
previous	8	4	0

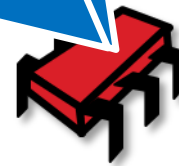
Update() - If (`curr.seq == my seq`), update current value. Else copy current to previous with new seq before updating current.

Single Scanner Algorithm

11

Due to *Scan()*,
scan-seq == #7

(curr.seq < my seq) →
copy



scan -seq	#7
-----------	----

current	10	23	5
curr.seq	#6	#7	#6
previous	8	42	0

Update() - If (curr.seq == my seq), update current value. Else copy current to previous with new seq before updating current.

Single Scanner Algorithm

Increment scan-seq,
and read (#8)

(curr.s

(curr.seq == my seq) →
collect previous(5)

(4)
#8)



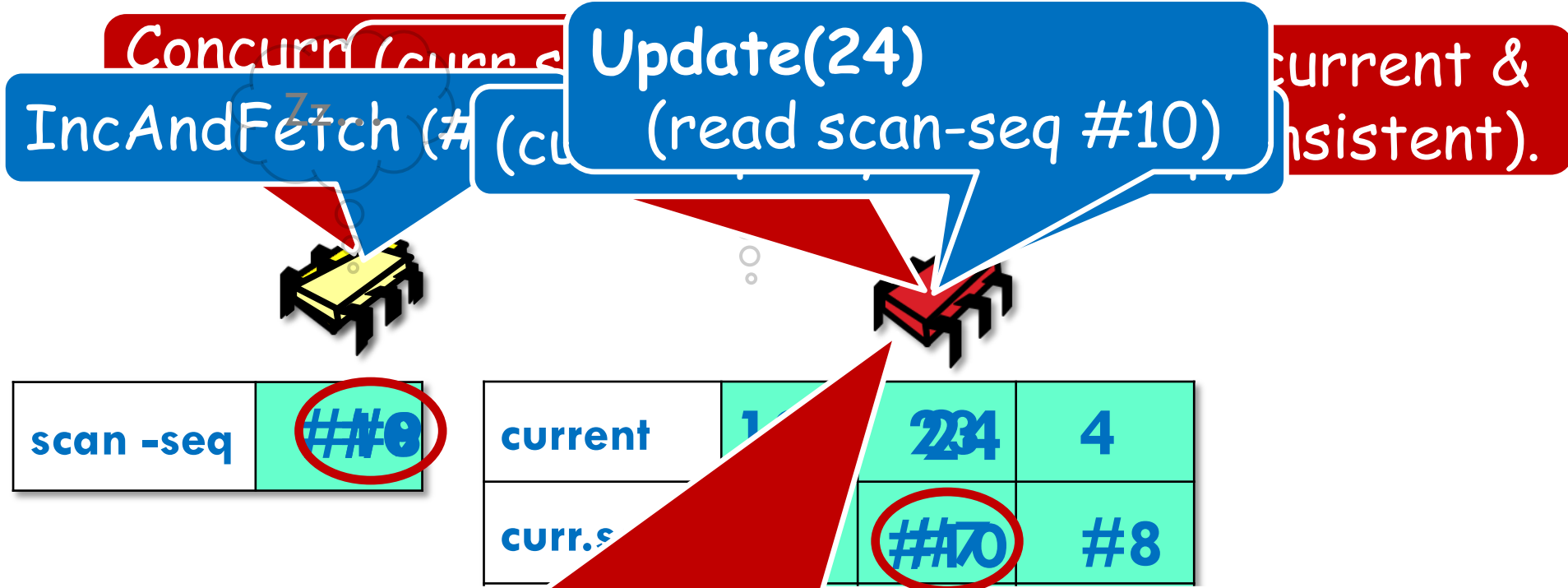
scan -seq #8

current	10	23	4
curr.seq	#6	#7	#8
previous	8	22	5

Scan() - Increment scan-seq (a write). Go over the local-counters, and collect values with $\text{curr.seq} < \text{my seq}$

Multi Scanner - Lock Free

13



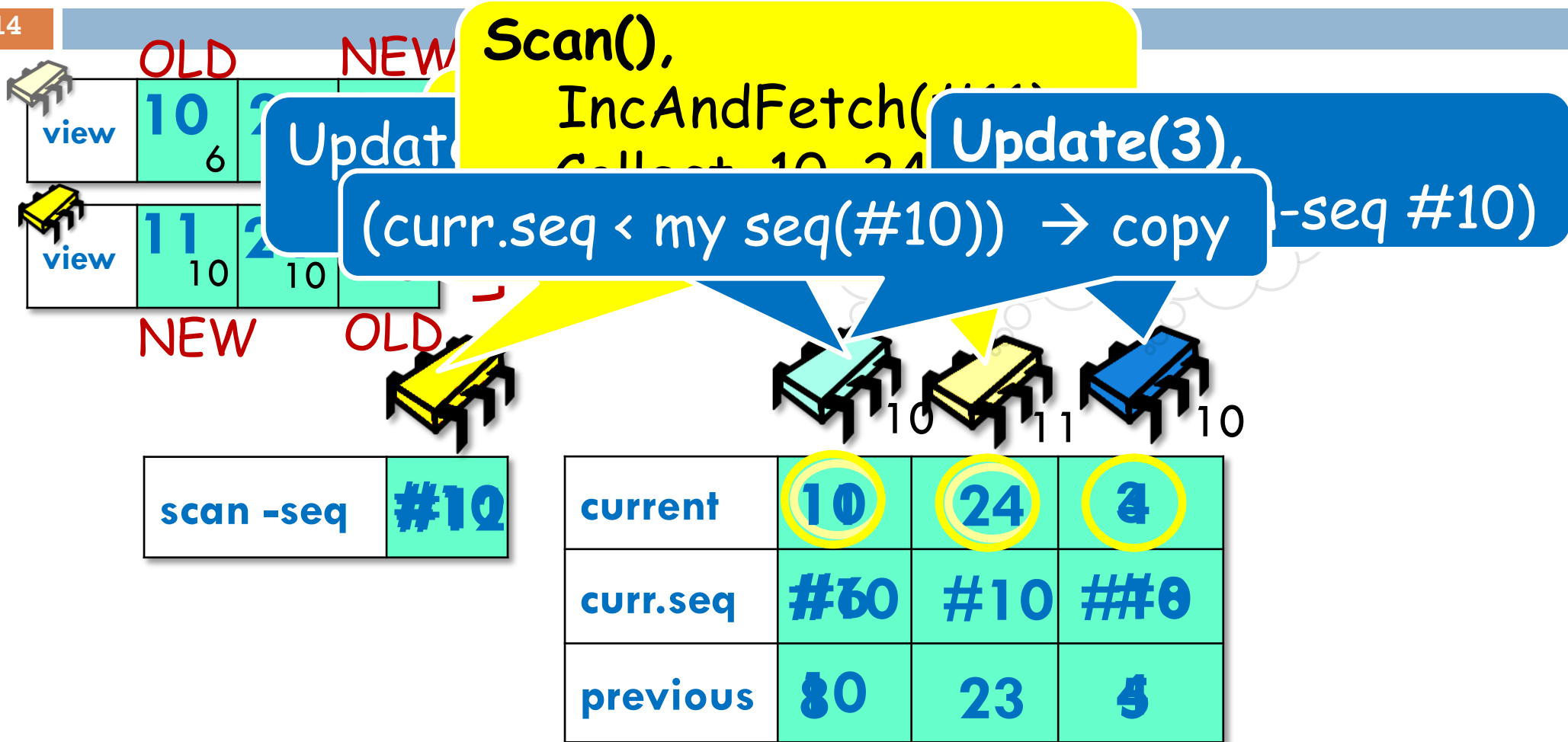
Solution: Do not increment scan-seq, adopt interrupting sequence (#10).
→ (Lock Free)

Scan() - Increment scan-seq using **CAS**. Collect values with `curr.seq < my seq`.

If (`curr.seq > my seq`) abort and restart (without incrementing the scan-seq).



Counter Example - Multi Scanner

14



Multi Scanner - Wait Free

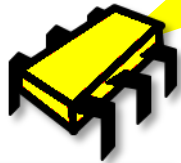
15

 view	10 6	24 10	3 10
 view	11 10	24 10	4 8

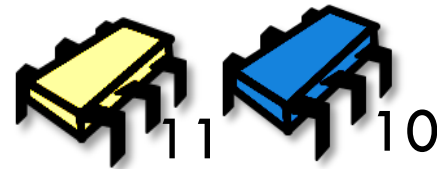
Scan(),
IncAnd
Collect 11, 24, 4

Try to CAS on last-view

CAS on last-view.
From [#9, 37] to [#12, 39]



scan -seq	#12
-----------	-----



current	11	24	4
curr.seq	#10	#10	#80
pro			

Last-view	##192	339
-----------	------------------	----------------

BORROW last-view result,
since happened in our interval.
→ wait-free

Complexity

16

- *Update()* - Wait-free $O(1)$ read/write operations
- *Scan()* - Wait-free Amortized $O(1)$ CAS operations

Performance Evaluation

17

Compare

- New Lock-free Size implementation(**NLF**)
- New wait-free (**NWF**)
- Java Concurrency Package (**JDK**)
- Using [RianyShavitTouitou] Snapshot (**RST**)

On

- Sun UltraSPARC T2 (64x)
- Azul Vega2 7200 series (~200x)

Performance Evaluation

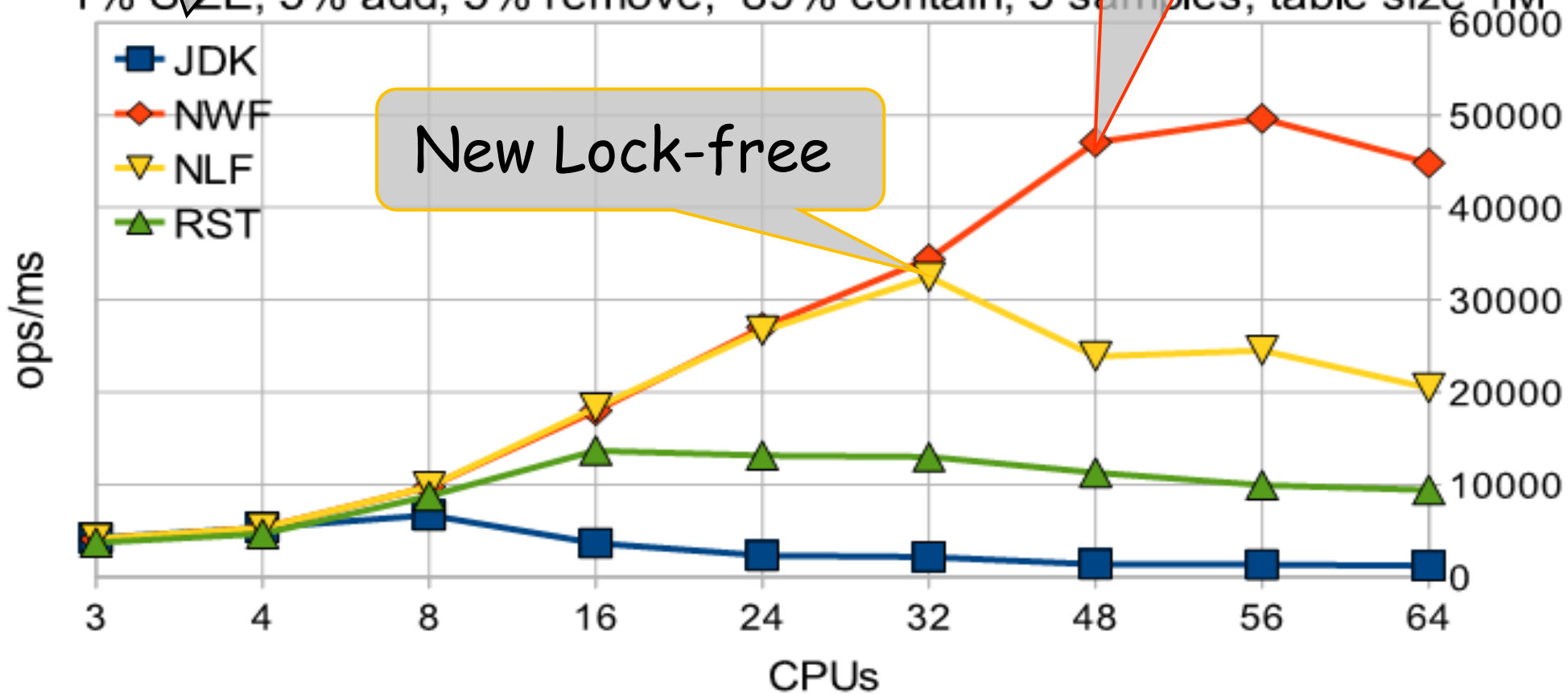
18

1% *Size()*

New Wait-free

SPARC T2 - Throughput as function CPUs

1% SIZE; 5% add; 5% remove; 89% contain; 5 samples; table size 1M



New Lock-free

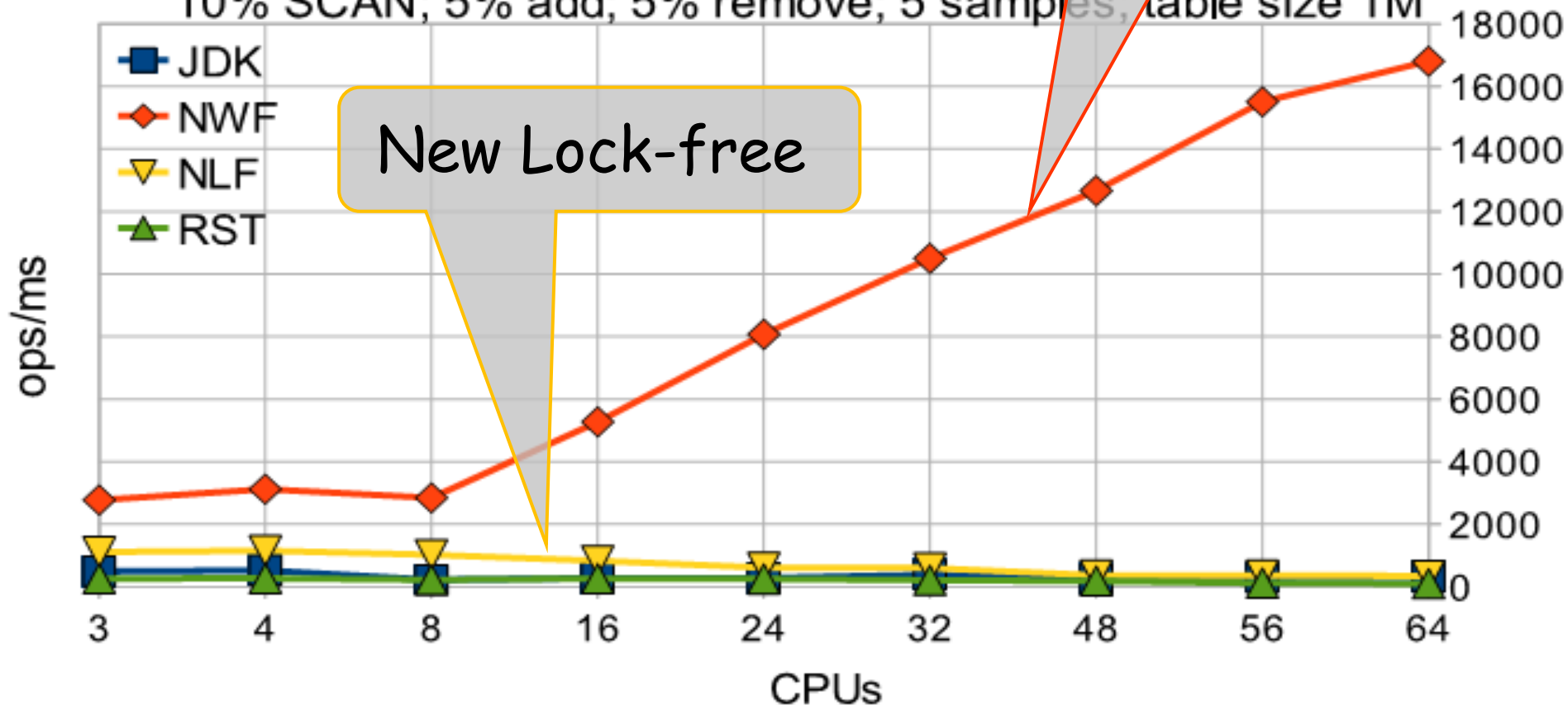
Performance Evaluation

10% *Size()*

New Wait-free

SPARC T2 - SCAN Throughput as function CPUs

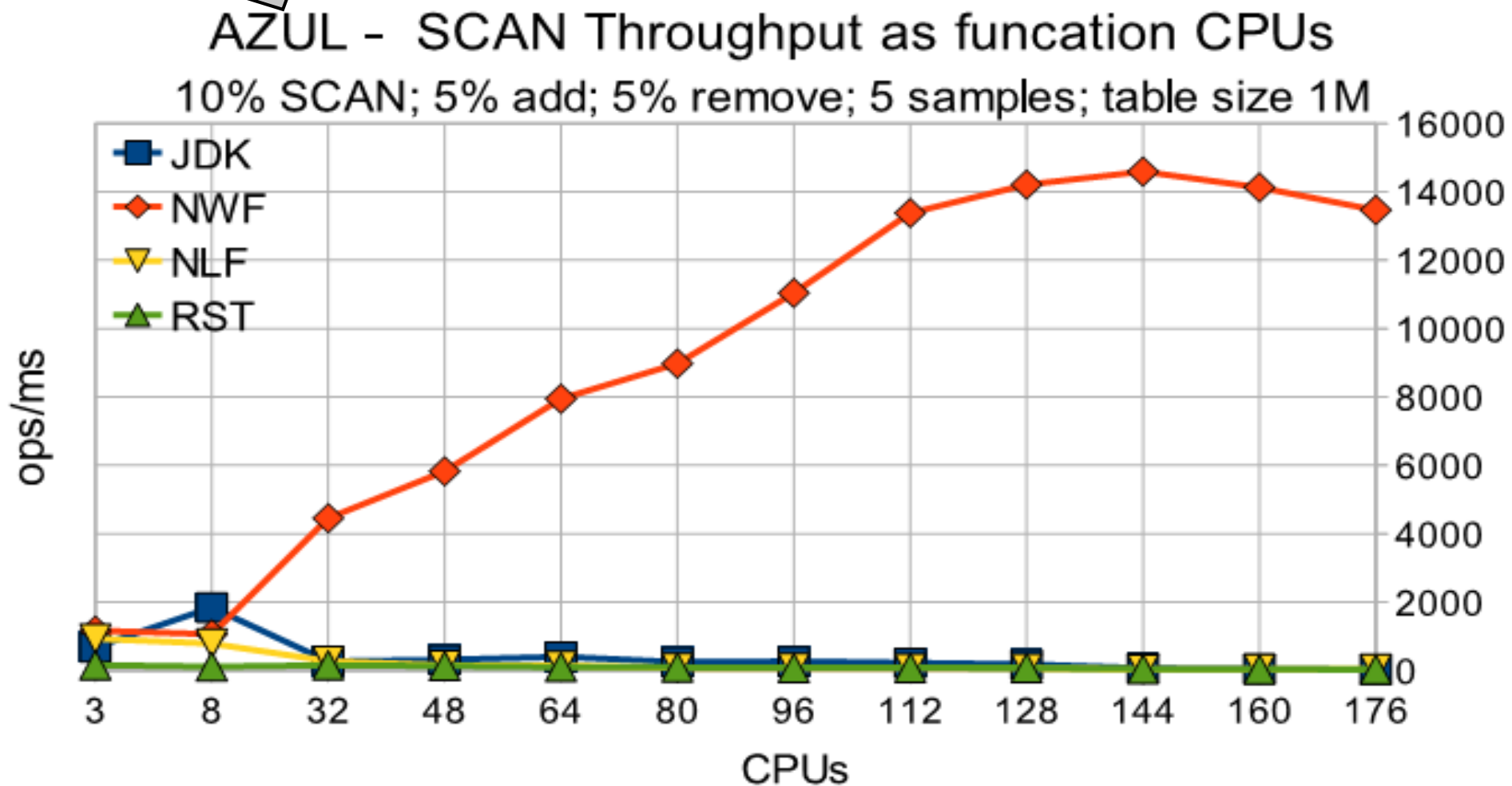
10% SCAN; 5% add; 5% remove; 5 samples; table size 1M



New Lock-free

Performance Evaluation

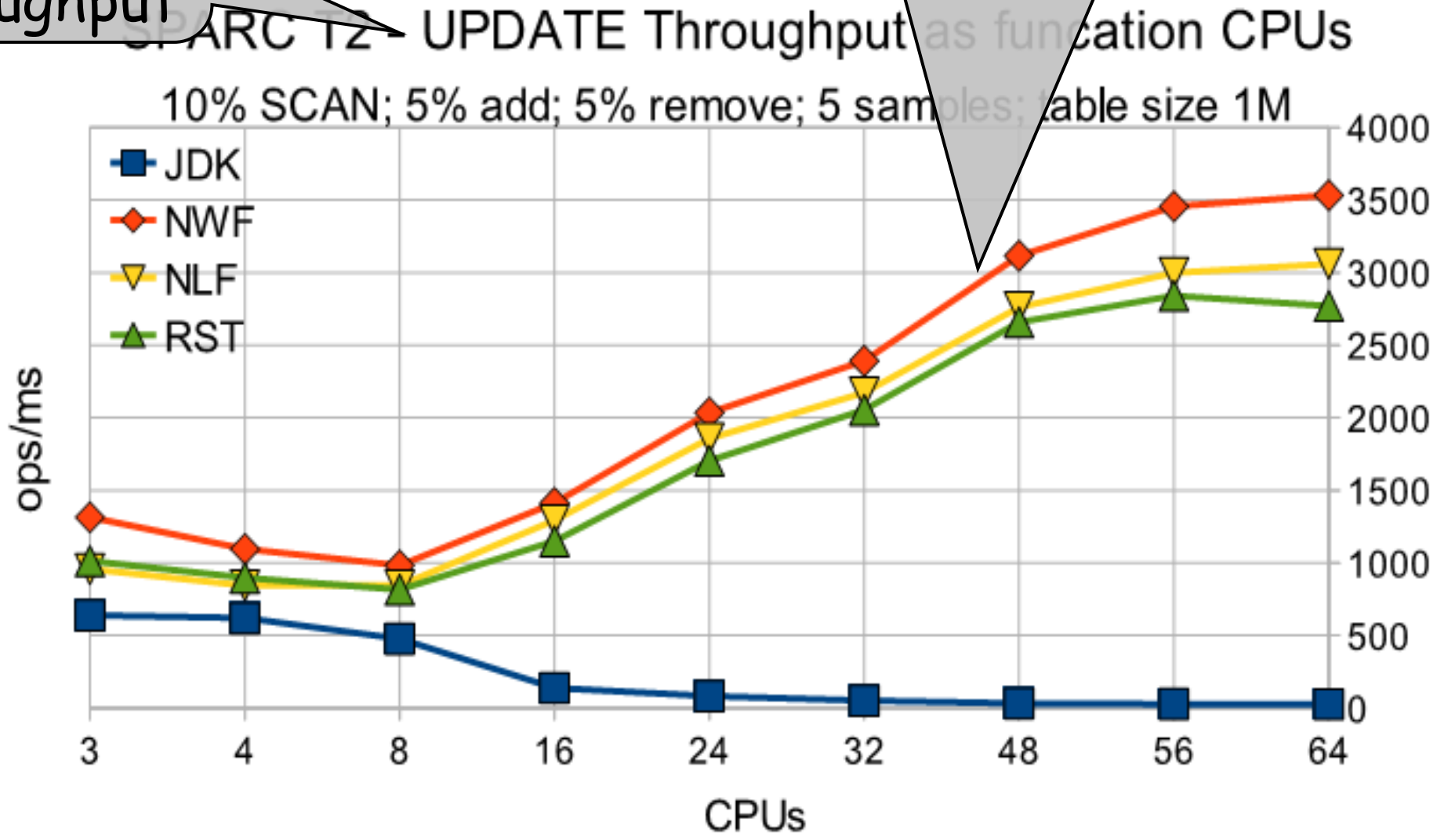
Same behavior
on Azul



Performance Evaluation

measure *Update()* throughput

Size() is responsible for difference in throughput



Conclusion

22

- Global operations *Size()*, *IsEmpty()* can be both accurate and scalable.
- Fundamental **atomic-snapshot** theory problem finds practical application

Future Work

23

Future Work:

- Uses for atomic-snapshots in TM, concurrent-applications, ...
- Extending the Interruption Snapshots algorithm to a view that is more than a single word.

Thanks for your time

Questions