

Read-Copy-Update for HelenOS

<http://d3s.mff.cuni.cz>



Martin Děcký

decky@d3s.mff.cuni.cz



CHARLES UNIVERSITY IN PRAGUE
faculty of mathematics and physics



HelenOS

- **HelenOS**

- Microkernel multiserer operating system
- Relying on asynchronous IPC mechanism
 - Major motivation for scalable concurrent algorithms and data structures

- **Martin Děcký**

- Researcher in computer science (operating systems)
- Not an expert on concurrent algorithms
 - But very lucky to be able to cooperate with hugely talented people in this area



Asynchronous IPC

=

Communicating parties may access the communication facilities concurrently

Asynchronous IPC

=

Communicating parties may access the communication facilities concurrently

→ The state of the shared communication facilities needs to be protected by explicit synchronization means



Asynchronous IPC

=

Communicating parties have to access the communication facilities concurrently



Asynchronous IPC

=

Communicating parties have to access the communication facilities concurrently

← In order to counterweight the overhead of the communication by doing other useful work while waiting for a reply

Asynchronous IPC



**Communication facilities have to use
concurrency-friendly (non-blocking)
synchronization means**

Asynchronous IPC

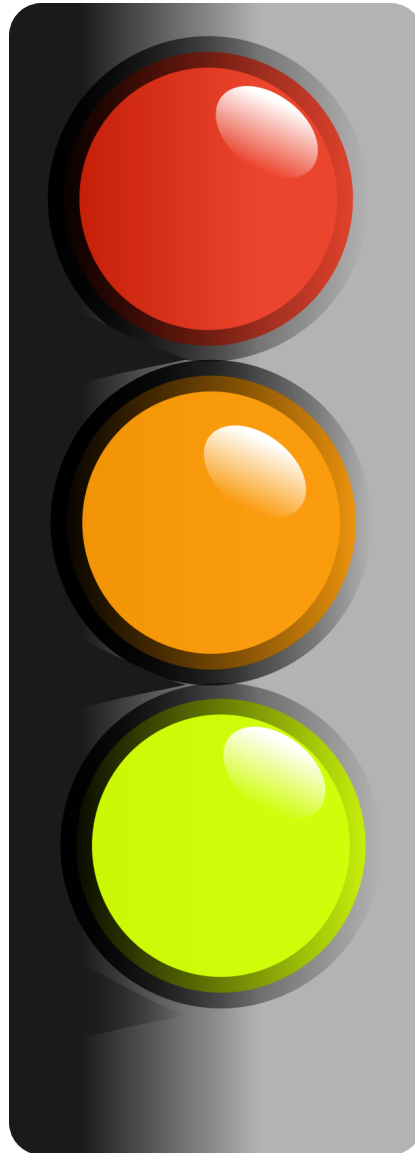


**Communication facilities have to use
concurrency-friendly (non-blocking)
synchronization means**

← In order to avoid limiting the achievable degree of concurrency

- For accessing shared data structures
- Mutual exclusion synchronization
 - Temporal separation of scheduling entities
 - Typical means
 - *Disabling preemption, Dekker's algorithm, direct use of atomic test-and-set operations, etc.*
 - Typical mechanisms
 - *Locks, semaphores, condition variables, etc.*
 - [+] Relatively intuitive semantics, well-known characteristics
 - [-] Overhead, restriction of concurrency, deadlocks

Mutual Exclusion Synchronization



- **Non-blocking synchronization**

- Replace temporal separation by sophisticated means that guarantee logical consistency
- Typical means
 - *Atomic writes, direct use of atomic read-modify-write operations, etc.*
- Typical mechanisms
 - *Transactional memory, hazard pointers, Read-Copy-Update, etc.*
- **[+]** Reasonable (almost no) overhead and restriction of concurrency in favorable cases, guarantee of progress
- **[-]** Less intuitive semantics, sometimes non-trivial characteristics, non-favorable cases, livelocks

Non-blocking Synchronization



HelenOS

16 Joseph Smith in & with Reynolds & Co. 1889				14 Dr Mrs Harry Vischer & Mrs Ella M Wendell 1905			
To Balance		\$7,717.44		July 19 A. S. Van Houton		July 5 Rent store	\$3.87
Feb 1 2600 th Middle	\$17.00	227.50		" Feb 48.989 Albee		31 " Office	17.5
" 20 993 rd Bush, Ontonagon, 338		327.94		" Ante G. Lauf 9.95	\$4.98	" " "	17.5
Mar " Freight on Bage		2.00		Aug 15 H. W. Pot 1072.5471	13.63	5 " Store	3.87
Sept 21 Disc. on D. Wiley Note		7.80		" Royal 38 E. Main	8.18	30 " Office	17.5
" 27 500 th Bush, Ontonagon, 444 1/2		222.00		" Pot 96.29.19 Mart		Oct 3 " Store	3.87
" " 318 th Church St. 10	\$12.20			" 38 E. Main	1.90	Nov 1 " Store	3.87
" " 5 th Lab. Chapel	7.00	38.00		Sept 7 Locity School Tax	71.28		
" " 10 th Bay Street	20	1.00		" 18 Sunday & School	58		
" " Old Mill, 9 th	1.60	1.02		" 23 of P. Frederici	74		
" " 80 th Church St. 10	21.00	69.36		Nov 4 Dft	123.89		
" " 40 th Bay Street	15.00	37.00			225		225.00
" " V. P. Barn	16.00	82.00					
" " 140 th Bage	24.00	273.58		1906		1905	
" " 780 th E. Middle	\$19.00	\$74.10		Feb 13 State & County Tax	\$16.21	Aug 3 Rent store	\$3.87
" " 400 th American	20.00	46.00		" 26 N. B. Becker	13.37	Dec 1 " Office	17.5
" " 800 th Bay Street	19.00	76.00		Mar 5 Dft	19.54	Jan 2 " Store	3.87
" " 5 th Bay Street	10 th	6.00				Jan 1 " "	3.87
" " 120 th Bage	30.00	226.70				" " Office	17.5
					225	27 " "	17.5
						Feb 19 " Store	3.87
						Mar 2 " Office	17.5
							225
				1906		1906	
				Mar 23 of P. Frederici	20	Mar 7 Rent store	\$3.87
				Apr 16 Water Tax	8.25	Apr 2 " Office	17.5
				June 13 A. S. Van Houton		May 1 " "	17.5
				" Pot 82.90 38 E. M	8.18	3 " Store	3.87
				July 4 " 4.56 247 81.5		Apr 14 " "	3.87
				" 38 E. Main St	4.20	June 1 " "	3.87
				10 May 29 Frederici	23.88	" " Office	17.5
				19		198 02	
				" Lockwood Note 200.00		31 37	
				" Lis. Bus	2.00	679 37	
				Sept 1 " Leach			
1896							
Apr 15	By check # 505	200.00					
May 22	" " " 589	150.00					
June 13	" " (A. C. Townsend)	90.00					
	" Cash	3.00					
19	" Lockwood Note 200.00	100.00					

● Wait-freedom

- Guaranteed system-wide progress and starvation-freedom (all operations are finitely bounded)
- Wait-freedom algorithms always exist [1], but the performance of general methods is usually inferior to blocking algorithms
- Wait-free queue by Kogan & Petrank [2]

● Lock-freedom

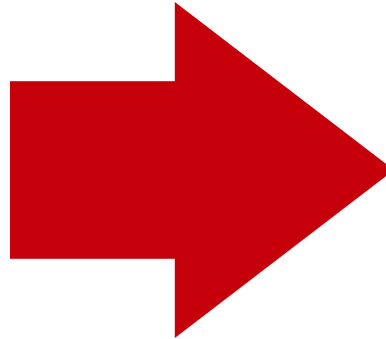
- Guaranteed system-wide progress, but individual threads can starve
- Four phases: Data operation, assisting obstruction, aborting obstruction, waiting

● Obstruction-freedom

- Guaranteed single thread progress if isolated for a bounded time (obstructing threads need to be suspended)

**Synchronization
means**

Individual instance of usage



**Synchronization
mechanism**

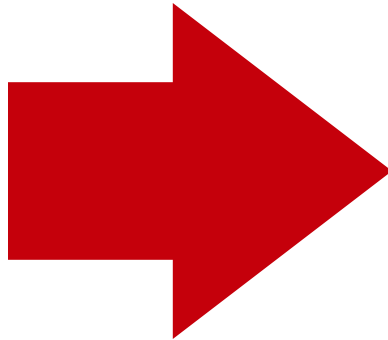
Generic reusable pattern



Synchronization means

Individual instance of usage

E.g. non-blocking list
implementation using
atomic pointer writes



Synchronization mechanism

Generic reusable pattern

E.g. non-blocking list
implementation using
Read-Copy-Update

- **Non-blocking synchronization mechanism**

- Targeting synchronization of read-mostly pointer-based data structures with immutable values
 - Favorable case: R/W ratio of $\sim 10:1$ (but even 1:1 is achievable)
 - Unlimited number of readers without blocking (not waiting for other readers or writers)
 - Little overhead on the reader side (smaller than taking an uncontended lock)
 - Readers have to tolerate “stale” data and late updates
 - Readers have to observe “safe” access patterns
 - Synchronization among writers out of scope of the mechanism
 - Optional provisions for asynchronous reclamation

- **Read-side critical section**

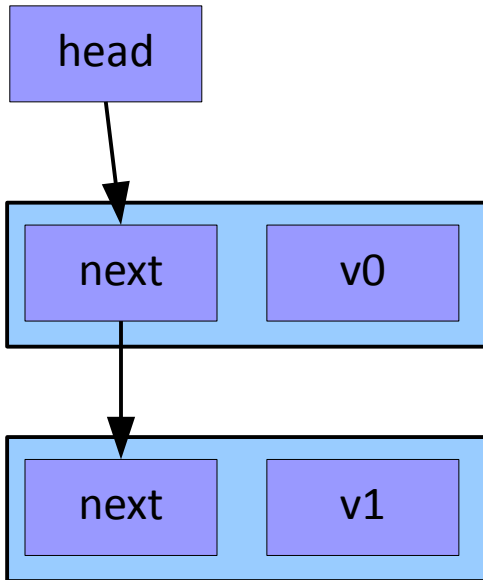
- Delimited by `read_lock()` and `read_unlock()` operations (non-blocking)
 - Protected data can be referenced only inside the critical section
- Safe **access()** methods for reading pointers
 - Avoiding unsafe compiler optimizations (reloading the pointer)
 - Not necessary for reading values
- **Quiescent state** (a thread outside a critical section)
- **Grace period** (all threads pass through a quiescent state)

- **Synchronous write-side update**
 - Atomically unlinking an old element
 - Calling a **synchronize()** operation
 - Blocks until a grace period elapses (all readers pass a quiescent state, no longer referencing the unlinked data)
 - Possibility to reclaim or free the unlinked data
 - Inserting a new element using safe **assign()** operation
 - Avoiding unsafe compiler optimizations and store reordering on weakly ordered architectures

Synchronous Update Example



I.

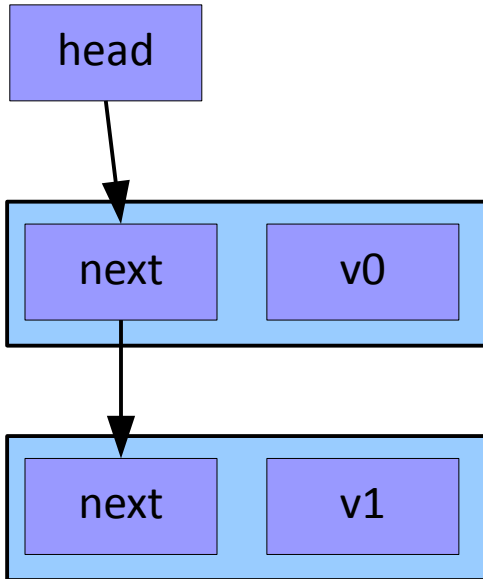


Atomic pointer update to remove the element with v0 from the list

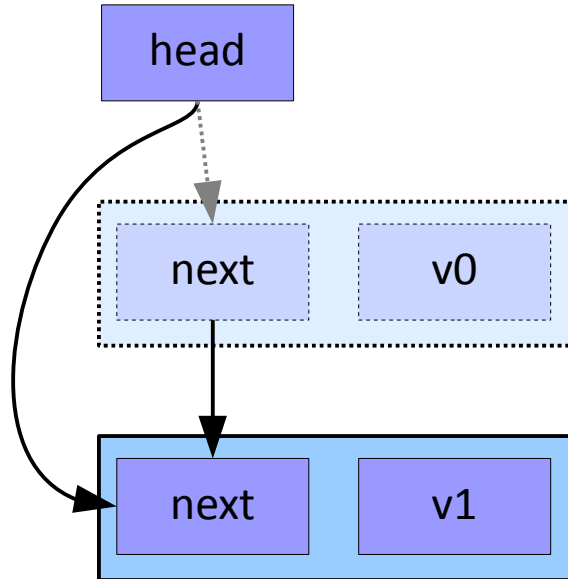
Synchronous Update Example



I.



II.



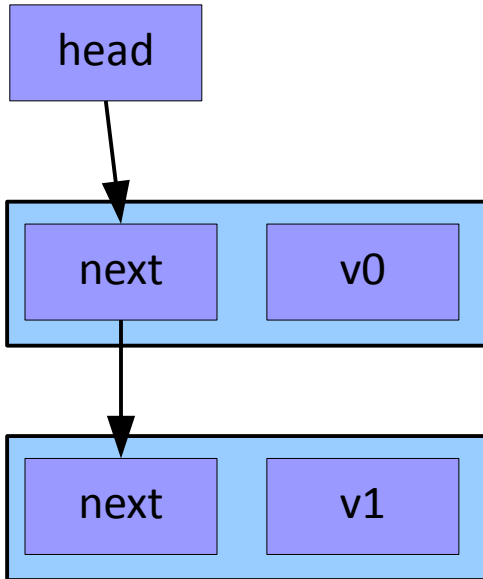
Blocking on **synchronize()**

During the grace period preexisting readers can still access the “stale” element with v0

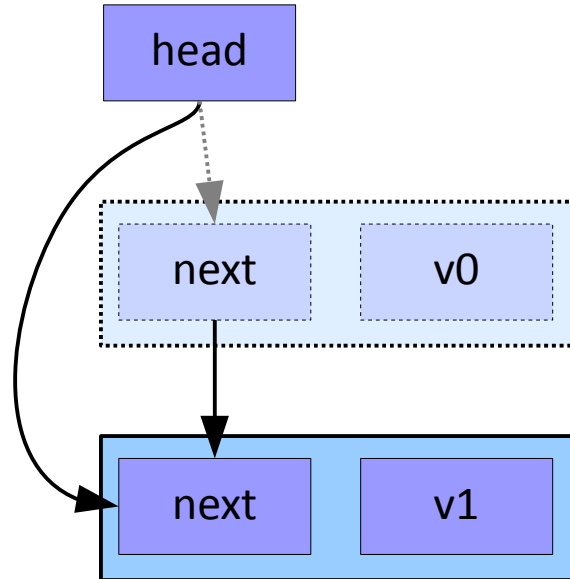
Synchronous Update Example



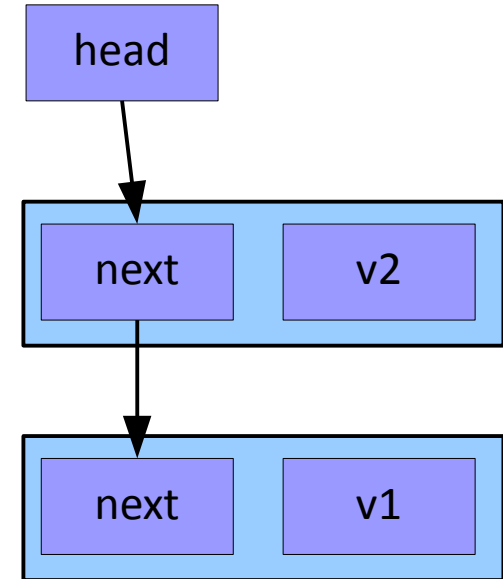
I.



II.



III.

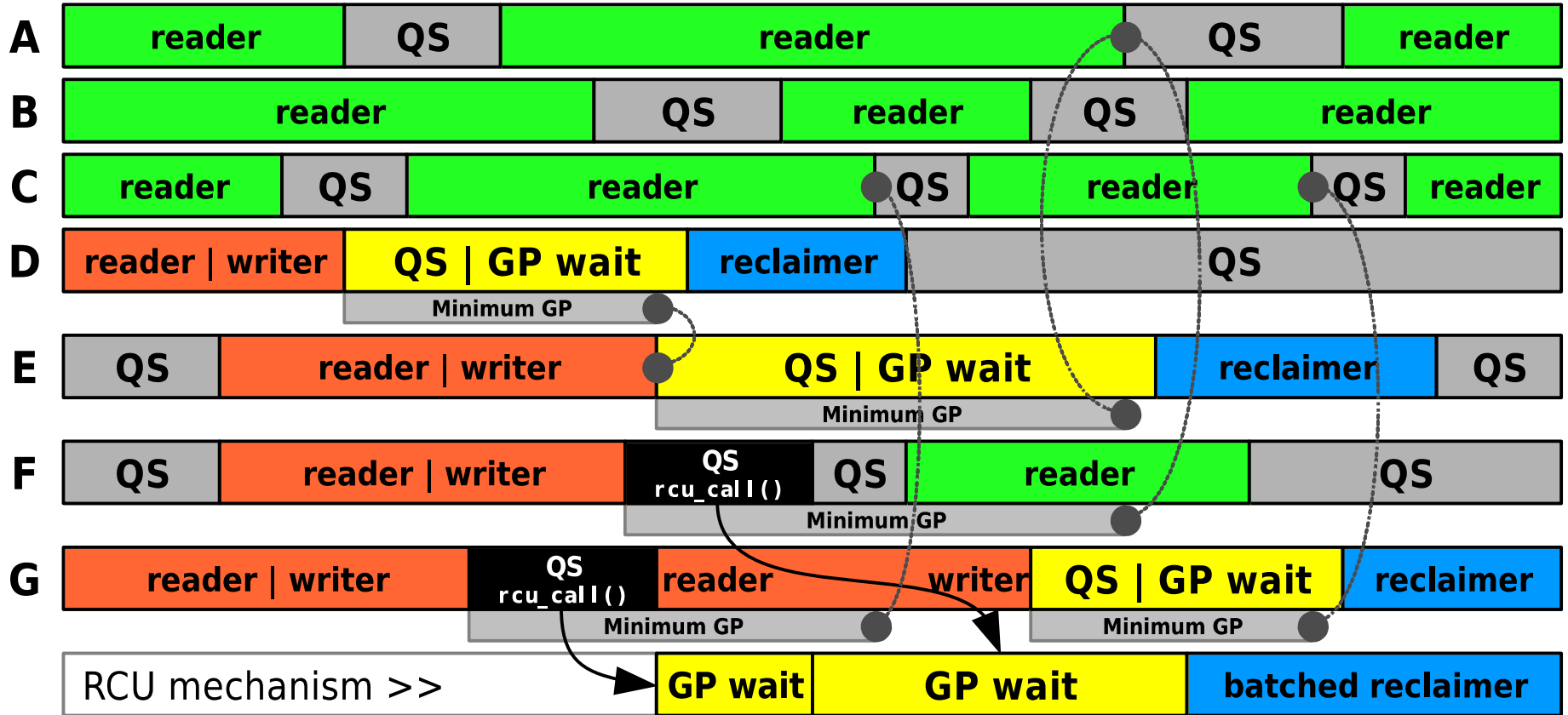


No reader can reference the element with v0 anymore – it can be reclaimed
New element with v2 can be atomically inserted

- **Asynchronous write-side update**
 - Using a **call()** operation
 - Non-blocking operation registering a callback
 - Callback is executed after a grace period elapses
 - Using a **barrier()** operation
 - Waiting for all queued asynchronous callbacks

- **Cornerstone of any RCU algorithm**
 - Implicit trade-off between precision and overhead
 - Any extension of a grace period is also a grace period
 - Long (imprecise) grace periods
 - Blocking synchronous writers for a longer time
 - Increasing memory usage due to unreclaimed elements
 - Short (precise) grace periods
 - Increasing overhead on the reader side (need for memory barriers, atomic operations, other heavy-weight operations, etc.)
 - Usual compromise
 - Identifying *naturally occurring quiescent states* for the given RCU algorithm
 - Context switches, exceptions (timer ticks), etc.

The Big Picture ...



- **Foundation for a scalable concurrent data structure**
- **Developing a microkernel-specific RCU algorithm**
 - Specific requirements, constraints and use cases
 - Last well-known RCU implementation for a microkernel in 2003 (K42)

● AP-RCU

- Non-intrusive, portable RCU algorithm
- Developed and implemented by Andrej Podzimek for UTS (OpenSolaris) [3] [4]

● AH-RCU

- Inspired by AP-RCU and several other RCU algorithms
- Developed and implemented by Adam Hraška for SPARTAN (HelenOS) [7]
- Foundation for the Concurrent Hash Table in HelenOS [8]
- Additional variants (preemptible AP-RCU, user space RCU)

- **The RCU algorithm must not impose design concepts of legacy systems on HelenOS**
 - E.g. a specific way how the timer interrupt handler is implemented
- **The kernel space RCU algorithm must support**
 - Read-side critical sections in interrupt and exception handlers
 - Asynchronous reclamation (**call()**) in interrupt and exception handlers
 - Read-side critical sections with preemption enabled (not affecting scheduling latency)

- **Concurrent Hash Table implementation**
 - Growing and shrinking
 - Interrupt and non-maskable interrupt tolerant
 - Suitable for a global page hash table
 - Concurrent reads with low overhead
 - Concurrent inserts and deletes

- **Basic characteristics**

- Kernel space algorithm
- Read-side critical sections are preemptible (without loss of performance)
 - Multiple read-side critical sections within a time slice
 - Expensive operations when a thread was preempted do not make much harm
- Support for asynchronous reclamation in interrupt and exception handlers
- No reliance on periodic timer

- **Grace period detection**

- Test if all CPUs passed a quiescent state
 - Sending an interprocessor interrupt (IPI) to each CPU
 - If the interrupt handler detects a nesting count of 0, it issues a memory barrier (representing a natural quiescent state)
 - Avoid sending IPI if context switch is detected
- Detect any preempted readers holding up the current grace period
 - Sleep and wait for the last preempted reader holding up the grace period to wake the detector thread

- **Advantages**

- Low overhead and preemptible read-side critical section, suitable for exception handlers
- No regular sampling

- **Disadvantages**

- Polling CPUs using interprocessor interrupts might be disruptive in large systems

● Basic characteristics

- Inspired by Triplett's relativistic hash table [5] and Michael's lock-free lists [6]
 - Hash collisions resolved using separate RCU-protected bucket lists
 - Buckets organized as lock-free lists without hazard pointers
 - RCU still protects against accessing invalid pointers and the ABA problem
 - Concurrent lookups and concurrent modifications
 - Tolerance for nested concurrent modifications from interrupt and exception handlers
- Growing and shrinking using background resizing by a factor of 2
 - Concurrent with lookups and updates
 - Requires four grace periods
- Deferred element freeing using RCU `call()`



Enough talk!

Enough talk!
Show me the (pseudo)code!



read_lock():

```
disable_preemption()
check_qs()
cpu.nesting_cnt++
```

read_unlock():

```
cpu.nesting_cnt--
check_qs()
enable_preemption()
```

check_qs():

```
if (cpu.nesting_cnt == 0) {
    if (cpu.last_seen_gp != cur_gp) {
        gp = cur_gp
        memory_barrier()
        cpu.last_seen_gp = gp
    }
}
```

read_lock():

```
disable_preemption()
check_qs()
cpu.nesting_cnt++
```

check_qs():

```
if (cpu.nesting_cnt == 0) {
    if (cpu.last_seen_gp != cur_gp) {
        gp = cur_gp
        memory_barrier()
        cpu.last_seen_gp = gp
    }
}
```

read_unlock():

```
cpu.nesting_cnt--
check_qs()
enable_preemption()
```

The first reader to notice the start of a new grace period on each CPU announces its quiescent state. Once all CPUs announce a quiescent state or perform a context switch (a naturally occurring quiescent state due to disabled preemption), the grace period ends.

Note: Writer forces a context switch on CPUs where no read-side critical section was not observed for a while.

Note: Except `memory_barrier()` only inexpensive operations.

Preemptible AP-RCU Reader Side

read_lock():

```
disable_preemption()
if (thread.nesting_cnt == 0)
    record_qs()
thread.nesting_cnt++
enable_preemption()
```

read_unlock():

```
disable_preemption()
if (thread.nesting_cnt-- == 0) {
    record_qs()
    if ((thread.was_preempted) ||
        (cpu.is_delaying_gp))
        signal_unlock()
}
enable_preemption()
```

record_qs():

```
if (cpu.last_seen_gp != cur_gp) {
    gp = cur_gp
    memory_barrier()
    cpu_last_seen_gp = gp
}
```

signal_unlock():

```
if (atomic_exchange(cpu.is_delaying_gp, false)
    == true)
    remaining_readers_semaphore.up()
if (atomic_exchange(thread.was_preempted, false)
    == true) {
    preempt_mutex.lock()
    preempted_list.remove(thread)
    if ((is_empty(cpu.cur_preempted)) &&
        (preempted_blocking_gp))
        remaining_readers_semaphore.up()
    preempt_mutex.unlock()
}
```



read_lock():

```
thread.nesting_cnt++  
compiler_barrier()
```

read_unlock():

```
compiler_barrier()  
thread.nesting_cnt--  
if (thread.nesting_cnt == was_preempted)  
    preempted_unlock()
```

preempted_unlock():

```
// avoid race between thread and interrupt handler  
if (atomic_exchange(thread.nesting_cnt, 0) == was_preempted) {  
    preempt_lock.lock()  
    preempted_list.remove(thread)  
    if ((is_empty(cpu.cur_preempted)) && (detection_waiting))  
        detection_semaphore.up()           // notify the detector thread about the grace period  
    preempt_lock.unlock()  
}
```

Note: Except **preempted_unlock()** only inexpensive operations.



synchronize():

```
memory_barrier()
mutex.lock()

cur_gp++ // start a new grace period
reader_cpus = [] // gather CPUs potentially in read-side CS
foreach cpu in cpus {
    if (!(cpu.idle) && (cpu.last_seen_gp != cur_gp)) {
        cpu.last_ctx_switch_cnt = cpu.ctx_switch_cnt
        reader_cpus += cpu
    }
}

wait(10ms) // longest acceptable grace period duration (tunable)

foreach cpu in reader_cpus { // enforce a quiescent state
    if (!(cpu.idle) && (cpu.last_seen_gp != cur_gp) &&
        (cpu.last_ctx_switch_cnt == cpu.ctx_switch_cnt))
        cpu.ctx_switch_force_wait()
}

mutex.unlock()
```

detector_thread:

```
forever {
    wait_for_callbacks()

    // run callbacks added before the current grace period
    execute_callbacks()

    // push callbacks registered since last processing to the queue
    advance_callbacks()

    wait_for_gp_end()
}
```


AH-RCU Writer Side (2)



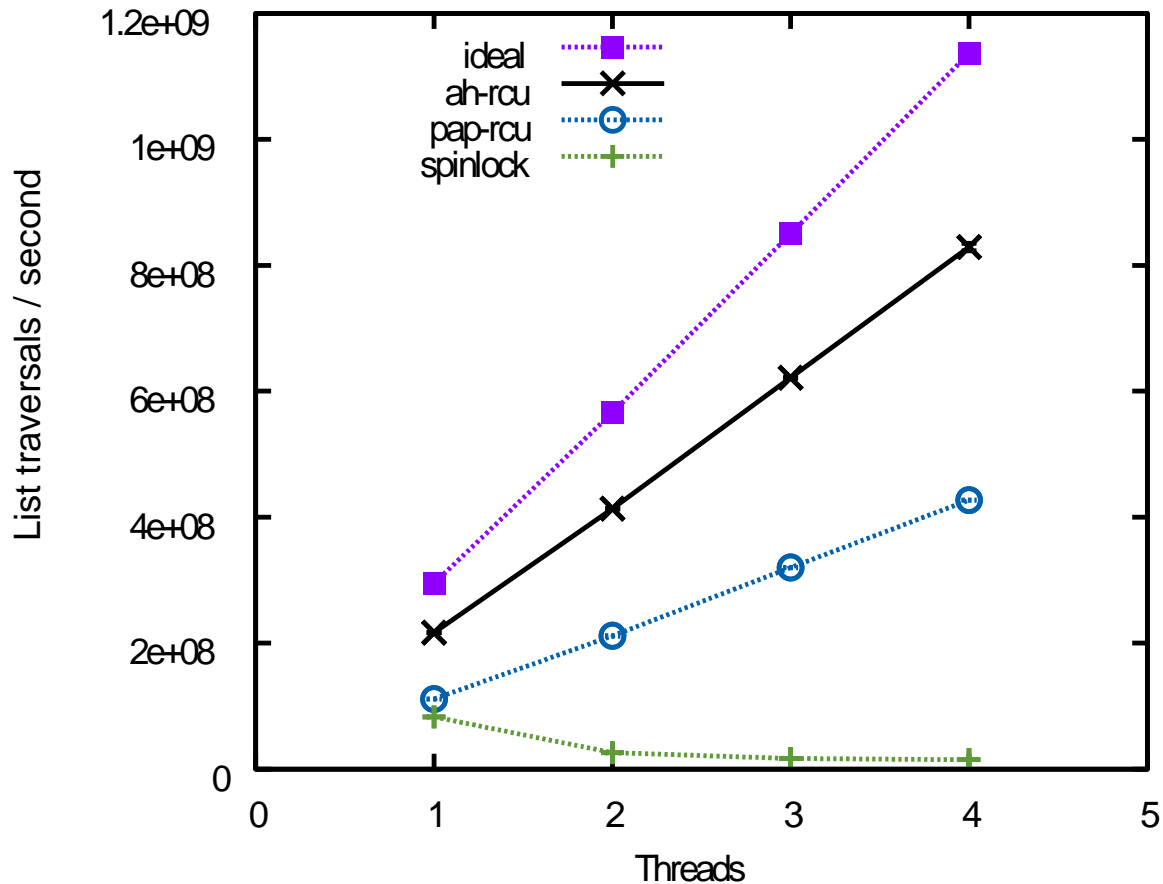
```
wait_for_gp_end():
    gp_mutex.lock()
    if (completed_gp != cur_gp) { // a grace period is already in progress
        wait_for_gp_end_signal()
        goto out
    } else { // start a new grace period
        preempt_lock.lock()
        cur_gp++
        preempt_lock.unlock()
    }
    gp_mutex.unlock()

    wait_for_readers()

    gp_mutex.lock()
    completed_gp = cur_gp

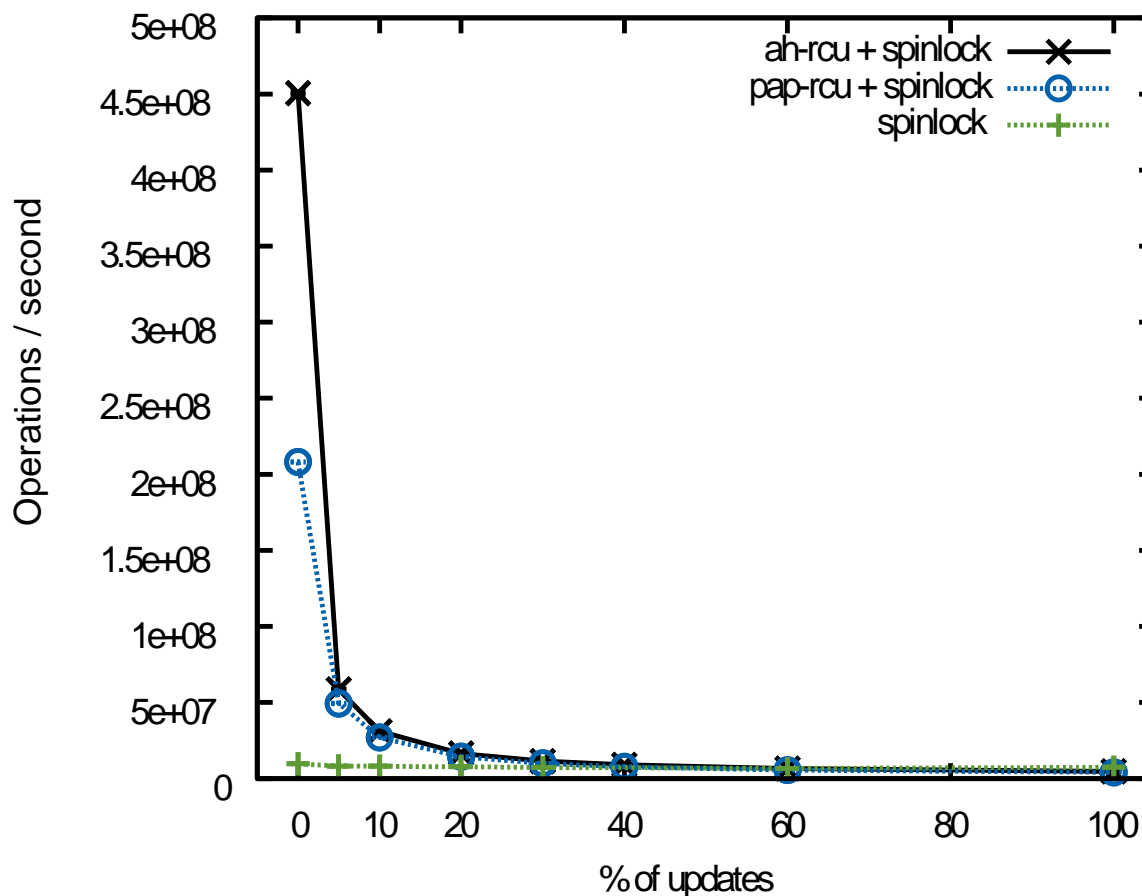
    out:
        gp_mutex.unlock()
```

Read-side critical section scalability: Traversal of a five-element list
The list is protected as a whole, it is only read, never modified.



Evaluation (2)

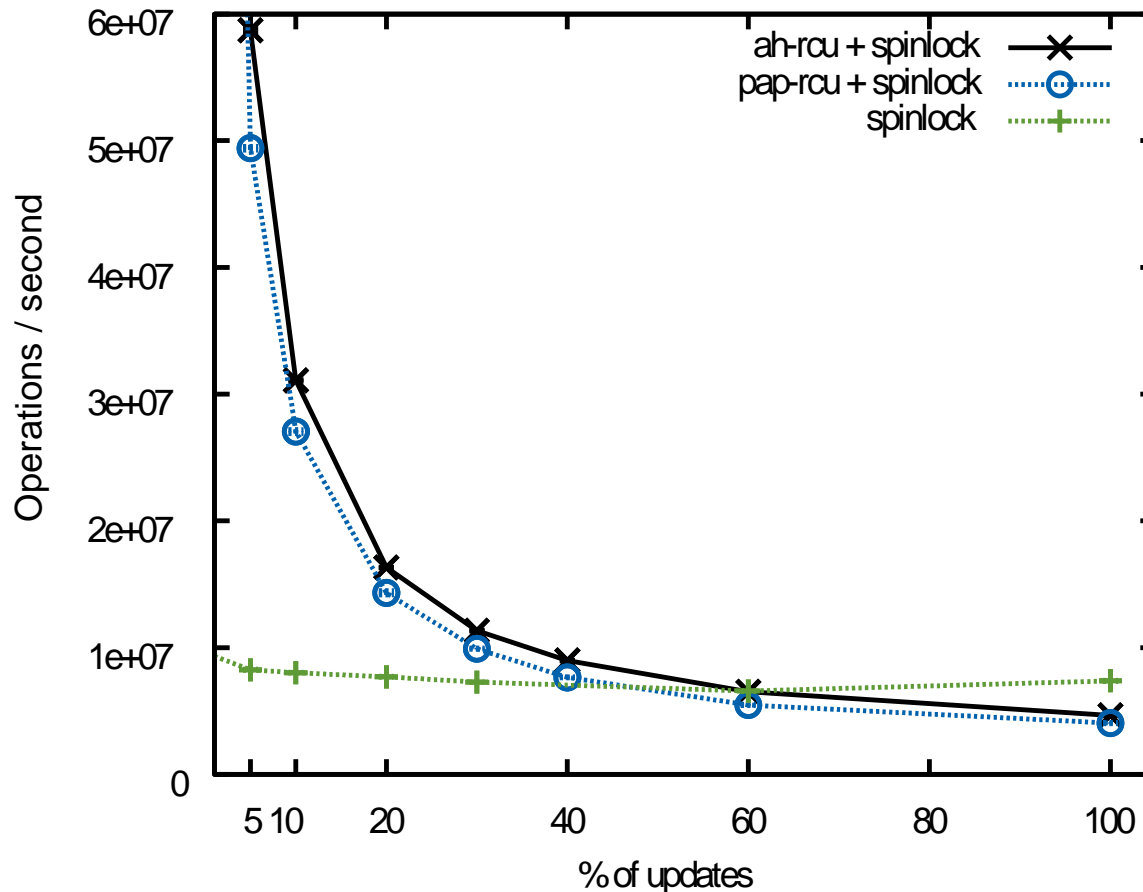
Write-side overhead: Different ratios of updates vs. lookups
Five-element list, four threads running in parallel. Updates are always synchronized by a spinlock.



Evaluation (3)

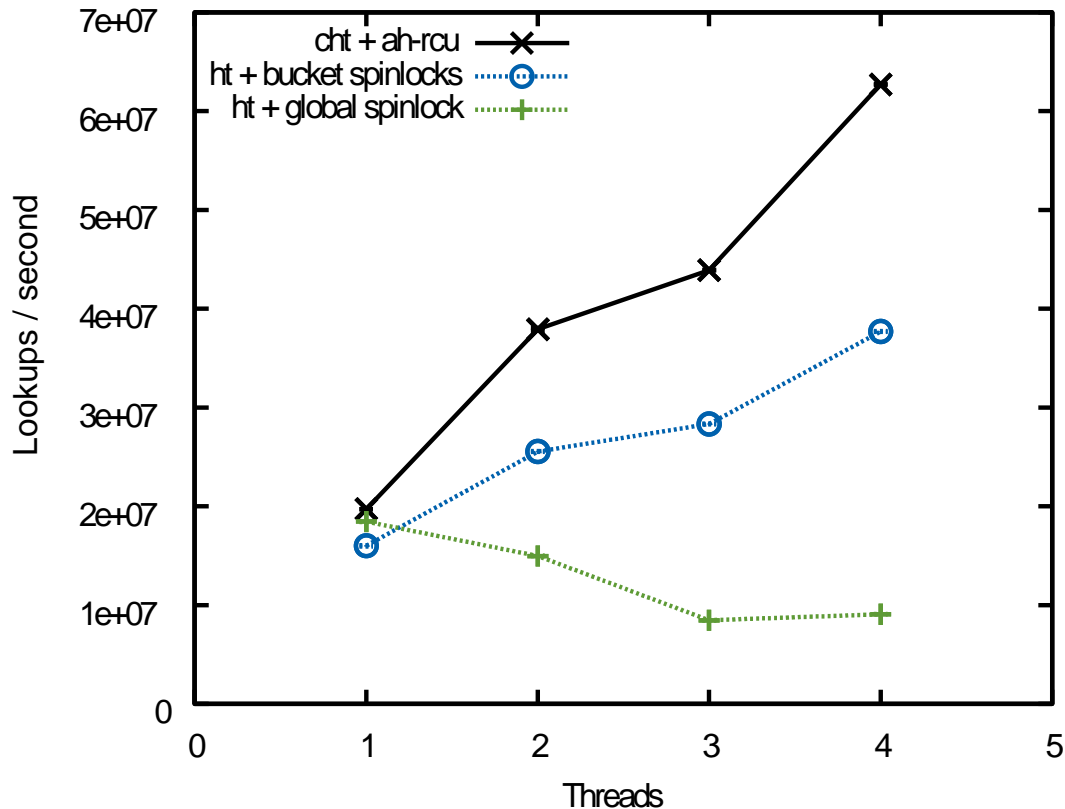


Read-side scalability vs. write-side overhead: Crossover point
Data points from previous figure with low fraction of updates are discarded.



Concurrent hash table lookup scalability

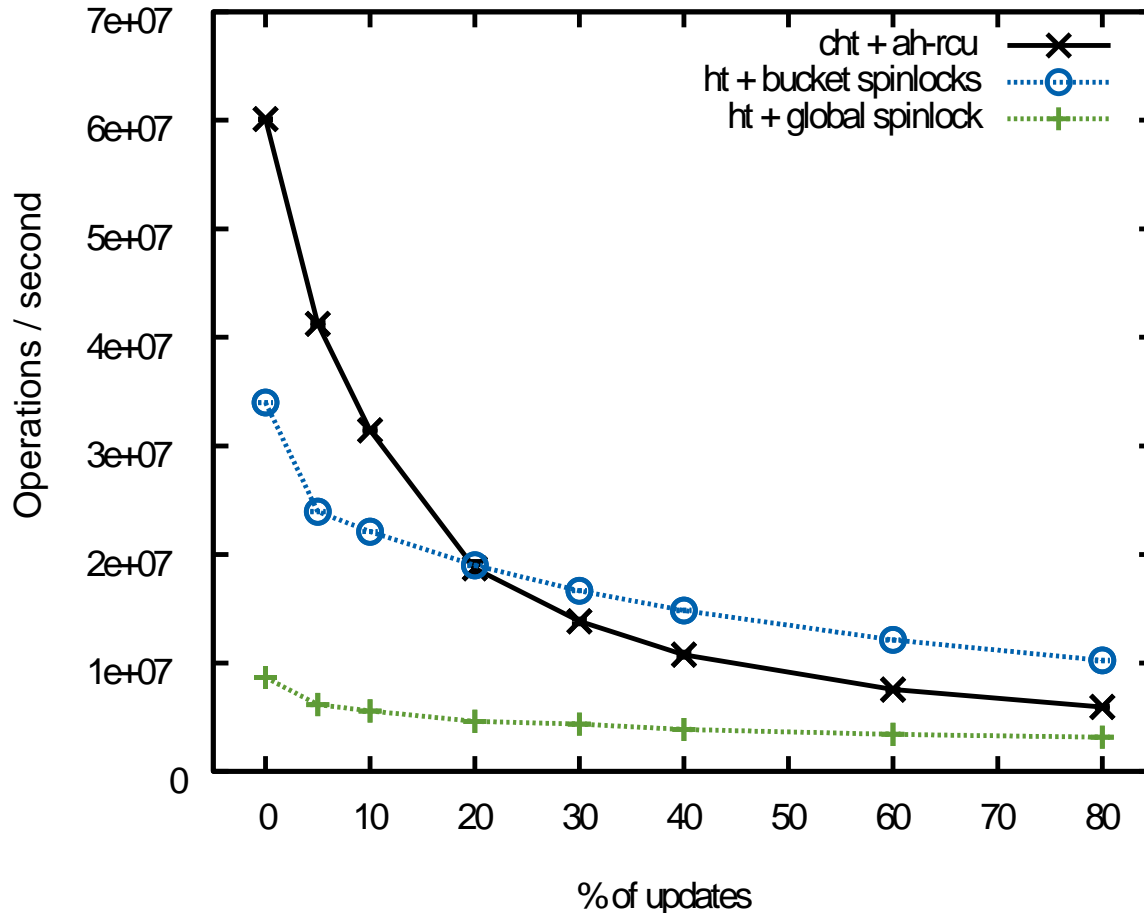
128 buckets, average load factor of 4 elements per bucket, 50 % of lookups for hitting keys, 50 % of lookups for missing keys (each thread used a separate list). The resize condition was checked, but never executed.



Evaluation (5)



Concurrent hash table update overhead: Different ratios of concurrent updates vs. lookups
Four threads running in parallel.



- **Novel scalable algorithms**
 - Preemptible AP-RCU for HelenOS
 - Preemptible AH-RCU for HelenOS
 - Resizeable Concurrent Hash Table for HelenOS
 - Suitable as a basic data structure for asynchronous HelenOS IPC
 - Suitable for other kernel uses (e.g. global page table)
- **Thorough evaluation**
 - Promising behavior



www.helenos.org

- [1] Herlihy M. P.: *Impossibility and universality results for wait-free synchronization*, in Proceedings of 7th Annual ACM Symposium on Principles of Distributed Computing, ACM, 1988
- [2] Kogan A., Petrank E.: *Wait-free queues with multiple enqueueers and dequeuers*, in Proceedings of 16th ACM Symposium on Principles and Practice of Parallel Programming, ACM, 2011
- [3] Podzimek A., Děcký M., Bulej L., Tůma P.: *A Non-Intrusive Read-Copy-Update for UTS*, in Proceedings of 18th IEEE International Conference on Parallel and Distributed Systems, IEEE, 2012,
<http://d3s.mff.cuni.cz/publications/download/PodzimekDeckyBulejTuma-ICPADS-2012.pdf>
- [4] <http://d3s.mff.cuni.cz/software/rcu/rcu.patch>
- [5] Triplett J., McKenney P. E., Walpole J.: *Resizable, scalable, concurrent hash tables via relativistic programming*, in Proceedings of the 2011 USENIX Annual Technical Conference, ACM, 2011
- [6] Michael M. M.: *High performance dynamic lock-free hash tables and list-based sets*, in Proceedings of 14th Annual ACM Symposium on Parallel Algorithms and Architectures, ACM, 2002
- [7] Hraška A.: *Read-Copy-Update for HelenOS*, master thesis, Charles University in Prague, 2013,
<http://www.helenos.org/doc/theses/ah-thesis.pdf>
- [8] <https://code.launchpad.net/~adam-hraska+lp/helenos/cht-bench>