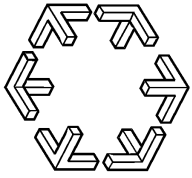# Memory Management for Lock-Free Concurrent Data-structures

Anders Gidenstam
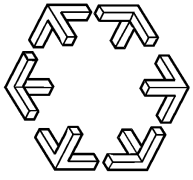
PostDoc,

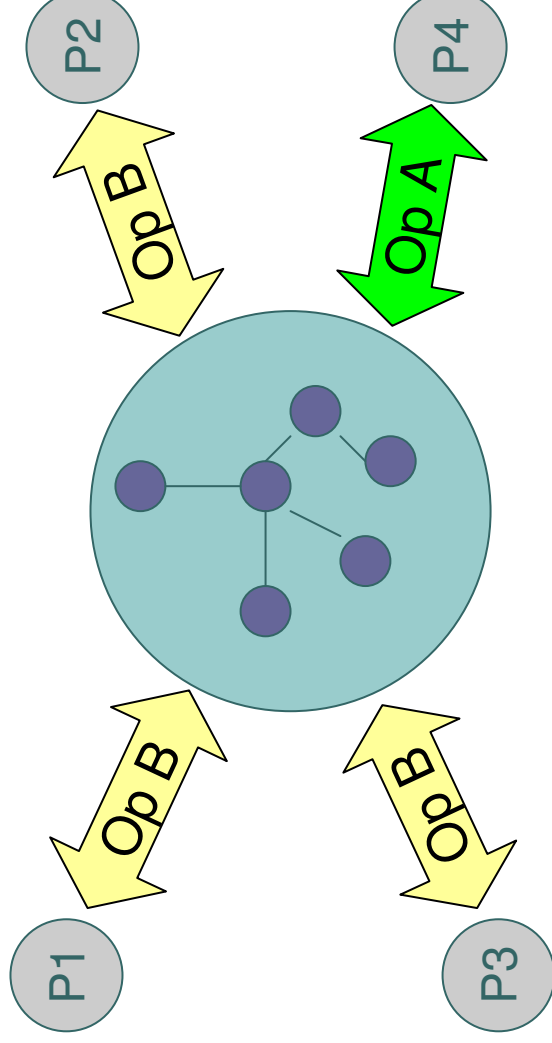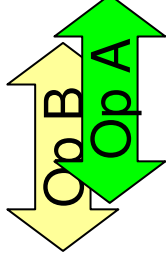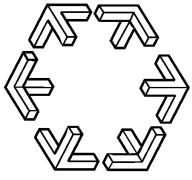AG1, Max-Planck-Institut für Informatik

# Outline

- Introduction to non-blocking algorithms
  - System model
  - Correctness criteria
  - Progress guarantees
  - Motivating example
- The lock-free memory reclamation problem
  - Solutions
  - LFMR [Gidenstam, Papatriantafilou, Sundell & Tsigas. I-SPAN 2005]
    - Idea
    - Properties
- What's out there: Some lock-free data-structures
- Current work

Anders Gidenstam, Dagstuhl
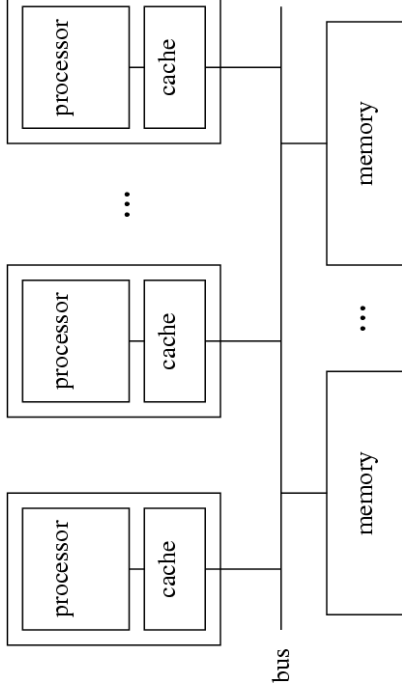
# Concurrent Shared Memory Data-structures

○ Object in shared memory
- Supports some set of operations (ADT)
- Concurrent access by many processes/threads

○ Useful to e.g.
- Exchange data between threads
- Coordinate thread activities

Op B
Op A

P1   P2   P3   P4

Op B   Op A   Op B   Op B
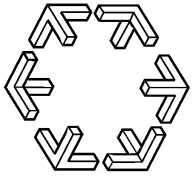
Anders Gidenstam, Dagstuhl

# System model



- Processes / threads
  - Asynchronous
    - Each executes a sequence of instructions

- Shared memory
  - Processes can read/write single memory words atomically
  - Hardware synchronization primitives/instructions
    - Compare-and-Swap(address, old, new)
      - Atomic read-modify-write (i.e. a critical section of one instruction)
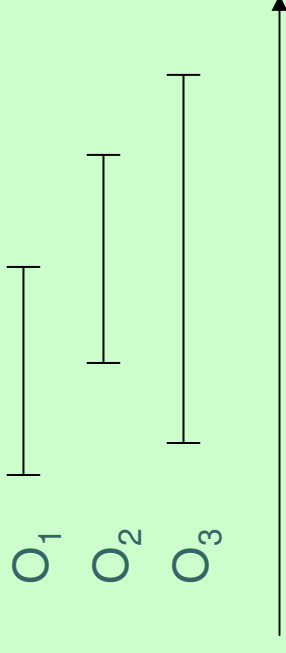    - Load-Linked(address) / Store-Conditional(address, new)

Anders Gidenstam, Dagstuhl

# Correctness of a concurrent object

○ Desired semantics of a shared data object

● Linearizability [Herlihy & Wing, 1990]

  • For each operation invocation there must be one single time instant during its duration where the operation appears to take effect.

  • The observed effects should be consistent with a sequential execution of the operations in that order.

$O_1$

$O_2$
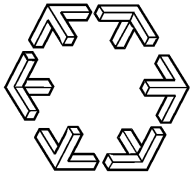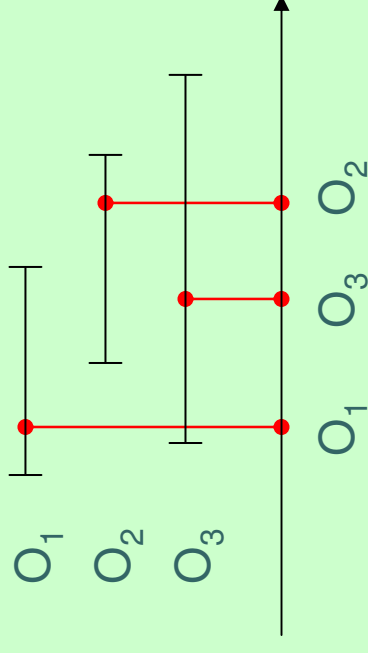
$O_3$

Anders Gidenstam, Dagstuhl

# Correctness of a concurrent object

- Desired semantics of a shared data object
  - Linearizability [Herlihy & Wing, 1990]
    - For each operation invocation there must be one single time instant during its duration where the operation appears to take effect.
    - The observed effects should be consistent with a sequential execution of the operations in that order.
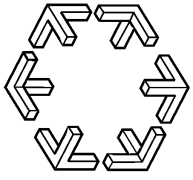
# Correctness of a concurrent object

○ Desired semantics of a shared data object

● Linearizability [Herlihy & Wing, 1990]

- For each operation invocation there must be one single time instant during its duration where the operation appears to take effect.

- The observed effects should be consistent with a sequential execution of the operations in that order.
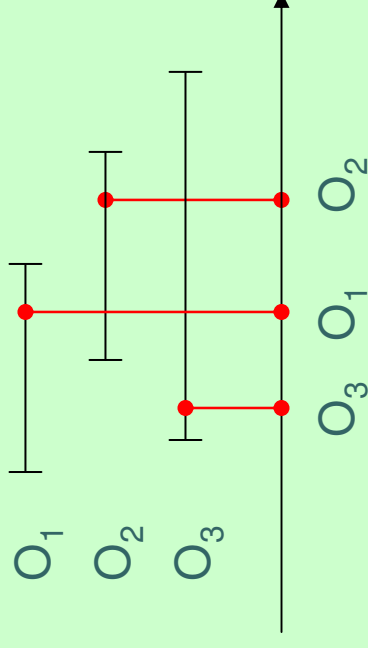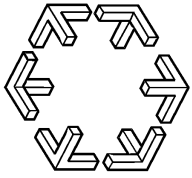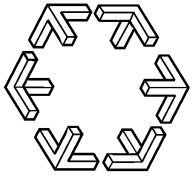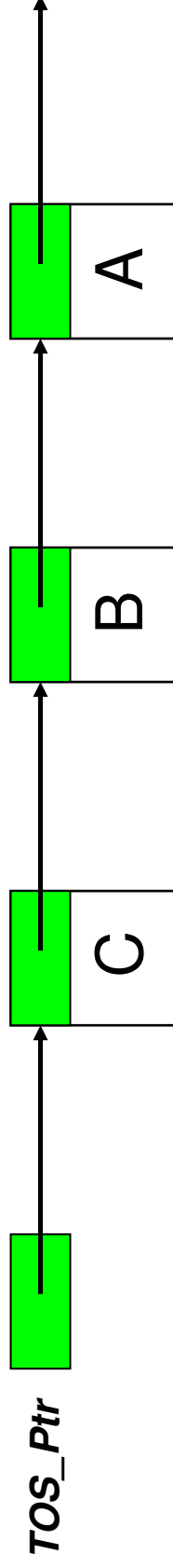
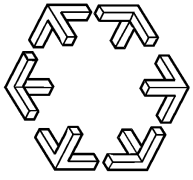13/03/2008

# Progress Guarantees

- **Traditional solution: Mutual exclusion**
  - Semaphores, mutexes, spin-locks, disabling interrupts
    - Protects critical sections
  - Drawbacks: blocking, deadlock, priority-inversion, limits parallelism

- **Non-blocking synchronization**
  - Wait-free synchronization [Lamport, 1977]
    - Every operation finishes in a finite number of its own steps.
  - **Lock-free synchronization** [Lamport, 1977]
    - At least one operation in a set of concurrent operation always makes progress.
  - Obstruction-free synchronization [Herlihy et. al. 2003]
    - Any operation that eventually executes in isolation is guaranteed to make progress.

13/03/2008

# Example: Lock-free Stack

- Operations
  - Push(item)
  - Pop: item

- Linked list based algorithm [IBM 1983]

*TOS_Ptr*

| | | |
|---|---|---|
| C | B | A |

Anders Gidenstam, Dagstuhl

# Example: Lock-free Stack

- Operations
  - Push(item)
  - Pop: item

  1. Read *TOS_Ptr*
  2. Prepare the new node

- Linked list based algorithm [IBM 1983]

*TOS_Ptr*

D

C

B

A

Anders Gidenstam, Dagstuhl
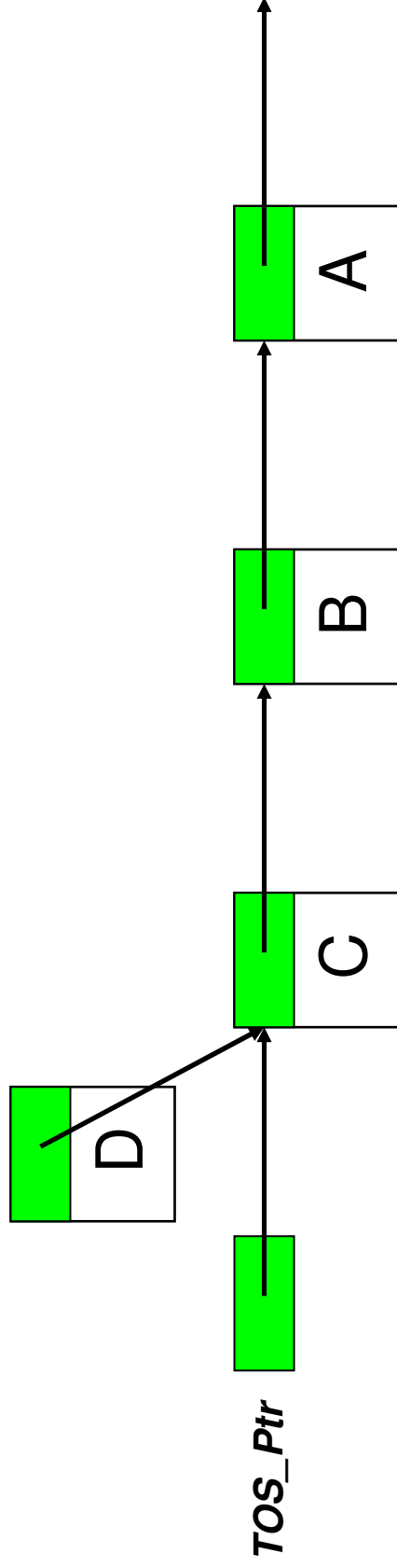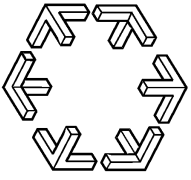
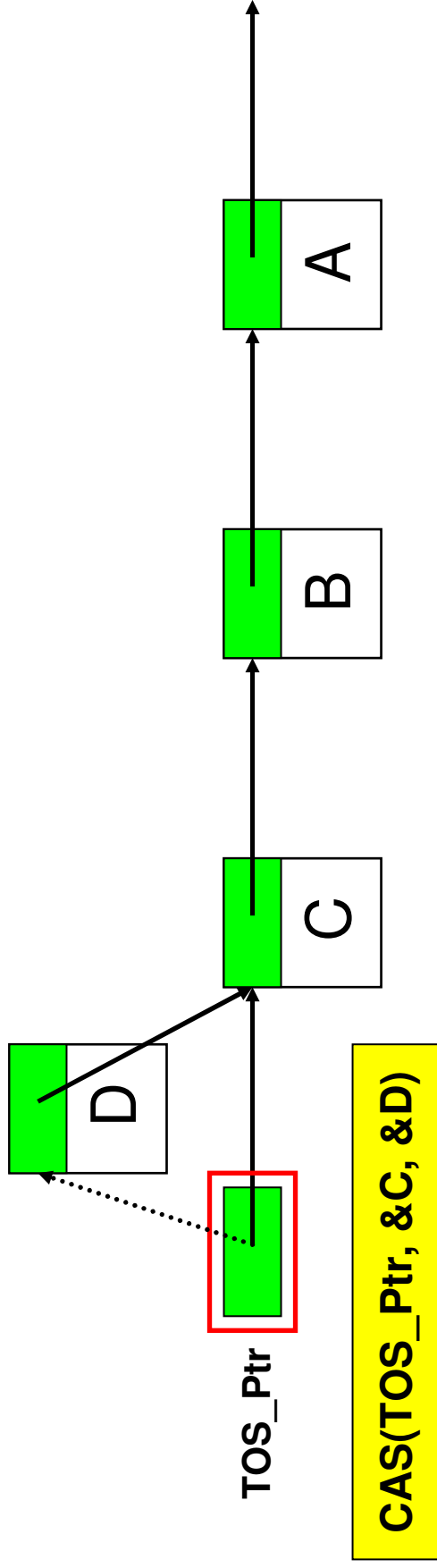13/03/2008

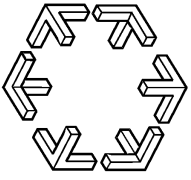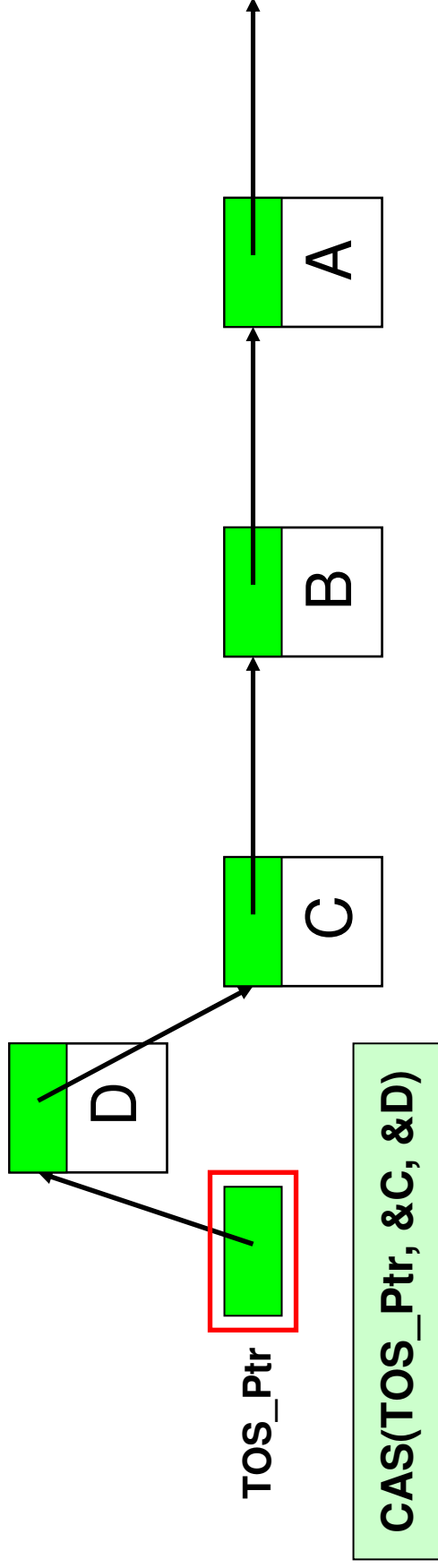# Example: Lock-free Stack

- Operations
  - **Push(item)**
  - Pop: item
- Linked list based algorithm [IBM 1983]

1. Read *TOS_Ptr*
2. Prepare the new node
3. Try to update *TOS_Ptr* with Compare&Swap

**TOS_Ptr**

**CAS(TOS_Ptr, &C, &D)**

A

B

C

D

# Example: Lock-free Stack

○ Operations

- Push(item)
- Pop: item

1. Read *TOS_Ptr*
2. Prepare the new node
3. Try to update *TOS_Ptr* with Compare&Swap
4. If successful then done else retry from 1.

○ Linked list based algorithm [IBM 1983]



A

B

C

D

**TOS_Ptr**

**CAS(TOS_Ptr, &C, &D)**

Anders Gidenstam, Dagstuhl

# Example: Lock-free Stack

- Operations
  - Push(item)
  - Pop: item

  1. Read *TOS_Ptr*
  2. Read *TOS_Ptr->next*

- Linked list based algorithm [IBM 1983]

*TOS_Ptr* → [ D ] → [ C ] → [ B ] →

Anders Gidenstam, Dagstuhl

13/03/2008
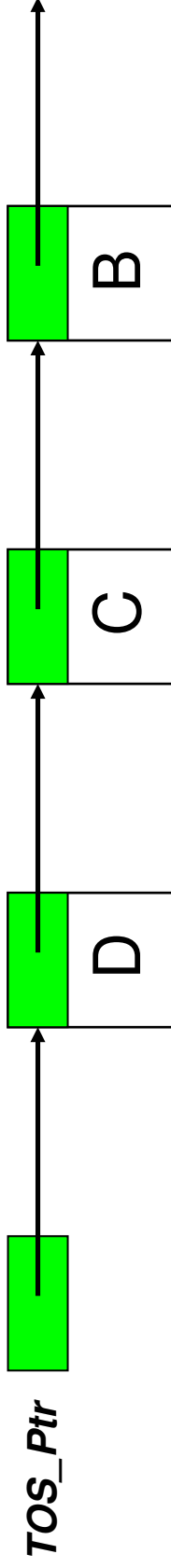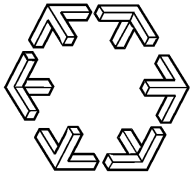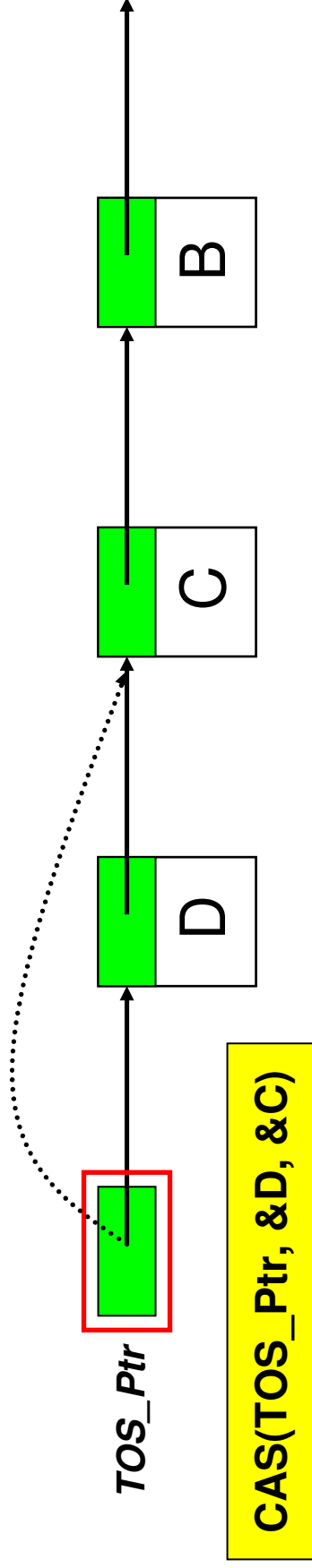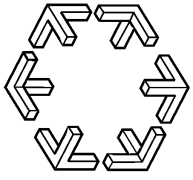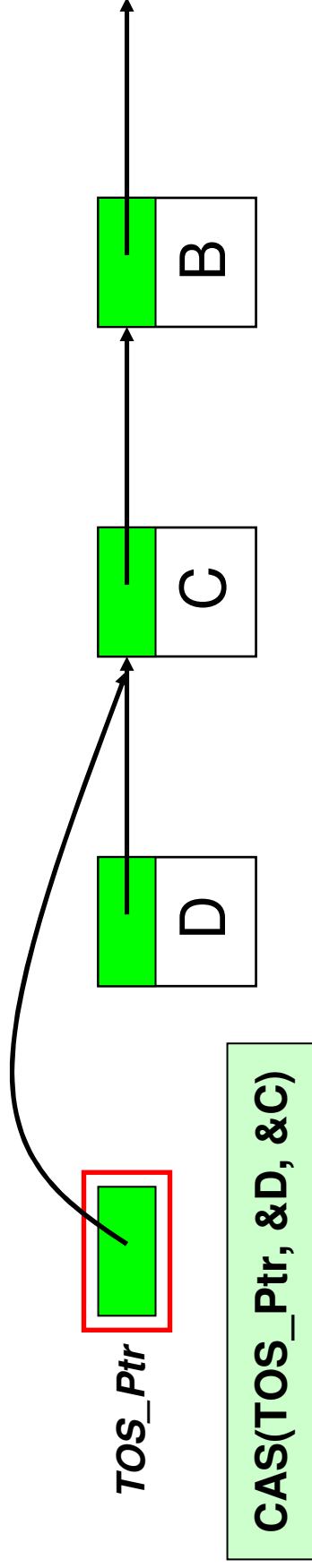
# Example: Lock-free Stack

- Operations
  - Push(item)
  - Pop: item

- Linked list based algorithm [IBM 1983]

Operations:

1. Read *TOS_Ptr*
2. Read *TOS_Ptr->next*
3. Try to update *TOS_Ptr* with Compare&Swap

*TOS_Ptr*

| | | | |
|---|---|---|---|
| D | C | B | |

**CAS(TOS_Ptr, &D, &C)**

13/03/2008

# Example: Lock-free Stack

- Operations
  - Push(item)
  - **Pop: item**

  1. Read *TOS_Ptr*
  2. Read *TOS_Ptr->next*
  3. Try to update *TOS_Ptr* with Compare&Swap
  4. If successful then delete the node and return the item else retry from 1.

- Linked list based algorithm [IBM 1983]



*TOS_Ptr*

**CAS(TOS_Ptr, &D, &C)**

Anders Gidenstam, Dagstuhl

# Example: Lock-free Stack

- Operations
  - Push(item)
  - Pop: item

  1. Read *TOS_Ptr*
  2. Read *TOS_Ptr->next*
  3. Try to update *TOS_Ptr* with Compare&Swap
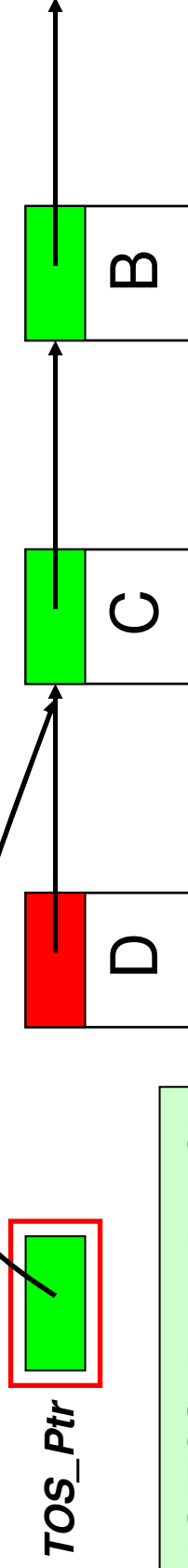  4. If successful then delete the node and return the item else retry from 1.

- Linked list based algorithm [IBM 1983]

**But is it safe to delete node D now?**

D

C

B

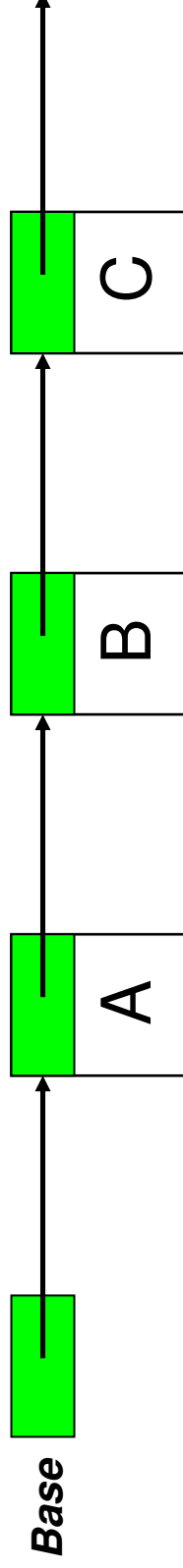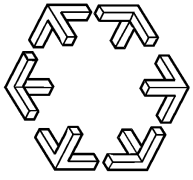*TOS_Ptr*

**CAS(TOS_Ptr, &D, &C)**

# The Lock-Free Memory Reclamation Problem

- Concurrent shared data structures with
  - Dynamic use of shared memory
  - Concurrent and overlapping operations by threads or processes

Can nodes be deleted and reused safely?

*Base* → [ ] → [A] → [B] → [C] →

13/03/2008

# The Lock-Free Memory Reclamation Problem

**Thread X**

*Local variables*

A B C

**Base**

# The Lock-Free Memory Reclamation Problem

**Thread X**

X has de-referenced the link (pointer) to B

**Base** → A → B → C →

13/03/2008

# The Lock-Free Memory Reclamation Problem

Another thread, Y, finds and deletes (removes) B from the active structure

*Thread X*

*Thread Y*

B

C

A

*Base*

Anders Gidenstam, Dagstuhl

13/03/2008

# The Lock-Free Memory Reclamation Problem

Property I: A node accessible via a private reference (i.e. *dereferenced*) should not be reclaimed

*Thread X*

*Thread Y*

Thread Y wants to reclaim(/free) B

C

?

A

*Base*

Anders Gidenstam, Dagstuhl

13/03/2008

# The Lock-Free Memory Reclamation Problem

Property II: Links in a dereferenced node should always be dereferencable.

The nodes B and C are deleted from the active structure.

**Thread X**

B

C

?

A

D

Base

13/03/2008

# The Lock-Free Memory Reclamation Problem

Solutions providing Property I but not Property II:

○ Hazard Pointers [Michael 2002]

○ Pass the Buck [Herlihy 2002]

**Hazard pointers** (of thread X) **Global**

*Thread X*

*Thread Y*

Thread Y to reclaim(/free) B

**?**

A

C

*Base*

Anders Gidenstam, Dagstuhl

13/03/2008

23

# The Lock-Free Memory Reclamation Problem

Reference counting can guarantee

- Property I
- Property II

But it needs to be lock-free!

**Thread X**

**Base**

A  B  C  D

Anders Gidenstam, Dagstuhl

13/03/2008

# Lock-free reference counting solutions

- [Valois et al, 1995]   Reference-count field MUST remain writable forever.
- LFRC [Detlefs et al, 2001]   Needs double word CAS.
- SLFRC [Herlihy et al, 2002/2005] Pure reference counting (RC).
- LFMR [Gidenstam et al, 2005]  RC + application guidance.

- Issues with pure reference counting
  - A slow thread with a private reference might prevent reclamation
  - Cyclic garbage



*Slow*

B  C  A

Anders Gidenstam, Dagstuhl

# Our approach – The basic idea

Joint work w. M. Papatriantafilou, H. Sundell & P. Tsigas

○ Combine the best of

● Hazard pointers [Michael 2002]

  ○ Tracks references from threads

  ○ Fast de-reference

  ○ Upper bound on the number of unreclaimed deleted nodes

  ○ Compatible with standard memory allocators

**Hazard pointers** (of thread X) **Global**

*Thread X*

B

Deletion list

# Our approach – The basic idea

Joint work w. M. Papatriantafilou, H. Sundell & P. Tsigas

o Combine the best of

- Hazard pointers [Michael 2002]
  - Tracks references from threads
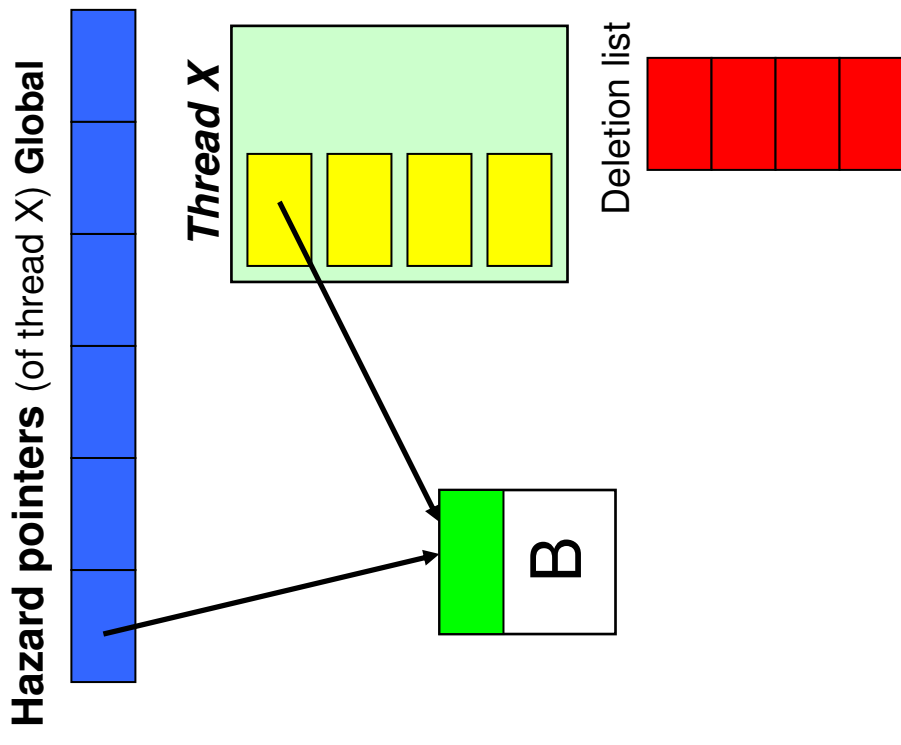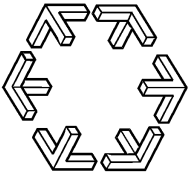  - Fast de-reference
  - Upper bound on the number of unreclaimed deleted nodes
  - Compatible with standard memory allocators

- Reference counting
  - Tracks references from links in shared memory
  - Manages links within dynamic nodes
  - Safe to traverse links (also) in deleted nodes

**Hazard pointers** (of thread X) **Global**

*Thread X*

Deletion list

| 1 | |
|---|---|
| | B |

| 1 | |
|---|---|
| | A |

*Base*

Anders Gidenstam, Dagstuhl

# The basic idea

o API
- DeRefLink
- ReleaseRef
- CompareAndSwapRef
- StoreRef
- NewNode
- DeleteNode

**Hazard pointers** (Thread X)

Deletion list

*Thread X*

*Base* → A → B → C

13/03/2008

# The basic idea

- API
  - **DeRefLink(Base)**
  - ReleaseRef
  - CompareAndSwapRef
  - StoreRef
  - NewNode
  - DeleteNode

**Hazard pointers** (Thread X)

Deletion list

*Thread X*

R

*Base*

1 — A

1 — B

1 — C

Anders Gidenstam, Dagstuhl

# The basic idea

- API
  - DeRefLink
  - **ReleaseRef(R)**
  - CompareAndSwapRef
  - StoreRef
  - NewNode
  - DeleteNode

**Hazard pointers** (Thread X)

Deletion list

*Thread X*

R

*Base*

A  |1|
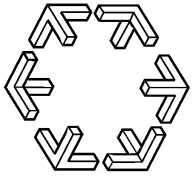
B  |1|

C  |1|

Anders Gidenstam, Dagstuhl

# The basic idea

- API
  - DeRefLink
  - ReleaseRef
  - CompareAndSwapRef
  - StoreRef
  - **NewNode**
  - DeleteNode

**Hazard pointers** (Thread X)

Deletion list

**Thread X**

new

*Base*

| 1 | A | | 1 | B | | 1 | C |

| 0 | D |

Anders Gidenstam, Dagstuhl

# The basic idea

o API
- DeRefLink
- ReleaseRef
- CompareAndSwapRef
- **StoreRef(new.next, R)**
- NewNode
- DeleteNode

**Hazard pointers** (Thread X)

Deletion list

*Thread X*

R

C

new

Base | A | B | C | D

Anders Gidenstam, Dagstuhl

# The basic idea

○ API
- DeRefLink
- ReleaseRef
- CompareAndSwapRef(prev.next, old, new)
- StoreRef
- NewNode
- DeleteNode

**Hazard pointers** (Thread X)

| | | | A | B | |
|---|---|---|---|---|---|

Deletion list

**Thread X**

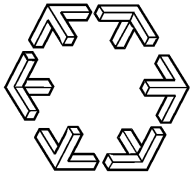| prev | A |
|---|---|
| | |
| old | B |
| new | |

Base

A — 1
D — 1
B — 0
C — 2

# The basic idea

- API
  - DeRefLink
  - ReleaseRef
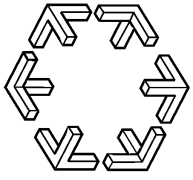  - CompareAndSwapRef
  - StoreRef
  - NewNode
  - DeleteNode(old)

**Hazard pointers** (Thread X)
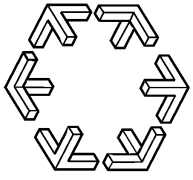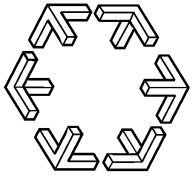
Deletion list

**Thread X**

| prev | A |
| old | |
| new | |

*Base*

| 1 | A |
| 1 | D |
| 0 | B |
| 2 | C |

Anders Gidenstam, Dagstuhl

# Bound on #unreclaimed nodes

**Hazard pointers** (Thread Y)

Deletion list
of Thread X

D — 1
C — 0
B — 0

**Theorem:** The maximum number of deleted but not yet reclaimed nodes in the system is never more than

$$N^2 \cdot (k + l_{max} + a + 1)$$

where

$N$ is the number of threads in the system,
$k$ is the number of hazard pointers per thread,
$l_{max}$ is the maximum number of links a node can contain and
$a$ is the maximum number of links in live nodes that may transiently point to a deleted node during an operation.

Anders Gidenstam, Dagstuhl

# Breaking chains of garbage

- Clean-up deleted nodes
  - Update links to point to live nodes
  - Performed on nodes in
    - Own deletion list
    - All deletion lists

**Hazard pointers** (Thread Y)

Deletion list

***Thread X***

***Base***

A | 1
D | 1
B | 0
C | 1
E | 2

13/03/2008

# Breaking chains of garbage

- Clean-up deleted nodes
  - Update links to point to live nodes
  - Performed on nodes in
    - Own deletion list
    - All deletion lists

**Hazard pointers** (Thread Y)

Deletion list

*Thread X*

Reclaimable

C

0

E

3

D

1

B

0

A

1

*Base*

Anders Gidenstam, Dagstuhl

# Bound on #unreclaimed nodes

- A deleted node is unreclaimable if
  - A hazard pointer points to it
    - Limited #hazard pointers: $N \cdot k$
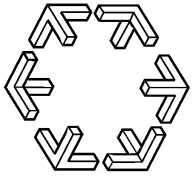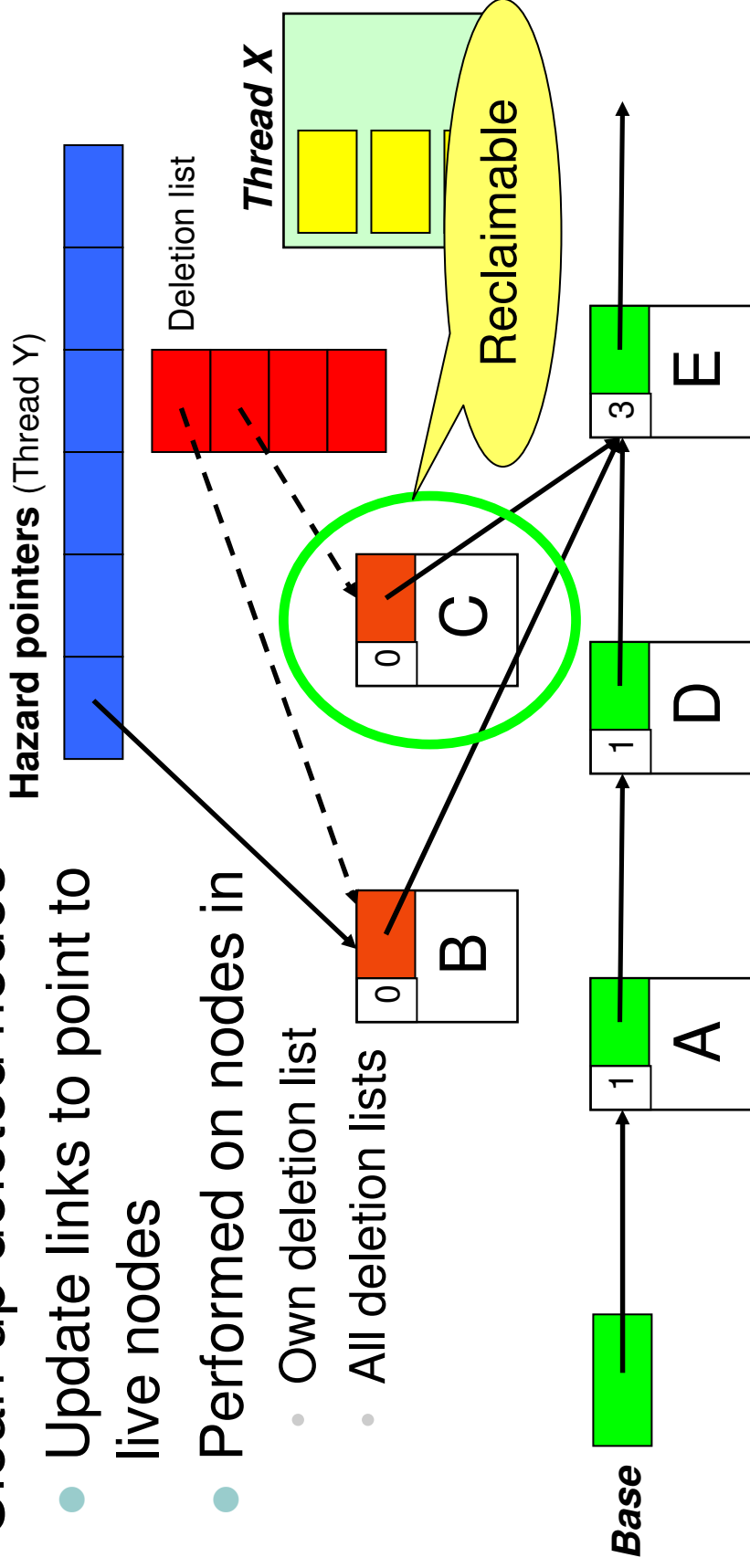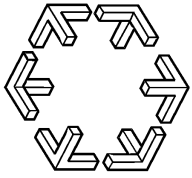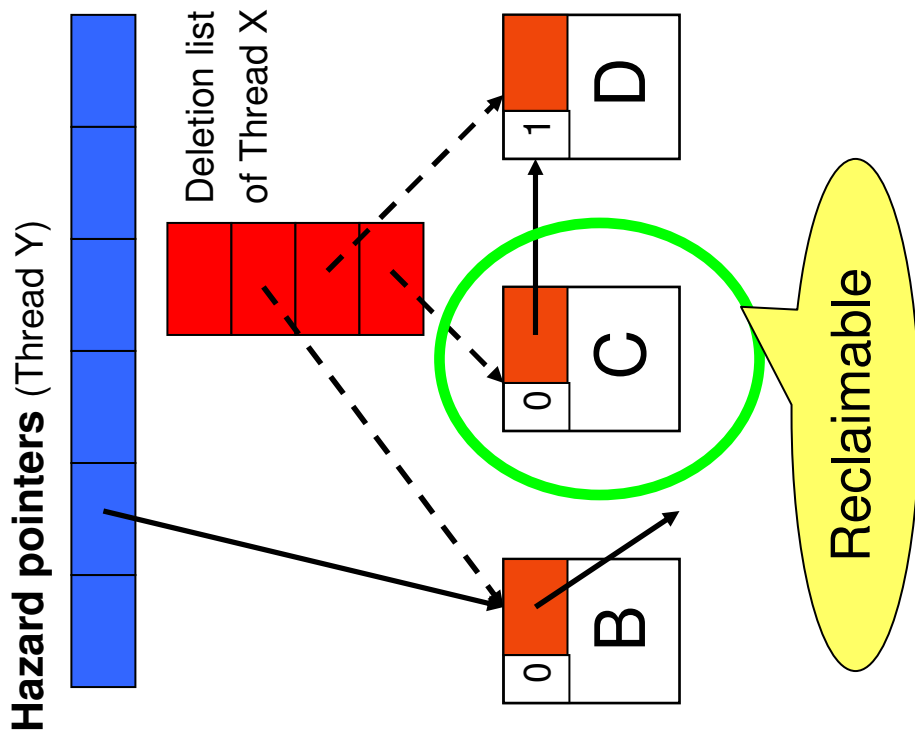  - Its reference count is nonzero
    - Limited #links in live nodes pointing to deleted nodes: $N \cdot a$
    - The links in most deleted nodes can be cleaned by any process.
      Exception: each thread can "hide" one node during a Delete operation.
      - #links in any node is bounded
        $$\Rightarrow N \cdot l_{max}$$
  - It is being cleaned by another thread: $N \cdot 1$

- $\Rightarrow$ The maximum size of a thread's deletion list is bounded by
  $$N \cdot (k + l_{max} + a + 1)$$

- $\Rightarrow$ The total number of unreclaimable deleted nodes is bounded by
  $$N^2 \cdot (k + l_{max} + a + 1)$$

Hazard pointers (Thread Y)

Deletion list of Thread X

D   1

C   0

B   0

Reclaimable

Anders Gidenstam, Dagstuhl

38

# Experimental evaluation

- Lock-free deque [Sundell and Tsigas 2004]
  (deque – double-ended queue)
  - The algorithm needs traversal of deleted nodes
  - Time for 10000 random operations/thread
- Tested memory reclamation schemes
  - Reference counting, Valois et al.
  - LFRM (a.k.a. the new algorithm)
- Systems
  - 4 processor Xeon PC / Linux (UMA)
  - 8 processor SGI Origin 2000 / IRIX (NUMA)

Anders Gidenstam, Dagstuhl

# Experimental evaluation



Lock-Free Deque - Linux Xeon, 4 Processors

VALOIS ET AL
NEW ALGORITHM

Anders Gidenstam, Dagstuhl

13/03/2008

# Experimental evaluation



Lock-Free Deque - SGI Origin 2000, 8 Processors

VALOIS ET AL ——
NEW ALGORITHM ---×---

Anders Gidenstam, Dagstuhl

13/03/2008

# Some lock-free data-structures

**Stacks**

| | Memory man. |
|---|---|
| [IBM 1983],[Treiber,1986] | version #s / HP |

**Queues**

| | |
|---|---|
| [Valois, 1994] | RC |
| [MS, 1996] | RC / HP |
| [TZ, 2001] | bounded size |
| [Hoffman et al, 2007] | version #s / RC |

**Double Ended Queues / Deques**

| | |
|---|---|
| [M, ?] | ? |
| [ST, 2004] | RC |

**Priority Queues**

| | |
|---|---|
| [Barnes, 1994] | bounded size |
| [ST 2003] | RC |

**Sets/Dictionaries**

| | |
|---|---|
| [Michael, 2002] | HP |
| [TS 2004] | RC |

**Linked lists**

| | Memory man. |
|---|---|
| [Valois, 1995] | RC |
| [Harris, 2001] | RC |

**Doubly linked lists**

| | |
|---|---|
| [TS, 2007] | RC |

**Hash tables**

| | |
|---|---|
| [Michael, 2002] | HP |
| [Shalev Shavit 2006] | ? |

**Trees (Binary, Red-Black)**

| | |
|---|---|
| [Fraser, 2004] | M-CAS / STM / RC |

# Current Work

- Memory management
  - Unified and easy to use interfaces to the memory management algorithms [NBAda library]
    - Easier to implement lock-free data-structures
    - Data-structure user could choose memory management method.

- Lock-free data-structures
  - Develop new (Red-Black trees with 1-CAS?)
  - Software library: NBAda
    - Ought to be here: http//www.mpi-inf.mpg.de/~andersg/ For now: Ask me.