# Lock-free algorithms for Kotlin coroutines

It is all about scalability

Presented at SPTCC 2017

/Roman Elizarov @ JetBrains

# Speaker: Roman Elizarov

- 16+ years experience
- Previously developed high-perf trading software @ Devexperts
- Teach concurrent & distributed programming @ St. Petersburg ITMO University
- Chief judge @ Northeastern European Region of ACM ICPC
- Now work on Kotlin @ JetBrains

# Agenda

- Kotlin coroutines overview & motivation for lock-free algorithms

- Lock-free doubly linked list

- Lock-free multi-word compare-and-swap

- Combining them to get more complex atomic operations (without STM)

# Kotlin basic facts

- Kotlin is a JVM language developed by JetBrains
- General purpose and statically-typed
- Object-oriented and functional paradigms
- Open source under Apache 2.0
- Reached version 1.0 in 2016
  - Compatibility commitment
  - Now at version 1.1
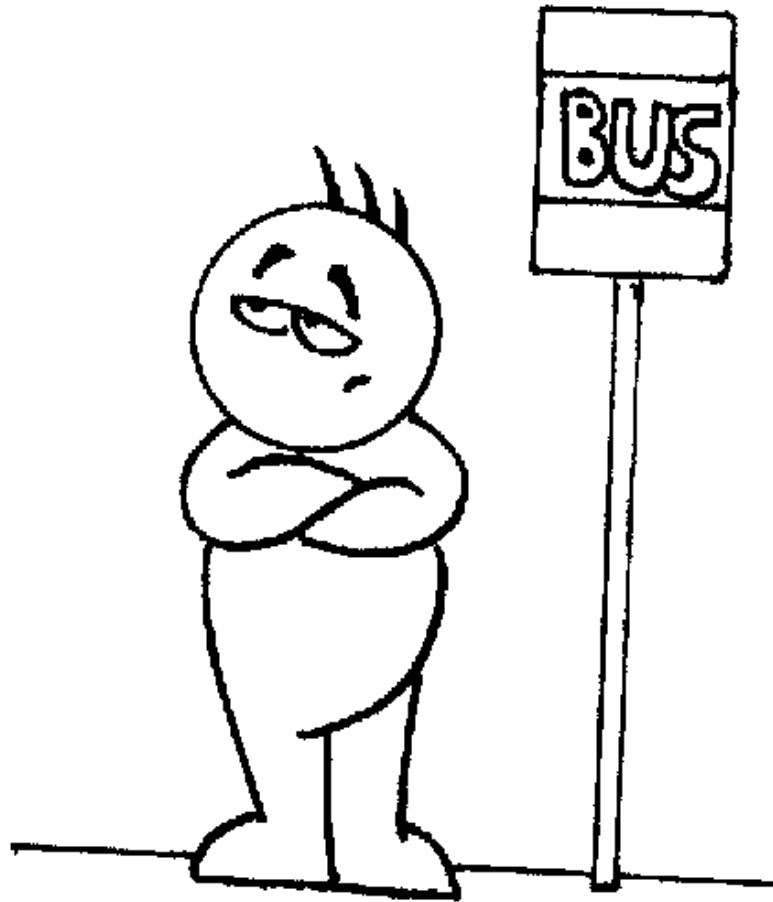- Officially supported by Google on Android

# Kotlin is …

- Modern
- Concise
- Safe
- Extensible
- Pragmatic
- Fun to work with!

# Kotlin is pragmatic

... and easy to learn

# Coroutines

Asynchronous programming made easy

How do we write code that waits for something most of the time?

# Blocking threads

```kotlin
fun postItem(item: Item) {
    val token = requestToken()
    val post = submitPost(token, item)
    processPost(post)
}
```

# Callbacks

```kotlin
fun postItem(item: Item) {
    requestToken { token ->
        submitPost(token, item) { post ->
            processPost(post)
        }
    }
}
```

# Futures/Promises/Rx

```kotlin
fun postItem(item: Item) {
    requestToken()
        .thenCompose { token ->
            submitPost(token, item)
        }
        .thenAccept { post ->
            processPost(post)
        }
}
```

# Coroutines

```kotlin
fun postItem(item: Item) {
    launch(CommonPool) {
        val token = requestToken()
        val post = submitPost(token, item)
        processPost(post)
    }
}
```

# CSP & Actor models

- A *style* of programming for modern systems
- Lots of concurrent tasks / jobs
  - Waiting most of the time
  - Communicating all the time

Share data by communicating

# Kotlin coroutines primitives

- Jobs/Deferreds (futures)
  - join/await
- Channels
  - send & receive
  - synchronous & buffered channels
- Select/alternatives
  - Atomically wait on multiple events
- Cancellation
- Parent-child hierarchies

# Implementation challenges

- Coroutines are like light-weight threads
- All the *low-level* scheduling & communication mechanisms have to *scale* to lots of coroutines

# Lock-free algorithms

# Building blocks

- Single-word CAS (that's all we have on JVM)

- Automatic memory management (GC)

- *Practical* lock-free algorithms
  - Lock-Free and Practical **Doubly Linked List**-Based Deques Using Single-Word Compare-and-Swap by Sundell and Tsigas
  - A Practical **Multi-Word Compare-and-Swap** Operation by Timothy L. Harris, Keir Fraser and Ian A. Pratt.

# Doubly linked list

**next** links form logical list contents
**prev** links are auxiliary



sentinel

sentinel

Use same node in practice

# Insert

PushRight (like in queue)

# Doubly linked list (insert 0)

create & init

# Doubly linked list (insert 1)



CAS

Retry insert on CAS
failure

# Doubly linked list (insert 2)

"finish insert"



CAS

Ignore CAS failure

# Remove

PopLeft (like in queue)
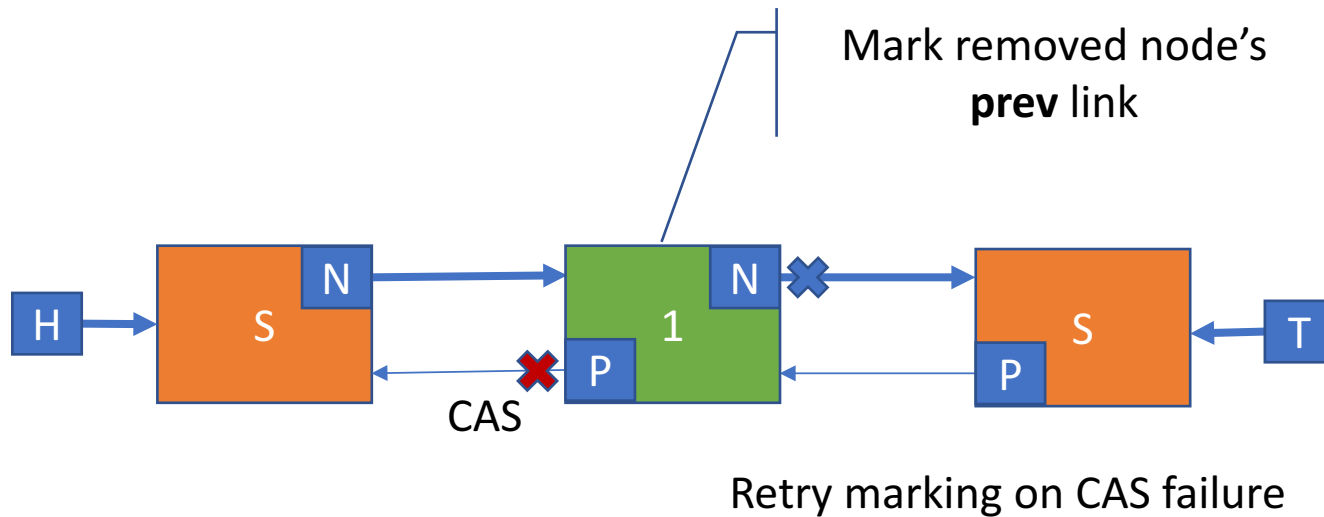
# Doubly linked list (remove 1)



Retry remove on CAS failure

Use *wrapper* object for mark in practice

*Don't use* **AtomicMarkableReference**

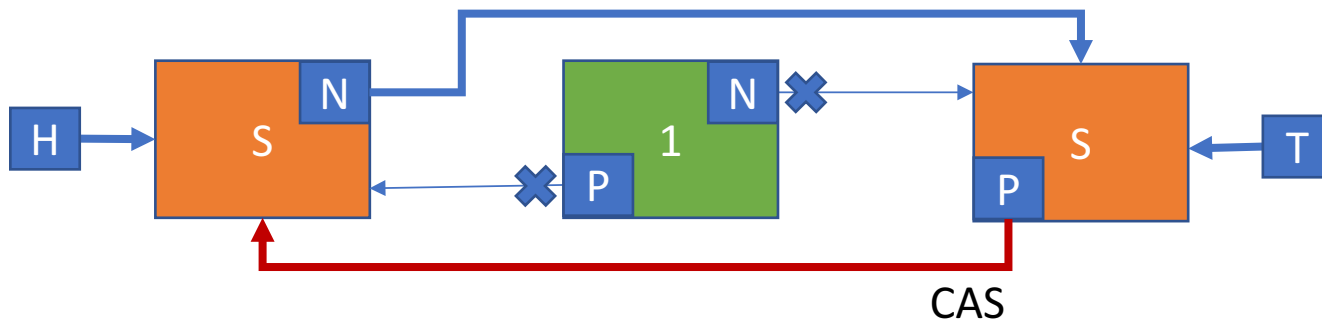Cache wrappers in pointed-to nodes

# Doubly linked list (remove 2)

"finish remove"



Mark removed node's **prev** link

CAS

Retry marking on CAS failure

# Doubly linked list (remove 3)

"help remove" – fixup **next** links

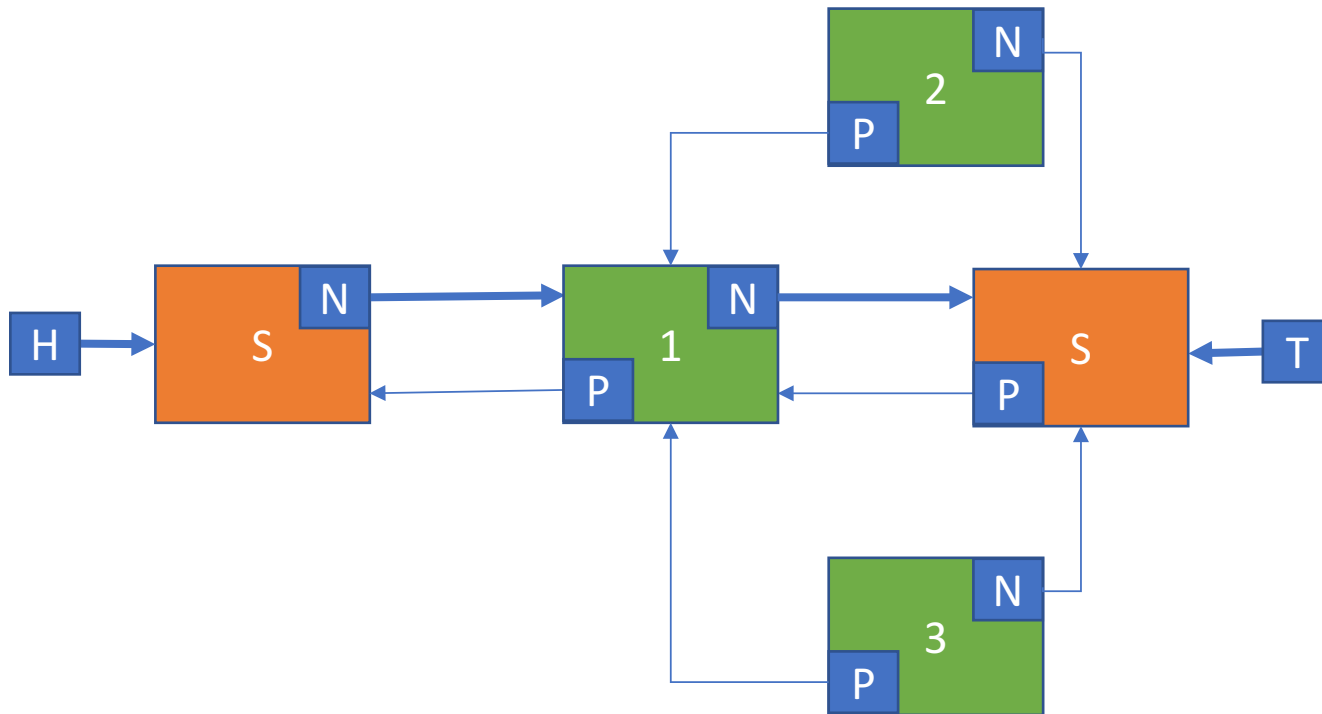# Doubly linked list (remove 4)

"correct prev" – fixup **prev** links

# State transitions

# Node states

correct prev

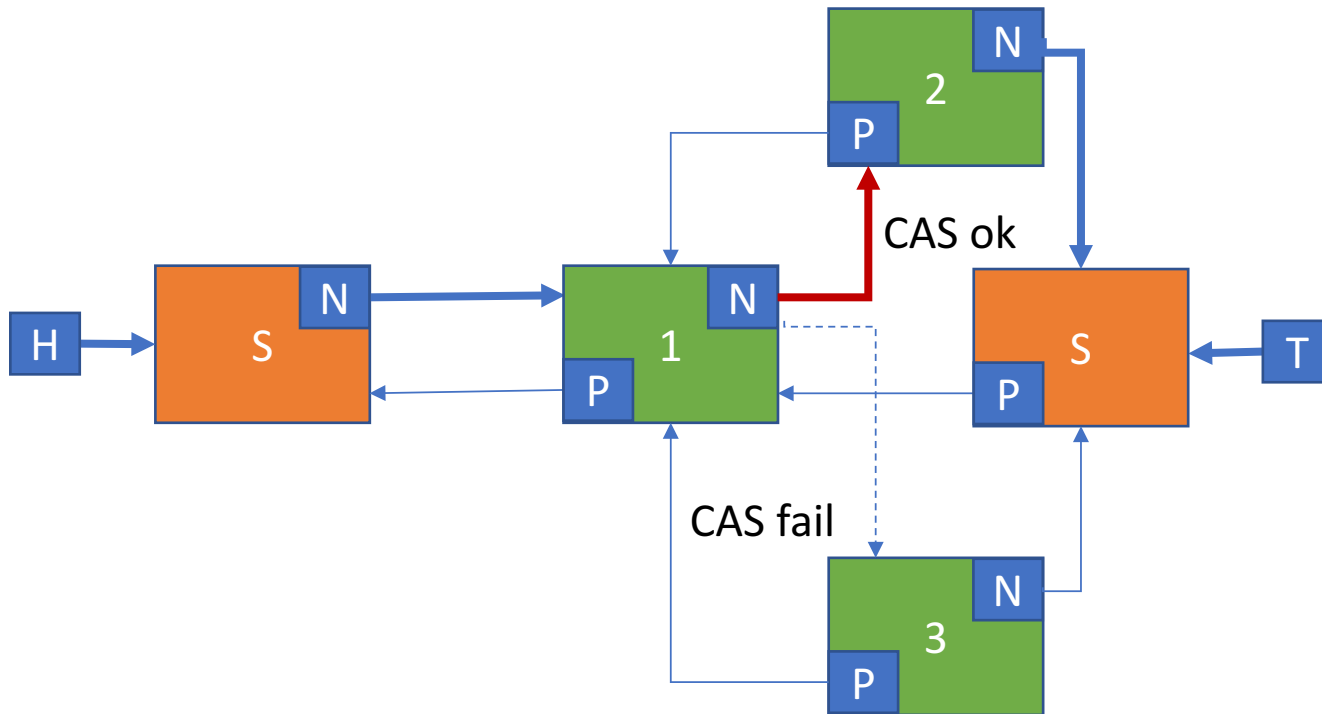**1** Init
next: Ok
prev: Ok
prev.next: --
next.prev: --

**2** Insert 1
next: Ok
prev: Ok
prev.next: me
next.prev: --

**3** Insert 2
next: Ok
prev: Ok
prev.next: me
next.prev: me

**4** Remove 1
next: Rem
prev: Ok
prev.next: me
next.prev: me

**5** Remove 2
next: Rem
prev: Rem
prev.next: me
next.prev: me

**6** Remove 3
next: Rem
prev: Rem
prev.next: ++
next.prev: me

**7** Remove 4
next: Rem
prev: Rem
prev.next: ++
next.prev: ++

help remove

correct prev

# Helping

# Concurrent insert

# Concurrent insert (0)

# Concurrent insert (1)

# Concurrent insert (2)

| I2 |
|---|
| I3 |



detect wrong prev
(t.prev.next != t)

# Concurrent insert (3)



correct prev

# Concurrent insert (4)
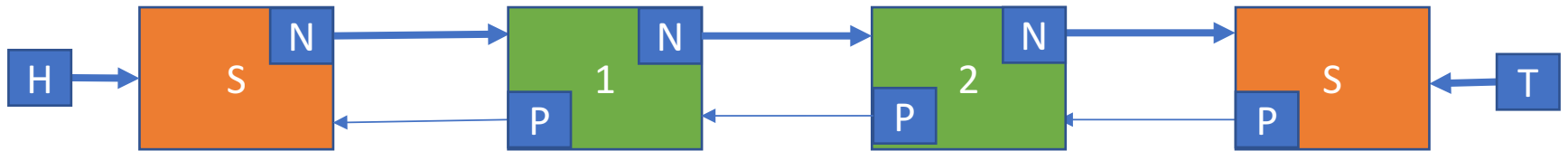
I2
I3

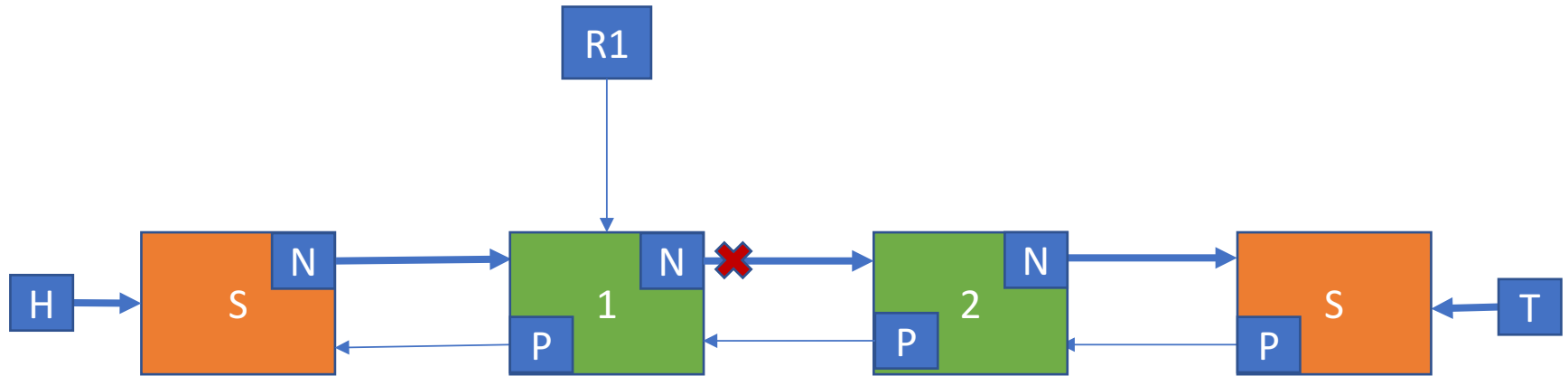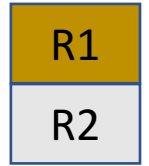H → S (N) → 1 (N) (P) → 2 (N) (P) → S (P) ← T
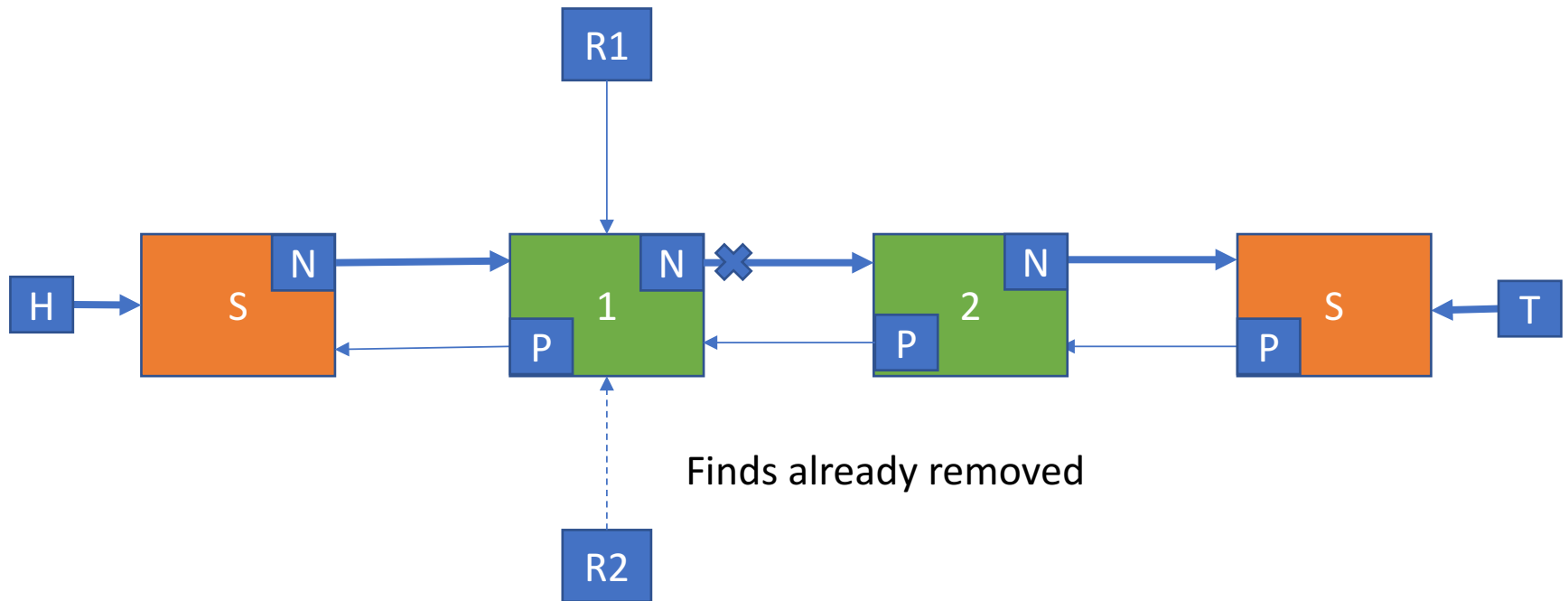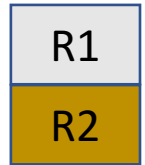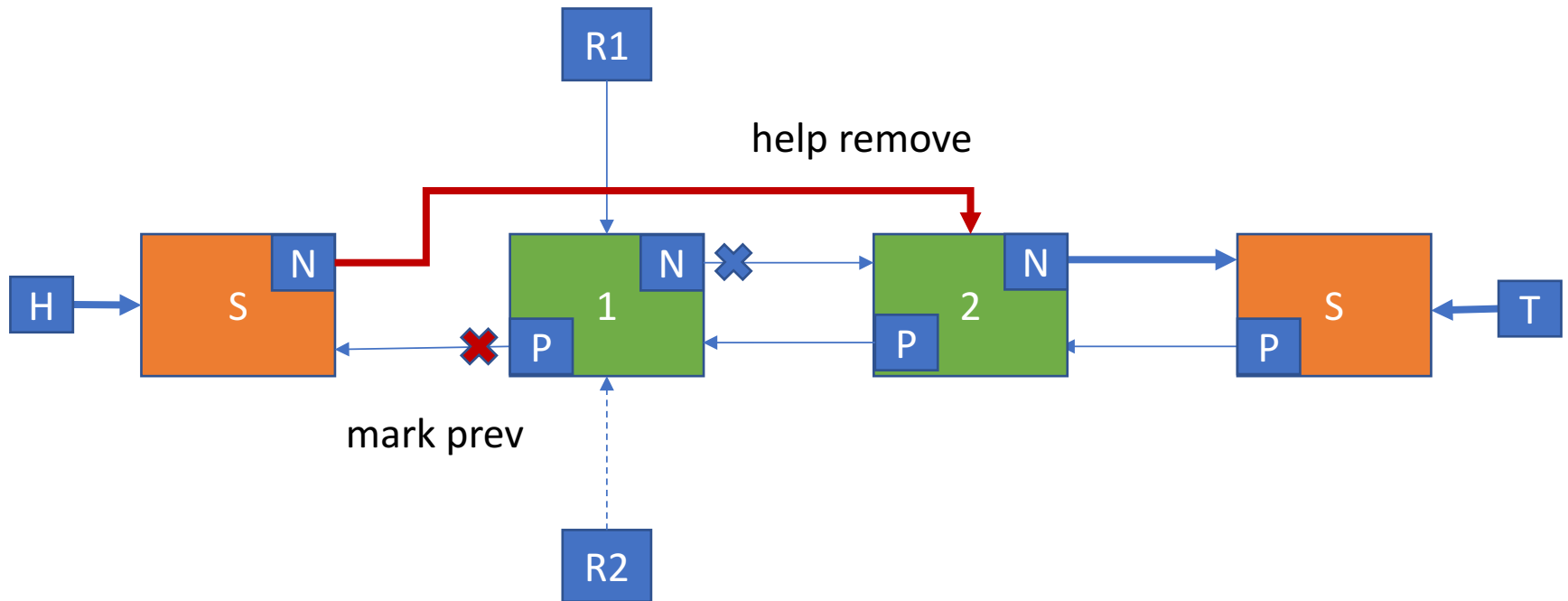
3 (N) (P)

reinit & repeat

# Concurrent remove

# Concurrent remove (0)

| |
|---|
| R1 |
| R2 |

# Concurrent remove (1)

R1

R2

R1

H  S  N  1  N  ❌  2  N  S  P  T

# Concurrent remove (2)



Finds already removed

# Concurrent remove (3)

R1

R2

R1

help remove

H

S
N

1
N

2
N

S
P
T

P

P

P

mark prev

R2

# Concurrent remove (4)



Retry with corrected next

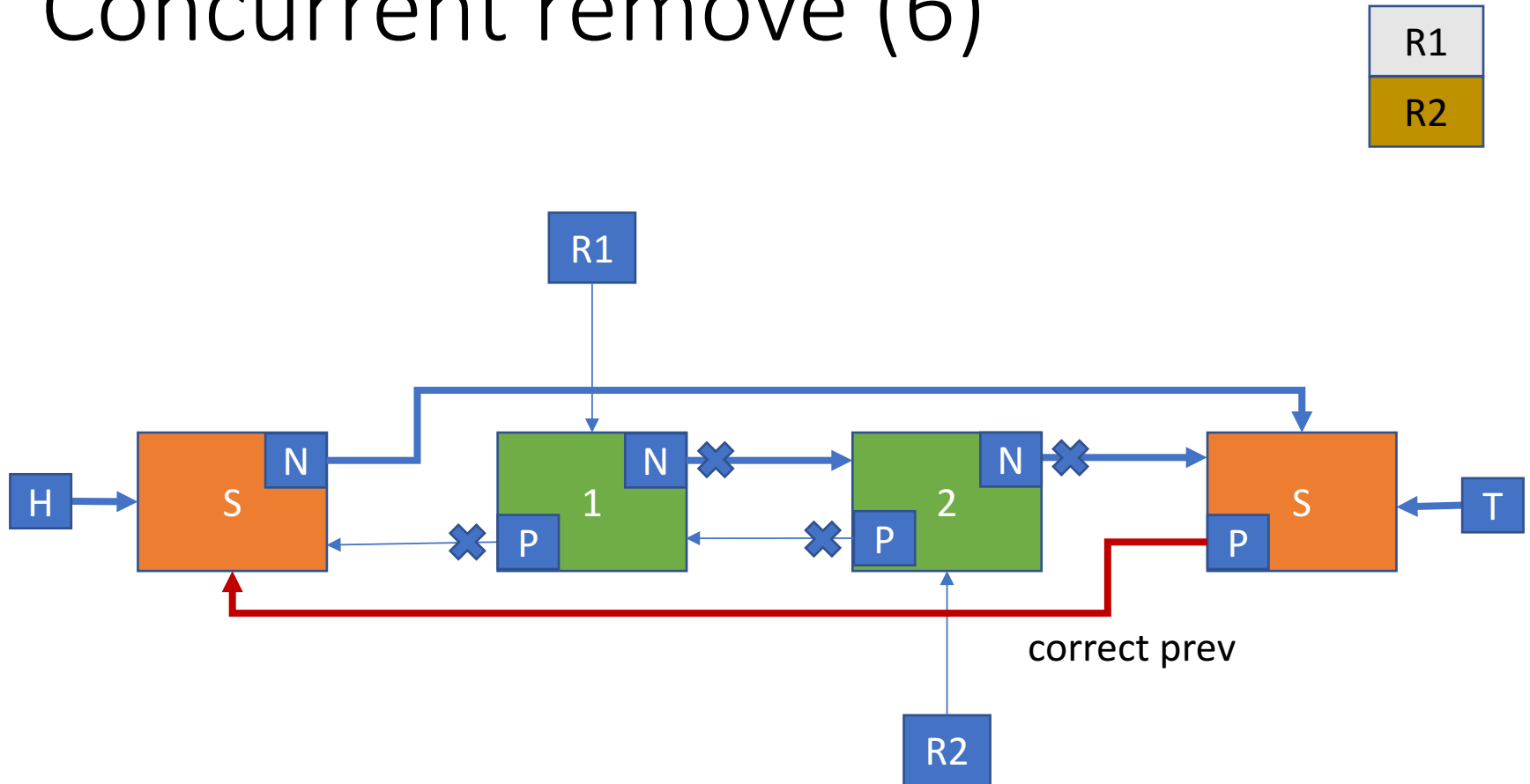# Concurrent remove (5)

# Concurrent remove (6)



correct prev

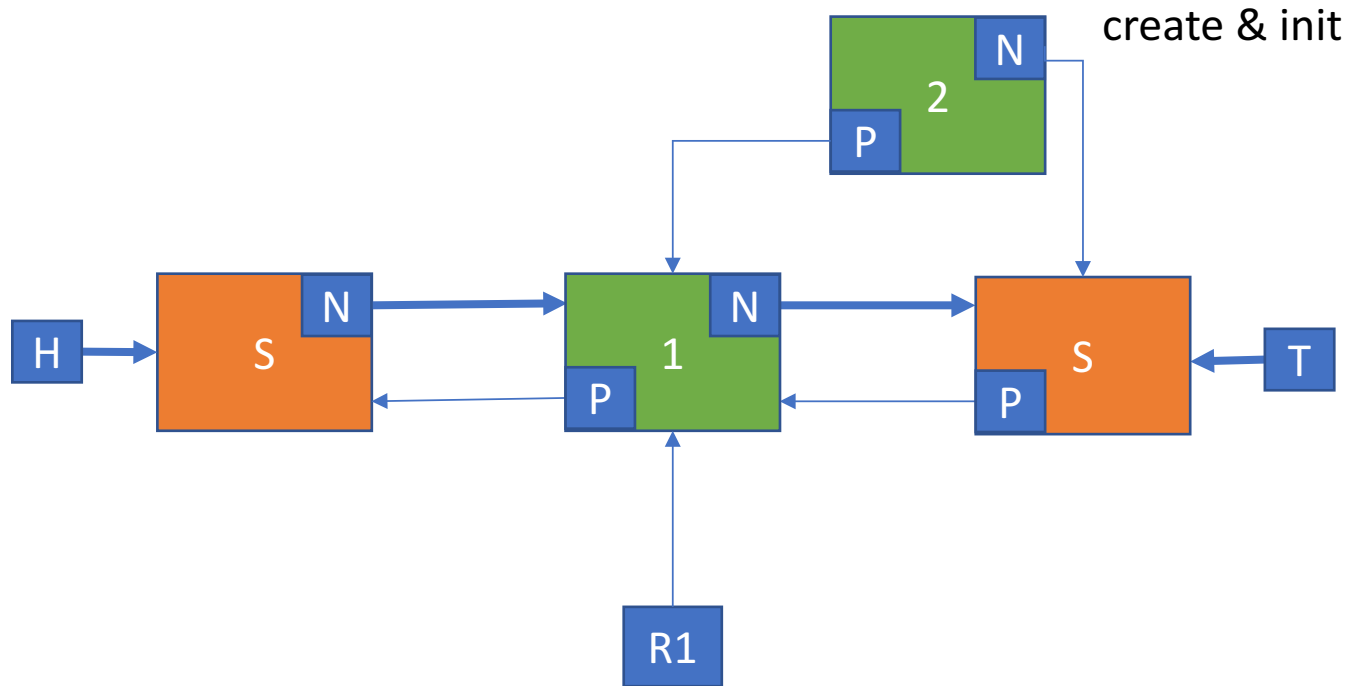# Concurrent remove & insert

When remove wins

# Concurrent remove & insert (0)

R1

I2

create & init

N

2

P

N

S

H

S

N

1

P

P

S

T

R1

# Concurrent remove & insert (1)

# Concurrent remove & insert (2)

| R1 |
|----|
| I2 |



CAS fail

# Concurrent remove & insert (3)

R1

I2

N

2

P

H

N

S

N

1

P

S

P

T

R1

detect wrong prev
(t.prev.next -- removed)
do "correct prev"

# Concurrent remove & insert (4)

R1
I2

fixup next

mark prev

R1

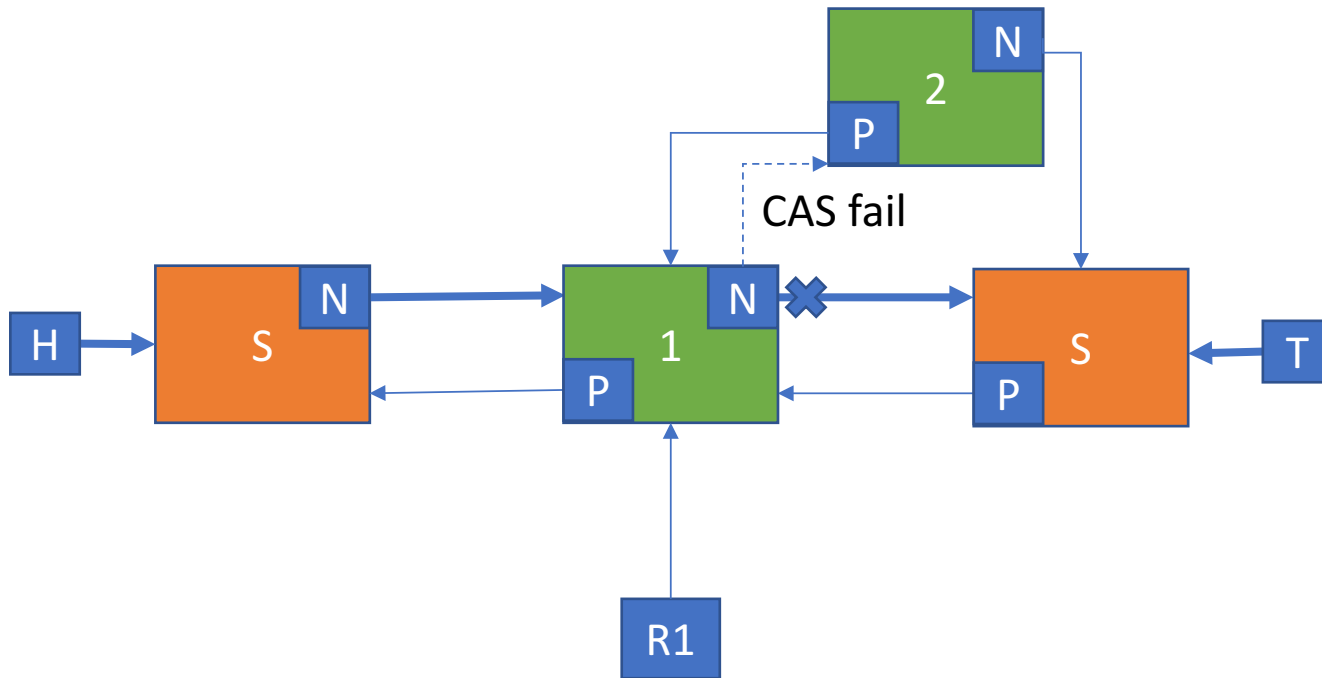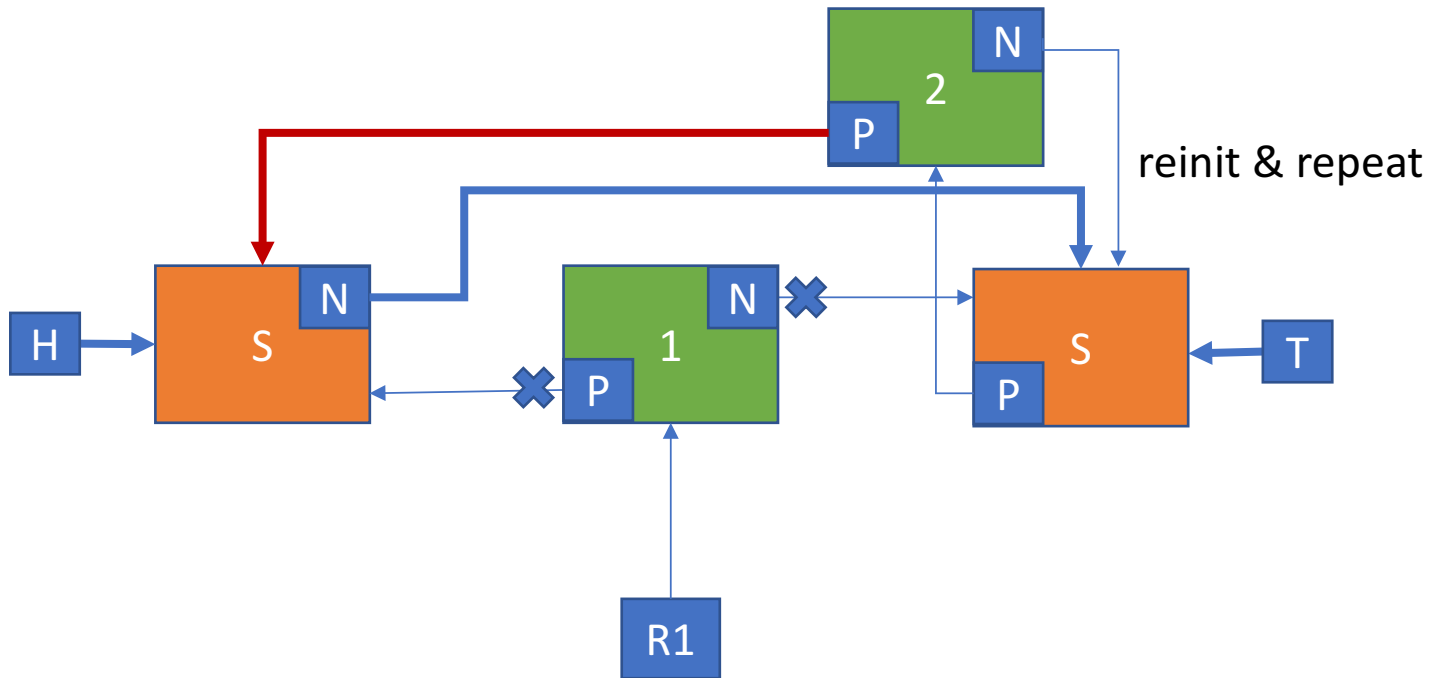# Concurrent remove & insert (5)



update prev

# Concurrent remove & insert (6)

# Concurrent remove & insert

When insert wins

# Concurrent remove & insert (0)

| R1 |
| --- |
| I2 |

create & init

# Concurrent remove & insert (1)

# Concurrent remove & insert (2)

| R1 |
|----|
| I2 |

will succeed marking on remove retry

# Concurrent remove & insert (3)

# Concurrent remove & insert (4)

R1

I2



2

N

P

correct prev

N

S

H

1

N

P

S

P

T

R1

Remove is over!

# Concurrent remove & insert (5)



correct prev

# Takeaways

- A kind of algo you need a paper for
- Hard to improve w/o writing another paper
- **Good news:** stress tests uncover most impl bugs
- **Bad news:** when stress test fails, you up to long hours
- **More bad news:** hard to find bugs that violate lock-freedomness of algorithm

# Summary: what we can do

- Insert items (at the end of the queue)
- Remove items (at the front of the queue)
- Traverse the list
- Remove items at arbitrary locations
  - In O(1)

# Linearizability

- Insert last
  - Linearizes at CAS of **next**

- Remove first / arbitrary
  - Success – at CAS of **next**
  - Fail – at read of **head.next**
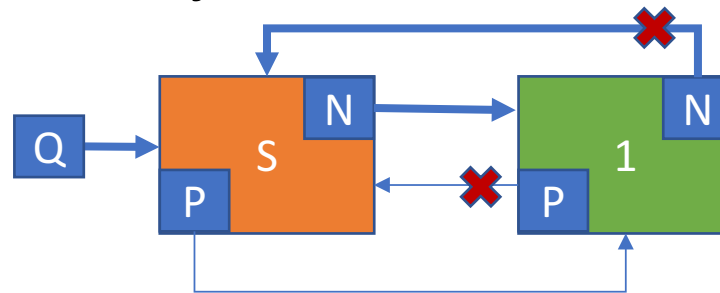
# More about algorithm

- Sundell & Tsigas algo supports *deque* operations
  - Can PushLeft & PopRight
- PopLeft is simple – read **head.next** & remove
- But cannot linearize them all at cas points
  - PushLeft, PushRight, PopRight -  Ok
  - PopLeft linearizes at **head.next** read (!!!)

# Summary of impl notes

- Use GC (drop all memory management details)
- Merge head & tail into a single sentinel node
  - Empty list is just one object (prev & next onto itself)
  - One item += one object



- Reuse "remove mark" objects
  - One-element lists reuse of ptrs to sentinel all the time
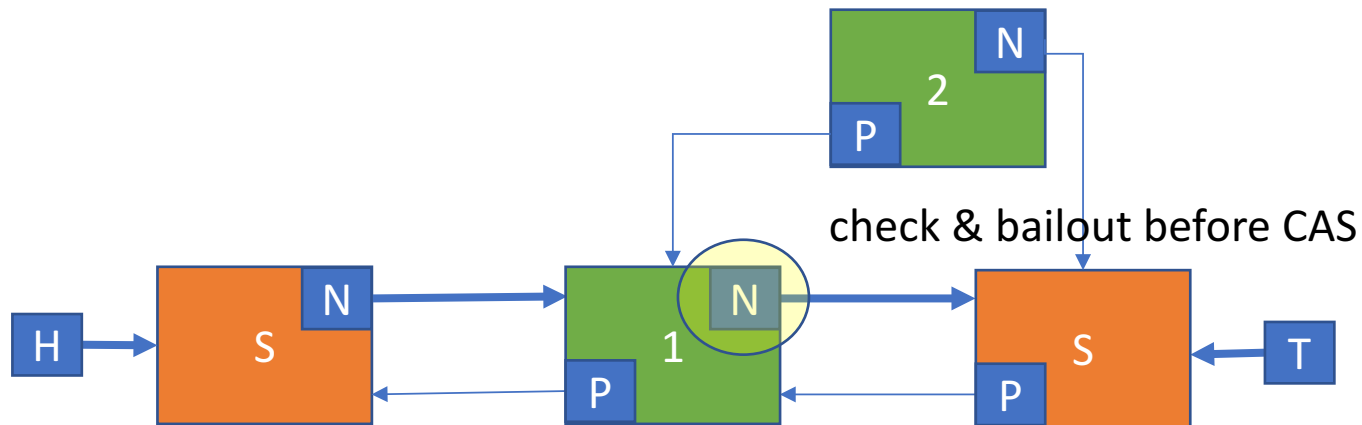- Encapsulate!

# Mods

More complex *atomic* operations

# Basic mods (1)

- Insert item conditionally on prev tail value



check & bailout before CAS

# Basic mods (2)

- Remove head conditionally on prev head value



check & bailout before CAS

# Practical use-case: synchronous channels

**(1)**
```
val channel = Channel<Int>()
```

**(2)**
```
// coroutine #1
for (x in 1..5) {
    channel.send(x * x)
}
```

**(3)**
```
// coroutine #2
repeat(5) {
    println(channel.receive())
}
```

# Senders wait

Incoming receivers

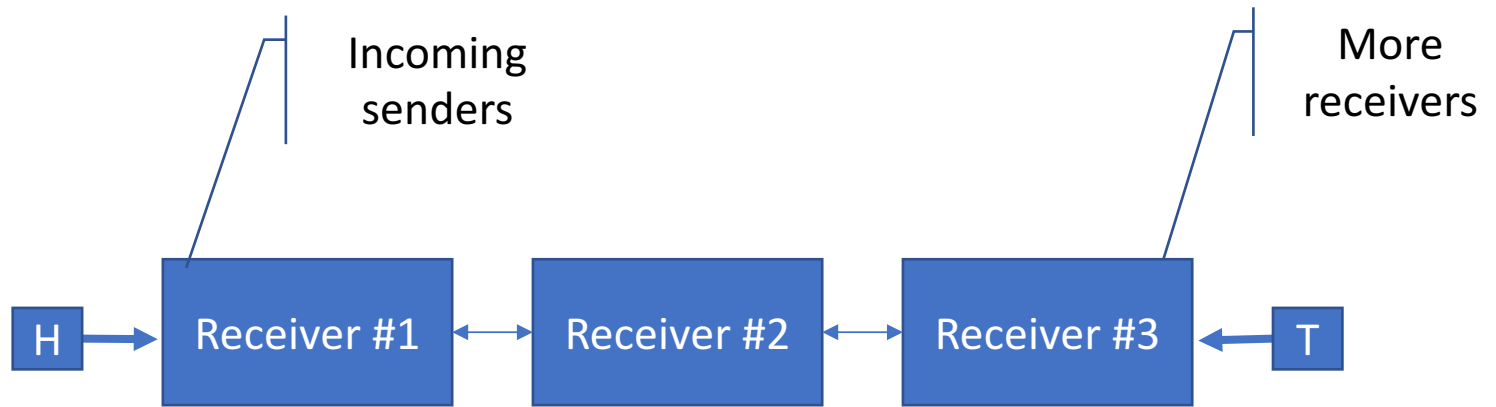More senders

```
H → Sender #1 ↔ Sender #2 ↔ Sender #3 ← T
```

Receiver removes first if it is a sender node

Sender inserts last if it is not a receiver node

# Receivers wait



Sender removes first if it is a receiver node

Receiver inserts last if it is not a sender node

# Send function sketch

```kotlin
fun send(element: T) {
    while (true) {
        // try to add sender, unless prev is receiver
        if (enqueueSend(element)) break
        // try to remove first receiver
        val receiver = removeFirstReceiver()
        if (receiver != null) {
            receiver.resume(element) // resume receiver
            break
        }
    }
}
```

**(1)** **(2)** **(3)** **(4)**

# Channel use-case recap

- Uses insert/remove ops conditional on tail/head node

- Can abort (cancel) wait to receive/send at any time by using remove
  - Full removal -- no garbage is left

- Pretty efficient in practice
  - One item lists – one "garbage" object

# Multi-word compare and swap (CASN)

Build even bigger atomic operations

# Use-case: select expression

```
val channel1 = Channel<Int>()
val channel2 = Channel<Int>()



select {
    channel1.onReceive { e -> ... }
    channel2.onReceive { e -> ... }
}
```
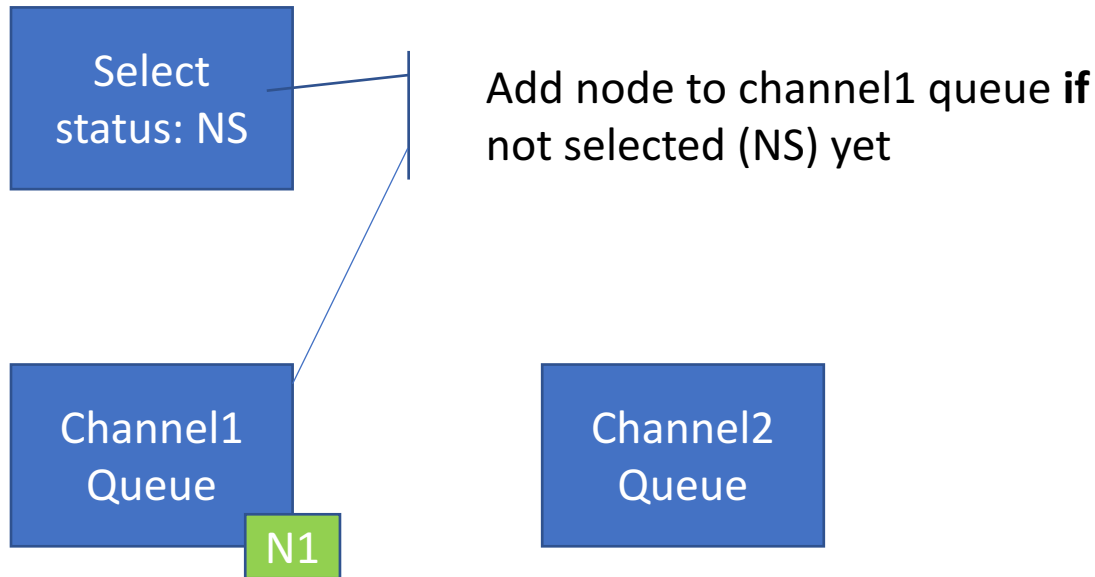
# Impl summary: register (1)
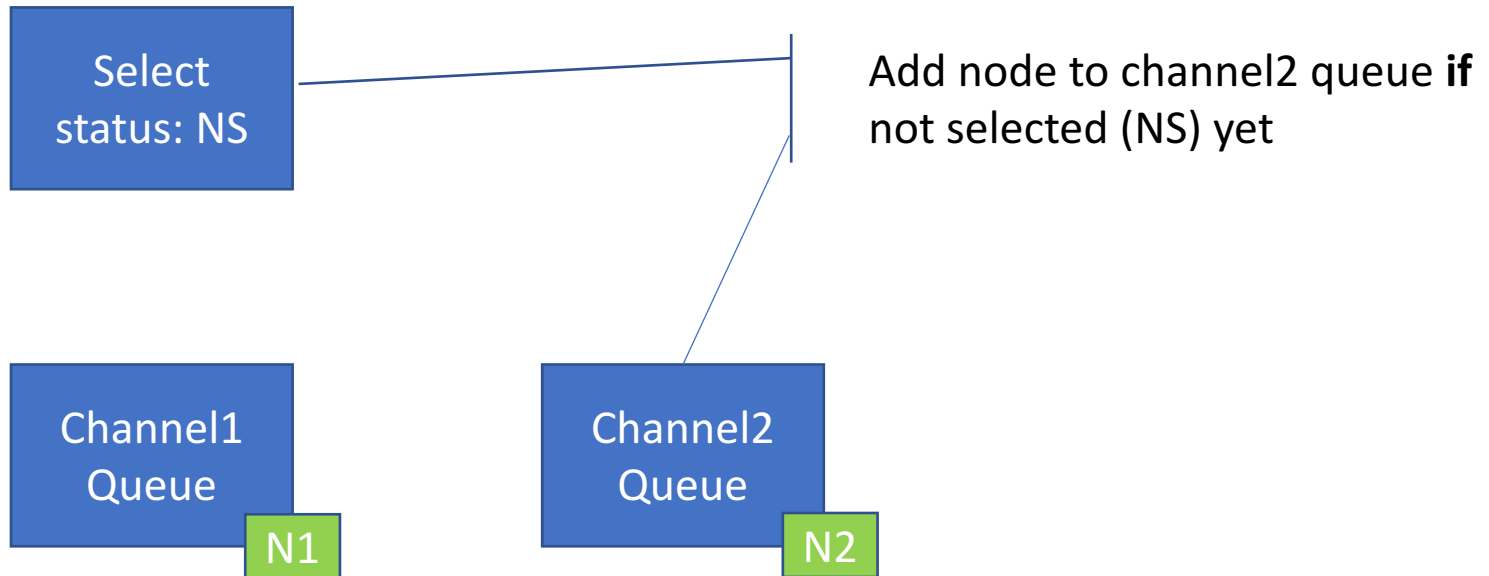
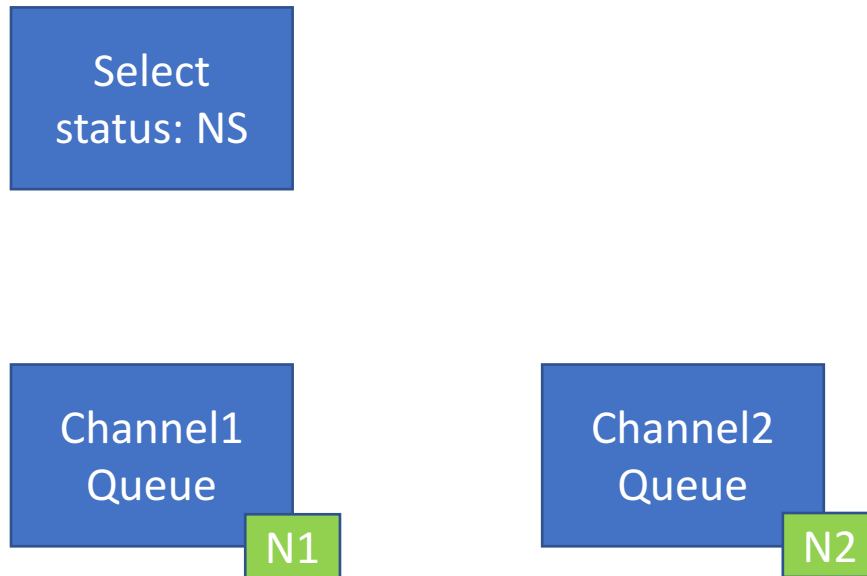**Select status: NS**

1. Not selected
2. Selected

**Channel1 Queue**
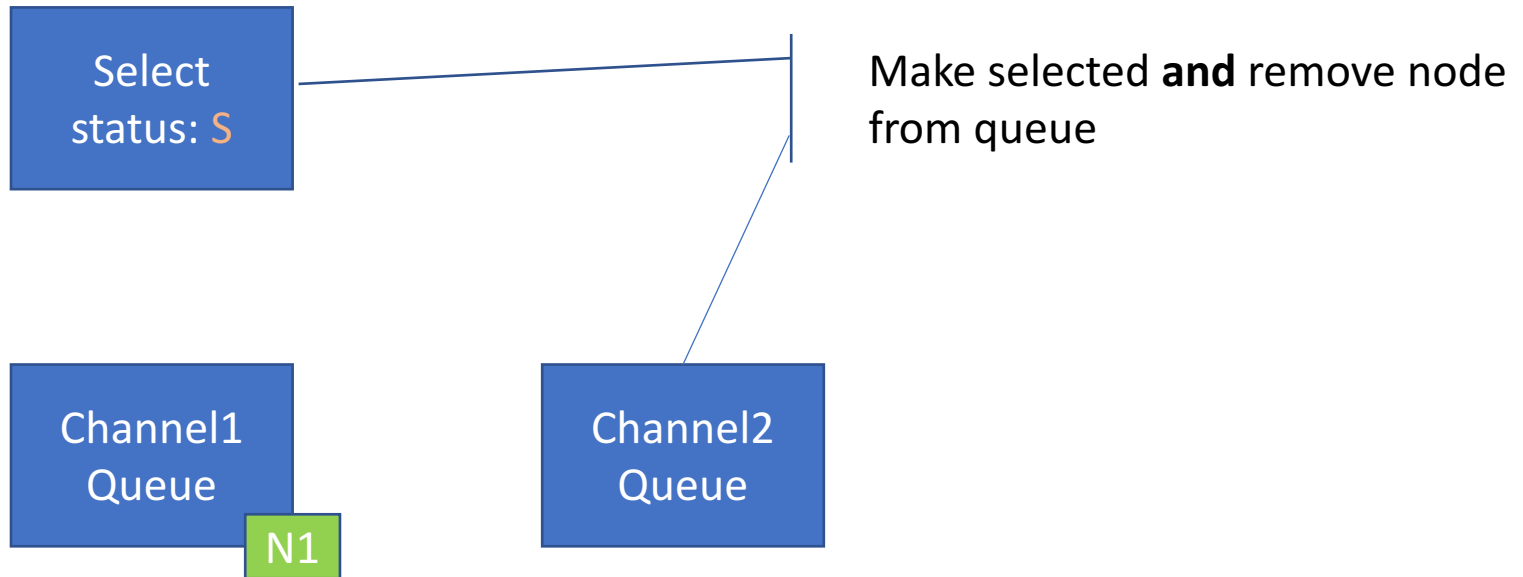
**Channel2 Queue**

# Impl summary: register (2)

Select
status: NS

Add node to channel1 queue **if**
not selected (NS) yet

Channel1
Queue

N1

Channel2
Queue

# Impl summary: register (3)

Select
status: NS

Add node to channel2 queue **if** not selected (NS) yet

Channel1
Queue

N1

Channel2
Queue

N2

# Impl summary: wait

Select
status: NS

Channel1
Queue

N1

Channel2
Queue

N2

# Impl summary: select (resume)

Select
status: S

Make selected **and** remove node
from queue

Channel1
Queue

N1

Channel2
Queue

# Impl summary: clean up rest

Select
status: S

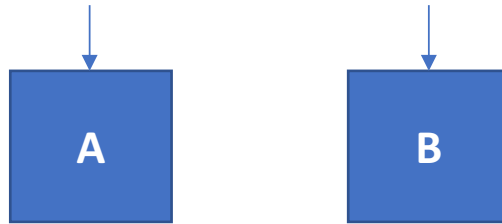Remove non-selected waiters
from queue

Channel1
Queue

Channel2
Queue

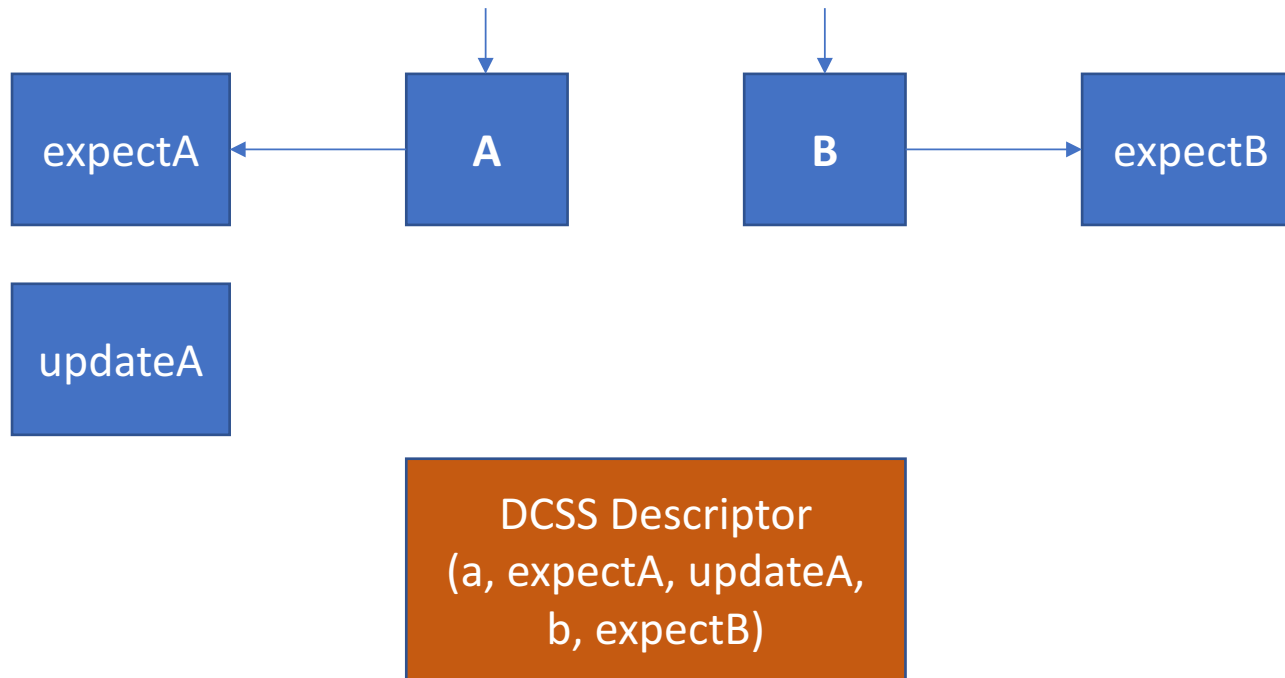# Double-Compare Single-Swap (DCSS)

Building block for CASN

# DCSS spec in pseudo-code
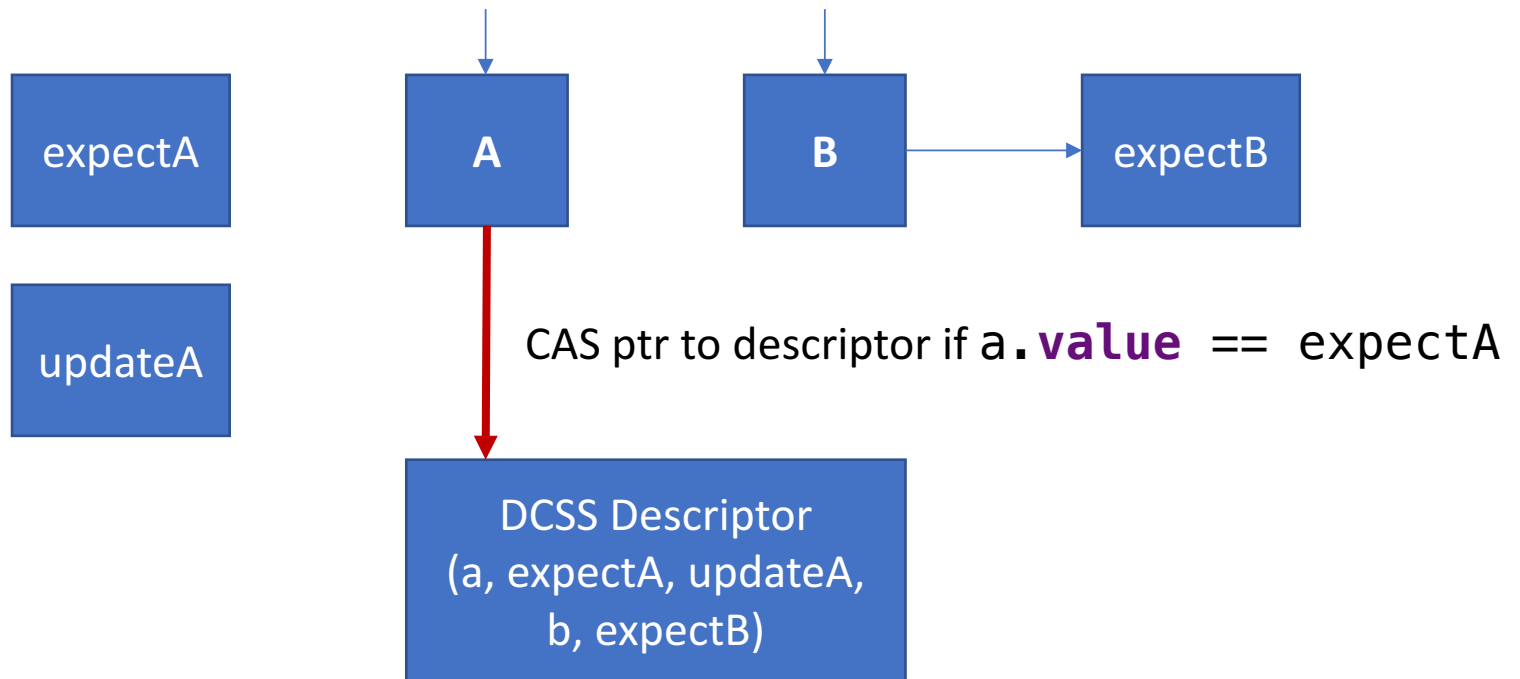


```
fun <A,B> dcss(
    a: Ref<A>, expectA: A, updateA: A,
    b: Ref<B>, expectB: B) =
        atomic {
            if (a.value == expectA && b.value == expectB) {
                a.value = updateA
            }
        }
```
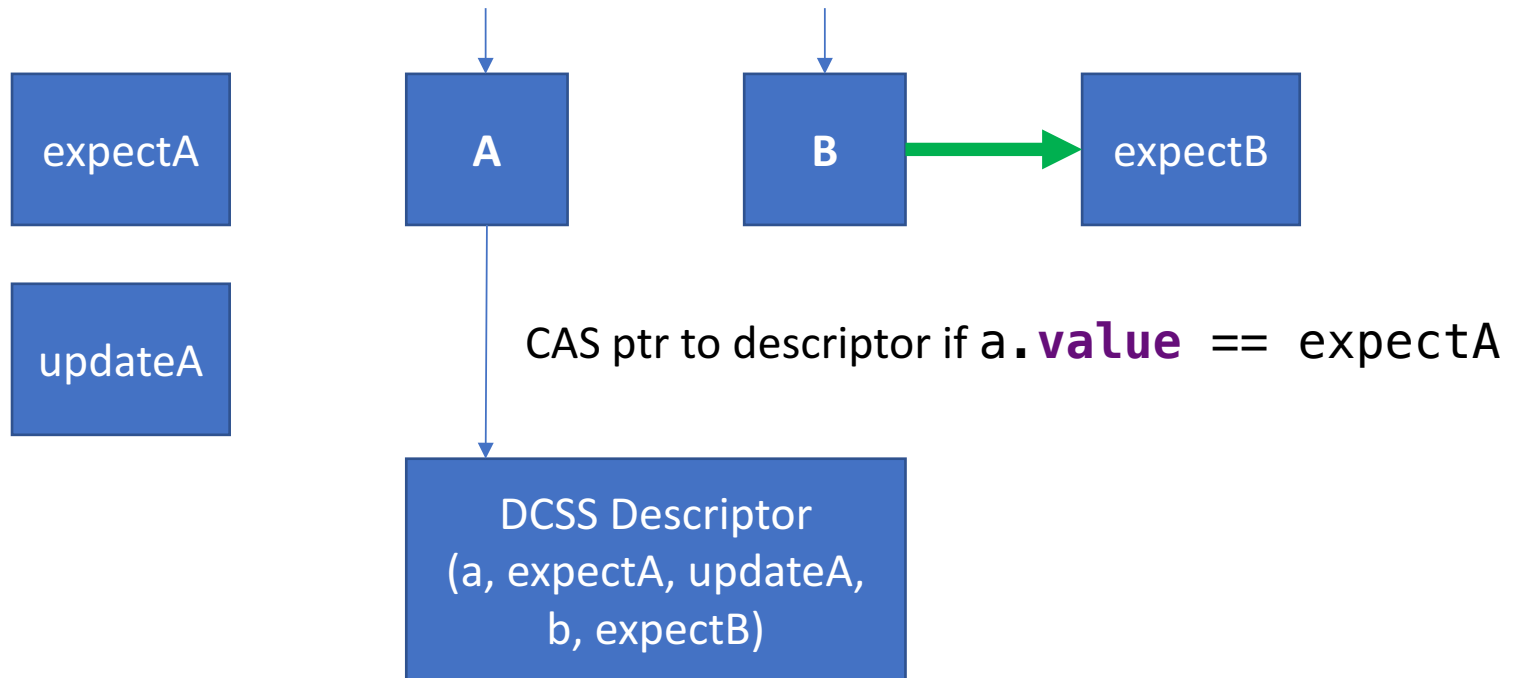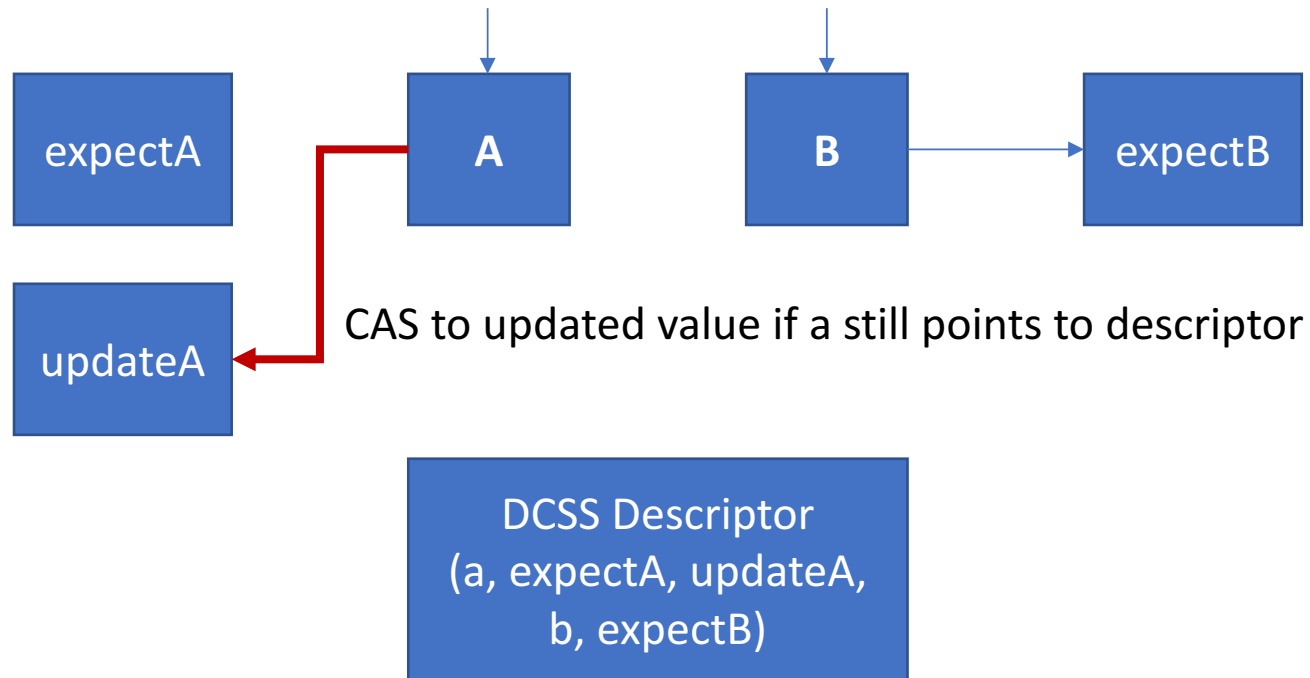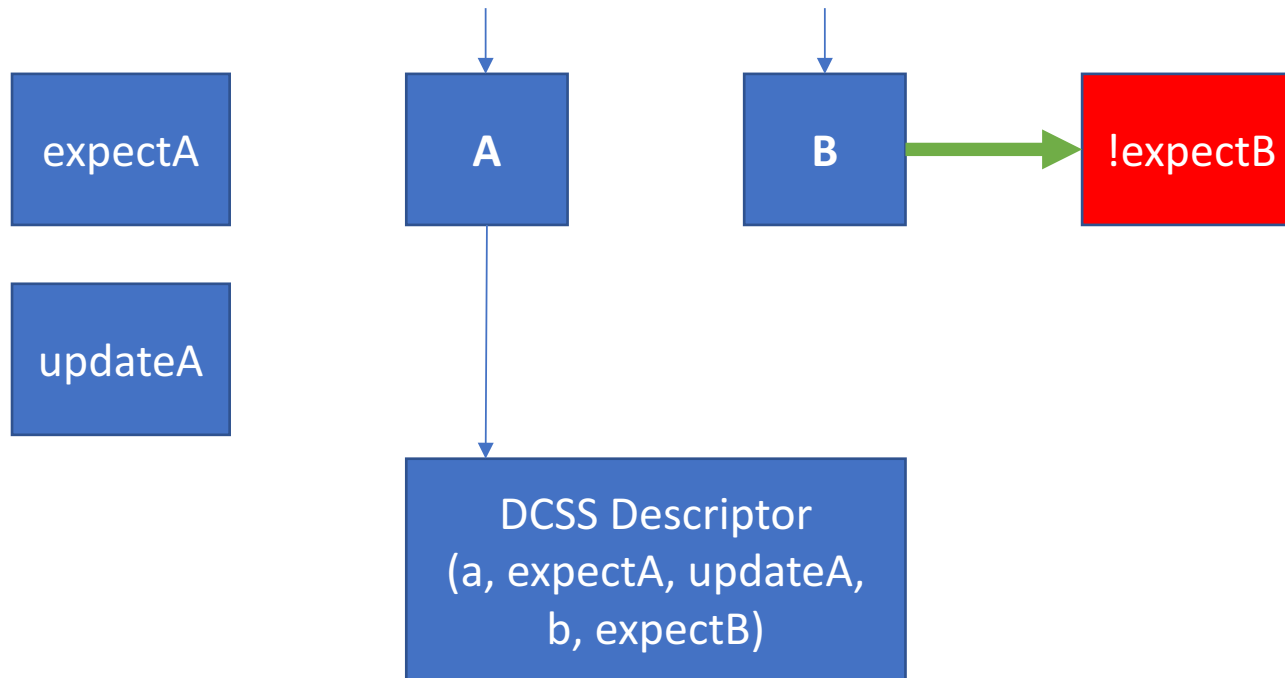
# DCSS: init descriptor
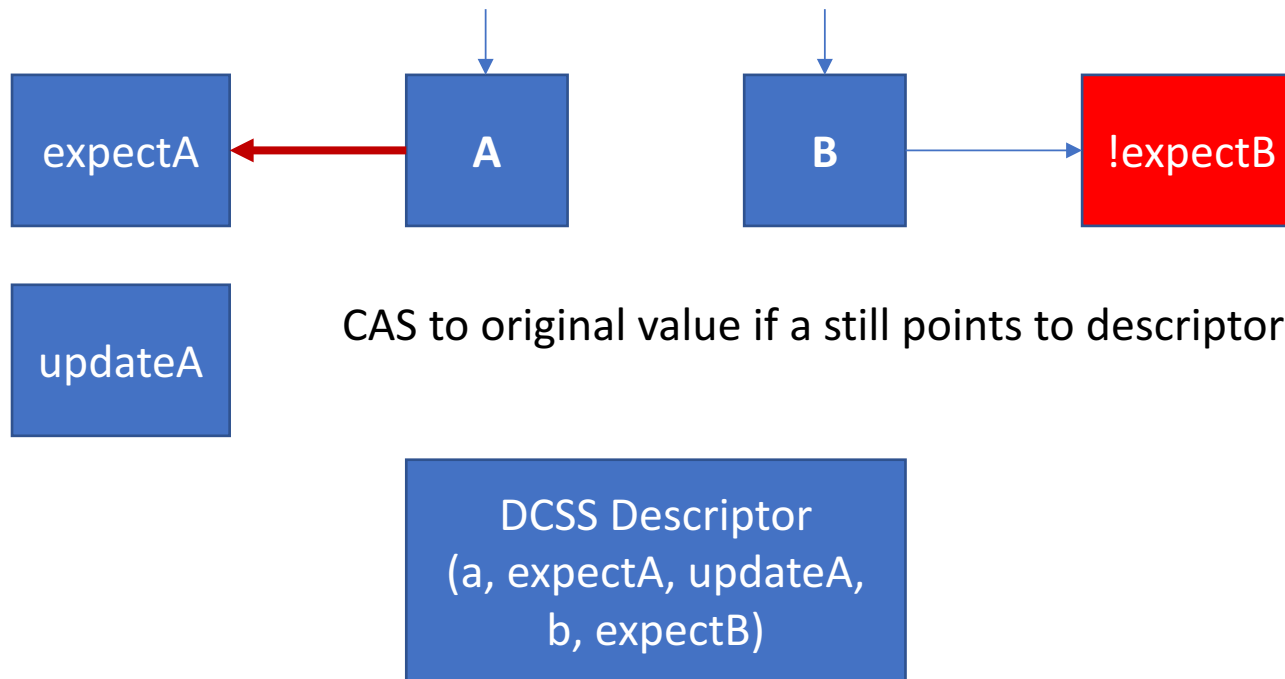
# DCSS: prepare



expectA

updateA

A

B → expectB

CAS ptr to descriptor if `a.value == expectA`

DCSS Descriptor
(a, expectA, updateA,
b, expectB)

# DCSS: read b.value



CAS ptr to descriptor if `a.value == expectA`

Boxes: expectA, updateA, A, B, expectB

DCSS Descriptor
(a, expectA, updateA,
b, expectB)

# DCSS: complete (when success)

expectA | A | B | expectB

CAS to updated value if a still points to descriptor

updateA

DCSS Descriptor
(a, expectA, updateA,
b, expectB)

# DCSS: complete (alternative)

expectA

updateA

A

B

!expectB

DCSS Descriptor
(a, expectA, updateA,
b, expectB)

# DCSS: complete (when fail)



CAS to original value if a still points to descriptor

# DCSS: States

Any other thread encountering descriptor helps complete

**1** Init
A: ???
(desc created)

*prep ok* →

**2** A: desc
A was expectA

*success* →

**3** A: updateA
B was expectB

*prep fail* ↓

**4** A: ???
A was !expectA

*fail* ↓

**5** A: expectA
B was !expectB

Originator cannot learn what was the outcome

one tread

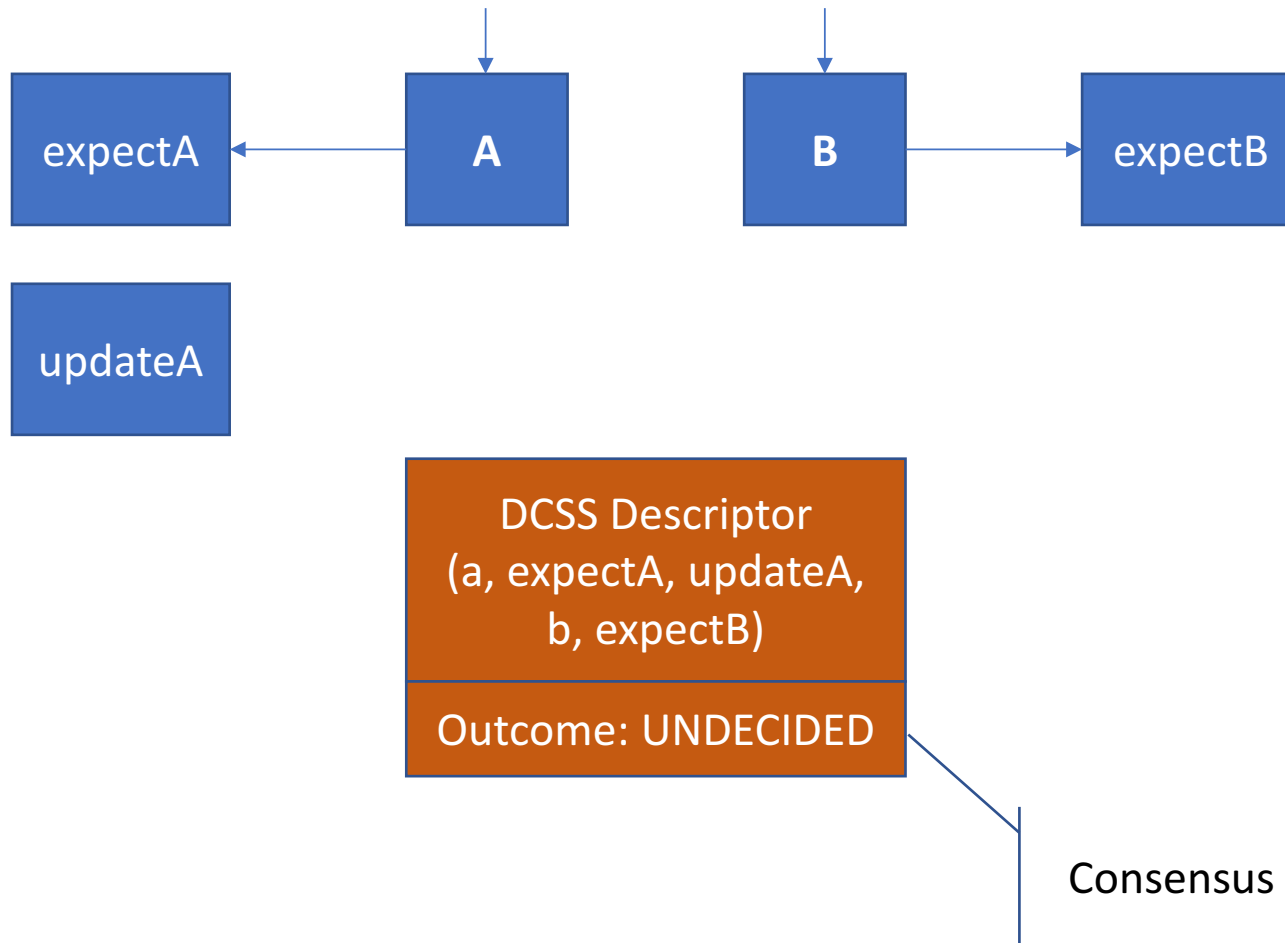Lock-free algorithm without loops!

# Caveats

- A & B locations must be totally ordered
  - or risk stack-overflow while helping
- One way to look at it: Restricted DCSS (RDCSS)

# DCSS Mod: learn outcome
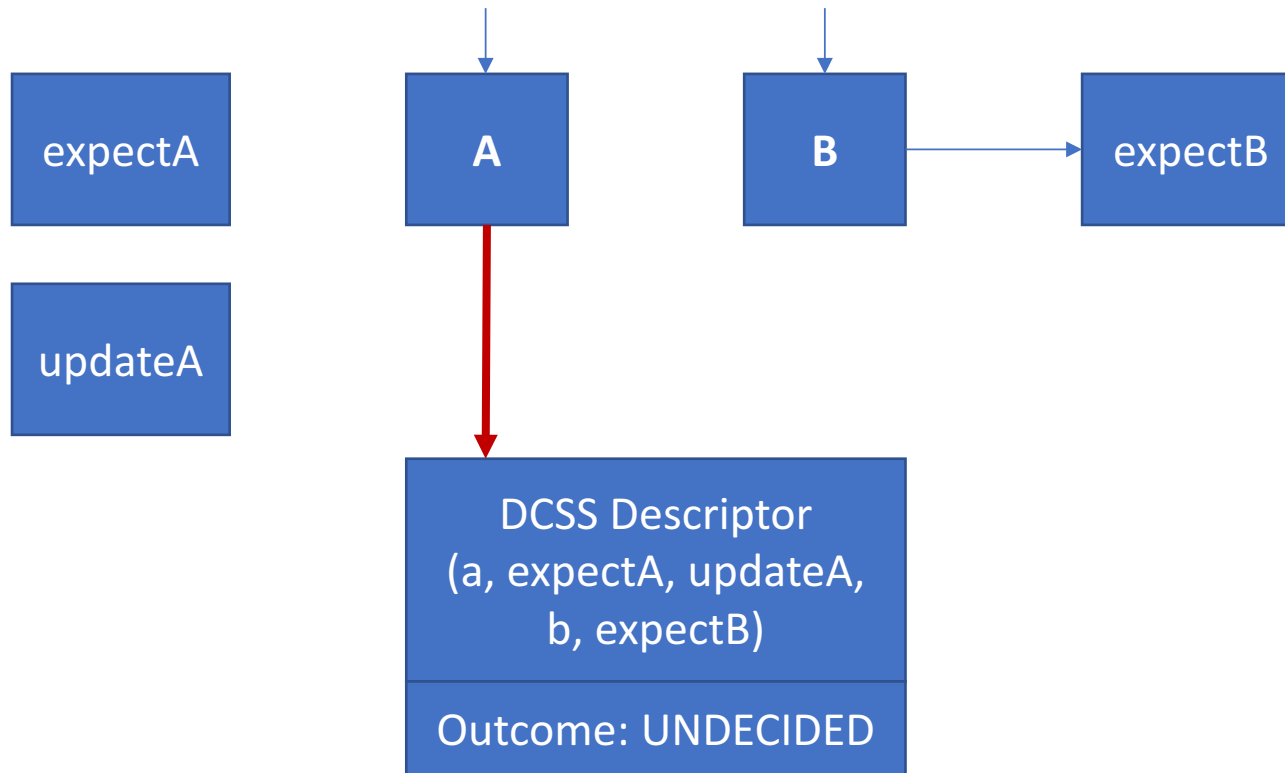


```
fun <A,B> dcssMod(
    a: Ref<A>, expectA: A, updateA: A,
    b: Ref<B>, expectB: B): Boolean =
        atomic {
            if (a.value == expectA && b.value == expectB) {
                a.value = updateA
                true
            } else
                false
        }
```
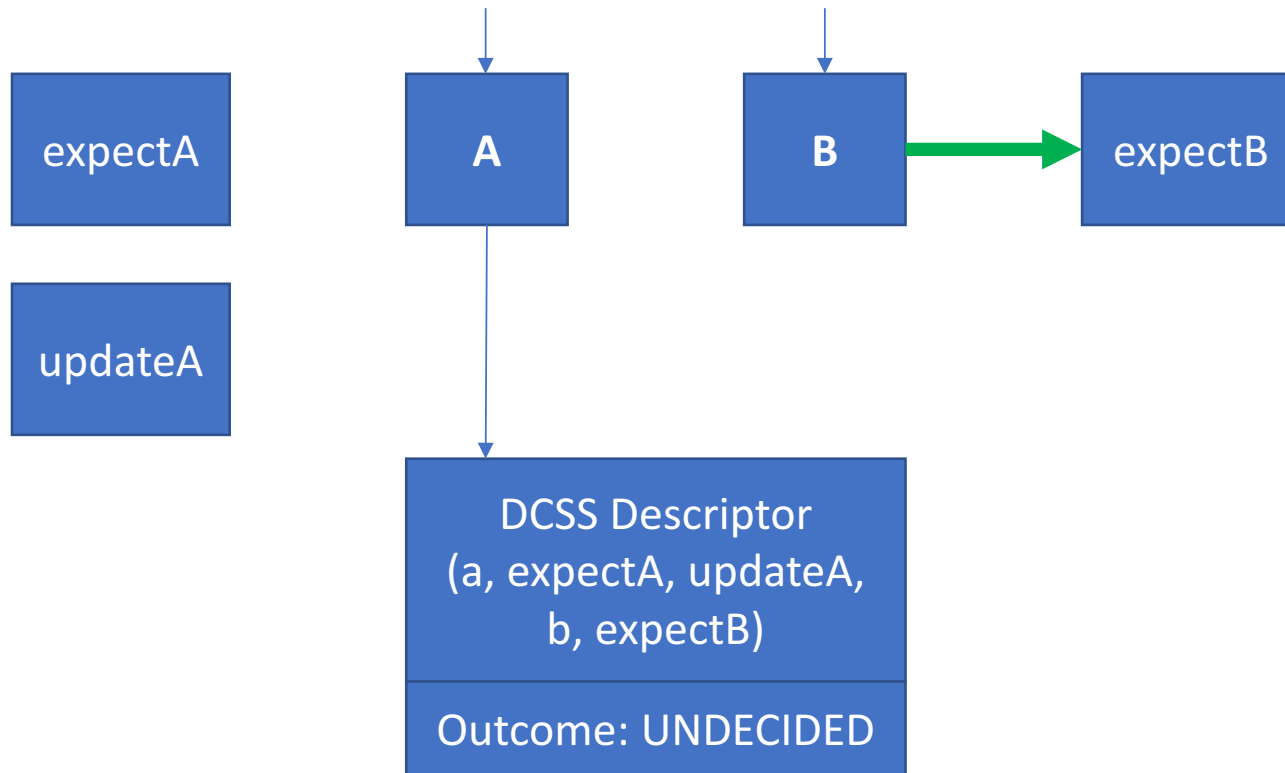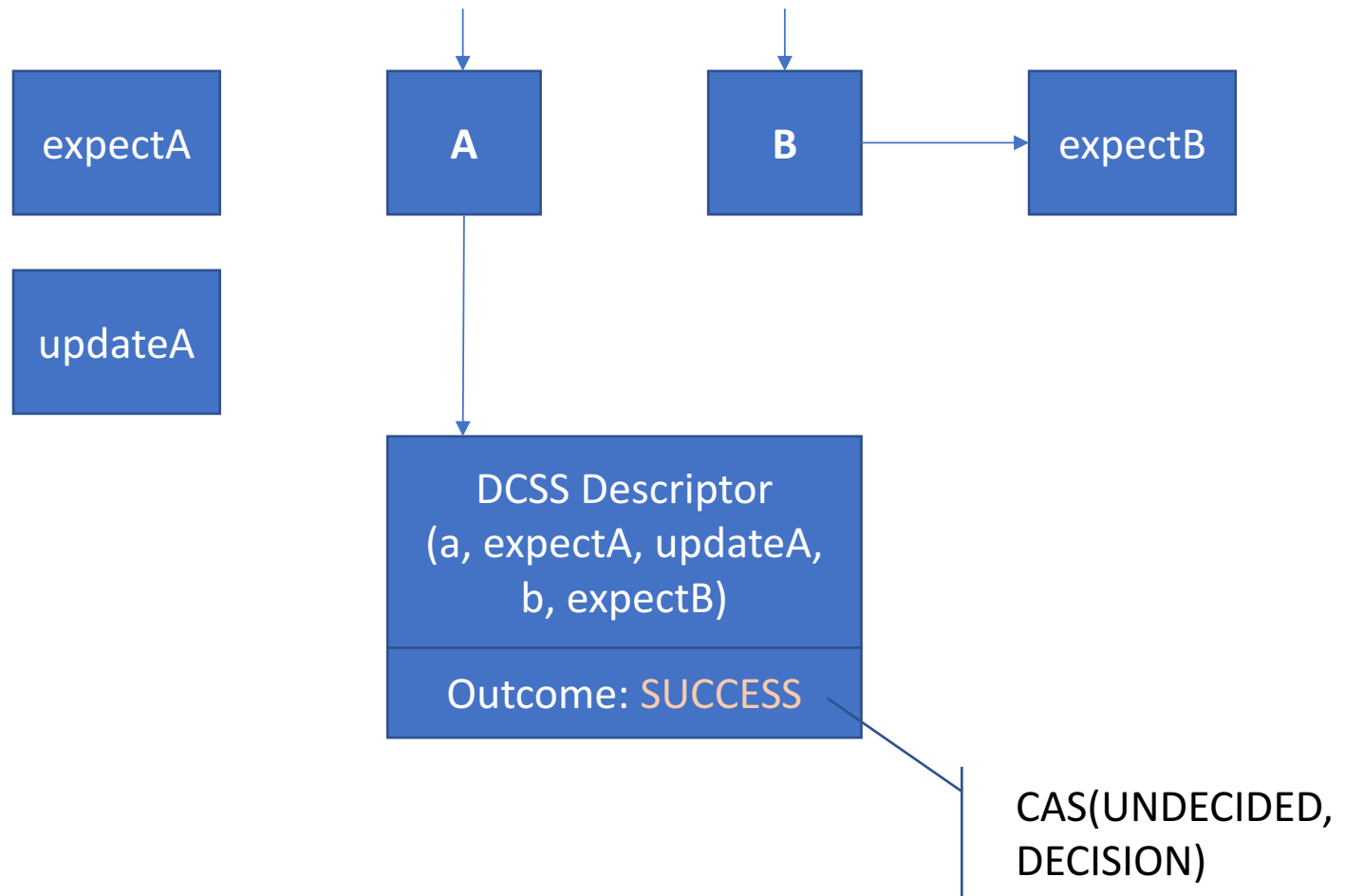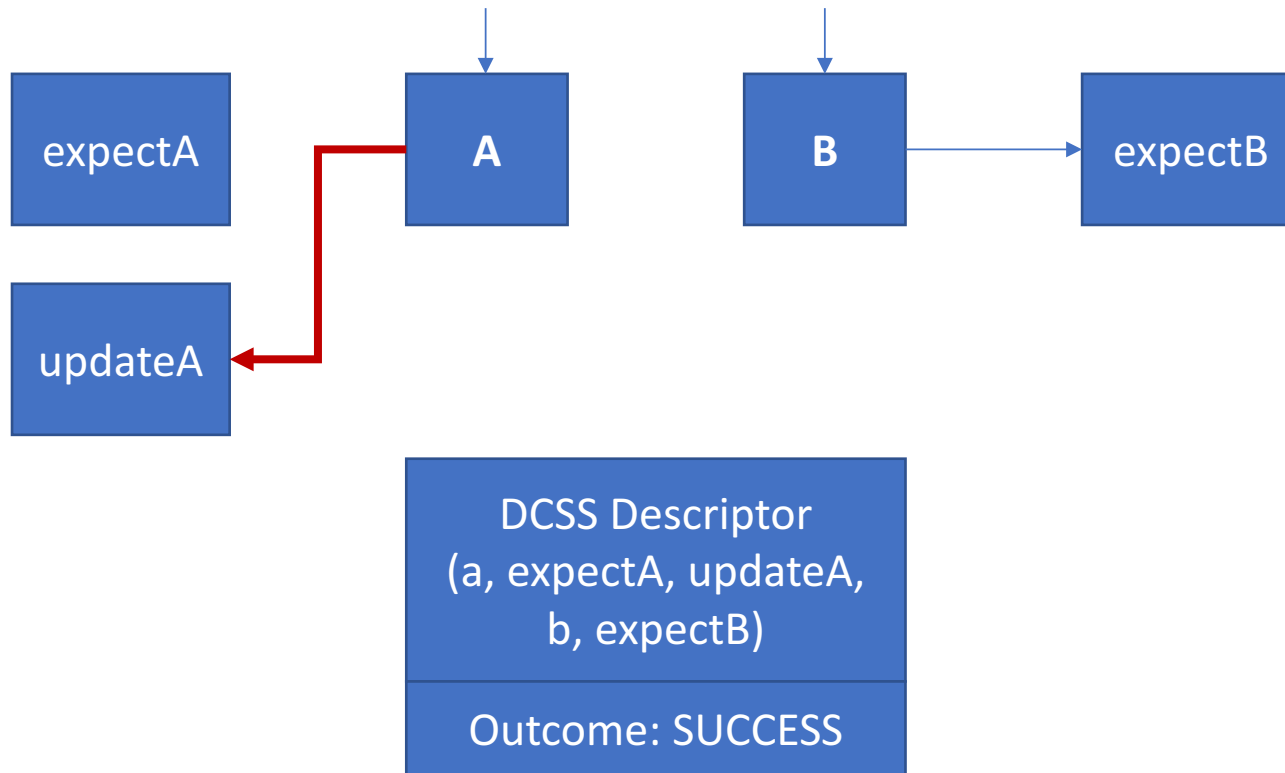
# DCSS Mod: init descriptor

# DCSS Mod: prepare

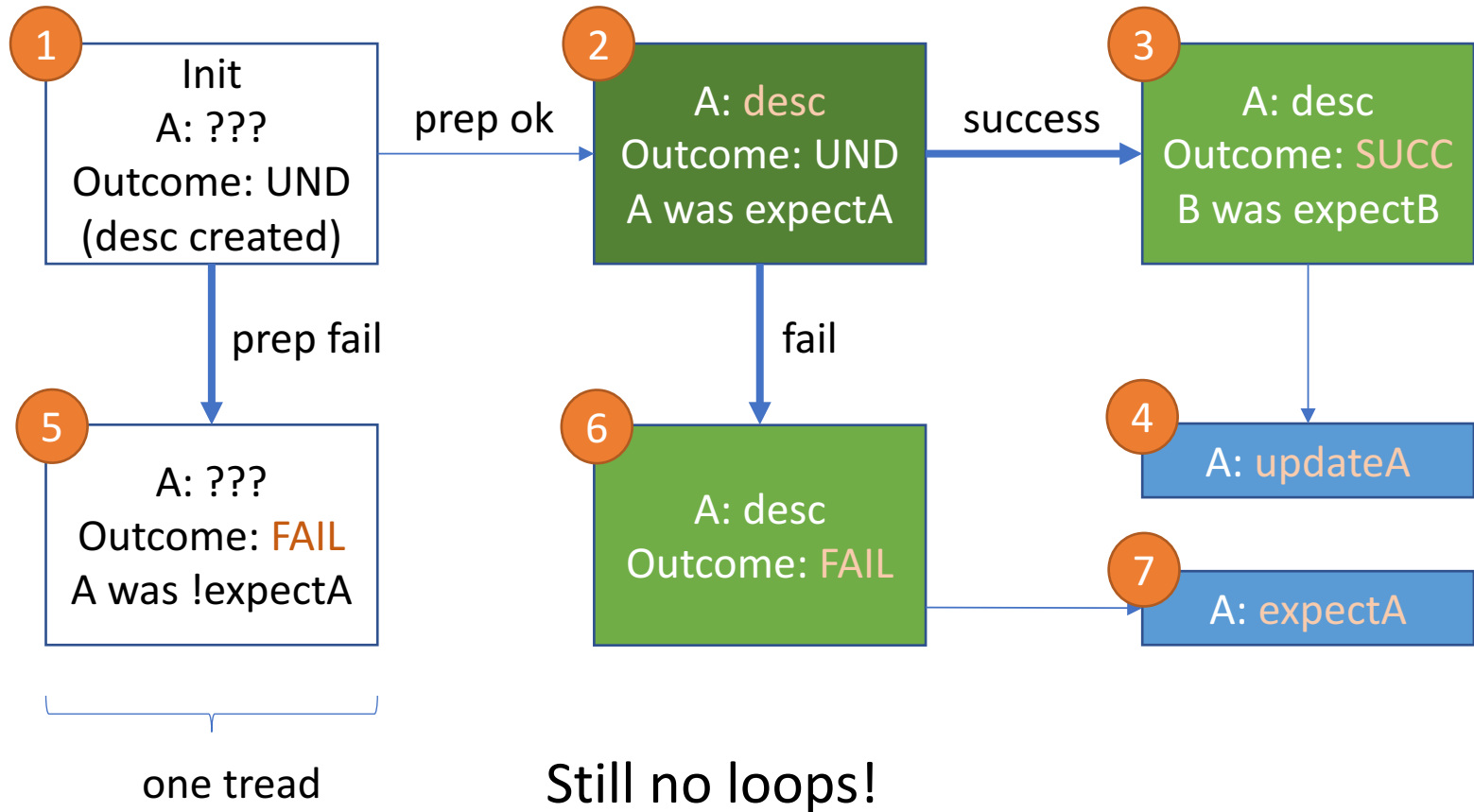# DCSS Mod: read b.value

# DCSS Mod: reach consensus

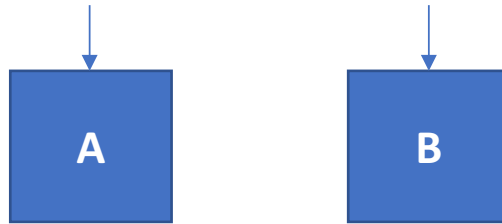# DCSS Mod: complete

# DCSS Mod: States

# Compare-And-Swap N-words (CASN)
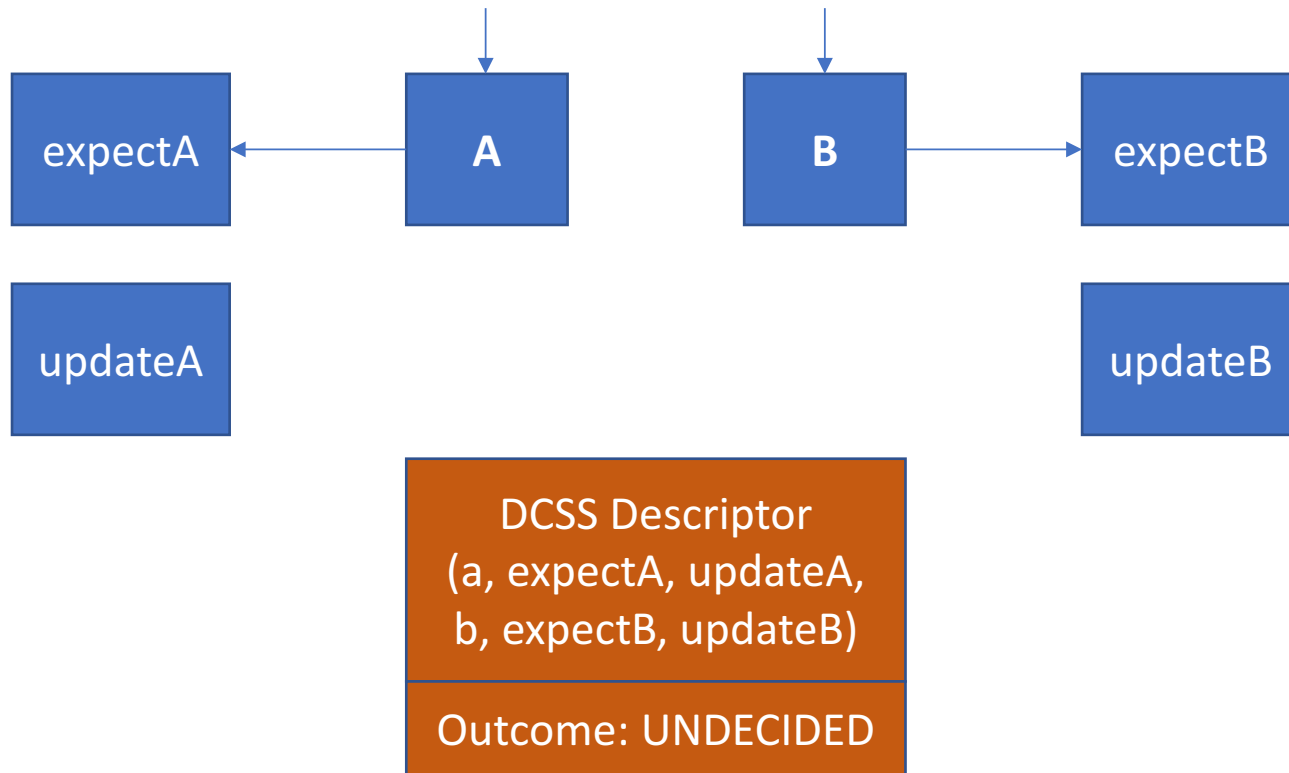
The ultimate atomic update

# CASN spec in pseudo-code

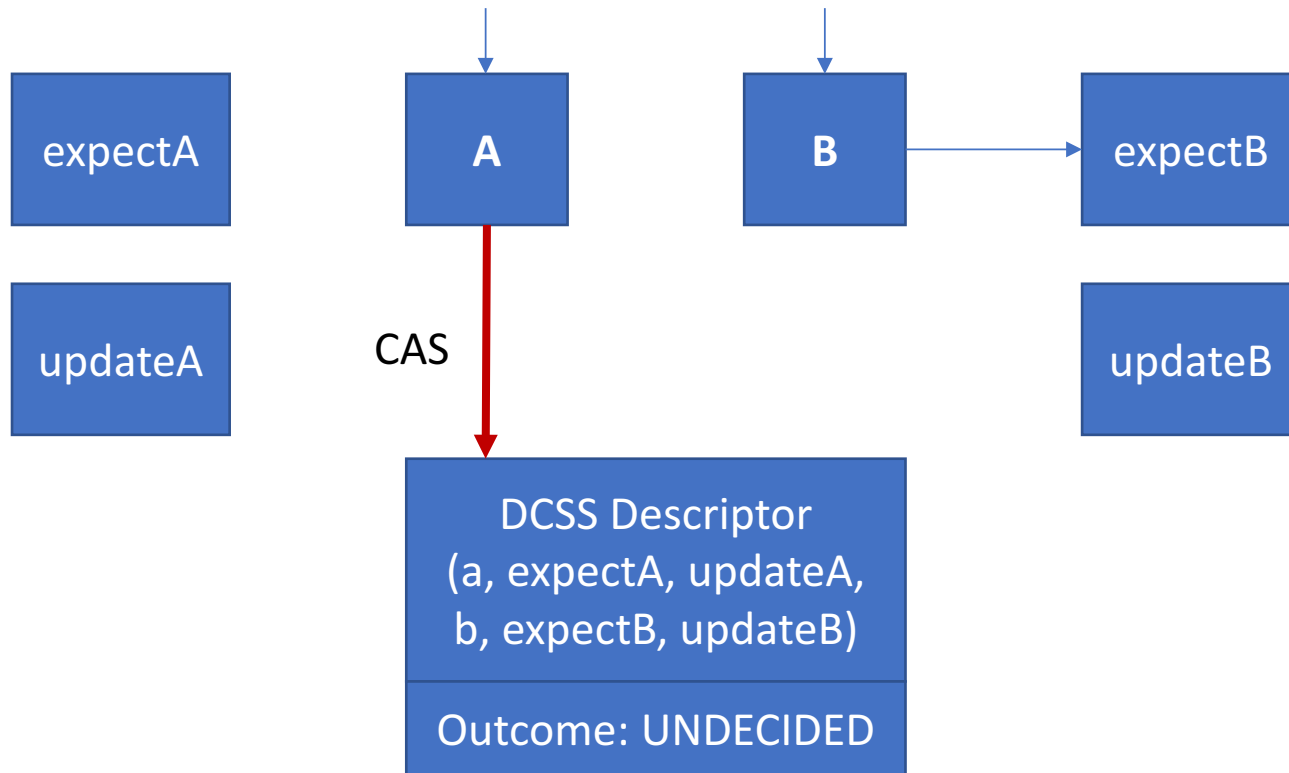For two words, for simplicity



```
fun <A,B> cas2(
    a: Ref<A>, expectA: A, updateA: A,
    b: Ref<B>, expectB: B, updateB: B): Boolean =
        atomic {
            if (a.value == expectA && b.value == expectB) {
                a.value = updateA
                b.value = updateB
                true
            } else
                false
        }
```
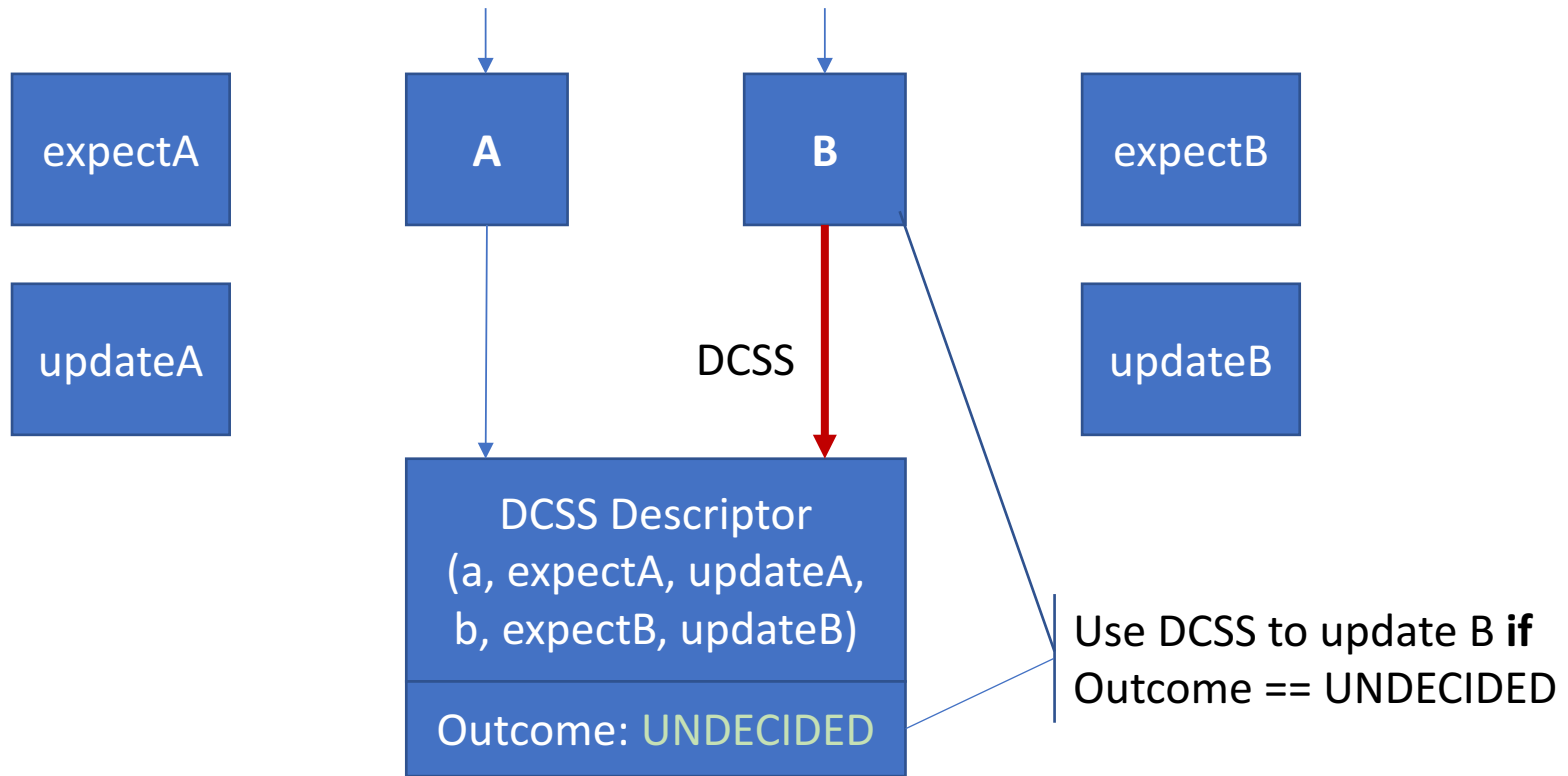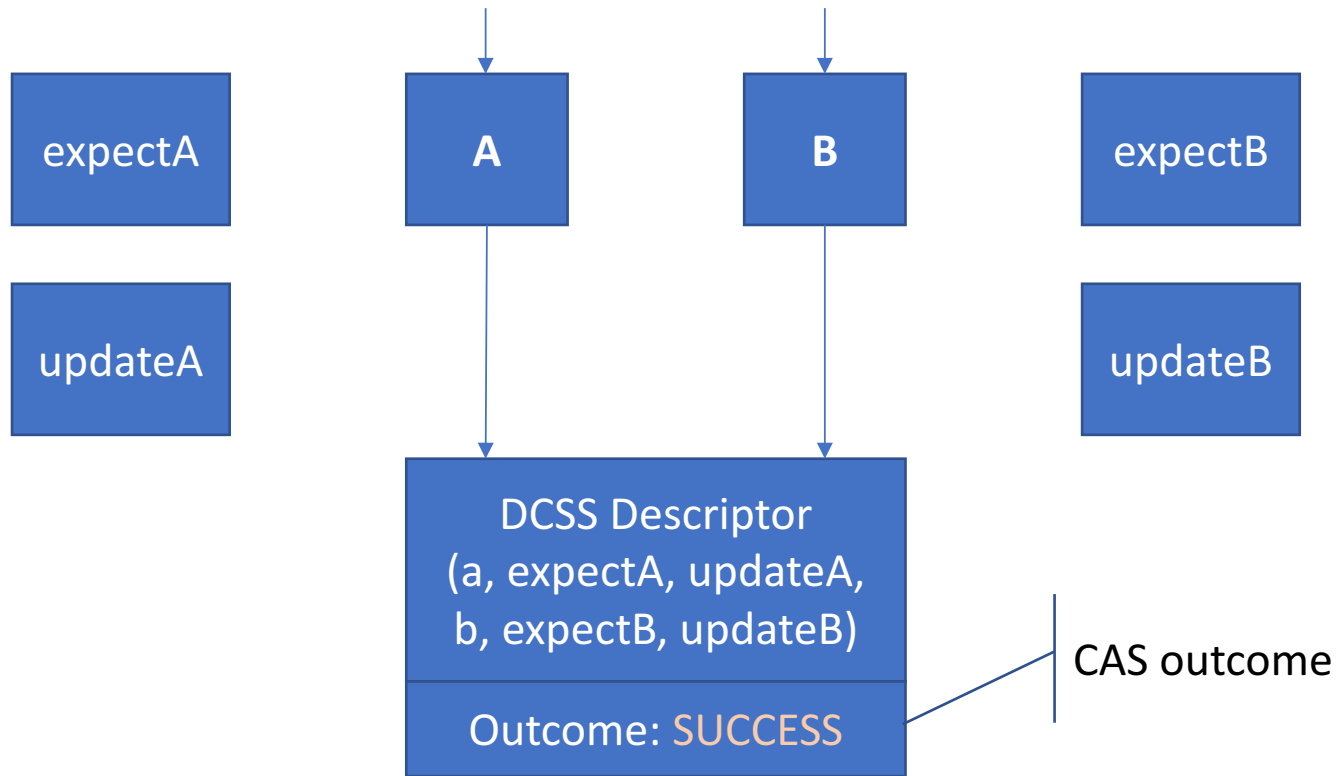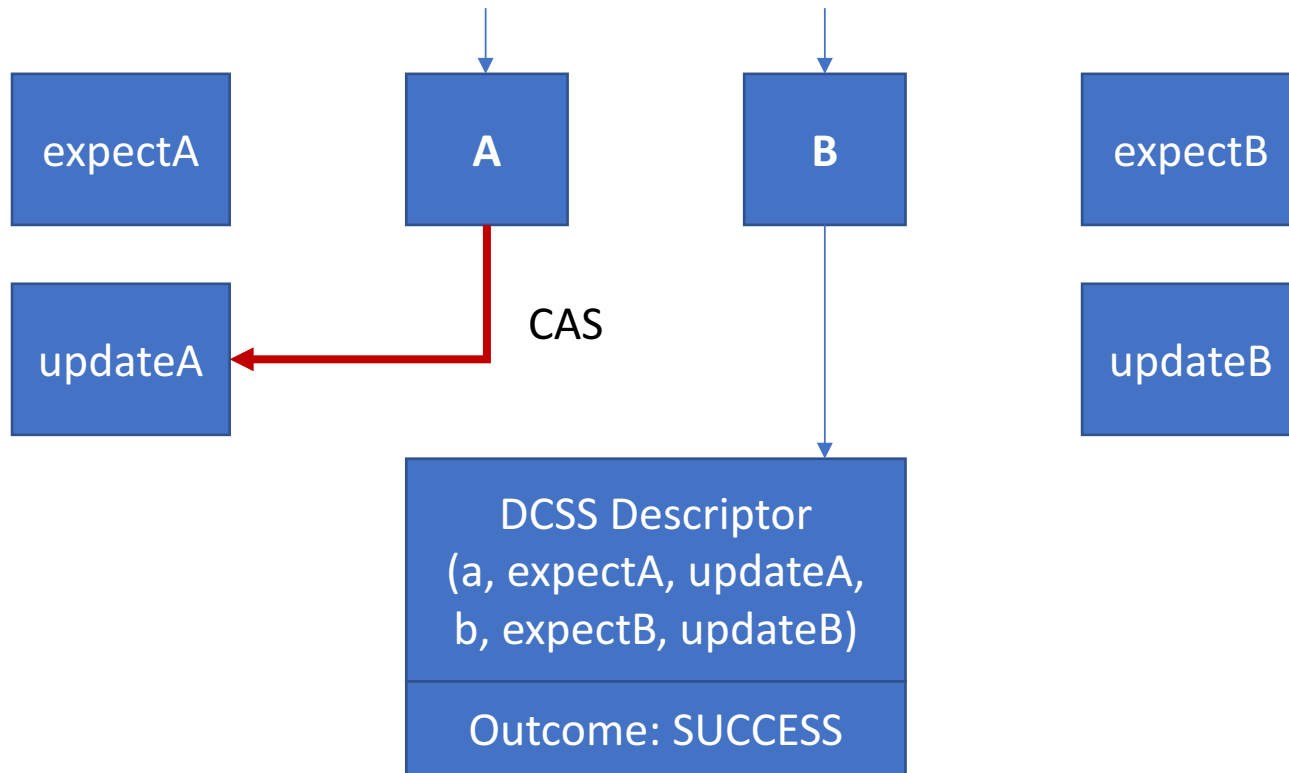
# CASN: init descriptor
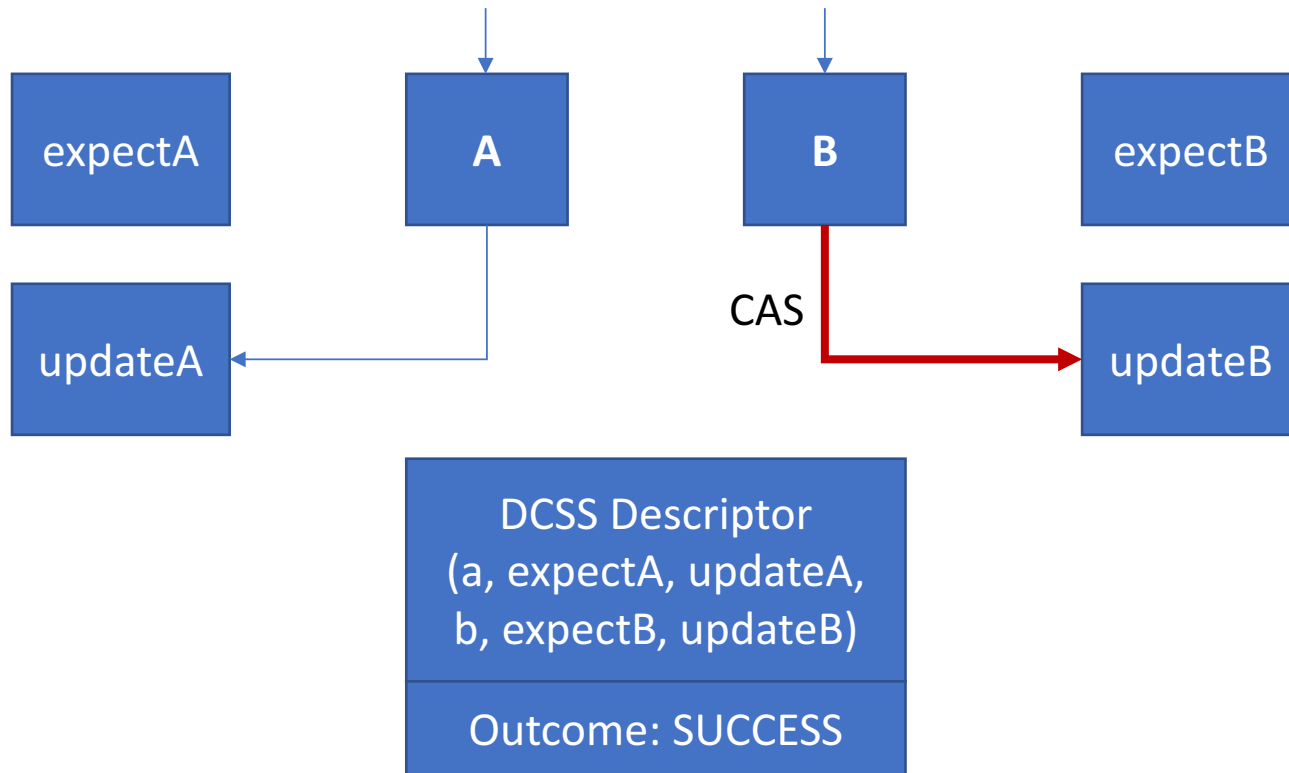
# CASN: prepare (1)

# CASN: prepare (2)

# CASN: decide

# CASN: complete (1)

# CASN: complete (2)

# CASN: States

**1**
A: ???
B: ???
O: UND

**2**
A: desc
B: ???
O: UND

**3**
A: desc
B: desc
O: UND

DCSS

**4**
A: desc
B: desc
O: SUCC

**5**
A: updateA
B: desc
O: SUCC

**6**
Init
A: updateA
B: updateB
O: SUCC

A != expectA

**7**
A: ???
B: ???
O: FAIL

one tread

B != expectB

**8**
A: desc
B: ???
O: FAIL

**9**
A: expectA
B: ???
O: FAIL

Prevents from
going back in this SM

descriptor is known to other (helping) threads
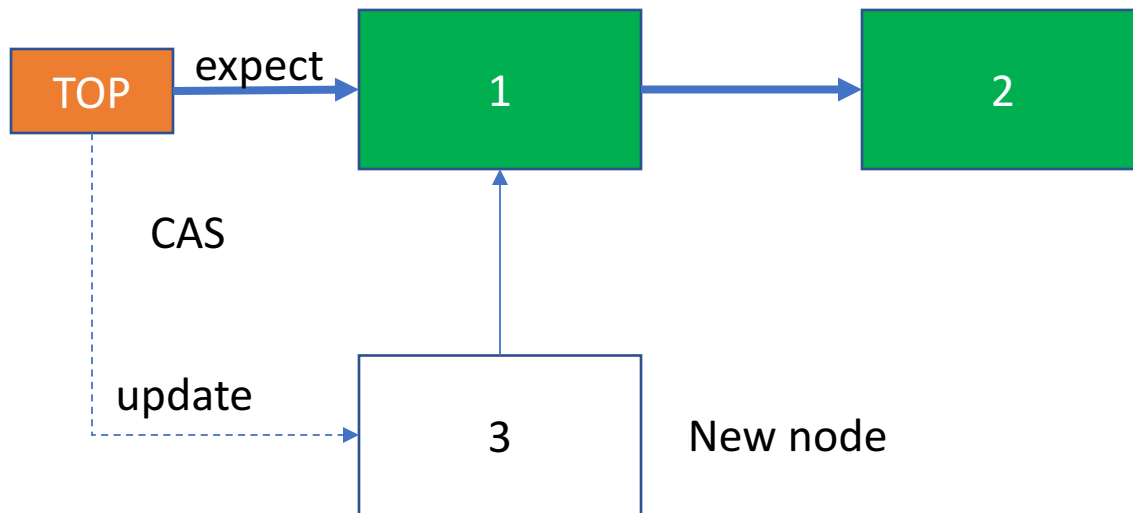
# Using it in practice

All the little things that matter

It is easy to combine multiple operations with DCSS/CASN that linearize on a **CAS** with a descriptor parameters that are **known in advance**
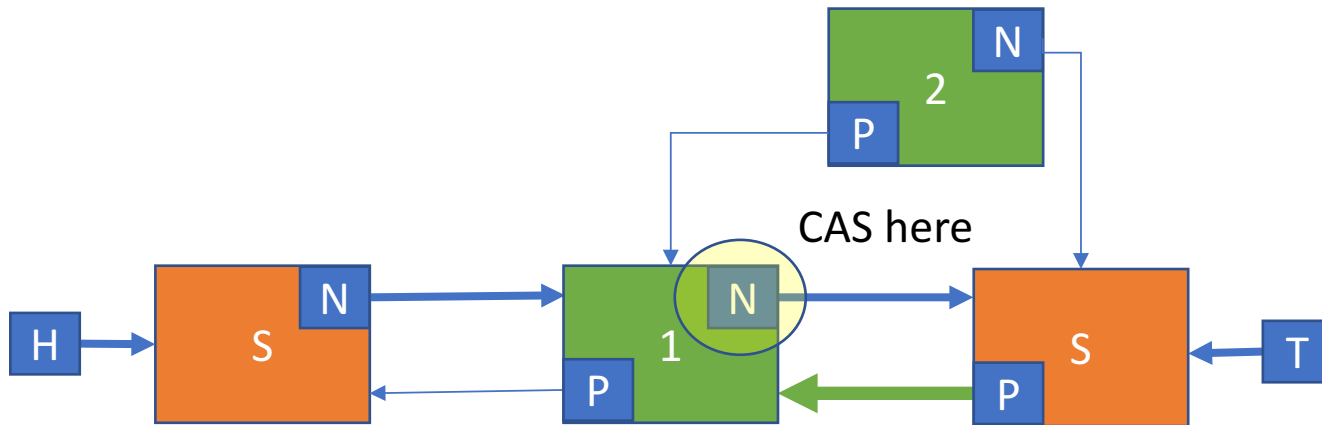
# Trivial example: Treiber stack

# Let's go deeper

Into unpublished territory

# Doubly linked list: insert last

# Doubly linked list: insert (0)



CAS here

DCSS here is needed (always!)

**Operation Descriptor**

A ref: ??? — **??? can fill in A before CAS & update on retry**
expectA: Sentinel — We know expected value for CAS in advance
updateA: Node #2 — We know updated value for CAS in advance
…
Outcome: UNDECIDED

# Doubly linked list: insert (1)

# Doubly linked list: insert (2)

Helpers are a bound to stumble upon the same descriptor

Competing inserts will complete (help) us first

CAS can only succeed on last node

**Operation Descriptor**

A ref: ???
expectA: Sentinel
updateA: Node #2
…
Outcome: UNDECIDED

**DCSS Descriptor**

affected node: #1
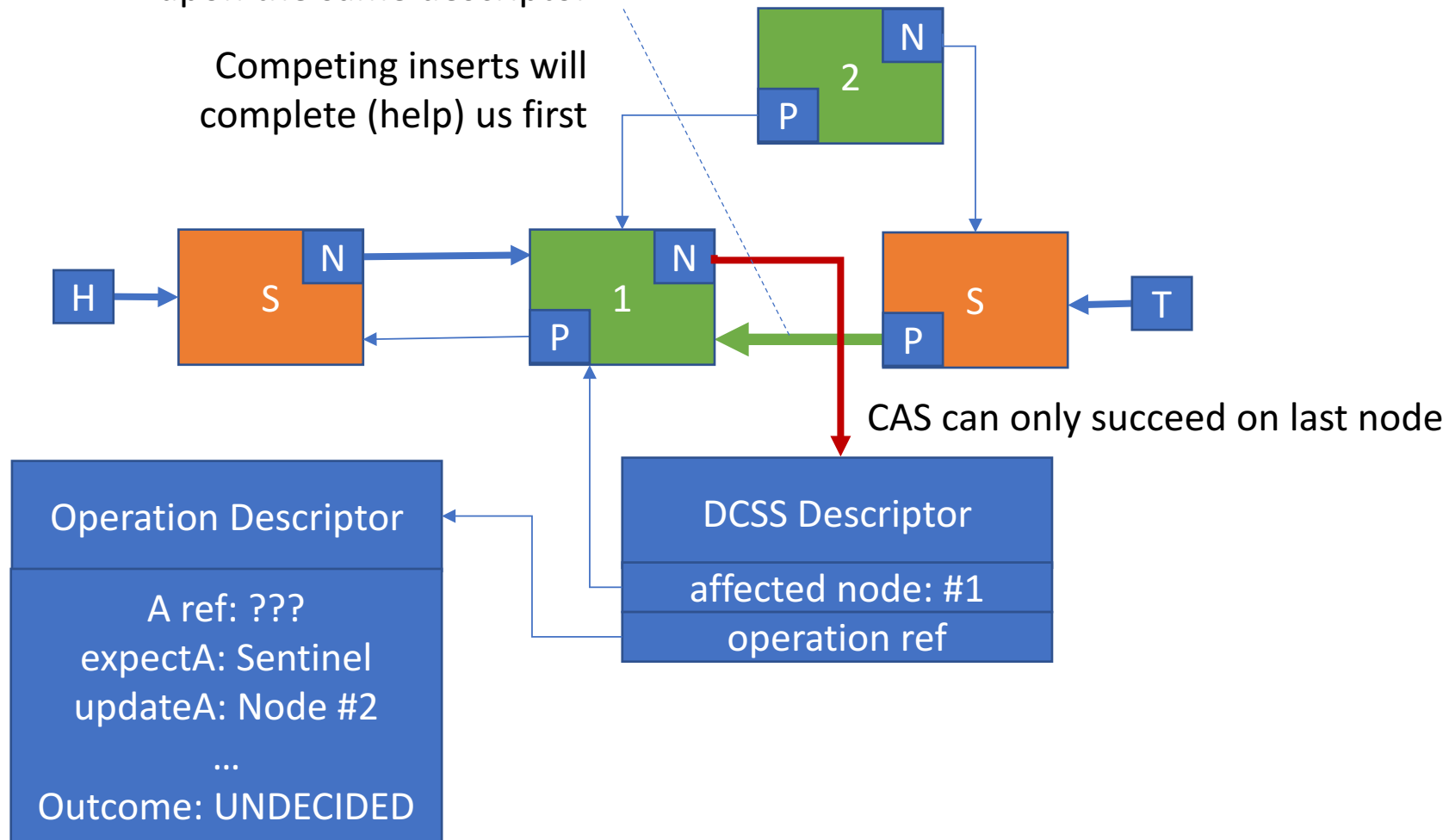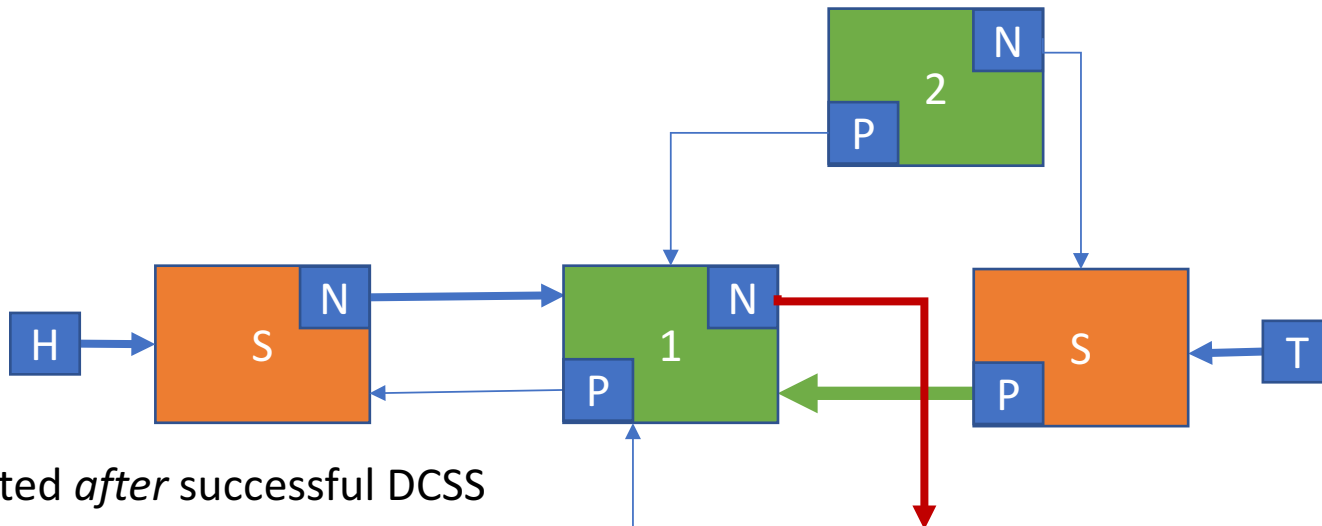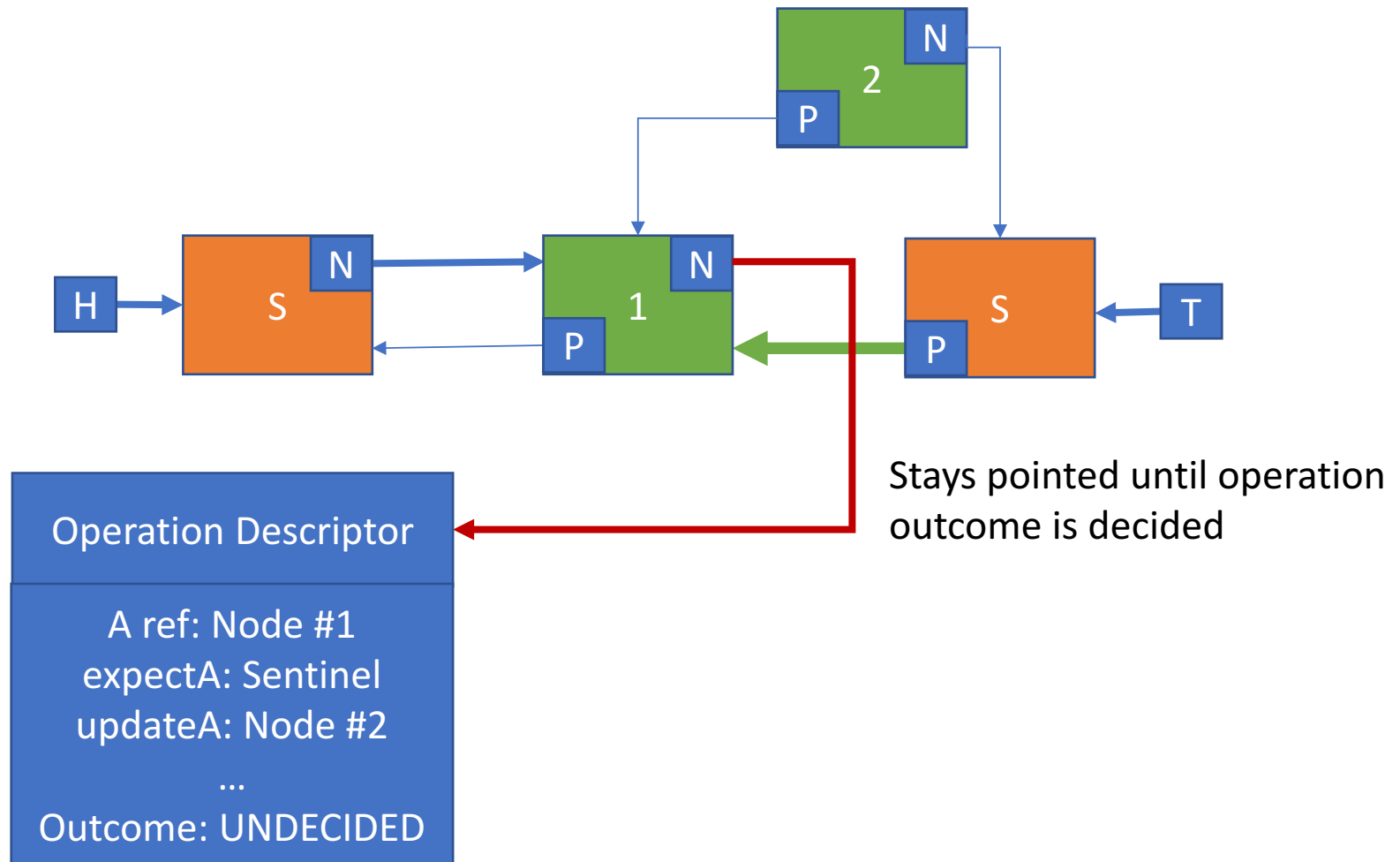operation ref

# Doubly linked list: insert (3)



desc is updated *after* successful DCSS

# Doubly linked list: insert (4)



Stays pointed until operation outcome is decided

**Operation Descriptor**

A ref: Node #1
expectA: Sentinel
updateA: Node #2
…
Outcome: UNDECIDED

# Doubly linked list: remove first

# Doubly linked list: remove (0)



CAS here

R1

H → S [N] → [N 1 P] → [N 2 P] → S [P] ← T

**Operation Descriptor**

A ref: ???
expectA: ???
updateA: Rem[???]
…
Outcome: UNDECIDED

Both not known in advance
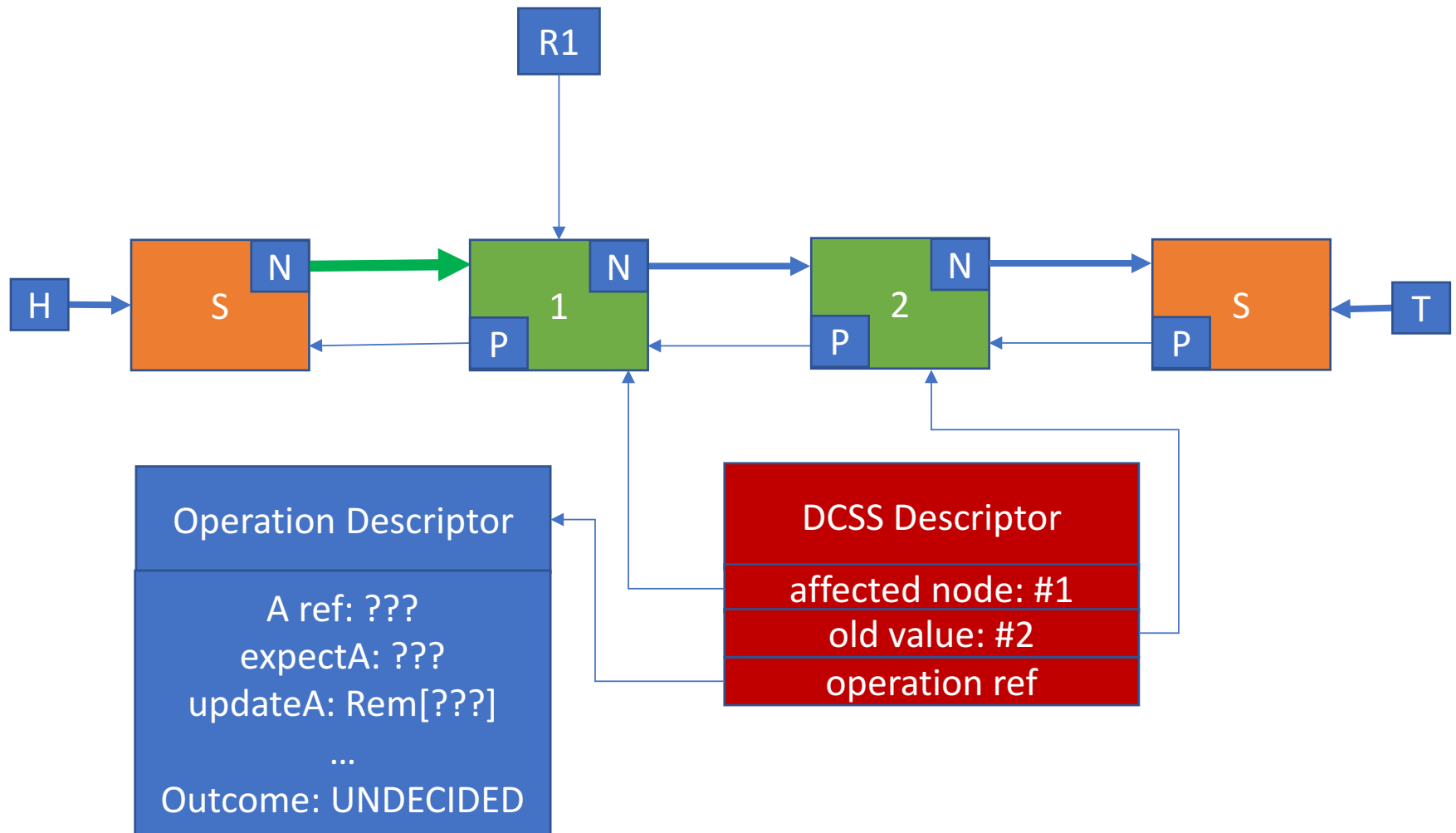
Deterministic f(expectA)

# Doubly linked list: remove (1)

# Doubly linked list: remove (2)
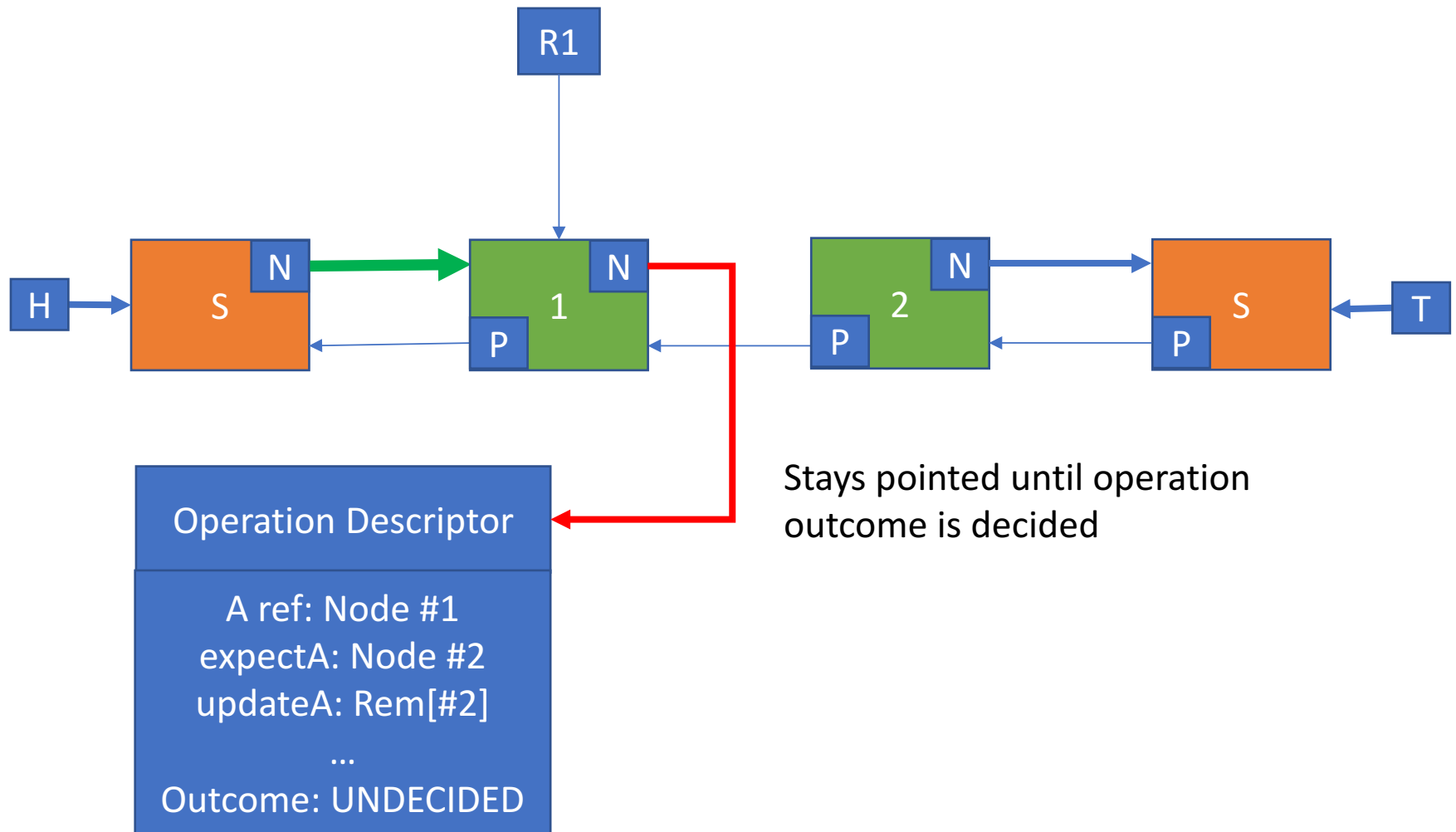
# Doubly linked list: remove (3)



desc is updated *after* successful DCSS

# Doubly linked list: remove (4)



R1

H → S | N
1 | N | P
2 | N | P
S | T

Stays pointed until operation outcome is decided

**Operation Descriptor**

A ref: Node #1
expectA: Node #2
updateA: Rem[#2]
…
Outcome: UNDECIDED

# Closing notes

- All we care about is **CAS** that linearizes operation
- Subsequent updates are *helper* moves
  - Invoke regular help/correct functions
- Perfect algorithm to combine with optional Hardware Transactional Memory (HTM)

Let's enjoy what we've accomplished

# References

- Kotlin language
  - http://kotlinlang.org
- Kotlin coroutines support library
  - http://github.com/kotlin/kotlinx.coroutines

# Thank you

## Any questions?

email me to **elizarov** at gmail

🐦 relizarov