

Consensus Numbers

魏恒峰

hfwei@nju.edu.cn

2017 年 12 月 14 日

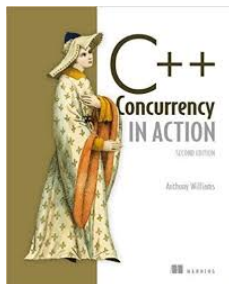
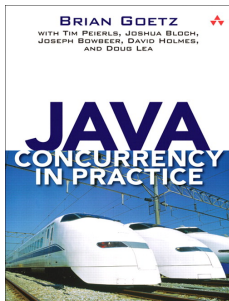




Concurrent Programming?



Synchronization



Synchronization Primitives

Lock
Semaphore

BlockingQueue
ConcurrentMap

Phaser
Barrier

Lock
Semaphore

BlockingQueue
ConcurrentMap

Phaser
Barrier

```
class AtomicInteger: // java.util.concurrent
    get()             // read
    set(int newVal)    // write

    getAndIncrement()
    getAndDecrement()
    getAndSet(int newVal)

    compareAndSet(int expect, int update)
```

Lock
Semaphore

BlockingQueue
ConcurrentMap

Phaser
Barrier

```
class AtomicInteger: // java.util.concurrent
    get()             // read
    set(int newVal)   // write

    getAndIncrement()
    getAndDecrement()
    getAndSet(int newVal)

    compareAndSet(int expect, int update) // atomically
```




(compareAndSet compareAndSwap)



(compareAndSet compareAndSwap)

```
CAS(int expect, int update) {  // atomically
    old = get()

    if (old == expect)
        set(update)
        return true

    return false
}
```

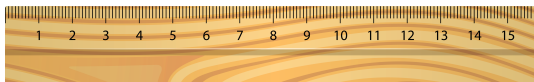


(Relative) *power* of various synchronization primitives

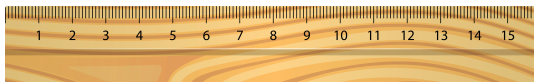
Power measured by the *consensus* problem



Q : Is it possible to solve the consensus problem using this synchronization primitive?



Consensus number: n -thread consensus problems



Consensus number: n -thread consensus problems



1, 2, n , \dots ∞

The beautiful idea of “Consensus Numbers”.

(Maurice Herlihy@TOPLAS'1991)

Wait-Free Synchronization

MAURICE HERLIHY

Digital Equipment Corporation

A *wait-free* implementation of a concurrent data object is one that guarantees that any process can complete any operation in a finite number of steps, regardless of the execution speeds of the other processes. The problem of constructing a wait-free implementation of one data object from another lies at the heart of much recent work in concurrent algorithms, concurrent data structures, and multiprocessor architectures. First, we introduce a simple and general technique, based on reduction to a consensus protocol, for proving statements of the form, “there is no wait-free implementation of X by Y .” We derive a hierarchy of objects such that no object at one level has a wait-free implementation in terms of objects at lower levels. In particular, we show that atomic read/write registers, which have been the focus of much recent attention, are at the bottom of the hierarchy: they cannot be used to construct wait-free implementations of many simple and familiar data types. Moreover, classical synchronization primitives such as *test&set* and *fetch&add*, while more powerful than *read* and *write*, are also computationally weak, as are the standard message-passing primitives. Second, nevertheless, we show that there do exist simple universal objects from which one can construct a wait-free implementation of any sequential object.

Consensus Number	Object
1	read/write registers
2	test&set, swap, fetch&add, queue, stack
\vdots	\vdots
$2n - 2$	n -register assignment
\vdots	\vdots
∞	memory-to-memory move and swap, augmented queue, compare&swap, fetch&cons, sticky byte

∞ : CAS

1 : Atomic read/write register

2 : Queue, getAndSet

Consensus



“It looks like we have a consensus.”

Propose



Decide

Propose



Decide

Definition (The Consensus Problem)

Agreement All (*non-faulty*) threads must *agree* on the same value.

Validity The common decision value must be the value *proposed* by some thread.

Wait-free Each (*non-faulty*) thread must decide in a *finite* number of steps.



Wait-free (very informally):

... *regardless of* the execution speeds of other threads



Wait-free (very informally):

- ... *regardless of* the execution speeds of other threads
- ... *eventually* finishes when it runs *solo*



We next focus on the *binary* consensus problem.

protocol

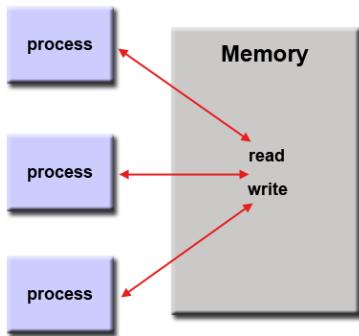
Definition (Consensus Protocol)

(Informally), a consensus protocol is a system of n processes that communicate through a set of shared objects.

Propose

Communicate

Decide



Processes communicate through shared objects.

(redraw)

consensus object (fig here)

```
1  public interface Consensus<T> {  
2      T decide(T value);  
3  }
```

consensus protocol

```
1 public abstract class ConsensusProtocol<T>
2     implements Consensus<T> {
3     protected T[] proposed = (T[]) new Object[N];
4
5     void propose(T value) {
6         proposed[ThreadID.get()] = value;
7     }
8
9     public abstract T decide(T value);
10 }
```

implement X using Y

```
1 public class CASConsensus<T>
2     extends ConsensusProtocol<T> {
3     private final int FIRST = -1;
4     private AtomicInteger r = new AtomicInteger(FIRST);
5
6     @Override
7     public T decide(T value) {
8         propose(value);
9
10        int i = ThreadID.get();
11        if (r.compareAndSet(FIRST, i))    // I won
12            return proposed[i];
13        else                               // I lose
14            return proposed[r.get()];
15    }
16 }
```

Theorem (Computational Power of CAS)

*A register providing **compareAndSet()** and **get()** methods can solve the consensus problem for **any number** of threads.*

Definition (Solving Consensus)

“ X solves n -process consensus” if there exists a consensus protocol

$$P_1, \dots, P_n; W, X$$

- ▶ W is a set of (atomic) read/write registers;
- ▶ W and X may be initialized to any state.

Definition (Consensus Number)

The **consensus number** for X is the largest n for which X solves n -thread consensus.

If no largest n exists, the consensus number is said to be **infinite**.

Theorem (Consensus Number of CAS)

A register providing *compareAndSet()* and *get()* methods has *infinite consensus number*.

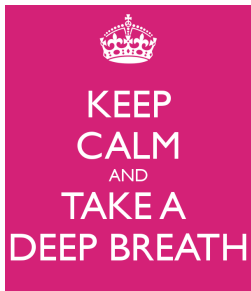


Lemma (Y Implements X)

Theorem (Consensus Number as ...)

in the following, main results (table)

beautiful ideas and proofs



Protocol in terms of set:

Protocol

$$\mathcal{P} = \{\text{All executions of this protocol}\}$$

Execution

$$e = \sigma_0 \xrightarrow{o_1} \sigma_1 \xrightarrow{o_2} \cdots \xrightarrow{o_{n-1}} \sigma_n$$

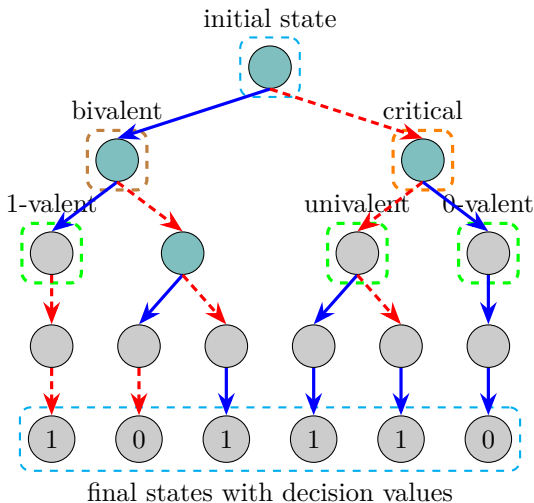
State

σ_i : states of individual threads + states of shared objects

Operation

o_i : method calls to a shared object

Modeling \mathcal{P} as a Computation “Tree” (Binary Consensus for 2 Threads)



Theorem (Bivalent Initial State)

Every 2-thread binary consensus protocol has a bivalent initial state.

Theorem (Bivalent Initial State)

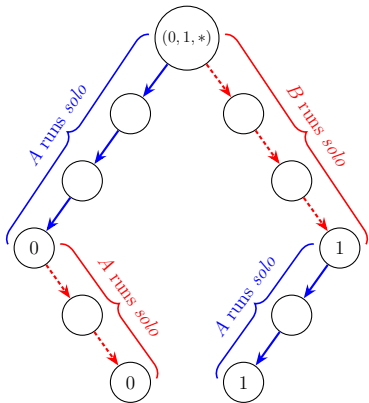
Every 2-thread binary consensus protocol has a bivalent initial state.

$$(A, B, O) : (0, 0, *) \quad (1, 1, *) \quad (0, 1, *) \quad (1, 0, *)$$

Theorem (Bivalent Initial State)

Every 2-thread binary consensus protocol has a bivalent initial state.

$$(A, B, O) : (0, 0, *) \quad (1, 1, *) \quad (0, 1, *) \quad (1, 0, *)$$



Lemma (Bivalent Initial State)

Every n -thread (binary) consensus protocol has a bivalent initial state.

Lemma (Bivalent Initial State)

Every n -thread (binary) consensus protocol has a bivalent initial state.

Theorem (Existence of Critical States)

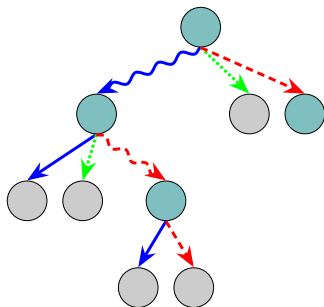
Every wait-free consensus protocol has a critical state.

Lemma (Bivalent Initial State)

Every n -thread (binary) consensus protocol has a bivalent initial state.

Theorem (Existence of Critical States)

Every wait-free consensus protocol has a critical state.







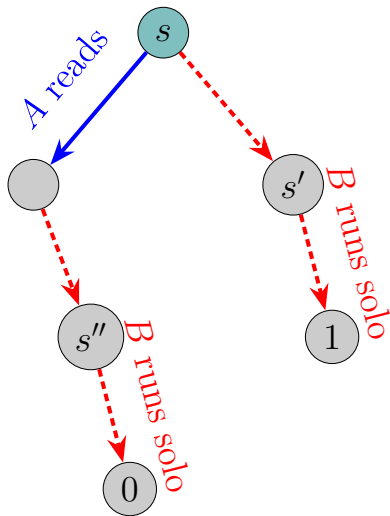
Theorem (Consensus Number of Atomic Registers)

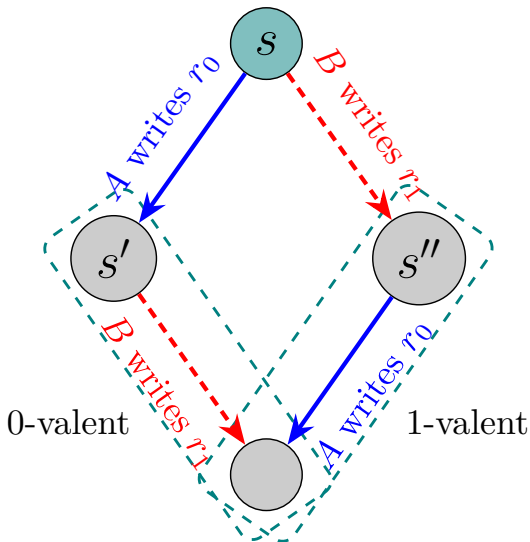
Atomic registers have consensus number 1.

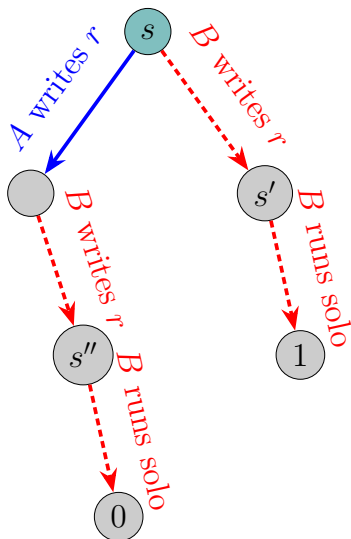
Proof.

- (1) Run the protocol until it reaches a critical state s .
- (2) What is the next?









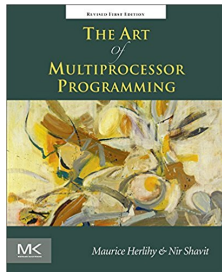
Theorem (“Weakness” of Atomic Read/Write Registers)

*It is **impossible** to construct a wait-free implementation of **any object with consensus number greater than 1** using atomic read/write registers.*

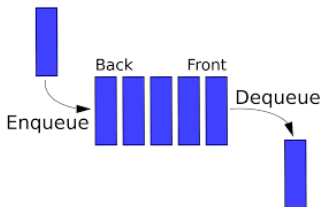
Theorem (“Weakness” of Atomic Read/Write Registers)

It is *impossible* to construct a wait-free implementation of *any object with consensus number greater than 1* using atomic read/write registers.

“... is perhaps one of the most striking impossibility results in Computer Science.”

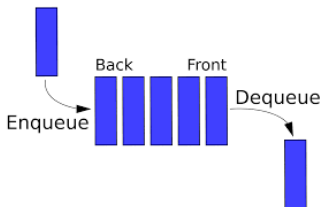






Theorem (Consensus Number of Queue)

FIFO queues have consensus number 2.



Theorem (Consensus Number of Queue)

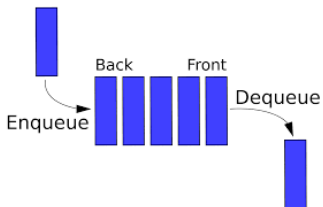
FIFO queues have consensus number 2.

Proof.

$$\geq 2$$

$$< 3$$





Theorem (Consensus Number of Queue)

FIFO queues have consensus number 2.

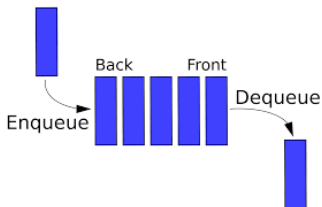
Proof.

$$\geq 2$$

$$< 3$$

By protocol for 2-thread consensus.





Theorem (Consensus Number of Queue)

FIFO queues have consensus number 2.

Proof.

$$\geq 2$$

By protocol for 2-thread consensus.

$$< 3$$

By the valency argument.



Theorem (“Power” of FIFO Queue)

The FIFO queue class can solve 2-thread consensus.

Theorem (“Power” of FIFO Queue)

The FIFO queue class can solve 2-thread consensus.

```
1 public class QueueConsensus<T>
2     extends ConsensusProtocol<T> {
3     private static final int WIN = 0;
4     private static final int LOSE = 1;
5
6     Queue queue;
7
8     // initialize queue with two items
9     public QueueConsensus() {
10         queue = new Queue();
11         queue.enq(WIN);
12         queue.enq(LOSE);
13     }
```

```
1 public class QueueConsensus<T>
2     extends ConsensusProtocol<T> {
3
4     // figure out which thread was first
5     public T decide(T val) {
6         propose(val);
7
8         int status = queue.deq();
9         int i = ThreadID.get();
10
11        if (status == WIN)
12            return proposed[i];
13        else
14            return proposed[1-i];
15    }
16 }
```

Theorem (“Weakness” of FIFO Queue)

The FIFO queue class cannot solve 3-thread consensus.

Proof.

By the valency argument and case analysis. □

Theorem (Consensus Number 2)

Many similar data types, such as FIFO queue, stack, set, list, and priority queue, all have consensus number 2.

Theorem (Consensus Number 2)

Many similar data types, such as FIFO queue, stack, set, list, and priority queue, all have consensus number 2.

Theorem (Consensus Number 2 (More))



The Idea of “Consensus Number”. (Maurice Herlihy@TOPLAS'1991)

Wait-Free Synchronization

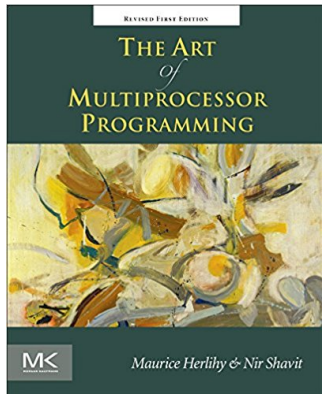
MAURICE HERLIHY

Digital Equipment Corporation

A *wait-free* implementation of a concurrent data object is one that guarantees that any process can complete any operation in a finite number of steps, regardless of the execution speeds of the other processes. The problem of constructing a wait-free implementation of one data object from another lies at the heart of much recent work in concurrent algorithms, concurrent data structures, and multiprocessor architectures. First, we introduce a simple and general technique, based on reduction to a consensus protocol, for proving statements of the form, “there is no wait-free implementation of X by Y .” We derive a hierarchy of objects such that no object at one level has a wait-free implementation in terms of objects at lower levels. In particular, we show that atomic read/write registers, which have been the focus of much recent attention, are at the bottom of the hierarchy: they cannot be used to construct wait-free implementations of many simple and familiar data types. Moreover, classical synchronization primitives such as *test&set* and *fetch&add*, while more powerful than *read* and *write*, are also computationally weak, as are the standard message-passing primitives. Second, nevertheless, we show that there do exist simple universal objects from which one can construct a wait-free implementation of any sequential object.

“The Art of Multiprocessor Programming”.

(Maurice Herlihy & Nir Shavit@2008)



The Idea of “Valency Argument”. (FLP@JACM'1985)

Impossibility of Distributed Consensus with One Faulty Process

MICHAEL J. FISCHER

Yale University, New Haven, Connecticut

NANCY A. LYNCH

Massachusetts Institute of Technology, Cambridge, Massachusetts

AND

MICHAEL S. PATERSON

University of Warwick, Coventry, England

Abstract. The consensus problem involves an asynchronous system of processes, some of which may be unreliable. The problem is for the reliable processes to agree on a binary value. In this paper, it is shown that every protocol for this problem has the possibility of nontermination, even with only one faulty process. By way of contrast, solutions are known for the synchronous case, the “Byzantine Generals” problem.

Thank
You!