ETH zürich

Parallel Programming Exercise 14





The second

Outline

- Post-Discussion: Assignment 13
- Pre-Discussion: Assignment 14
- Theory
- Exam Tasks & Tips





Post-Discussion Assignment 13

Contraction of the Party of

Sequential Consistency

For each of the following histories, indicate if they are sequentially consistent or not

registers: r and s queue: q

--|r.write(1)|-----A: A: q.enq(5)-----|r.read():0|-----B: q.enq(3)В: -----|r.read():1|-----A: void C: B: void A: q.deq() B: q.deq() --|s.write(1)|-----A: A: 3 -----|r.read():0|-----В: B: 3 -----|r.read():1|-----C:

A: --|s.write(1)|-----B: -----|r.read():1|---|r.read():0|-----



The second second

Sequential Consistency

For each of the following histories, indicate if they are sequentially consistent or not

registers: r and s queue: q







The second second second

Linearizability

Infer the object type from the supported operations, registers are initially zero, stacks/queues initially empty.

- A: s.push(1)
- A: void
- B: s.push(2)
- B: void
- B: s.pop()
- A: s.pop()
- B: 1
- A: 2

A: --|s.write(1)|-----B: -----|r.read():0|-----



Cart de la carte d

Linearizability

Infer the object type from the supported operations, registers are initially zero, stacks/queues initially empty.



A: --|s.write(1)|-----B: -----|r.read():0|-----



Equivalence

For all threads T: H|T = G|T

- A: r.read()
- A: 0
- C: r.write(1)

C: void

- B: r.read()
- B: 1

is equivalent to

- B: r.read()
- B: 1

The section of

- C: r.write(1)
- C: void
- A: r.read()
- A: 0



Incomplete Histories

When histories are obtained from a program trace, the history might be incomplete.

This can be dealt with in two ways.

Why do we need both ways? Give an example where discarding all pending invocations will lead to a non-linearizable history, but adding a response will lead to a linearizable history.





Sequential Consistency vs Linearizability

- A: write(1)
- A: void
- B: read()
- B: 0

SC but not linearizable





Pre-Discussion Assignment 14



Consensus

- Code snippets: argue about consistency (result same for each thread), validity (result was proposed by a thread) and wait-freedom
- Implement consensus protocol with wait-free FIFO queue
- Equivalence between consensus and binary consensus for two threads





Theory

13



Recap: Consensus Protocols



A few moments later... (a finite number of steps)



all the story was

Which other scenarios are allowed?



Consistent Result





This is illegal!

Consensus result needs to be consistent: the same on all threads.

a light of an a light to the



Valid Result





This is illegal!

Consensus result needs to be valid: proposed by some thread.

and the state of the second



Wait-Free



I cannot finish because I am waiting for the other thread.



This is illegal!

all the sectors and

Consensus needs to be wait-free: All threads finish after a finite number of steps, independent of other threads.



Consensus: Motivation

Example: Databases (of social networks)





Consensus: Motivation

Example: Databases (of social networks)





Consensus: Motivation

Example: Databases (of social networks)





Consensus Number

"The **consensus number** of a concurrent object is defined to be the maximum number of processes in the system which can reach consensus by the given object in a wait-free implementation. "

atomic read/write registers, mutex: 1 TAS, wait-free queue & stack: 2 CAS: ∞

Note: wait-free queue & stack cannot be implemented with atomic read/write registers





Atomic read/write registers: Consensus Number 1

Proof: There is no wait-free implementation of n-thread consensus (n > 1) with atomic read/write registers.



Proof simplification

Proof that is does not work for 2 threads to show that it does not work for n threads.

Why can we make this assumption?

Assume our consensus protocol is correct (consistent, valid, wait-free).

Assume n - 2 threads die / get descheduled.

Since our consensus is **wait-free**, it should still work for the two remaining threads.

***SPCL

Simplification: Binary Consensus

- Instead of proposing an integer, every thread now proposes either 0 or 1
- Equivalent to "normal" consensus for two threads
 - How can we proof this?

binary_decide(bit b) {
 return int_decide(b)
}

We can implement binary consensus using normal consensus.

```
int_decide(int d) {
    prop[id] = d; //prop is shared
    other = (id + 1)%2;
    int win = bin_decide(id);
    return prop[win];
}
```

We can implement binary consensus using normal consensus (id in {0,1} and unique).



State Diagrams of Two-thread Consensus Protocols

Cycles among states cannot exist in a wait-free algorithm: The state "looks" the same each time we visit, so we are trapped forever in the loop and not wait-free.

Start state, both threads (A and B) have not yet executed the first instruction of the consensus protocol.

Each state has at most two **successors**: Either A or B execute an instruction.



Anatomy of a State (in two-thread consensus)



and the second se



Anatomy of a State



A DESCRIPTION OF THE PARTY OF T



Critical States



A CONTRACTOR TO THE





ANT A COMPANY



Critical State Existence Proof

Lemma: Every consensus protocol has a critical state.

Proof: From (bivalent) start state, let the treads only move to other bivalent states.

- If it runs forever the protocol is not wait free.
- If it reaches a position where no moves are possible this state is critical.



Impossibility Proof Setup – Critical State



So what actions can a thread perform in his "move"?

A State of the second

Either read or write a shared register! – Let's see why.



32

Impossibility Proof Setup – Possible actions of a thread



The second of the

Impossibility Proof Setup – Possible actions of a thread



State and and the state

Many cases... let's make tables



Many Cases to check

		First Action		n			Is binary
		A: r1.read()	A: r1.write()	A: r1.write()	A: r2.write()		consensus
Second Action	B: r1.read()						of those?
	B: r2.read()						
	B: r1.write()						Can we simplify
	B: r2.write()						somehow?

Del Charles

		Second Action						
		A: r1.read()	A: r2.read()	A: r1.write()	A: r2.write()			
	B: r1.read()							
First	B: r2.read()	narahla Lat	's look at the	cases wher	o A roads			
Action	B: r1.write()	iagabie Let		cases when	CATCOUS			
	B: r2.write()							

Let's say A always moves first, otherwise, switch names.

Similarly, we can call the register A reads r1 in both cases.



Contraction and

Impossibility Proof Case I: A reads





What did we just prove?



and the second second


Contraction of the

Impossibility Proof Case I': B reads





What did we just prove?

		First Action	
		A: r1.read()	A: r1.write()
Second Action	B: r1.read()	No, Case I	No, Case I'
	B: r2.read()	No, Case I	No, Case I'
	B: r1.write()	No, Case I	O
	B: r2.write()	No, Case I	

State Same and

Is binary consensus possible for any of those?



Impossibility Proof Case II: A and B write to different registers





What did we just prove?

		First Action	
		A: r1.read()	A: r1.write()
	B: r1.read()	No, Case I	No, Case I'
Second Action	B: r2.read()	No, Case I	No, Case I'
	B: r1.write()	No, Case I	?
	B: r2.write()	No, Case I	No, Case II

States and the second

Is binary consensus possible for any of those?



Impossibility Proof Case III: A and B write to the same register

The second





That's all

		First A	Action		
		A: r1.read()	A: r1.write()		la hina
Second Action	B: r1.read()	No, Case I	No, Case I'	IS DIR CONSE	is bina consen
	B: r2.read()	No, Case I	No, Case l'		possible fo
	B: r1.write()	No, Case I	No, Case III	~ -	of thos
	B: r2.write()	No, Case I	No, Case II		No





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





What did we prove?

• Atomic read/write registers have consensus number 1





Consensus: TAS consensus number

Proof: TAS has consensus number 2



Proof outline

- First proof that TAS has consensus number n >= 2 by construction
- Then proof that TAS has consensus number n < 3 by contradiction</p>
- ➔ n needs to be 2



Implementing two thread consensus with TAS

 Assume you have a machine with atomic registers and an atomic test-and-set operation with the following semantics (mem[s] is initially 0):

```
boolean TAS(memref s) {
    if (mem[s] == 0) {
        mem[s] = 1;
        return true;
    }
    return false;
```

Implement a wait-free two-process consensus protocol using TAS and atomic registers.



Implementing two thread consensus - Solution

Pseudo-Code for both threads

AtomicIntegerArray proposed = new AtomicIntegerArray(2); flag = 0;

```
int decide (int value) {
    int i = ThreadID.get();
    proposed.set(i, value);
    if (TAS(flag)) {
        return value;
    }
    else {
        return proposed.get((i + 1) % 2);
    }
}
```



TAS consensus number: proof outline that n < 3

- Assume there exist some correct wait-free consensus protocol with TAS for n > 2 threads.
- Proof that it does not work for 3 threads (which implies that it does not work for n > 3 threads)
- Go over all cases just like in the proof for atomic read/write registers
- Difference to proof for atomic read/write registers
 - Each node can have 3 children
 - Multivalent would be a better word for states that did not decide yet





Consensus: CAS consensus number

Proof: CAS has consensus number ∞



Proof by construction: CAS consensus number ∞

```
class CASConsensus {
  private final int FIRST = -1;
  private AtomicInteger r = new AtomicInteger(FIRST); // supports CAS
  private AtomicIntegerArray proposed; // suffices to be atomic register
  ... // constructor (allocate array proposed etc.)
  public Object decide (Object value) {
     int i = ThreadID.get();
     proposed.set(i, value);
     if (r.compareAndSet(FIRST, i)) // I won
       return proposed.get(i); // = value
     else
```

```
return proposed.get(r.get());
```



Consensus: What should you definitely know for the exam

- Know about properties: consistent, valid, wait-free
- Bivalent, univalent, critical state
- Argue about properties: why acyclic graph?
- Know consensus numbers
- Argue with consensus numbers





MPI (Message Passing Interface)



MPI: Motivation

- Locking is too slow
- Shared memory is not always viable
- What if we have multiple compute nodes and want to distribute work?



High Performance Computing: Supercomputers

- Connect many different computers (cluster) and let them communicate with each other over high speed interconnects (e.g. Infiniband)
- Example: Dragonfly-Topology





all the stars way

Infiniband demo

[eli@sbrinz1 ~]\$ ib_read_bw						

RDMA_Read BW TestDual-port: OFFDevice: mlx5_0Number of qps: 1Transport type: IBConnection type: RCUsing SRQ: OFFPCIe relax order:ONLock-free: OFFibv_wr* API: ONUsing DDP: OFFCQ Moderation: 1Mtu: 4096[B]Link type: IBOutstand reads: 16rdma_cm QPs: OFFData ex. method: Ethernet						
local address: LID 0x03 QPN 0x18bac PSN 0xdfad59 OUT 0x10 RKey 0xb83a00 VAddr 0x00400000529000 remote address: LID 0x03 QPN 0x18bab PSN 0x6c0903 OUT 0x10 RKey 0xb69600 VAddr 0x00400000529000						
#bytes #iterations BW peak[MiB/sec] BW average[MiB/sec] MsgRate[Mpps] 65536 1000 39185.09 39155.05 0.626481						



MPI

- Message Passing Interface defines the semantics
- There are many different MPI implementations: e.g. OpenMPI, MPICH
- → When you build a program using MPI, you can decide which implementation to use
- Different implementations can lead to different performance, but the semantic is the same as defined by the MPI standard
- \rightarrow One is also free to choose the number of MPI-ranks & processes



Contra and and the

MPI

```
MPI.Init(args);
int rank = MPI.COMM_WORLD.Rank();
int size = MPI.COMM_WORLD.Size();
```

// CODE...

MPI.Finalize();



MPI

comm_∎send(

Object buf,	//	the data object to be sent
<pre>int offset,</pre>	//	start offset from the data start iter
int count,	//	number of items to be sent
Datatype datatype,	//	data type of items
int dest,	//	destination process id
int tag	//	data id tag

comm.recv(

Qbject buf,	//	the object to write the data to
int offset,	//	
int count,	//	number of items to be receiv
Datatype datatype,	//	data type of items
int src,	//	souce process id or MPI_ANY_SOURCE
int tag	//	data id tag MPI_ANY_TAG



The second state

MPI barrier

```
public void mybarrier() {
    int rank = MPI.COMM_WORLD.Rank();
    int size = MPI.COMM_WORLD.Size();
    boolean[] buf = new boolean[1];
    if (rank == 0) {
       // Wait until every process sent a meaningless boolean
       for (int i = 1; i < size; i++) {</pre>
           MPI.COMM_WORLD.Recv(buf, 0, 1, MPI.BOOLEAN, i, 0);
       // All processes are ready
       // Send a signal to all processes so that they can continue
       for (int i = 1; i < size; i++) {</pre>
           MPI.COMM_WORLD.Send(buf, 0, 1, MPI.BOOLEAN, i, 0);
       }
    } else {
         MPI.COMM_WORLD.Send(buf, 0, 1, MPI.BOOLEAN, 0, 0); // send meaningless boolean
         MPI.COMM_WORLD.Recv(buf, 0, 1, MPI.BOOLEAN, 0, 0); // wait until process 0 gives signal
```



Collective Computation paradigms



ALCON STREET STR



Collective Computation paradigms

Distribution



A STATISTICS AND AND A STATISTICS





Exam tasks

62



Consensus

11. Welche der folgenden Methoden sind korrekte Implementierungen von Wait-free Consensus für die angegebene Zahl N an Threads? Begründen Sie Ihre Antworten. Which of the following Java Methods are valid implementations of wait-free consensus for the given number N of threads? Justify your answers.

The section



Consensus

11. Welche der folgenden Methoden sind korrekte Implementierungen von Wait-free Consensus für die angegebene Zahl N an Threads? Begründen Sie Ihre Antworten. Which of the following Java Methods are valid implementations of wait-free consensus for the given number N of threads? Justify your answers.

The sections

(a) $N = 1$	(2) (b) $N = 1$	(2)
<pre>1 int consensus(int x) { 2 return 0; 3 }</pre>	<pre>1 int consensus(int x) { 2 return x; 3 }</pre>	
Not valid	correct	



(2)

(c) N = 2

```
(2) (d) N = 2
```

a light and and the

```
volatile int decision;
volatile bool decided = false;
synchronized int consensus(int x) {
    if (!decided)
    decision = x;
    decided = true;
    return decision;
    }
```

```
1 AtomicInteger dec = 0;
2 int[] prop = new int[] {0,0};
3
4
5 int consensus(int x) {
6     prop[dec.Get()] = x;
7     t = dec.IncrementAndGet();
8     return prop[t-1];
9 }
```



(2)

(c) N = 2

(2) (d) N = 2

```
volatile int decision;
volatile bool decided = false;
synchronized int consensus(int x) {
    if (!decided)
    decision = x;
    decided = true;
    return decision;
    }
```

```
AtomicInteger dec = 0;
2 int[] prop = new int[] {0,0};
3
4
5 int consensus(int x) {
```

```
6 prop[dec.Get()] = x;
```

```
7 t = dec.IncrementAndGet();
```

```
8 return prop[t-1];
```

A CATA COMPANY

```
9 }
```

Not wait free

Not consistent

Clamenters.

The sections

No, as using locks can lead to threads blocking each other \rightarrow Not wait-free



and and and



Ald and and

Yes

What are the consensus numbers of the following three objects:

spcl.inf.ethz.ch

EHzürich

a) atomic registers that support only the operations Read() and Write(),
b) atomic registers that also support the

TestAndSet() operation, and

c) atomic registers that also support the CompareAndSet() operation?

What are the consensus numbers of the following three objects: a) atomic registers that support only the 1operations Read() and Write(), 1b) atomic registers that also support the 2TestAndSet() operation, and c) atomic registers that also support the ∞ CompareAndSet() operation?

The second of the P

spcl.inf.ethz.ch

EHzürich
What is the difference between MPI point-to-point and collective operations? Why are collective operations used?

Contra and and the

What is the difference between MPI point-to-point and collective operations? Why are collective operations used?

Point-to-Point: One sender to one receiver Collective operations: see previous slide about "collective computation paradigms"

We mainly use collective operations to get scalable performance

***SPCL

Which two collective operations can be combined to imlement the Allgather operation?



Distribution						
Broadcast			Scatter/Gather			
PO	A		A	PO A B C	D Scatter	A
P1		Broadcast	А	P1		В
P2			А	P2	Gather	С
P3			A	P3		D
AllGathe	r			Scan		
PO	A		A B C D	PO A		А
P1	В	Allgather	A B C D	P1 B	Scan	A+B
P2	С		A B C D	P2 C		A+B+C
P3	D		A B C D	P3 D		A+B+C+D
AllToAll						
PO	A0 A1 A2 A3		A0 B0 C0 D0			
P1	B0 B1 B2 B3	Alltoall	A1 B1 C1 D1			
P2	C0 C1 C2 C3		A2 B2 C2 D2			
P3	D0 D1 D2 D3		A3 B3 C3 D3			

Station and and

***SPCL

Which two collective operations can be combined to imlement the Allgather operation?





and the states

Gather and Broadcast

Ein MPI Programm wird von 8 Prozessen ausgeführt. Jeder Prozess kann zu jeder Zeit maximal an einem Nachrichtenaustausch teilnehmen. Eine Nachricht welche einen double Wert enthält kann innerhalb 1 ms zwischen zwei beliebigen Prozessen ausgetauscht werden. In den unten gezeichneten Schemata (A, B und C) wird jeder Nachrichtenaustausch als gerichtete Kante zwischen Sender und Empfänger dargestellt.

Das Ziel des MPI Programms ist es, einen double Wert von Prozess 0 zu jedem anderen Prozess zu senden. Beantworten Sie für jedes der Schemata unten folgende Fragen: Wird das Ziel erreicht? Was ist die minimale Zeitdauer, bis der letzte Prozess seine Aufgabe erfüllt hat. An MPI program is executed by 8 processes. Each process can only participate in one communication at a time (either as a sender or a receiver). A message containing a double can be sent between any two processes in 1 ms. In the communication schemes below (A, B, and C) each directed edge represents a message sent from a sender to a receiver process.

The goal is to send a double from process 0 to all other processes. For each of the schemes below, answer the following: Do they achieve the goal? What is the minimum time for the last process to finish its task?



We have point-topoint communication





Benjamin Scherling Catalan @bscherling · 10 months ago · edited 10 months ago



A: Yes, message passes through all processes. Minimum time = 3ms

B: Yes, all processes recieve message directly from 0. Minimum time = 7ms

C: No, since we start in process 0 we see that process 7 has no in-going edge and therefore can't recieve any message



Advice / Outline for exam preparation

- Revise topics to get an overview over everything
- Practice some "easy points" tasks until you are comfortable with them (always the same pattern)
 - Amdahl, Gustafson, Pipelining, Histories, State Diagrams, Fork-Join, etc.
- Practice harder tasks (some "creativity" needed)
 - Wait / notify, questions about properties of locks, barriers, code snippets, etc.
- Practice "Mixer" tasks
 - Very unpredictable
 - Revise theory; try to understand and make connections between different concepts

- Practice some old exams without a timer first
- Identify tasks where you struggle → specifically practice those tasks and skip easy tasks
- Try to solve under time contraints when you feel comfortable with the tasks and theory
- Try to conclude how much time to invest for each task based on the assigned points



Good theory overview

- https://luck-wildflower-a5b.notion.site/pprog-summary-leon-thomm
- PVW script
- Summaries on ComSol (<u>https://exams.vis.ethz.ch/</u>)
- The Art of Multithreading book



Viel Glück und Erfolg für die Prüfung!!

and the second second second