Distributed Systems CS6421

Clocks & Coordination

Prof. Tim Wood

Challenges Heterogeneity Openness Security Failure Handling Concurrency Quality of Service Scalability

Transparency

Clocks and Timing

Distributed systems often need to order events to help with consistency and coordination

Coordinating updates to a distributed file system Managing distributed locks Providing consistent updates in a distributed DB

Coordinating time?

How can we synchronize the clocks on two servers?

A <u>clock</u>: *8:03*

B

clock: 8:01

Cristian's Algorithm

Easy way to synchronize clock with a time server



Client sends a clock request to server

Measures the round trip time

Set clock to t + 1/2*RTT (8:01.505)

Cristian's Algorithm

What will affect accuracy?



Cristian's Algorithm

Suppose the minimum delay between A and B is X



Ordering

Sometimes we don't actually need clock time

We just care about the order of events!

What event happens before another event?

- e->e' means event e happens before event e'

Easy: we'll just use counters in each process and update them when events happen!

- Maybe not so easy...



An avant is **and** of the following:

- Action that occurs within a process
- Sending a message
- Receiving a message

What is true? What can't we know?

Tim Wood - The George Washington University

N)

7



11/hat in trua?

- a->b, b->g, c->d, e->f (events in same process)
- b->c, d->f (send is before receive)

What can't we know?

- e??a
- e??c

N)

2



independently, are those counters meaningful?

Tim Wood - The George Washington University

2

N



Each process maintains a counter, L

Increment counter when an event happens

When receiving a message, take max of own and sender's counter, then increment

Clock Comparison

Independent clocks p_1 1 2 3 g g p_2 p_3 g p_3 p_4 p_4 p_5 p_6 p_7 p_8 p_7 p_9 $p_$

if e->e', then: C(e) ??? C(e')



Clock Comparison

Is the opposite true? if L(e) < L(e') then do we know e->e'?



Clock Comparison

Is the opposite true? No! Lamport clocks don't actually let us compare two clocks to know how they are related :(



Lamport Clocks

Lamport clocks are better than nothing

- but only let us make limited guarantees about how things are ordered
- Ideally we want a clock value that indicates:
 - If an event happened before another event
 - If two events happened concurrently



vector clocks



Each process keeps an array of counters: (p1, p2, p3)

- When p_i has an event, increment V[p_i]
- Send full vector clock with message
- Update each entry to the maximum when receiving a clock

vector clocks



Now we can compare orderings! if V(e) < V(e') then e->e' - (a,b,c) < (d,e,f) if: $a \le d \& b \le e \& c \le f$ If neither V(e) < V(e') nor V(e') < V(e) then e and e' are concurrent events

PI	P2	P3
a: 1,0,0	c: 2, I ,0	e: 0,0, I
b: 2,0,0	d: 2,2,0	f: 2,2,2
g: 3,0,0		

Lamport vs Vector

Which clock is more useful when you can't see the timing diagram?

- Remember, your program will only see these counters!

ΡI	P2	P3
a:I	c:3	e:I
b:2	d:4	f:5
g:3		

PI	P2	P3
a: 1,0,0	c: 2, I ,0	e: 0,0, I
b: 2,0,0	d: 2,2,0	f: 2,2,2
g: 3,0,0		

VC Worksheet



What are the vector clocks at each event? Assume all processes start with (0,0,0)

How to Compare VC?



How does g compare to d?

Vector Clocks

Allow us to compare clocks to determine a **partial** ordering of events

Example usage: versioning a document being edited by multiple users. How do you know the order edits were applied and who had what version when they edited?

Is there a drawback to vector clocks compared to Lamport clocks?

Clock Worksheet

Do the worksheet in groups of 2-3 students

When you finish, do this on the back:

- Draw the timeline for the four processes with vector clocks shown in problem 3. Compare your answer with another group.

P1	P2	P3	P4
a: 1,0,0,0	e: 1,1,0,0	i: 0,0,1,0	I: 0,0,0,1
b: 2,0,0,0	f: 1,2,0,1	j: 0,0,2,2	m: 0,0,0,2
c: 3,0,0,0	g: 1,3,0,1	k: 0,0,3,2	n: 0,0,0,3
d: 4,2,0,1	h: 1,4,3,2		



We can apply the vector clock concept to versioning a piece of data

- This is used in many distributed data stores (DynamoDB, Riak)

When a piece of data is updated:

- Tag it with the **actor** who is modifying it and the version #
- Treat the (actor: version) pairs like a vector clock

The version vectors can be used to determine a causal ordering of updates

Also can detect concurrent updates

Need to have a policy for resolving conflicts

- If two versions are concurrent, they are "siblings", return both!

Alice tells everyone to meet on Wednesday

- Dave and Cathy discuss and decide on Thursday
- Ben and Dave exchange emails and decide Tuesday
- Alice wants to know the final meeting time, but Dave is offline and Ben and Cathy disagree... what to do?



Alice tells everyone to meet on Wednesday

Dave and Cathy discuss and decide on Thursday

Ben and Dave exchange emails and decide Tuesday

Alice	Bob	Cathy	Dave
Wednesday	Wednesday	Wednesday	Wednesday
	Tuesday	Thursday	Thursday
			Tuesday

Alice tells everyone to meet on Wednesday

Dave and Cathy discuss and decide on Thursday

Ben and Dave exchange emails and decide Tuesday

Alice	Bob	Cathy	Dave
Wednesday A:1	Wednesday A:1	Wednesday A:1	Wednesday A:1
		Thursday A:1, C:1, D:1	Thursday A:1, C:1, D:1
	Tuesday A:1, B:1, C:1, D:2		Tuesday A:1, B:1, C:1, D:2

The result ends on the order of:

- Dave and Cathy discuss and decide on Thursday
- Ben and Dave exchange emails and decide Tuesday



Resolving Conflicts

What if we have?



What are the conflicts?

Resolving Conflicts

What if we have?



How to resolve Alice vs the rest?

- The Tuesday vs Thursday debate is not a real conflict since we can order them based on their version vectors

We need a policy for resolving the conflicts

- Random
- Priority based
- User resolved

Dependencies

Vector clocks also help understand the dependency between different events and processes



Multi-Tier Backup

Consider a multi-tier web app backup system

- Some tiers have a disk that must be protected
- All writes to protected disks must be replicated to a backup
- Can only send responses to a client once writes have been successfully backed up!



Tracking Dependencies

Use Vector Clocks to track pending writes

- One entry per protected disk: $<D_1, D_2, ..., D_n >$



Node \boldsymbol{i} increments \boldsymbol{D}_i on each write

Use vector clocks to determine a causal ordering



Tracking Dependencies

Use Vector Clocks to track pending writes

- One entry per protected disk: $<D_1, D_2, ..., D_n >$



Node \boldsymbol{i} increments \boldsymbol{D}_i on each write

Use vector clocks to determine a causal ordering



Ordered Asynchrony

Allowing processing to proceed asynchronously provides major performance advantage!

- But need vector clocks to determine ordering and dependencies



Tim Wood - The George Washington University

Time and Clocks

Synchronizing clocks is difficult

But often, knowing an order of events is more important than knowing the "wall clock" time!

Lamport and Vector Clocks provide ways of determining a consistent ordering of events

- But some events might be treated as concurrent!

The concept of vector clocks or version vectors is commonly used in real distributed systems

Distributed Coordination

(Distributed) Locking

We need mutual exclusion to protect data

- How does this limit scalability?

Among processes and threads:

- Mutexes and Semaphores

Among distributed servers?

Centralized or decentralized?

Centralized Approach

Simplest approach: put one node in charge

Other nodes ask coordinator for each lock

- Block until they are granted the lock
- Send release message when done

Coordinator can decide what order to grant lock

Do we get:

- Mutual exclusion?
- Progress?
- Resilience to failures?
- Balanced load?



Distributed Approach

Use Lamport Clocks to order lock requests across nodes

Send Lock message with clock

- Wait for OKs from all nodes

When receiving Lock msg:

- Send OK if not interested
- If I want the lock:
 - Send OK if request's clock is smaller
 - Else, put request in queue

When done with a lock:

- Send OK to anybody in queue



Ring Approach

- Nodes are ordered in a ring
- One node has a token
- If you have the token, you have the lock
- If you don't need it... pass it on



Token Ring

Can be slow...

- Will we make progress?



Comparison

Messages per lock/release

- Centralized:
- Distributed:
- Token Ring:

Delay before entry

- Centralized:
- Distributed:
- Token Ring:

Problems

- Centralized:
- Distributed:
- Token Ring:

Are the distributed approaches better in any way?

Comparison

Messages per lock/release

- Centralized: 3
- Distributed: 2(n-1)
- Token Ring: ???

Delay before entry

- Centralized: 2
- Distributed: 2(n-1) in parallel
- Token Ring: 0 to n-1 in sequence

Problems

- Centralized: Coordinator crashes
- Distributed: anybody crashes
- Token Ring: lost token, crashes

Are the distributed approaches better in any way?

Distributed Systems are Hard

Going from centralized to distributed can be..

Slower

- If everyone needs to do more work

More error prone

- 10 nodes are 10x more likely to have a failure than one

Much more complicated

- If you need a complex protocol
- If nodes need to know about all others

Distributed Architectures

Purely distributed / decentralized architectures are difficult to run correctly and efficiently



Decentralized



Centralized

Elections

Appoint a central coordinator

- But allow them to be replaced in a safe, distributed way

Must be able to handle simultaneous elections

- Reach a consistent result

Who should win?



Bully Algorithm

The biggest (ID) wins

Any process P can initiate an election

P sends **Election** messages to all process with higher Ids and awaits **OK** messages

If it receives an OK, it drops out and waits for an **I won**

If a process receives an **Election** msg, it returns an **OK**...



Bully Algorithm

The biggest (ID) wins

Any process P can initiate an election

P sends **Election** messages to all process with higher Ids and awaits **OK** messages

If it receives an OK, it drops out and waits for an **I won**

If a process receives an **Election** msg, it returns an **OK**...

What next?





Bully Algorithm

The biggest (ID) wins

Any process P can initiate an election

P sends **Election** messages to all process with higher Ids and awaits **OK** messages

If it receives an OK, it drops out and waits for an **I won**

If a process receives an **Election** msg, it returns an **OK** and starts an election

If no **OK** messages, P becomes leader and sends **I won** to all process with lower lds

If a process receives a **I won**, it treats sender as the leader



P4

P5

Ring Algorithm

Any other ideas?



Ring Algorithm

Initiator sends an **Election** message around the ring

Add your ID to the message

When Initiator receives message again, it announces the winner

What happens if multiple elections occur at the same time?



Ring Algorithm

Initiator sends an **Election** message around the ring

Add your ID to the message

When Initiator receives message again, it announces the winner

What happens if multiple elections occur at the same time?



Comparison

Number of messages sent to elect a leader:

Bully Algorithm

- Worst case: lowest ID node initiates election
 - Triggers n-1 elections at every other node = $O(n^2)$ messages
- Best case: Immediate election after n-2 messages

Ring Algorithm

- Always 2(n-1) messages
- Around the ring, then notify all

Elections + Centralized Locking

Elect a leader

Let them make all the decisions about locks

- What kinds of failures can we handle?
 - Leader/non-leader?
 - Locked/unlocked?
 - During election?

