CS162 Operating Systems and Systems Programming Lecture 23

Distributed Decision Making (Con't), Networking and TCP/IP

November 18th, 2020 Prof. John Kubiatowicz http://cs162.eecs.Berkeley.edu

Recall: How do entities communicate? A Protocol!



- A protocol is an agreement on how to communicate, including:
 - Syntax: how a communication is specified & structured
 » Format, order messages are sent and received
 - Semantics: what a communication means
 - » Actions taken when transmitting, receiving, or when a timer expires
- · Described formally by a state machine
 - Often represented as a message transaction diagram
 - Can be a partitioned state machine: two parties synchronizing duplicate sub-state machines between them
 - Stability in the face of failures!

```
11/18/20
```

Kubiatowicz CS162 © UCB Fall 2020

Lec 23.2

Recall: Distributed Applications

- How do you actually program a distributed application?
 Need to synchronize multiple threads, running on different machines
 - » No shared memory, so cannot use test&set



- One Abstraction: send/receive messages
 » Already atomic: no receiver gets portion of a message and two receivers cannot get same message
- Interface:
 - Mailbox (mbox): temporary holding area for messages
 » Includes both destination location and queue
 - Send(message,mbox)
 - » Send message to remote mailbox identified by mbox
 - Receive(buffer,mbox)
 - » Wait until mbox has message, copy into buffer, and return
 - » If threads sleeping on this mbox, wake up one of them

Using Messages: Send/Receive behavior

- When should send (message, mbox) return?
 - When receiver gets message? (i.e. ack received)
 - When message is safely buffered on destination?
 - Right away, if message is buffered on source node?
- · Actually two questions here:
 - When can the sender be sure that receiver actually received the message?
 - When can sender reuse the memory containing message?
- Mailbox provides 1-way communication from T1 \rightarrow T2
 - $-T1 \rightarrow buffer \rightarrow T2$
 - Very similar to producer/consumer
 - » Send = V, Receive = P
 - » However, can't tell if sender/receiver is local or not!

11/18/20



Lec 23.7

General's Paradox (con't) **Two-Phase Commit** Since we can't solve the General's Paradox Can messages over an unreliable network be used to guarantee two entities do something simultaneously? (i.e. simultaneous action), let's solve a related problem - Remarkably, "no", even if all messages get through Distributed transaction: Two or more machines agree to do ¹¹ am ok> something, or not do it, atomically - No constraints on time, just that it will eventually happen! but what it Two-Phase Commit protocol: Developed by Turing award Don't an this ack? winner Jim Gray - (first Berkeley CS PhD, 1969) - Many important DataBase breakthroughs also from Jim Gray - No way to be sure last message gets through! Jim Grav - In real life, use radio for simultaneous (out of band) communication • So, clearly, we need something other than simultaneity! 11/18/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 23.9 11/18/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 23.10 **Two-Phase Commit Protocol 2PC Algorithm** · Persistent stable log on each machine: keep track of whether commit has One coordinator happened N workers (replicas) - If a machine crashes, when it wakes up it first checks its log to recover state of · High level algorithm description: world at time of crash Prepare Phase: - Coordinator asks all workers if they can commit - The global coordinator requests that all participants will promise to commit or - If all workers reply "VOTE-COMMIT", then coordinator broadcasts "GLOBAL-COMMIT" rollback the transaction Otherwise coordinator broadcasts "GLOBAL-ABORT" - Participants record promise in log, then acknowledge - Workers obey the GLOBAL messages - If anyone votes to abort, coordinator writes "Abort" in its log and tells everyone to abort; each records "Abort" in log · Use a persistent, stable log on each machine to keep track of what you are Commit Phase doing - After all participants respond that they are prepared, then the coordinator writes - If a machine crashes, when it wakes up it first checks its log to recover state of "Commit" to its log world at time of crash - Then asks all nodes to commit; they respond with ACK - After receive ACKs, coordinator writes "Got Commit" to log · Log used to guarantee that all machines either commit or don't 11/18/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 23.11 11/18/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 23.12

	 Two-Phase Commit: Setup One machine (<i>coordinator</i>) initiates the protocol It asks <i>every</i> machine to vote on transaction Two possible votes: Commit Abort Commit transaction only if unanimous approval 		v	Two-Phase Commit: Preparing Vorker Agrees to Commit Machine has guaranteed that it will accept transaction Must be recorded in log so machine will remember this decision and restarts Vorker Agrees to Abort Machine has guaranteed that it will never accept this transact Must be recorded in log so machine will remember this decision and restarts	nit: Preparing cept transaction vill remember this decision if it fails ver accept this transaction vill remember this decision if it fails	
11/18/20	Kubiatowicz CS162 © UCB Fall 2020	Lec 23.13	11/18/20	Kubiatowicz CS162 © UCB Fall 2020	Lec 23.14	
	 Two-Phase Commit: Finishing Commit Transaction Coordinator learns all machines have agreed to commit Record decision to commit in local log Apply transaction, inform voters Abort Transaction Coordinator learns at least on machine has voted to abort Record decision to abort in local log Do not apply transaction, inform voters 			Two-Phase Commit: Finishing Commit Transaction • Coordinator learns all machines have agree to conflict • Record decision to commit in local log to the conflict • Apply transaction, inform voters to the conflict • Apply transaction to the conflict • Coordinator learns at lease to the conflict • Apply transaction to the conflict • Coordinator learns at lease to the conflict • Coordinator learns at lease to the conflict • Do not apply transaction to the conflict		
11/18/20	Kubiatowicz CS162 © UCB Fall 2020	Lec 23.15	11/18/20	Kubiatowicz CS162 © UCB Fall 2020	Lec 23.16	



Lec 23.19

11/18/20

Kubiatowicz CS162 © UCB Fall 2020





Blocking for Coordinator to Recover

Distributed Decision Making Discussion (1/2) A worker waiting for global decision can ask fellow workers about their Why is distributed decision making desirable? state - Fault Tolerance! - If another worker is in ABORT or A group of machines can come to a decision even if one or more of them fail COMMIT state then coordinator INIT during the process must have sent GLOBAL-* Recv: VOTE-REQ Recv: VOTE-REQ » Simple failure mode called "failstop" (different modes later) Send: VOTE-ABORT » Thus, worker can safely Send: VOTE-COMMIT abort or commit, respectively - After decision made, result recorded in multiple places READY Why is 2PC not subject to the General's paradox? Recv: GLOBAL-ABORT Recv: GLOBAL-COMMIT - If another worker is still in - Because 2PC is about all nodes eventually coming to the same decision - not INIT state then both workers ABORT сомміт necessarily at the same time! can decide to abort Allowing us to reboot and continue allows time for collecting and collating decisions - If all workers are in ready, need to BLOCK (don't know if coordinator wanted to abort or commit) 11/18/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 23.29 11/18/20 Kubiatowicz CS162 © UCB Fall 2020 Lec 23.30 Alternatives to 2PC Distributed Decision Making Discussion (2/2) Three-Phase Commit: One more phase, allows nodes to fail or block and still Undesirable feature of Two-Phase Commit: Blocking make progress. - One machine can be stalled until another site recovers: • PAXOS: An alternative used by Google and others that does not have 2PC » Site B writes "prepared to commit" record to its log, sends a "yes" vote to the blocking problem coordinator (site A) and crashes - Develop by Leslie Lamport (Turing Award Winner) » Site A crashes - No fixed leader, can choose new leader on fly, deal with failure » Site B wakes up, check its log, and realizes that it has voted "yes" on the update. - Some think this is extremely complex! It sends a message to site A asking what happened. At this point, B cannot decide to abort, because update may have committed RAFT: PAXOS alternative from John Osterhout (Stanford) - Simpler to describe complete protocol » B is blocked until A comes back - A blocked site holds resources (locks on updated items, pages pinned in · What happens if one or more of the nodes is malicious? memory, etc) until learns fate of update - Malicious: attempting to compromise the decision making - Use a more hardened decision making process: **Byzantine Agreement and Block Chains**

11/18/20

Lec 23.31



Is a BlockChain a Distributed Decision Making Algorithm?



- · BlockChain: a chain of blocks connected by hashes to root block
 - The Hash Pointers are unforgeable (assumption)
 - The Chain has no branches except perhaps for heads
 - Blocks are considered "authentic" part of chain when they have authenticity info in them
- · How is the head chosen?
 - Some consensus algorithm
 - In many BlockChain algorithms (e.g. BitCoin, Ethereum), the head is chosen by solving hard problem
 - » This is the job of "miners" who try to find "nonce" info that makes hash over block have specified number of zero bits in it
 - » The result is a "Proof of Work" (POW)
 - » Selected blocks above (green) have POW in them and can be included in chains



11/18/20

Kubiatowicz CS162 © UCB Fall 2020

```
Lec 23.35
```

11/18/20



Is a Blockchain a Distributed Decision







Packet Forwarding

- Upon receiving a packet, a router
 - read the IP destination address of the packet
 - consults its forwarding table \rightarrow output port
 - forwards packet to corresponding output port
- · Default route (for subnets without explicit entries)
 - Forward to more authoritative router



11/18/20



Lec 23.49

11/18/20

IP Addresses vs. MAC Addresses

- Why not use MAC addresses for routing?
 Doesn't scale
- Analogy
 - MAC address \rightarrow SSN
 - IP address \rightarrow (unreadable) home address
- · MAC address: uniquely associated with device for the entire lifetime of the device
- · IP address: changes as the device location changes
 - Your notebook IP address at school is different from home



IP Addresses vs. MAC Addresses

- · Why does packet forwarding using IP addr. scale?
- Because IP addresses can be aggregated
 - E.g., all IP addresses at UC Berkeley start with 0xA9E5, i.e., any address of form 0xA9E5**** belongs to Berkeley
 - Thus, a router in NY needs to keep a single entry for all hosts at Berkeley
 - If we were using MAC addresses the NY router would need to maintain an entry for every Berkeley host!!
- · Analogy:
 - Give this letter to person with SSN: 123-45-6789 vs.
 - Give this letter to "John Smith, 123 First Street, LA, US"



Lec 23.51

Setting up Routing Tables

- How do you set up routing tables?
 - Internet has no centralized state!
 - » No single machine knows entire topology
 - » Topology constantly changing (faults, reconfiguration, etc.)
 - Need dynamic algorithm that acquires routing tables
 - » Ideally, have one entry per subnet or portion of address
 - » Could have "default" routes that send packets for unknown subnets to a different router that has more information
- Possible algorithm for acquiring routing table
 - Routing table has "cost" for each entry
 - » Includes number of hops to destination, congestion, etc.
 - » Entries for unknown subnets have infinite cost
 - Neighbors periodically exchange routing tables
 - » If neighbor knows cheaper route to a subnet, replace your entry with neighbors entry (+1 for hop to neighbor)
- In reality:
 - Internet has networks of many different scales
 - Different algorithms run at different scales



- » Get person to use incorrect IP address!
- Attempt to avoid substitution attacks:
 - » Query includes random number which must be returned
- In July 2008, hole in DNS security located!
 - Dan Kaminsky (security researcher) discovered an attack that broke DNS globally
 - » One person in an ISP convinced to load particular web page, then all users of that ISP end up pointing at wrong address
 - High profile, highly advertised need for patching DNS
 - » Big press release, lots of mystery
 - » Security researchers told no speculation until patches applied

Lec 23.55

11/18/20

- Routing is limited to within a physical link (wire) or perhaps through a switch
- Our goal in the following is to show how to construct a secure, ordered, ٠ message service routed to anywhere:

Physical Reality: Packets	Abstraction: Messages		
Limited Size	Arbitrary Size		
Unordered (sometimes)	Ordered		
Unreliable	Reliable		
Machine-to-machine	Process-to-process		
Only on local area net	Routed anywhere		
Asynchronous	Synchronous		
Insecure	Secure		

11/18/20

Internet Architecture: The Five Layers

- · Lower three layers implemented everywhere
- · Top two layers implemented only at hosts

Internet Architecture: Five Layers

- · Communication goes down to physical network
- · Then from network peer to peer
- Then up to relevant layer



Summary (2/2)

- Internet Protocol (IP): Datagram packet delivery
 - Used to route messages through routes across globe
 - 32-bit addresses, 16-bit ports
- DNS: System for mapping from names⇒IP addresses
 - Hierarchical mapping from authoritative domains
 - Recent flaws discovered
- Next time: TCP: Reliable byte stream between two processes on different machines over Internet (read, write, flush)
 - Uses window-based acknowledgement protocol
 - Congestion-avoidance dynamically adapts sender window to account for congestion in network

- 1	4	14	0	10	\sim
	ь	/ 1	0	12	υ

Kubiatowicz CS162 © UCB Fall 2020

Lec 23.61