

Distributed Systems



2PC and 3PC

Continuing our consistency saga

2

- Recall from prior lectures:
 - ▣ Cloud-scale performance centers on replication
 - ▣ Consistency of replication depends on our ability to talk about notions of time.
 - Lets us use terminology like “If B accesses service S after A does, then B receives a response that is at least as current as the state on which A’s response was based.”
 - Lamport: Don’t use real clocks, use logical clocks
 - We have seen two forms, logical clocks and *vector* clocks

Next steps?

3

- We'll create a second kind of building block
 - ▣ Two-phase commit
 - ▣ It's cousin, three-phase commit
- These commit protocols (or a similar pattern) arise often in distributed systems that replicate data
- Closely tied to “consensus” or “agreement” on events, and event order, and hence replication

The Two-Phase Commit Problem

4

- The problem first was encountered in *distributed* database systems
- Suppose a database system is updating some complicated data structures that include parts residing on more than one machine
- As they execute, a “transaction” is built up in which participants join as they are contacted

... so what's the “problem”?

5

- Suppose that the transaction is interrupted by a crash before it finishes
 - Perhaps, it was initiated by a leader process L
 - By now, we've done some work at P and Q, but a crash causes P to reboot and “forget” the work L had started
 - Implicitly assumes that P might be keeping the pending work in memory rather than in a safe place like on disk
 - But this is very common, to speed things up
 - Forced writes to a disk are very slow compared to in-memory logging of information, and “persistent” RAM memory is costly
 - How can Q learn that it needs to back out?

The basic idea

6

- We make a rule that P and Q (and other participants) treat pending work as transient
 - ▣ You can safely crash and restart and discard it
 - ▣ If such a sequence occurs, we call it a “forced abort”
- Transactional systems often treat commit and abort as a special kind of keyword

A transaction

7

- L executes:

Begin

{

Read some stuff, get some **locks**

Do some **updates** at P, Q, R...

}

Commit

- If something goes wrong, executes “Abort”

Transaction...

8

- Begins, has some kind of system-assigned id
- Acquires pending state
 - ▣ Updates it did at various places it visited
 - ▣ Read and Write locks it acquired
- If something goes horribly wrong, can Abort
- Otherwise if all went well, can *request* a Commit
 - ▣ But commit can fail. This is where the 2PC and 3PC algorithms are used

The Two-Phase Commit (2PC) problem

9

- Leader L has a set of places { P, Q, ... } it visited
 - ▣ Each place may have some pending state for this xtn
 - ▣ Takes form of pending updates or locks held

- Phase 1 starts
- L asks “OK to commit?” and P, Q ... must reply
 - ▣ Each participant takes local actions to decide if it can vote in favor of commit
 - May need to set up a persistent data structure, record that 2PC is underway and save info needed to perform desired action if a commit occurs
 - “No” if something has caused them to discard the state of this transaction (lost updates, broken locks)
 - “No “ usually occurs if a member crashes and then restarts
 - ▣ No reply treated as “No” (handles failed members)

2PC (2)

10

- If a member replies “Yes,” this means it has moved to a state we call *prepared-to-commit*
 - Up to then it could just abort in a unilateral way; i.e., if data or locks were lost due to a crash/restart (or a timeout)
 - Once it says “I’m prepared to commit”, must not lose locks or data.
 - probably needs to force data to disk at this stage
 - Many systems push data to disk in background so all they need to do is update a single bit on disk: “prepared=true” but this disk-write is still considered costly event!

- Then can reply “Yes”

2PC (3)

11

- L waits and eventually has replies from {P, Q, ... }
 - ▣ L uses timer to limit wait time (duration of first phase)
- L tallies replies
 - ▣ “No” if someone replies no, or if a timeout occurs
 - ▣ “Yes” only if that participant actually replied “yes” and hence is now in the prepared-to-commit state
- Phase 1 ends

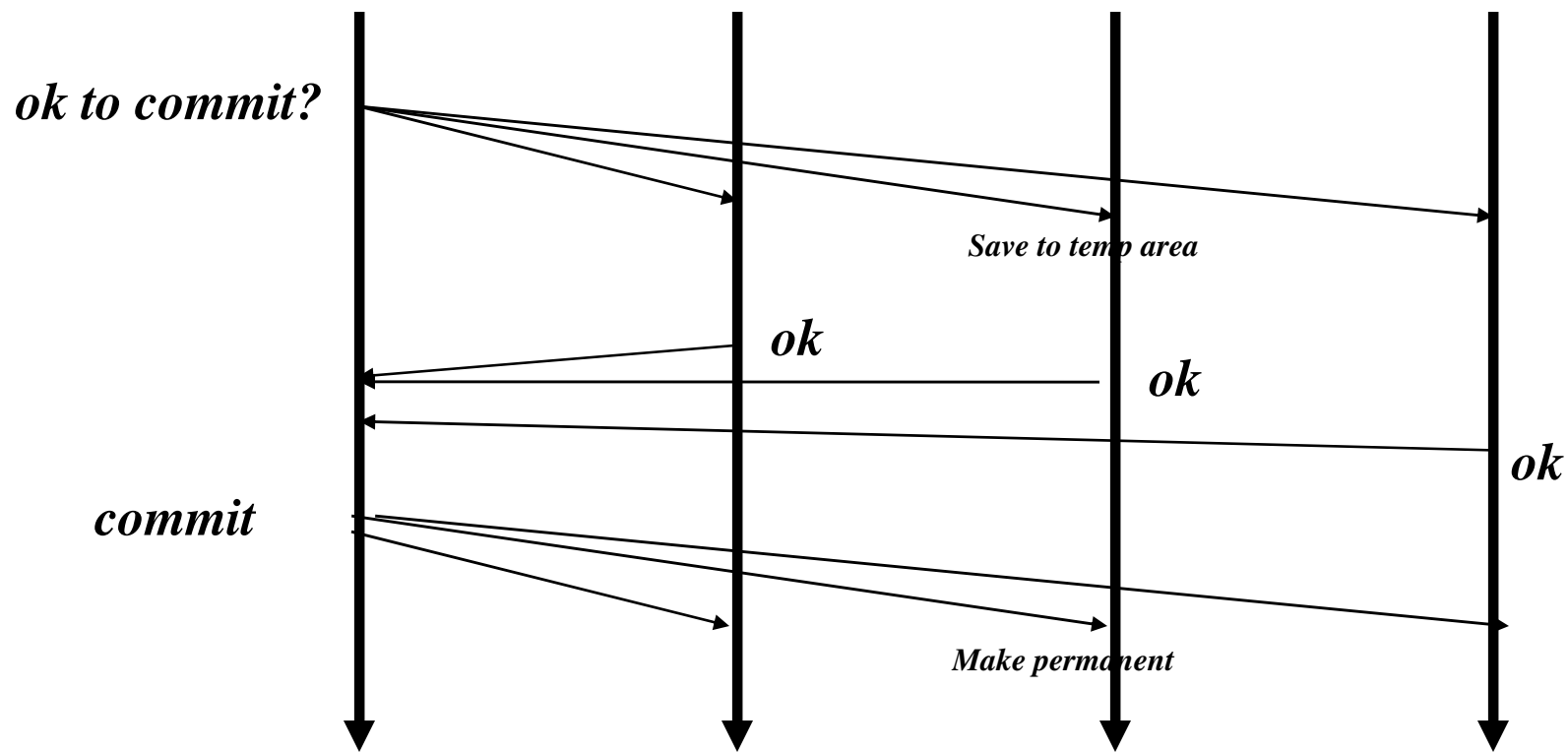
- Phase 2 starts: If all participants are prepared to commit, L sends (multicasts) a “Commit” message. Else L must send “Abort”
 - ▣ Notice that L could mistakenly abort.
 - E.g., timer in first phase goes off before all replies received
 - This is ok

2PC (4)

12

- If participant is prepared to commit, it waits for outcome to be known. On receipt from L:
 - ▣ If leader decided to Commit, participant “finalizes” the state by making updates permanent
 - ▣ If leader decided to Abort, participant discards any updates
 - ▣ Then can release locks

2PC protocol illustrated



2PC basic skeleton (no failures)

14

Coordinator:

multicast: *ok to commit?*

collect replies

all *ok* → multicast *commit*

else → multicast *abort*

Participant:

ok-to-commit →

save to temp area, reply *ok*

commit → make change permanent

abort → delete temp area

Failure cases to consider

15

- Two possible worries
 - ▣ Some participant might fail at some step of the protocol
 - ▣ The leader might fail at some step of the protocol

- Notice how a participant moves from “participating” to “prepared-to-commit” to “committed/aborted”

- Leader moves from “doing work” to “inquiry” to “committed/aborted”

- Must ensure protocol terminates with the desired all-or-nothing semantics

Can think about cross-product of states

16

- This is common in distributed protocols
 - ▣ We need to look at each member, and each state it can be in
 - ▣ The system state is a vector (S_L, S_P, S_Q, \dots)
 - ▣ Since each can be in 3 states there are 3^N possible scenarios we need to think about!

- Many protocols are actually written in a state-diagram form, but we'll use English today

Handling participant failures

17

- Suppose L stays healthy and only participants fail
- If a participant failed before voting, leader just aborts the protocol
- The participant might later recover and needs a way to find out what happened
 - ▣ If failure causes it to forget the txn, no problem
 - ▣ For cases where a participant may know about the txn and want to learn the outcome, we just keep a log of outcomes and it can look this txn up by its ID to find out
 - ▣ Writing to this log is a role of the leader (and slows it down)

Handling participant failures (2)

18

- If participant votes “Yes” and hence is prepared, but then fails
- In this case it won’t receive the Commit/Abort message
 - ▣ Solved because the leader logs the outcome
 - ▣ On recovery that participant notices that it is in *prepared-to-commit* state and consults the log
 - ▣ Must find the outcome there and must wait if it can’t find the outcome information
 - Important because apps often limit processing of new requests while in *prepared-to-commit* state

Participant recovery after failure

19

- If participant in initial *participating* state, it can always unilaterally abort
- If participant in *prepared-to-commit* state, (i.e., had voted “Yes”), it must learn the outcome and can’t terminate the txn until it does
 - ▣ E.g., must keep holding any pending updates and locks
 - ▣ Can’t release them without knowing outcome
 - ▣ Obtains outcome from L, or from the outcomes log
- If participant in *commit/abort* state, needs to complete commit/abort action even if repeatedly disrupted by failures while doing so
 - ▣ Action must be idempotent (e.g., copying a file)

2PC extended to handle participant failures

20

Coordinator:

multicast: *ok to commit?*

collect replies

all *ok* → log “commit” to “outcomes” table

multicast *commit*

else → multicast *abort*

collect *acks*

garbage-collect protocol outcome info

Participant:

ok-to-commit →

save to temp area, reply *ok*

commit → make change permanent

abort → delete temp area

After failure:

for each pending protocol “session”

contact coordinator to learn outcome

If leader tracks protocol outcome until all participants are known to have completed commit or abort actions → EXTRA PHASE needed to collect acks from participants

Handling coordinator failure

21

- Suppose a participant P votes “Yes” but then leader L seems to vanish
 - ▣ Maybe it died... maybe became disconnected from the system (partitioning failure)
 - ▣ P is “stuck”. We say that it is “blocked”

- Can P deduce the state?
 - ▣ If log (stored not on L) reports outcome, P can make progress
 - As long as we follow rule that L logs outcome before telling anyone, it is safe to commit in this case
 - ▣ What if the log doesn't know the outcome or there is no log?
 - Could ask other participants

Handling coordinator failure

22

- When participant P enters *prepared-to-commit* state sets timer
 - ▣ On timeout, seeks to complete protocol on its own
 - ▣ If P was told the list of participants when L contacted it for its vote, P could poll them
 - ▣ E.g. P asks Q, R, S... “what state are you in?”

- Suppose someone says “*commit*” or “*abort*”?
 - ▣ Now P can just commit or abort!
 - ▣ Repeats second phase of original protocol (so all can commit/abort)

- But what if N-1 say “*prepared-to-commit*” and 1 is inaccessible?

P remains blocked in this case

23

- L plus one member, perhaps S, might know outcome
- P is unable to determine what L decided
- Worse possible situation: L is both leader and also participant and hence a single failure leaves the other participants blocked!

2PC extended – on leader/net failure, participants try to terminate protocol without blocking

24

Coordinator:

multicast: *ok to commit?*

collect replies

all *ok* → log “commit” to “outcomes” table
wait until safe on persistent store
multicast *commit*

else → multicast *abort*

collect *acks*

After failure:

for each pending protocol “session” in outcomes table
send outcome (commit or abort)
wait for acks

Periodically:

query each process: *terminated protocols?*
determine *fully terminated protocol sessions*
2PC to garbage-collect protocol outcome info

Participant: first time msg received

ok-to-commit →

save to temp area, reply *ok*

commit → log outcome, make change permanent

abort → log outcome, delete temp area

Message is a duplicate (recovering coordinator)

send *ack*

After failure:

for each pending protocol “session”
contact coordinator to learn outcome

After timeout in *prepare-to-commit* state:

query other participants about state

outcome can be deduced →

run coordinator-recovery protocol

outcome uncertain →

must wait

2PC -- version from previous page

25

- Gains higher availability at cost of more communication
 - ▣ Participants sometimes can terminate even if coordinator down
 - ▣ Must garbage collect outcomes held for sessions that have terminated at all participants
- Can still block!
 - ▣ If failure of both coordinator and a participant occurs during decision stage

Skeen & Stonebraker: 3PC

26

- Skeen proposed a 3PC protocol, that adds one step (and omits any log service)
- With 3PC the leader runs 3 rounds of communication:
 - ▣ “Are you able to commit”? Participants reply “Yes/No”
 - ▣ “Abort” or “Prepare to commit”. They reply “OK”
 - ▣ “Commit”
- Notice that Abort happens in round 2 but Commit only can happen in round 3
- Ensures state of system can be deduced by subset of processes, provided they can communicate reliably

State space gets even larger!

27

- Now we need to think of 5^N states
 - But Skeen points out that many can't occur
 - For example we can't see a mix of processes that are in the Commit and Abort state
 - We could see some in "Participating" and some in "Yes"
 - We could see some in "Yes" and some in "Prepared"
 - We could see some in "Prepared" and some in "Commit"
 - But by pushing "Commit" and "Abort" into different rounds we reduce uncertainty

3PC recovery is complex

28

- Skeen shows how, on recovery, we can poll the system state
- Any (or all) processes can do this
- Can always deduce a safe outcome... provided that we have
 - ▣ Fail-stop failures (processes fail only by crashing)
 - ▣ Failures are *accurately detectable* operational processes
- 3PC, without any log service, and with accurate failure detection is guaranteed to be non-blocking

3PC outline (garbage collection not shown)

29

Coordinator:

multicast: *ok to commit?*

collect replies

all *ok* → log “precommit”

multicast *precommit*

else → multicast *abort*

collect *acks* from non-failed participants

all *ack* → log “commit”

multicast *commit*

collect *acks*

garbage-collect protocol outcome info

Participant: logs state on each message

ok-to-commit →

save to temp area, reply *ok*

precommit →

enter precommit state, *acknowledge*

commit →

make change permanent

abort →

delete temp area

After failure:

collect participant state information

any precommit or committed →

push forward to commit

else →

push back to abort

3PC

30

- 3PC, without any log service, and with accurate failure detection is guaranteed to be non-blocking
- However, in actual systems
 - ▣ inaccurate failure detections and net partitions are very possible
 - ▣ hence, in reality 3PC is blocking

Failure detection in a network

31

- Many think of Skeen's 3PC as a practical protocol

- But to really use 3PC we would need a perfect failure detection service that never makes mistakes
 - ▣ It always says "P has failed" if, in fact, P has failed
 - ▣ And it never says "P has failed" if P is actually up

- Is it possible to build such a failure service?

Notions of failure

32

- This leads us to think about failure “models”

Best: “Fail-stop” with trusted notifications

Many things can fail in a distributed system

- ▣ Network can drop packets, or the O/S can do so
- ▣ Links can break causing a network partition that isolates one or more nodes
- ▣ Processes can fail by halting suddenly
- ▣ A clock could malfunction, causing timers to fire incorrectly
- ▣ A machine could freeze up for a while, then resume
- ▣ Processes can corrupt their memory and behave badly without actually crashing
- ▣ A process could be taken over by a virus and might behave in a malicious way that deliberately disrupts our system

Worst: Byzantine

What do “real” systems do?

33

- Linux and Windows use timers for failure detection
 - ▣ These can fire even if the remote side is healthy
 - ▣ So we get “inaccurate” failure detections
 - ▣ Of course many kinds of crashes can be sensed accurately so for those, we get trusted notifications

- Some applications depend on TCP, but TCP itself uses timers and so has the same problem

Aside: Byzantine case

34

- Much debate around this
 - ▣ Do we design systems assuming a Byzantine failure model?

- Since programs are buggy (always), it can be appealing to just use a Byzantine model. A bug gives random corrupt behavior... like a mild attack

- But Byzantine model is hard to work with and can be costly (you often must “outvote” the bad process)

Failure detection in a network

35

- Return to our use case
- 2PC and 3PC are normally used in standard Linux or Windows systems with timers to detect failure
 - ▣ Hence we get *inaccurate* failure sensing with possible mistakes (e.g. P thinks L is faulty but L is fine)
 - ▣ **3PC is thus also blocking** in this case, although less likely to block than 2PC
 - ▣ Can prove that any commit protocol would have blocking states with inaccurate failure detection

World-Wide Failure Sensing



36

- Vogels wrote a paper in which he argued that we really could do much better than using just timers
 - ▣ In a cloud computing setting, the cloud management system often “forces” slow nodes to crash and restart
 - Used as a kind of all-around fixer-upper
 - Also helpful for elasticity and automated management

The Postman Always Rings Twice

37

- Vogels suggests that there are many reasons a machine might timeout and yet not be faulty
- Suppose the mailman wants to see you...
 - He rings and waits a few seconds
 - Nobody comes to the door... should he assume you've died?

- Hopefully not



Causes of delay

38

- Scheduling can be sluggish
- A node might get a burst of messages that overflow its input sockets and triggers message loss, or network could have some kind of malfunction in its routers/links
- A machine might become overloaded and slow because too many virtual machines were mapped on it
- An application might run wild and page heavily

Vogels suggests?



39

- He recommended that we add some kind of failure monitoring service as a standard network component
- Instead of relying on timeout, even protocols like remote procedure call (RPC) and TCP would ask the service and it would tell them
- It could do a bit of sleuthing first... e.g. ask the O/S on that machine for information... check the network...

Why clouds *don't* do this



40

- Hamilton: In the cloud our focus tends to be on keeping the “majority” of the system running
 - ▣ No matter what the excuse it might have, if some node is slow it makes more sense to move on
 - ▣ Keeping the cloud up, as a whole, is way more valuable than waiting for some slow node to catch up
 - ▣ End-user experience is what counts!

- So the cloud is casual about killing things
- ... and avoids services like “failure sensing” since they could become bottlenecks



Also, most software is buggy!

41

- A mix of “Bohrbugs” and “Heisenbugs”
 - ▣ Bohrbugs: Boring and easy to fix. Like Bohr model of the atom
 - ▣ Heisenbugs: They seem to hide when you try to pin them down (caused by concurrency and problems that corrupt a data structure that won’t be visited for a while). Hard to fix because crash seems unrelated to bug
- Studies show that pretty much all programs retain bugs over their full lifetime.
 - ▣ So if something is acting strange, it may be failing!

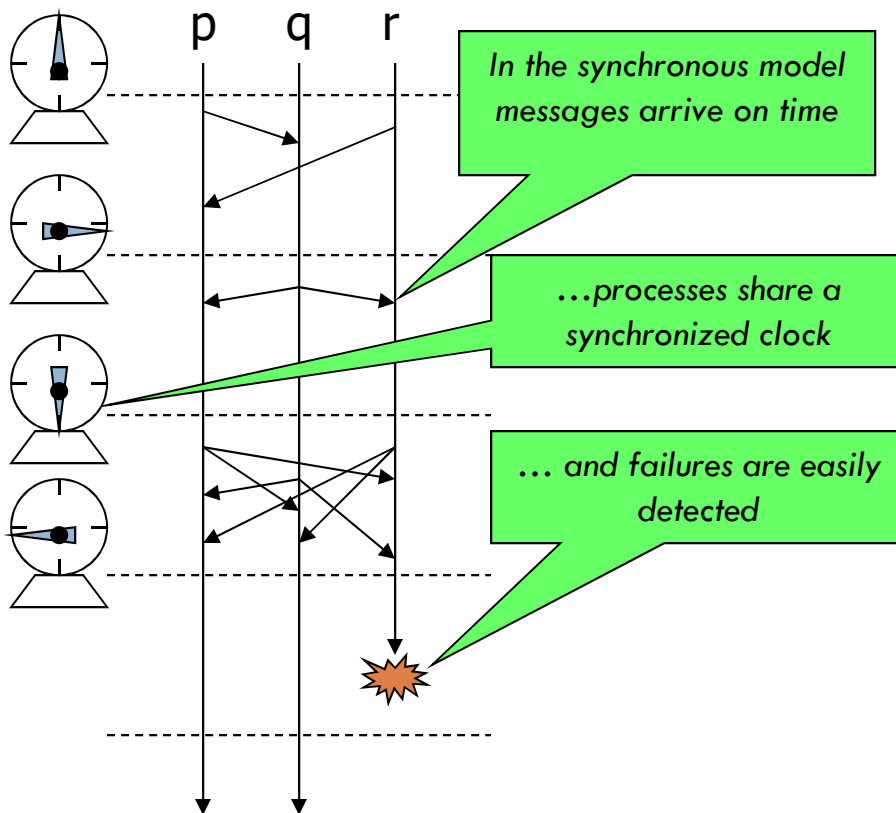
Worst of all... timing is flakey

42

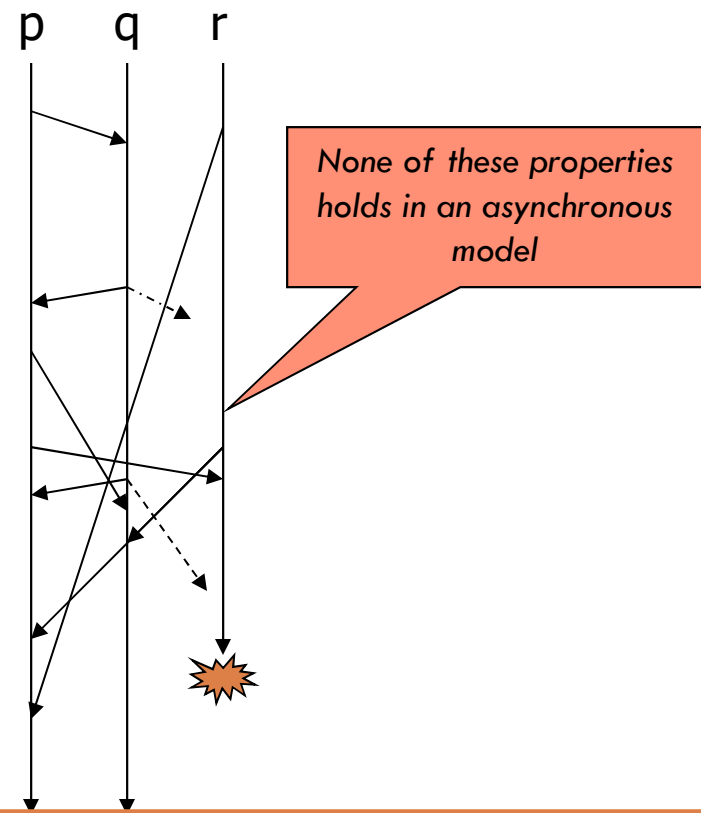
- At cloud-scale, with millions of nodes, we can't trust timers at all
- Too many things can cause problems that manifest as timing faults or timeouts
- There are some famous models... and none is ideal for describing real systems

Synchronous and Asynchronous Executions

43



Failures are very easy to detect.



Impossible to detect difference between a crashed node and network delay.

Reality: neither model applies

44

- Real distributed systems aren't synchronous
 - ▣ Although a flight control computer can come close
- Nor are they asynchronous
 - ▣ Software often treats them as asynchronous
 - ▣ In reality, clocks work well... so in practice we often use time cautiously and make assumptions about upper bounds on message delays
- For our purposes we usually start with an asynchronous model
 - ▣ Subsequently enrich it with sources of time when useful.
 - ▣ We sometimes assume a “public key” system. This lets us sign or encrypt data where need arises

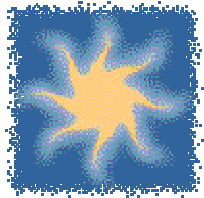
Thought problem

45

- Harry and Sally will meet for lunch. They'll eat in the cafeteria unless both are sure that the weather is good
 - ▣ Sally's cubicle is in the basement, so Harry will send email
 - ▣ Both have lots of meetings, and might not read email. So she'll acknowledge his message.
 - ▣ They'll meet inside if one or the other is away from their desk and misses the email.
- Harry sees sun. Sends email. Sally acks. Can they meet outside?

Meeting for lunch

46



Harry

Sally

H: Sally, the weather is beautiful! Let's meet at the sandwich stand outside.

S: I can hardly wait. I've been in this dungeon studying and haven't seen the sun in weeks!

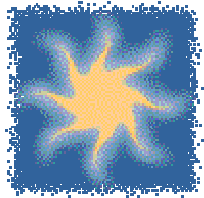
They eat *inside!* Harry reasons:

47

- “Sally sent an acknowledgement but doesn’t know if I read it
- “If I didn’t get her acknowledgement I’ll assume she didn’t get my email
- “In that case, I’ll go to the cafeteria
- “She’s uncertain, so she’ll meet me there

Harry had better send an Ack

48



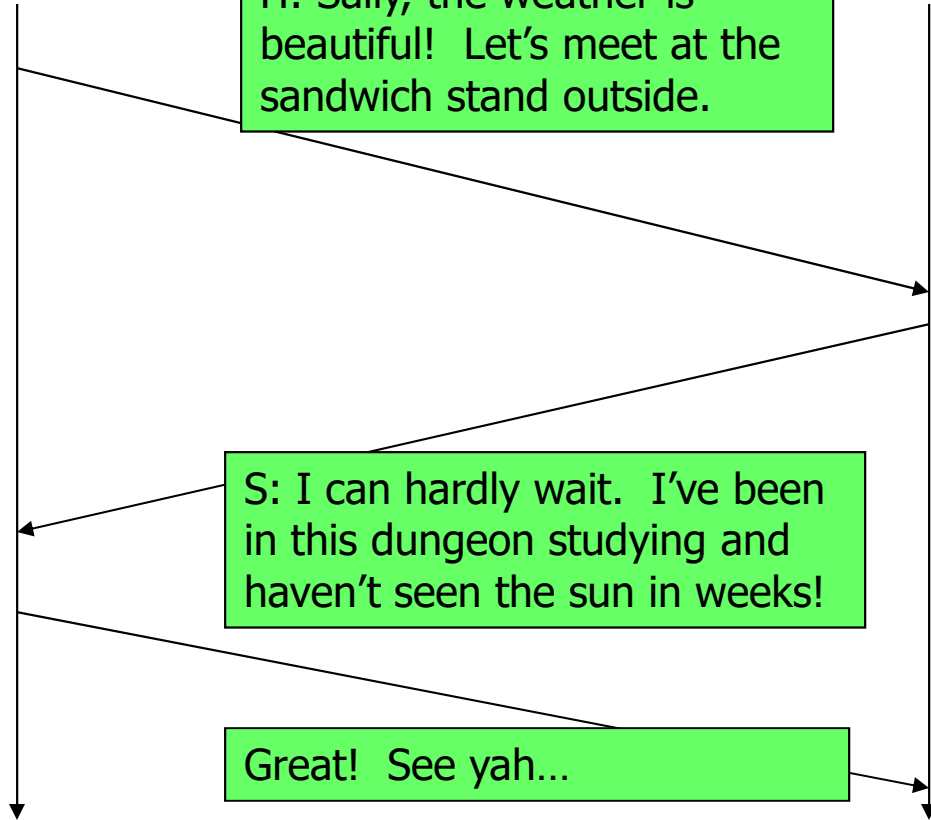
Harry

Sally

H: Sally, the weather is beautiful! Let's meet at the sandwich stand outside.

S: I can hardly wait. I've been in this dungeon studying and haven't seen the sun in weeks!

Great! See yah...



Why didn't this help?

49

- Sally got the ack... but she realizes that Harry won't be sure she got it
- Being unsure, he's in the same state as before
- So he'll go to the cafeteria, being dull and logical. And so she meets him there.

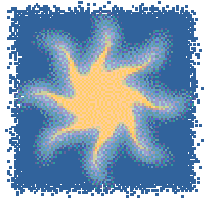
New and improved protocol

50

- Sally sends an ack. Harry acks the ack. Sally acks the ack of the ack....
- Suppose that noon arrives and Sally has sent her 117'th ack.
 - ▣ Should she assume that lunch is outside in the sun, or inside in the cafeteria?

How Harry and Sally's romance ended

51



Harry

Sally

H: Sally, the weather is beautiful! Let's meet at the sandwich stand outside.

S: I can hardly wait. I've been in this dungeon studying and haven't seen the sun in weeks!

Great! See yah...

Yup...

Got that...

...

Oops, too late for lunch

Maybe tomorrow?

Moral of the story?

52

- ~~❑ Logicians are dull people and have miserable lives.~~
- ❑ The real world demands leaps of faith: pure logic isn't enough.
- ❑ For our computing systems, this creates a puzzle, since software normally behaves logically!

How do real people meet for lunch?

53

- They send one email, then go outside
 - ▣ Sally doesn't need an ack in the original protocol to ensure Harry hasn't failed in some way; she just assumes he got her reply
 - ▣ Mishaps happen, now and then, but we deal with those.
 - ▣ In fact we know perfectly well that we can't achieve perfect agreement, and we cope with that
 - ▣ In some sense a high probability of meeting outside for lunch is just fine and we don't insist on more

Things we just can't do (in distributed systems)

54

- We can't detect failures in a trustworthy, consistent manner
- We can't reach a state of "common knowledge" concerning something not agreed upon in the first place
- We can't guarantee agreement on things (election of a leader, update to a replicated variable) in a way certain to tolerate failures

Back to 2PC and 3PC

55

- Summary of the state of the world
 - 3PC would be better than 2PC in a perfect world
 - In the real world, 3PC is more costly (extra round) but blocks just the same (inaccurate failure detection)
 - Failure detection tools could genuinely help but the trend in large data centers has been in the opposite direction