# Kyle Kingsbury

## Distsys Safety Tester

# In a nutshell

3 liveness (crash, delayed messages)

7 safety (dups, last write, split-brain, aborted read, G1c)

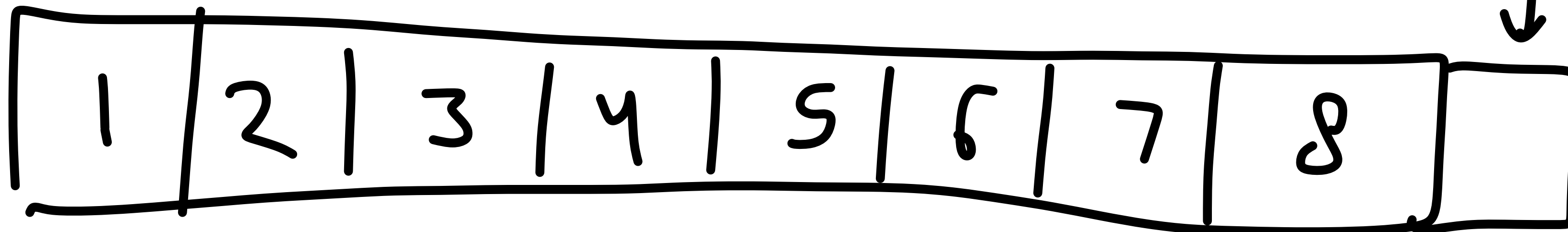2 ambiguous (G0, G1c, indefinite error code)

# Redpanda

# Fundamentals

# Distributed Log, ala kafka

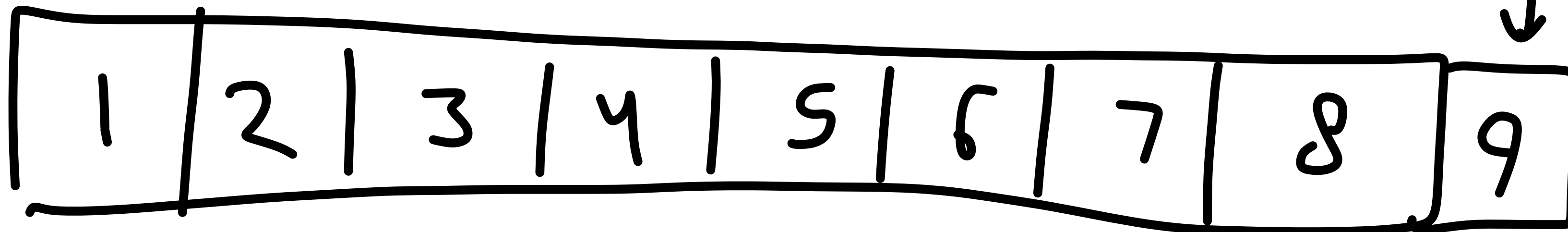- API compatible
- Aiming for better performance
- No Zookeeper

# A    Partition

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | |

Producer.send( 9 )

position

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |

consumer.poll()

position

1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

consumer.poll() → [1,2,3]

position

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |

consumer.poll()

position

1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

consumer.poll() → [4, 5, 6, 7, 8, 9]

position
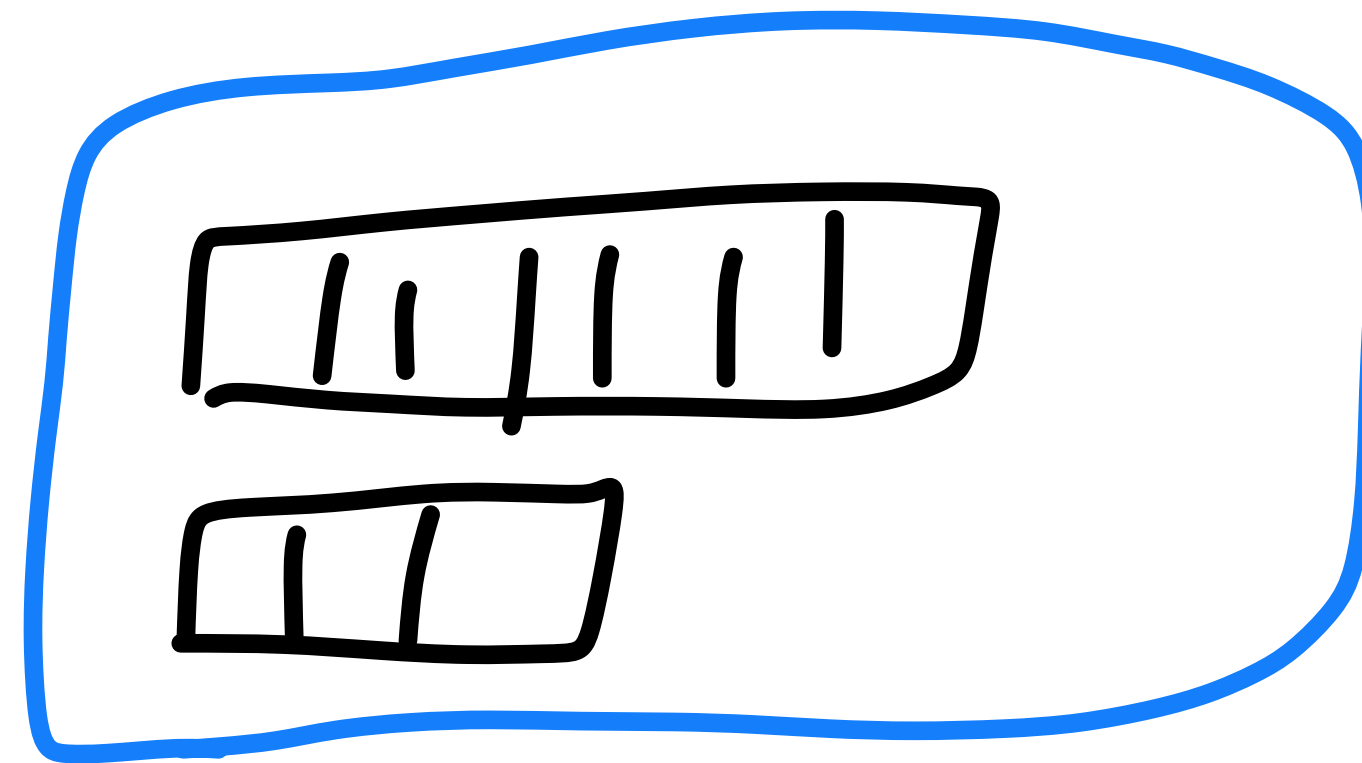
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |

Consumer.poll() → [ ]

A Topic

# Redpanda

rainfall $\rightarrow$



harvests $\rightarrow$

$C_1$ $C_2$ $C_3$

Consumer Group

poll

$C_1$    $C_2$    $C_3$

Consumer Group

Commit

$C_1$ $C_2$ $C_3$

Consumer Group

poll

$C_1$ $C_2$ $C_3$

Consumer Group

Consumer Group

$C_1$  $C_2$  $C_3$

poll

$C_1$  $C_2$  $C_3$

Consumer Group

Commit

$C_1$ $C_2$ $C_3$

Consumer Group

Poll

$C_1$ $C_2$ $C_3$

Consumer Group

Test Design

Jepsen (JVM)

n1

n2

n3

Redpanda Cluster

Queue

# History

invoke    subscribe    [ k ]

ok    subscribe    [ k ]

# History

invoke   subscribe   [ k ]

ok   subscribe   ( k )

invoke   poll   [ poll ]

ok   poll   [ poll   { $k_1$   [[0:a] [1:b]]} ]

# History

invoke    subscribe    [k]

ok    subscribe    [k]

invoke    poll    [poll]

ok    poll    [poll  {k_1  [[0:a]  [1:b]]}]

invoke    send    [send  k_2 : c]

ok    send    [send  k_2  [3 : c]]

# History

ok txn [[poll {$k_1$ [[0:a] [1:b]]}],
[send $k_2$ [3:c]]]

Checker

fail $[[send \ x \ a]]$     OK $[[poll \ \{x([[],a])\}]]$

fail $[[$send $x$ $@$$]]$    OK $[[$poll $\{x([\exists, @])\}]]$

Aborted Read!

(G1a)

$[send \ k, \ [3, 6]]$



$k_1:$

6

"Version order"

$[pall \ \{k_1 \ ([0 \ a] \ [2 \ b] \ [3 \ c])\}]$



$k_1:$

| a |  | b | b/c |
|---|---|---|---|

Ok! We might not obscure all

$k_1$:

| a | | b | b/c |

$k_1:$

duplicate!

a      b   b/c

$k_1$:

inconsistent offset!

"what affect is a?"

$k_1:$

a | | b | b/c

⊙

"What value got offset 2?"

$k_1$:

| a | | b | b/c |
|---|---|---|---|
| 0 | 1 | 2 | |

"Which value came just after a?"

$k_1$:

a | | | ⓑ | b/c

0     1     2

$k_i:$ | a | b | c |

$$[\text{pall } \{k_1 : [[0, a]]\}], \quad [\text{pall } \{k_1 : [[2, c]]\}]$$

$$k_i : \boxed{a \mid b \mid c}$$

$[\text{pall } \{k_1 : [[0, a]]\}]$ , $[\text{pall } \{k_1 : [[2, c]]\}]$

skip!

$k_i :$

| a | b | c |
|---|---|---|

$$[\text{pall } \{k_1 : [[1, b]]\}] \quad , \quad [\text{pall } \{k_1 : [[0, a]]\}]$$

non-monotonic!

$$k_1 : \boxed{a \mid b \mid c}$$

$[\text{send } k, b]$ , $[\text{send } k, a]$

non-monotonic!

$k_i$: | a | b | c |

| a | b | c | d | e | f | g | h | i | ... |

| a | b | c | d | e | f | g | h | i | ... |

Polled by $c_1$

Polled by $c_2$

Unseen

| a | b | c | d | e | f | g | h | i | ... |

Polled by $c_1$

Polled by $c_2$

$T_1$

$T_2$

$([\text{send } k_1 \quad b])$ $([\text{poll } \{k_1 [... [1, b]...]\}])$

$K_1:$

| a | b | c |
|---|---|---|

$$T_1 \xrightarrow{\text{write-read}} T_2$$

$$\left(\left[\text{send } k_1 \; \textcircled{b}\right]\right) \quad \left(\left[\text{poll } \{k_1 \; [\ldots [1, \textcircled{b}] \ldots]\}\right]\right)$$

$K_1$:

| a | b | c |
|---|---|---|

$$T_1$$
$$[[\text{send } k, \ b]]$$

$$T_2$$
$$[[\text{send } k, \ c]]$$

$K_1$:

| a | b | c |
|---|---|---|

$T_1$ — Write-write → $T_2$

$\Big[\big[\text{send } k_1 \ b\big]\Big]$    $\Big[\big[\text{send } k_1 \ c\big]\Big]$

$k_1:$

| a | b | c |
|---|---|---|

$T_1$ GO $T_2$

ww

ww

Write Cycle

$$T_1 \xrightarrow{ww} T_2$$

$$Glc$$

$$T_1 \xleftarrow{wr} T_2$$

Circular
Information
Flow

Faults

- Crash
- Pause
- Clock skew
- Network Partition

- Membership change

# Results

# #1 Duplicate Writes by Default

producer.send(a)

# #1 Duplicate Writes by Default

producer.send(a)

.
.
.

(retry)

ok!

| a | | |

| (a) | b | c | (a) |

# #1 Duplicate Writes by Default

$$enable.idempotence = true$$

$$(default \quad in \quad 3.0.0 \quad client)$$

# #1 Duplicate Writes by Default

producer.send(a)



.
.
.
.

(retry)

ok!

| a | b | c | | | |

no dup!

# #1 Duplicate Writes by Default

Except... we still saw dups
w/ any crash, pause, partition...

# #1 Duplicate Writes by Default

Partly a client bug!

Redpanda server needs:

- enable_idempotence true
- id_allocator_replication 3

# #1 Duplicate Writes by Default

Should be fixed by enabling

idempotence by default

Planned for 22.1.1

# #3039 Duplicate Writes Not by Default

We <u>still</u> saw duplicate messages

| a | b | c | b | c | d |

# #3039 Duplicate Writes Not by Default

send (a)
:
(retry)

a

a | b

Out Of Order Sequence Exception

# #3039 Duplicate Writes Not by Default

send (a)
    ⋮
(retry)

send (a) → `| a |`

(retry) → `| a | b |`

**Out Of Order Sequence Exception**

(retry) → `| (a) | b | (a) |`

# #3039 Duplicate Writes Not by Default

- 21.10.2 : Fixed by adding deduplication

- 21.10.3 : Performance improvements

# #3335 Assert Failure Deallocating Partitions

With membership changes...

Assert failure... allocated partitions >
allocation capacity

Still investigating

# 3336 Assert failure in response.partition_index

w/process crashes...

Assert failure: ... Response $ current partition ids have to be the same. Current response 0, update 1

Patched → 22.1.1

# #3003 Incansistent Offsets

w/ process crashes or

network partitions...

# #3003 Inconsistent Offsets

Caused by Raft state machine applying uncommitted log entries

→ Fixed in 21.10.3 by waiting for commit pointer to advance first

# #6 Lost / Stale Messages

W/ process pause

producer

| a | b | c | d ) |

? ?

Consumer

# #6 Lost /Stale Messages

Consumer.poll()
→ OK! Messages: [ ]

Consumer.poll()
→ OK! Messages: [ ]

Consumer.poll()
→ OK! Messages: [ ]

# #6 Lost/Stale Messages



21.11.2 queue assign acks=all retries=0 aor=earliest membership,partition unseen

# #6 Lost/Stale Messages

Still on disk... just not being delivered...

Redpanda still investigating!

# #7 Aborted Read — NotLeaderOrFollowerException

send (a) $\longrightarrow$

NotLeaderOrFollowerException

# #7 Aborted Read — NotLeaderOrFollowerException

"This failed, right?"

"I think so..."

# #7 Aborted Read — Not Leader Or Follower Exception

Broker returns this error if a request could not be processed because the broker is not the leader or follower for a topic partition. This could be a transient exception during leader elections and reassignments. For Produce and other requests which are intended only for the leader, this exception indicates that the broker is not the current leader.

Kafka team:

"Actually, it's indeterminate!

Could have succeeded!"

# #7 Aborted Read — NotLeaderOrFollowerException

Just needs docs!

# #8 Write Cycles

$T_1$: send(a)

ok!

send(b)

ok!

Commit

ok!

$T_2$:

send(c)

ok!

Commit

ok!

# #8 Write Cycles

$T_1$: send(a)

send(b)

$T_2$:
send(c)

| a | c | b |
|---|---|---|

# #8   Write Cycles

$T_1$: send(a)

send(b)

ww →

← ww

$T_2$:
send(c)

| a | c | b |

# #8 Write Cycles

GO happens constantly in healthy

Redpanda & Kafka clusters...

Forbidden by Adya read-uncommitted!

# #8 Write Cycles

GO happens constantly in healthy

Redpanda & Kafka clusters...

Forbidden by Adya read-uncommitted!

# #8 Write Cycles



Commit order: X2 < X0

Since X2 is committed first, each partition will expose messages from X2 before X1.

Consumer processing order

p0: Mx-p0, X2: M1-p0, X2: M4-p0, X1: M1-p0, X1: M2-p0, X1: M3-p0, X1: M5-p0

p1: My-p1, X2: M2-p1, X2: M3-p1, X1: M4-p1, X1: M6-p1

# #8 Write Cycles

Timeline (top to bottom labels along time axis):
- X1: BEGIN
- X1: M1-p0
- X2: BEGIN
- X2: M1-p0
- X1: M2-p0
- X1: M3-p0
- X2: M2-p1
- X1: M4-p1
- Mx-p0
- X2: M3-p1
- My-p1
- X2: M4-p0
- X1: M5-p0
- X2: COMMIT
- X1: M6-p1
- X1: COMMIT

time →

Commit order: X2 < X0

Since X2 is committed first, each partition will expose messages from X2 before X1.

Consumer processing order

p0: Mx-p0, X2: M1-p0, X2: M4-p0, X1: M1-p0, X1: M2-p0, X1: M3-p0, X1: M5-p0

p1: My-p1, X2: M2-p1, X2: M3-p1, X1: M4-p1, X1: M6-p1

Liiiiiieeees

Since X2 is committed first each partition will expose messages from X2 before X1

"Transactional Messaging in Kafka"

we have no strong evidence that applications can benefit from the commit order option, we opted for not implementing it

"Exactly Once Delivery and Transactional Messaging in Kafka"

In theory, theory and practice are the same. In practice, they are not

Kafka protocol
permits GO
(tx interleaving)

Seattle

NYC

Seattle

NYC

2 x 72 ms

Seattle
NYC

Seattle

NYC

# Non-Monotonic Snapshot Isolation: scalable and strong consistency for geo-replicated transactional systems

Masoud Saeida Ardekani
Université Pierre-et-Marie Curie (UPMC-LIP6)
Paris, France
masoud.saeida-ardekani@lip6.fr

Pierre Sutra
University of Neuchâtel
Neuchâtel, Switzerland
pierre.sutra@unine.ch

Marc Shapiro
INRIA & UPMC-LIP6
Paris, France
http://lip6.fr/Marc.Shapiro/

*Abstract*—Modern cloud systems are geo-replicated to improve application latency and availability. Transactional consistency is essential for application developers; however, the corresponding concurrency control and commitment protocols are costly in a geo-replicated setting. To minimize this cost, we identify the following essential scalability properties: (i) only replicas updated by a transaction $T$ make steps to execute $T$; (ii) a read-only transaction never waits for concurrent transactions and always commits; (iii) a transaction may read object versions committed after it started; and (iv) two transactions synchronize with each other only if their writes conflict. We present Non-Monotonic Snapshot Isolation (NMSI), the first strong consistency criterion to allow implementations with all four properties. We also present a practical implementation of NMSI called Jessy, which we

Transactional YCSB (on 4 sites)
Varying the update/read-only proportion from 10%/90% to 30%/70%
Update Transactions: 1 Read/ 1 Write
Read-only Transactions: 2 Reads

1) all-or-nothing for multi-partitions writes

2) atomic consume-transform-produce loop

3) ~~consistent snapshots~~

# #8 Write Cycles

This is just how Kafka/ Redpanda txns work

→ Needs documentation

# #3036 Aborted Reads $ Circular Information Flow

poll()
OK

send(a)
ok

addOffsets
fail

poll
[a]

# #3036 Aborted Reads & Circular Information Flow

poll()
ok
send (a)
ok
addOffsets
(fail)

G1a: Aborted Read

Happened often in healthy clusters!

Poll
(a)

# #3036  Aborted Reads & Circular Information Flow



$T_1:$ send $(a)$ → [a]

ok

poll() → [a,b]

[a | b]

$T_2:$ send $(b)$ → ok

poll() → [a]

# #3036 Aborted Reads & Circular Information Flow



T1: send(a)

poll()

[a,b]

Wr

Wr

T2: send(b)

poll()

[a]

# #3036 Aborted Reads & Circular Information Flow

$T_1$ $\xrightarrow{wr}$ $T_2$

$T_1$ $\xleftarrow{wr}$ $T_2$

G1c: Circular Information Flow

Illegal in Adya Read Committed!

# #3036 Aborted Reads $ Circular Information Flow

- Off-by-one error allowed LSO to advance beyond committed offsets

- #3232 accidentally aborted more txns than necessary

Fixed in 21.10.2

# #10 Internal Non-monotonic Polls

```
beginTxn()
poll()
poll()
addOffsetsToTxn(...)
commitTxn()
```

# #10 Internal Non-monotonic Polls

```
beginTxn()
poll() → [a, b, c]
poll() → [b, c, d]
addOffsetsToTxn(...)
commitTxn()
```

# #10 Internal Non-monotonic Polls

- Caused by consumer rebalance

- Clients **must** register rebalance listeners to detect potential non-monotonicity in transactions

# #3616-a    Aborted   Read   w/ Invalid Txn State Exception

begin Txn()

send (a)

commitTxn() -> Invalid Txn State Exception

# 3616 -a   Aborted   Read   w/ Invalid Txn State Exception

beginTxn()

send (a)

commitTxn() → Invalid Txn State Exception

:
:

poll() → [a]

# #3616-a   Aborted   Read   w/ Invalid Txn State Exception

Pause, crash, or partition

led to G1a

# #3616-a  Aborted Read  w/ Invalid Txn State Exception



begin Txn()

send (a)

commit Txn()

(retry)

Broker    Coordinator

commits!

X

Coordinator

???

Invalid Txn State Exception

# #3616-a   Aborted   Read  w/ Invalid Txn State Exception

Fixed  by   returning   UnknownServerError

January  21,  2022

⮡ Released  in  21.11.15

# #3616-6 Last Transactional Writes

beginTxn()

send (a)         -> Offset 2

commitTxn() -> OK

# #3616 - 6   Last Transactional Writes

beginTxn()

send (a)        -> Offset 2

commitTxn() -> OK

---

poll() ->

| b | c |   | d |
|---|---|---|---|
| 0 | 1 | 2 | 3 |

# 3616-6   Lost Transactional Writes

beginTxn()

send(a)        -> Offset 2

commitTxn() -> OK

---

poll() ->

# #3616-6   Last Transactional Writes

Happened in healthy clusters

on  21.11.2

# #3616-6   Last Transactional Writes

```
lock.acquire()
if ( ! log.sync() )  {
    return;
}
do_x();
do_y();
lock.release()
```

Critical section

# 3616-6    Last Transactional Writes

```
lock.acquire()
if (! log.sync()) {
    return;
}
do_x();
do_y();
lock.release()
```

Critical section

Node loses, then regains leadership — other leader did things in our critical section!

# #3616-6  Last Transactional Writes

```
lock.acquire()
if (! log.sync()) {
    return;
}
term = _insync_term
do_x(term);
do_y(term);
lock.release()
```

Critical section

# 3616-6   Last Transactional Writes

Fixed   January 21, 2022

   ↳ Released in 21.11.15

# Discussion

| # | Issue | Fault | Resolved |
|---|---|---|---|
| 1 | Duplicate writer by default | P ∨ C ∨ Part | 22.1.1 |
| 3039 | Duplicate writer w/ idempotence | P ∨ C ∨ Part | 21.10.3 |
| 3335 | Assert crash deallocating partitions | M | |
| 3336 | Assert crash w/ partition IDs | C | 22.1.1 |
| 3003 | Inconsistent offsets | C ∨ Part | 21.10.3 |
| 6 | Lost/stale messages | P ∨ C ∨ Part | |
| 7 | Aborted read — NotLeaderOrFollower | M ∧ P | |
| 8 | Write cycles | ∅ | |
| 3036 | Aborted read/Circular information flow | ∅ | 21.10.2 |
| 10 | Internal non-monotonic polls | ∅ | |
| 3616-a | Aborted read — Invalid Txn state | P ∨ C | 21.11.15 |
| 3616-b | Lost txn writes | ∅ | 21.11.15 |

| # | Issue | Fault | Resolved |
|---|---|---|---|
| 1 | Duplicate writer by default | P ∨ C ∨ Part | 22.1.1 |
| 3039 | Duplicate writer w/ '' ' | P.. C ∨ Part | 21.10.3 |
| 3335 | Assert crash | | |
| 3336 | Assert crash | | 22.1.1 |
| 3003 | Inconsistent ot. | , art | 21.10.3 |
| 6 | Lost/stale messages | P ∨ C ∨ Part | |
| 7 | Aborted read — NotLeaderOrFollower | M ∧ P | |
| 8 | Write cycles | ∅ | |
| 3036 | Aborted read/Circular information flow | ∅ | 21.10.2 |
| 10 | Internal non-monotonic polls | ∅ | |
| 3616-a | Aborted read — Invalid Txn state | P ∨ C | 21.11.15 |
| 3616-b | Lost txn writes | ∅ | 21.11.15 |

Need Docs

| # | Issue | Fault | Resolved |
|---|-------|-------|----------|
| 1 | Duplicate writer by default | P ∨ C ∨ Part | 22.1.1 |
| 3039 | Duplicate write | ∿ ∨ Part | 21.10.3 |
| 3335 | Assert crash | | |
| 3336 | Assert crash | | 22.1.1 |
| 3003 | Inconsistent offset | ∿ ∨ Part | 21.10.3 |
| 6 | Lost/stale messages | P ∨ C ∨ Part | |
| 7 | Aborted read — NotLeaderOrFollower | M ∧ P | |
| 8 | Write cycles | ∅ | |
| 3036 | Aborted read/Circular information flow | ∅ | 21.10.2 |
| 10 | Internal non-monotonic polls | ∅ | |
| 3616-a | Aborted read — Invalid Txn State | P ∨ C | 21.11.15 |
| 3616-b | Lost txn writes | ∅ | 21.11.15 |

Effective data loss

Redpanda found many
of these bugs independently

#3039     #3003    #3036

Most frequent issues
resolved in 21.10.3 —
but some serious problems
in 21.11.2

No clack skew
issues

There may be more.....

Ease of operation

# Not fault tolerant by default

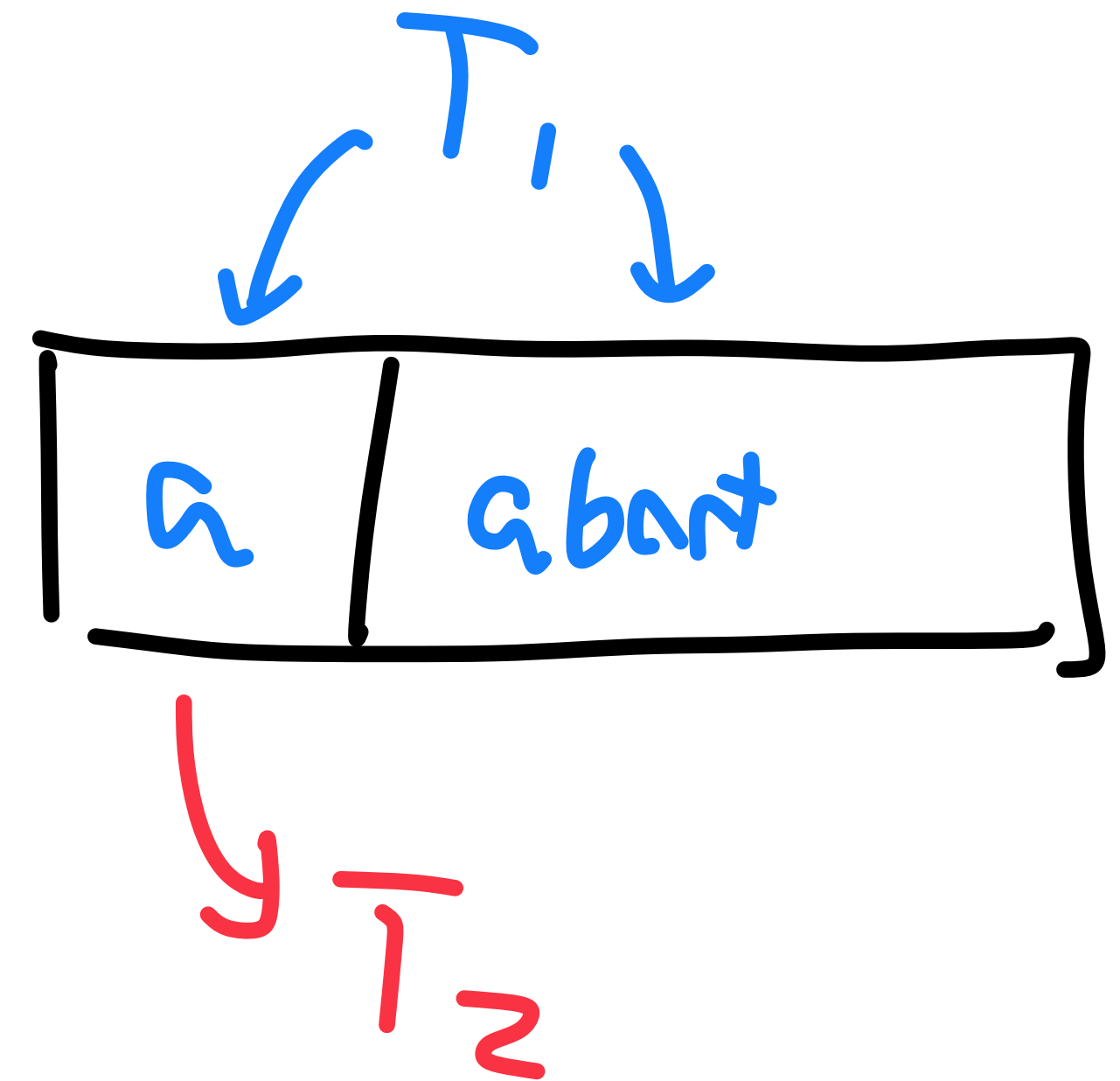Idempotence                    Txns
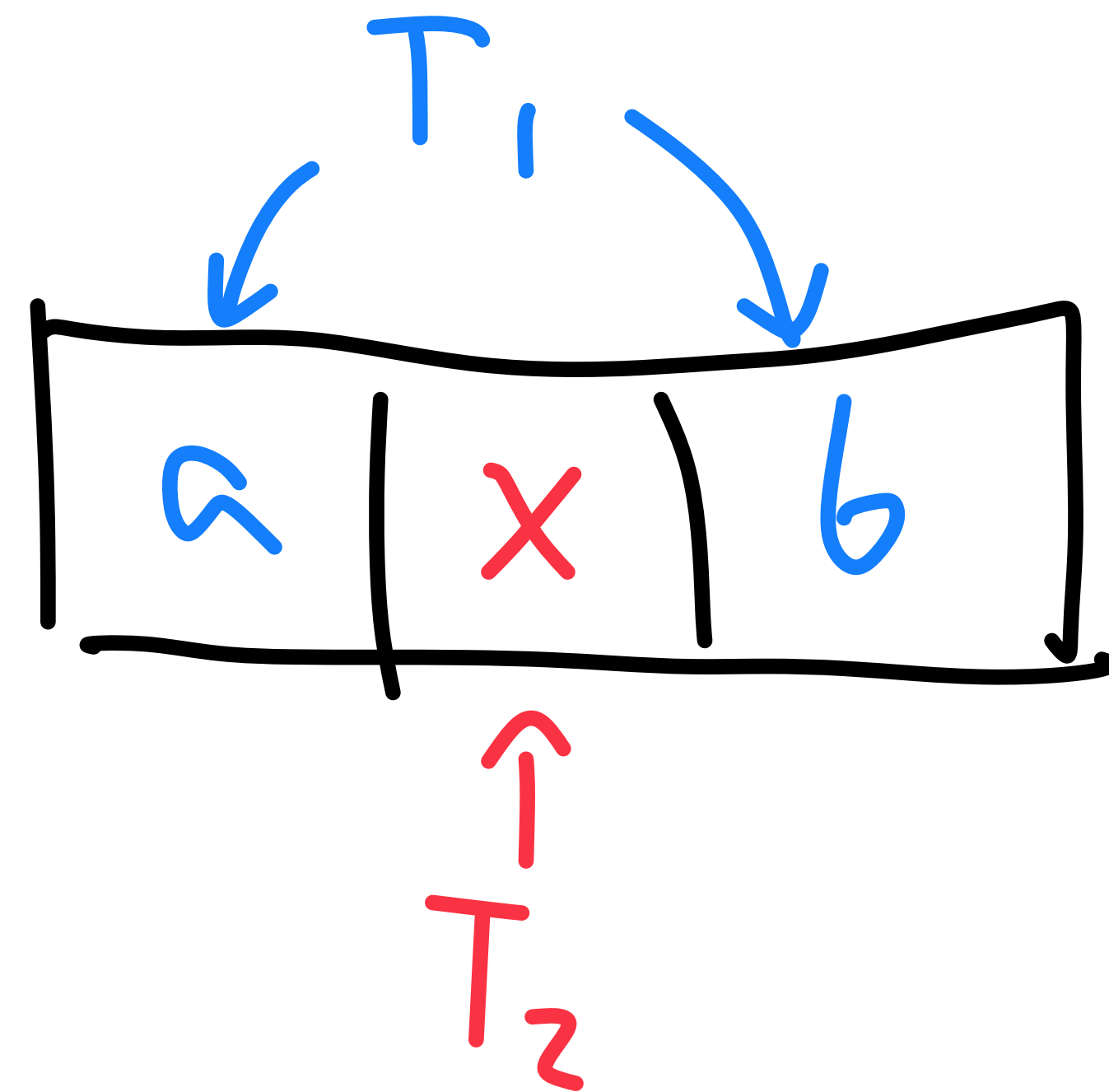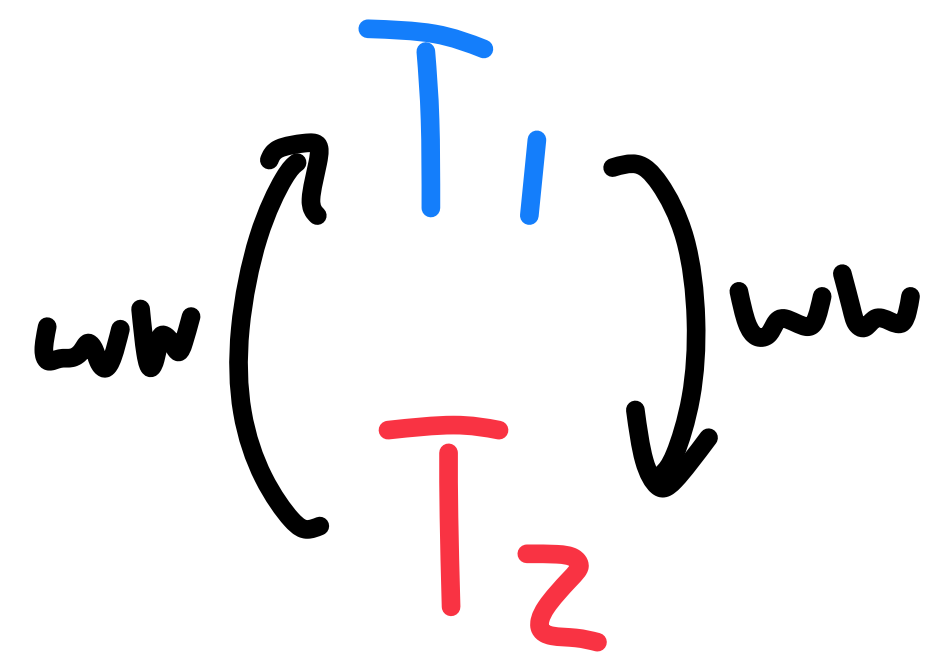
Default topics
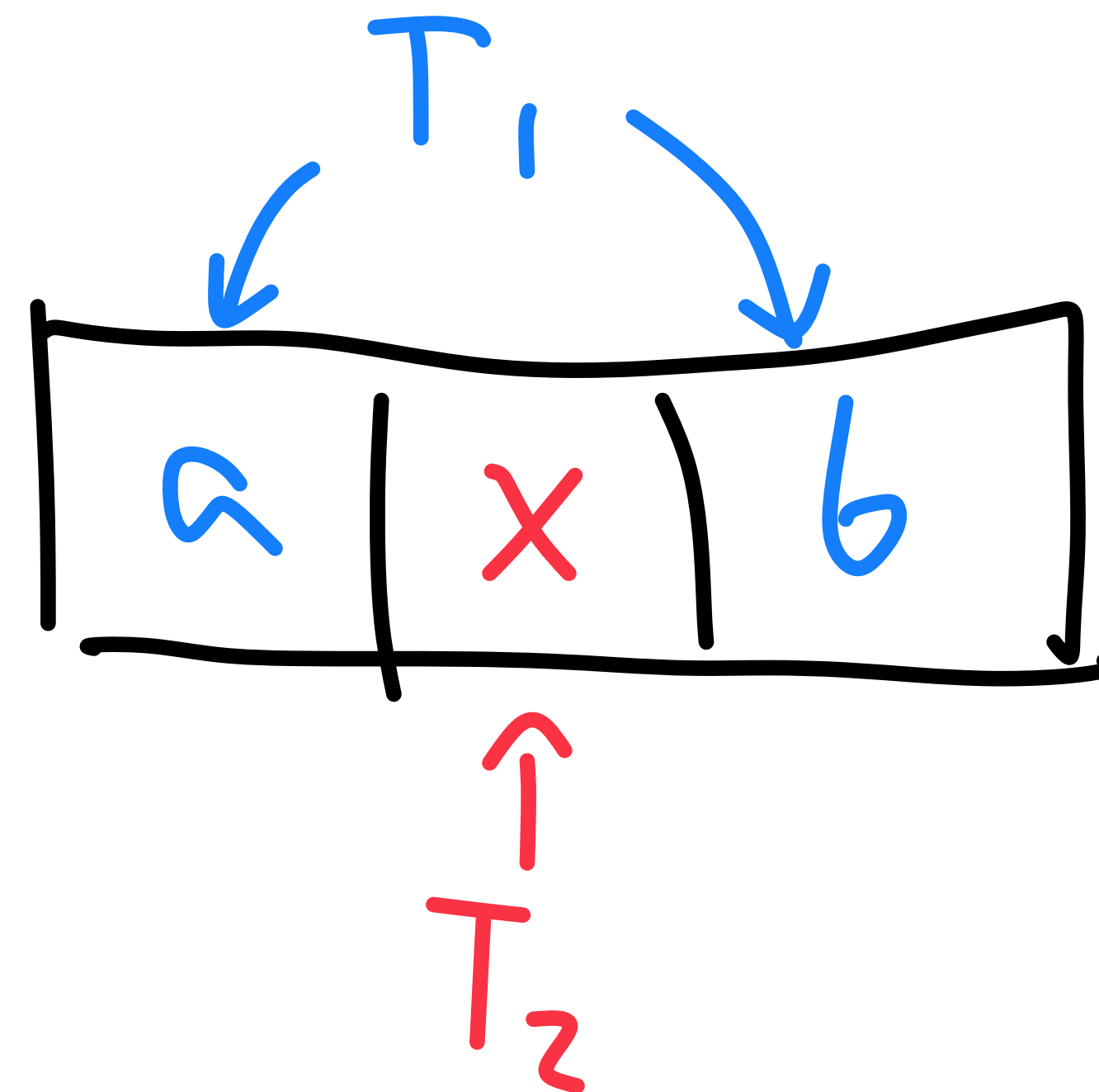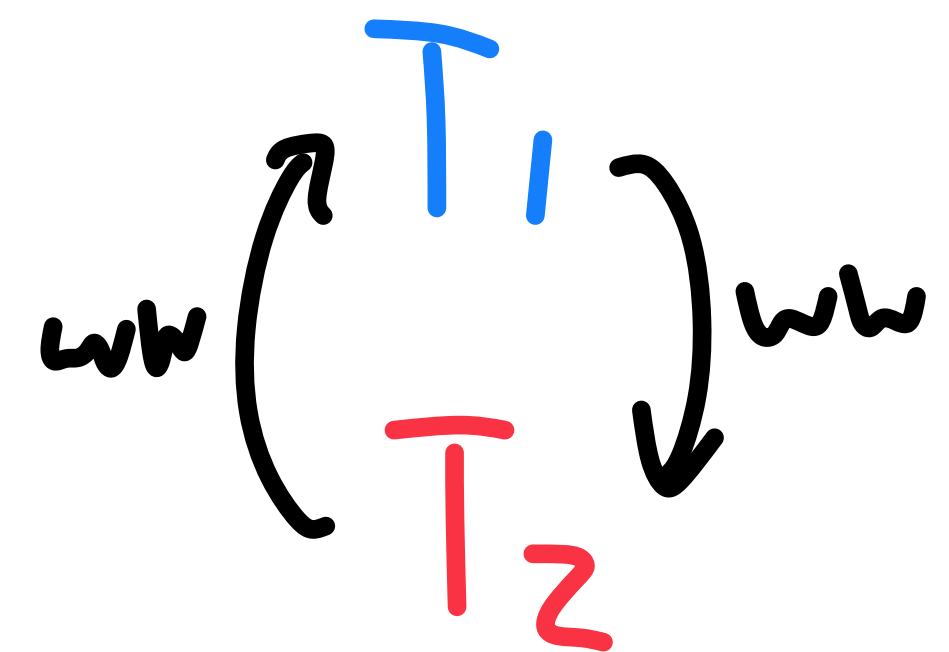
# Transactions

# G1a

$T_1$ (aborted)

$\downarrow wr$

$T_2$



As of 21.11.15, forbidden
by read_committed

GO.

$ww \left( \begin{matrix} T_1 \\ T_2 \end{matrix} \right) ww$

G0

$_{ww}\left(\begin{array}{c} T_1 \\ T_2 \end{array}\right)_{ww}$

| a | x | b |

G1c

$_{wr}\left(\begin{array}{c} T_1 \\ T_2 \end{array}\right)_{ww}$

| a | b |

# G0

$$ww \left( \begin{array}{c} T_1 \\ T_2 \end{array} \right) ww$$



(Ambiguous if this is 0/r)

# G1c

$$wr \left( \begin{array}{c} T_1 \\ T_2 \end{array} \right) ww$$

Should be fixed under
read_committed in 21.11.15!

G1c
w/all wr edges

$wr\left(\begin{array}{c}T_1\\T_2\end{array}\right)wr$

# Didn't check exactly-once semantics

1. Multiple palls per txn

2. Allowed more than one transactional.id to consume a partition

# All bets are off!?

1. Multiple palls per txn

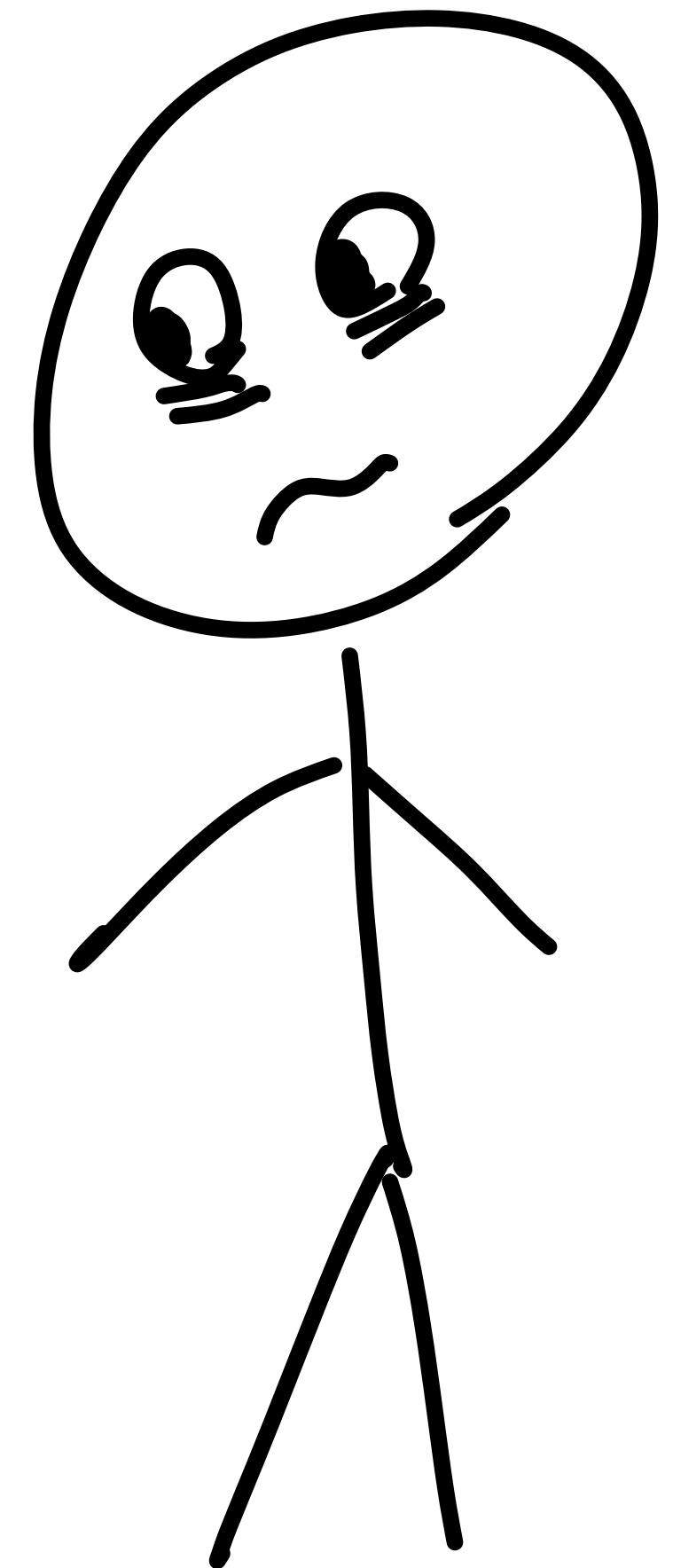2. Allowed more than one transactional.id to consume a partition

# Txn Documentation

- Absent
- Incomplete
- Vague
- Misleading
- Wrong

- Largely explained in terms of implementation

- Go read 67 pages of design doc?

# Membership Changes

- No docs, but customers were learning how in Slack
- Improved APIs & some docs now
- Don't re-use node IPs!

# FUTURE WORK!

- Binding clients to nodes
- Adding partitions
- Exactly-once semantics
- Streams

Thanks!

- Camilo Aguilar
- Travis Bischel
- Bob Dever
- Juan Castillo
- Alexander Gallego
- Dhruv Gupta

- Michal Maslanka
- Denis Rystsov
- John Spray
- Coral Waters
- David Wang
- Noah Watkins
- Allison Daly

https://jepsen.io

https://redpanda.com