# Viewstamped Replication

10/19/18

# MIDTERM

Next Wednesday 10/24 at 7 - 9pm in CS 104

Covers all material up to and including Monday's lecture

# Viewstamped Replication

A way to implement replicated state machines

Goal: strong consistency across replicas

Similar to Paxos and RAFT, but less popular

# Viewstamped Replication

## Normal operation

*2f + 1 = 3 nodes*

*Can tolerate f = 1
node failing at once*

A

| status | normal |
| replica | 0 |
| view | 0 |
| op | 0 |
| commit | -1 |

<empty>

B

| status | normal |
| replica | 1 |
| view | 0 |
| op | 0 |
| commit | -1 |

<empty>

C

| status | normal |
| replica | 2 |
| view | 0 |
| op | 0 |
| commit | -1 |

<empty>

A

| status | normal |
|--------|--------|
| replica | 0 |
| view | 0 |
| op | 1 |
| commit | -1 |

<0, 1> x = 18

*<view, op>*

**Prepare**
view: 0
op: 1
commit: -1
<Request>

B

| status | normal |
|--------|--------|
| replica | 1 |
| view | 0 |
| op | 0 |
| commit | -1 |

<empty>

C

| status | normal |
|--------|--------|
| replica | 2 |
| view | 0 |
| op | 0 |
| commit | -1 |

<empty>

*Primary informs backups that op 1 is committed during the next* `Prepare`

A
| status | normal |
| replica | 0 |
| view | 0 |
| op | 1 |
| commit | 1 |

<0, 1> x = 18 ✓

*<view, op>*

✓ *committed*

B
| status | normal |
| replica | 1 |
| view | 0 |
| op | 1 |
| commit | -1 |

<0, 1> x = 18

C
| status | normal |
| replica | 2 |
| view | 0 |
| op | 1 |
| commit | -1 |

<0, 1> x = 18

A

status    normal
replica   0
view      0
op        2
commit    1

<0, 1> x = 18 ✔
<0, 2> x += 3

<view, op>

✔ committed

PrepareOK
view: 0
op: 2
replica: 1

PrepareOK
view: 0
op: 2
replica: 2

B

status    normal
replica   1
view      0
op        2
commit    1

<0, 1> x = 18 ✔
<0, 2> x += 3

C

status    normal
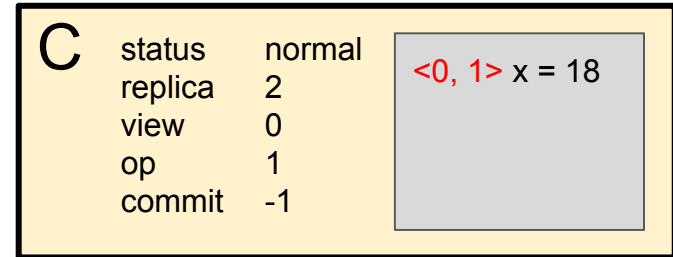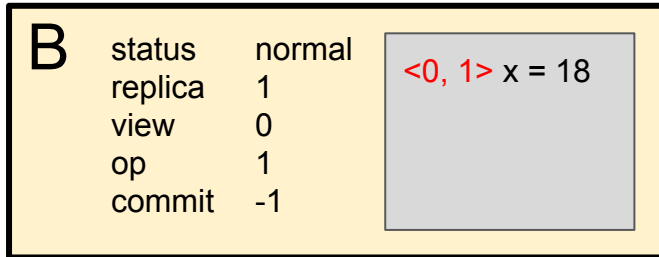replica   2
view      0
op        2
commit    1

<0, 1> x = 18 ✔
<0, 2> x += 3

*What if the next `Prepare` never comes?*

*Primary times out and sends a `Commit` message to each backup*

**A** status normal
replica 0
view 0
op 2
commit 2

<0, 1> x = 18 ✔
<0, 2> x += 3 ✔

*<view, op>*

✔*committed*

```
Commit
view: 0
commit: 2
```

**B** status normal
replica 1
view 0
op 2
commit 1

<0, 1> x = 18 ✔
<0, 2> x += 3

**C** status normal
replica 2
view 0
op 2
commit 1

<0, 1> x = 18 ✔
<0, 2> x += 3

# Why is waiting for *f* nodes enough?

Op is guaranteed to have been executed on *f + 1* nodes (majority)

# Overlapping quorums



Write quorum
contains *f + 1* nodes

x = 1

A

x = 1

B

C

# Overlapping quorums

# Overlapping quorums



Write quorum
contains *f + 1* nodes

x = 1

A

B
FAILED

C

x = ?

Client

# Overlapping quorums

x = 1

x = ?

**Client**

Write quorum
contains *f + 1* nodes

A

B

FAILED

C

Read quorum
contains *f + 1* nodes

# Overlapping quorums



x = 1

Write quorum
contains *f + 1* nodes

x = ?

**Client**

x = ?

A

B

FAILED

C

Read quorum
contains *f + 1* nodes

# Overlapping quorums



Write quorum contains $f + 1$ nodes

x = 1

x = 1

A

B

FAILED

C

Client

x?

Read quorum contains $f + 1$ nodes

# *Non*-overlapping quorums?

# Viewstamped replication
## View change

*Primary fails before sending* `Prepare` *to B*

A
status    normal
replica   0
view      0
op        3
commit    2

<0, 1> x = 18 ✔
<0, 2> x += 3 ✔
<0, 3> y = 100

*<view, op>*

✔ *committed*

**Prepare**
```
view: 0
op: 3
commit: 2
<Request>
```

B
status    normal
replica   1
view      0
op        2
commit    2

<0, 1> x = 18 ✔
<0, 2> x += 3 ✔

C
status    normal
replica   2
view      0
op        3
commit    2

<0, 1> x = 18 ✔
<0, 2> x += 3 ✔
<0, 3> y = 100

*Logs are out of sync*

<0, 1> x = 18 ✓
x += 3 ✓
y = 100

<view, op>

✓ *committed*

**B**
status normal
replica 1
view 0
op 2
commit 2

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓

**C**
status normal
replica 2
view 0
op 3
commit 2

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓
<0, 3> y = 100

*C times out on hearing from the primary and starts view change*

<view, op>

✓ *committed*

18 ✓
x += 3 ✓
y = 100

???

B  status    normal
   replica   1
   view      0
   op        2
   commit    2

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓

C  status    normal
   replica   2
   view      0
   op        3
   commit    2

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓
<0, 3> y = 100

**Start view change:**

*Status* = change
*Increment local view*
*Send* SVC *to all nodes*

18 ✓
x += 3 ✓
y = 100



*<view, op>*

✓*committed*

B | status | normal
| replica | 1
| view | 0
| op | 2
| commit | 2

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓

C | status | normal
| replica | 2
| view | 0
| op | 3
| commit | 2

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓
<0, 3> y = 100

**Start view change:**

  *Status* = change
  *Increment local view*
  *Send* SVC *to all nodes*

<view, op>

✓ committed

18 ✓
x += 3 ✓
y = 100

| B | status | normal |
|---|---|---|
| | replica | 1 |
| | view | 0 |
| | op | 2 |
| | commit | 2 |

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓

**StartViewChange**
view: 1
replica: 2

| C | status | change |
|---|---|---|
| | replica | 2 |
| | view | 1 |
| | op | 3 |
| | commit | 2 |

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓
<0, 3> y = 100

**Receive SVC where:**

*SVC.view > local view {*
    *Status* = `view change`
    *Advance local view*
    *Send* SVC *to other nodes*
*}*



18 ✓
x += 3 ✓
y = 100

*<view, op>*

✓*committed*

**B**
| status | normal |
| replica | 1 |
| view | 0 |
| op | 2 |
| commit | 2 |

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓

`StartViewChange`
view: 1
replica: 2

**C**
| status | change |
| replica | 2 |
| view | 1 |
| op | 3 |
| commit | 2 |

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓
<0, 3> y = 100

**Receive SVC where:**

*SVC.view > local view {*
    *Status =* `view change`
    *Advance local view*
    *Send* SVC *to other nodes*
*}*



<view, op>

✓*committed*

18 ✓
x += 3 ✓
y = 100

---

B
| status | change |
| replica | 1 |
| view | 1 |
| op | 2 |
| commit | 2 |

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓

**StartViewChange**
view: 1
replica: 1

→

C
| status | change |
| replica | 2 |
| view | 1 |
| op | 3 |
| commit | 2 |

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓
<0, 3> y = 100

**Receive f SVCs where:**

*SVC.view == local view {*
    *Send DVC to new primary*
*}*

<0, 1> x = 18 ✓
x += 3 ✓
y = 100

<view, op>

✓ *committed*

**B**
| status | change |
|--------|--------|
| replica | 1 |
| view | 1 |
| op | 2 |
| commit | 2 |

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓

**StartViewChange**
view: 1
replica: 1

**C**
| status | change |
|--------|--------|
| replica | 2 |
| view | 1 |
| op | 3 |
| commit | 2 |

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓
<0, 3> y = 100

**Receive f SVCs where:**

*SVC.view == local view {*
    *Send DVC to new primary*
*}*



$\langle view, op \rangle$

✓ *committed*

18 ✓
x += 3 ✓
y = 100

**DoViewChange**
replica: 2
view: 1
op: 3
commit: 2
<log>

B

| status | change |
|--------|--------|
| replica | 1 |
| view | 1 |
| op | 2 |
| commit | 2 |

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓

C

| status | change |
|--------|--------|
| replica | 2 |
| view | 1 |
| op | 3 |
| commit | 2 |

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓
<0, 3> y = 100

Logs are no longer out of sync!

With more nodes, we may receive multiple different logs

Pick the one with highest view and op number



*<view, op>*

✓ *committed*

<0, 1> x = 18 ✓
x += 3 ✓
y = 100

**B**

| status | change | |
|--------|--------|---|
| replica | 1 | <0, 1> x = 18 ✓ |
| view | 1 | <0, 2> x += 3 ✓ |
| op | 3 | <0, 3> y = 100 |
| commit | 2 | |

**C**

| status | change | |
|--------|--------|---|
| replica | 2 | <0, 1> x = 18 ✓ |
| view | 1 | <0, 2> x += 3 ✓ |
| op | 3 | <0, 3> y = 100 |
| commit | 2 | |

**Receive f DVCs:**

*Become new primary*
*Send* `StartView` *to others*

Why do we send the log here?

<0, 1> x = 18 ✓
x += 3 ✓
y = 100

*<view, op>*
✓*committed*

👑

B    status    normal
     replica   1
     view      1
     op        3
     commit    2

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓
<0, 3> y = 100
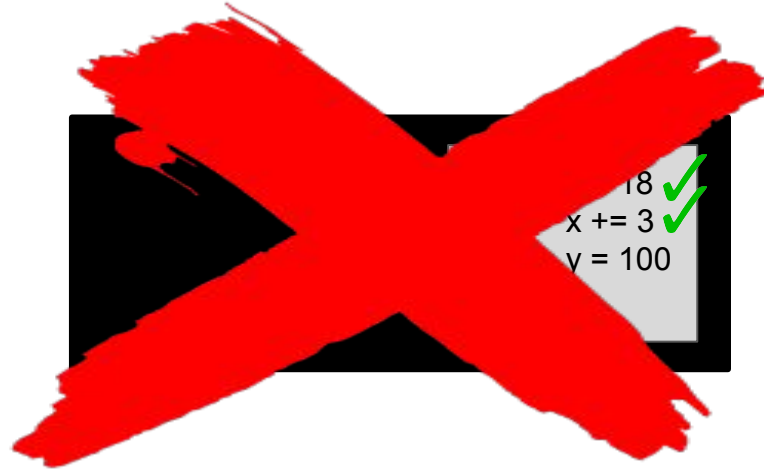
`StartView`
view: 1
replica: 1
op: 3
commit: 2
<log>

C    status    change
     replica   2
     view      1
     op        3
     commit    2

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓
<0, 3> y = 100

*Notice <0, 3> is uncommitted and from an old view...*

*Do we commit it?*

<0, 3> y = 100

18 ✓
x += 3 ✓
y = 100

<view, op>

✓ *committed*

**B**

| | |
|---|---|
| status | normal |
| replica | 1 |
| view | 1 |
| op | 3 |
| commit | 2 |

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓
<0, 3> y = 100

`PrepareOK`
view: 0
op: 3
replica: 2

**C**

| | |
|---|---|
| status | normal |
| replica | 2 |
| view | 1 |
| op | 3 |
| commit | 2 |

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓
<0, 3> y = 100

*Are uncommitted ops like <0, 3> guaranteed to survive into the new view?*

*What about committed ops? (e.g. <0, 1> and <0, 2>)*

<view, op>

✓ *committed*



**B**

| status | normal |
|--------|--------|
| replica | 1 |
| view | 1 |
| op | 3 |
| commit | 3 |

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓
<0, 3> y = 100 ✓

**C**

| status | normal |
|--------|--------|
| replica | 2 |
| view | 1 |
| op | 3 |
| commit | 2 |

<0, 1> x = 18 ✓
<0, 2> x += 3 ✓
<0, 3> y = 100

# Summary: view change in VR

New primary is pre-selected based on IP address (round-robin)

View change triggered by timeout, could be any node

Wait for $f$ SVC that matches our view number before sending DVC

Wait for $f$ DVC to start new view (primary)

- Why $f$ in both cases?

- Provided that at most $f$ servers fail, is *liveness* guaranteed?

# Additional reading for viewstamped replication

http://pmg.csail.mit.edu/papers/vr-revisited.pdf

https://blog.acolyer.org/2015/03/06/viewstamped-replication-revisited/

# Q & A