

# Viewstamped replication

10/20/17

# A note on assignment 2

Your tests need to pass *deterministically*

Use SyncMap, err on using too many (correctness > performance here)

You don't need maps of maps (bad design in general)

Due tonight!!

# MIDTERM

Next Friday 10/27 at 10am or 11am, you choose (90 minutes)

Covers all material up to and including today's class

# 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



<b>A</b>	status	normal	
	replica	0	
	view	0	
	op	0	
	commit	-1	

<b>B</b>	status	normal	
	replica	1	
	view	0	
	op	0	
	commit	-1	

<b>C</b>	status	normal	
	replica	2	
	view	0	
	op	0	
	commit	-1	



Client 136

Request  
op: x = 18  
cid: 136  
request num: 0



<b>A</b>	status	normal	<div style="border: 1px solid gray; background-color: #d3d3d3; padding: 5px; width: 100px; height: 80px; display: flex; align-items: center; justify-content: center;">&lt;empty&gt;</div>
	replica	0	
	view	0	
	op	0	
	commit	-1	

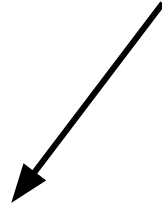
<b>B</b>	status	normal	<div style="border: 1px solid gray; background-color: #d3d3d3; padding: 5px; width: 100px; height: 80px; display: flex; align-items: center; justify-content: center;">&lt;empty&gt;</div>
	replica	1	
	view	0	
	op	0	
	commit	-1	

<b>C</b>	status	normal	<div style="border: 1px solid gray; background-color: #d3d3d3; padding: 5px; width: 100px; height: 80px; display: flex; align-items: center; justify-content: center;">&lt;empty&gt;</div>
	replica	2	
	view	0	
	op	0	
	commit	-1	

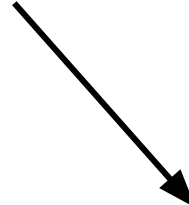


<b>A</b>	status	normal	<code>&lt;0, 1&gt; x = 18</code>
	replica	0	
	view	0	
	op	1	
	commit	-1	

*<view, op>*



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

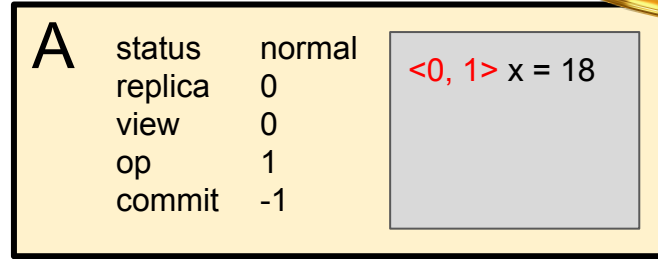


<b>B</b>	status	normal	<code>&lt;empty&gt;</code>
	replica	1	
	view	0	
	op	0	
	commit	-1	

<b>C</b>	status	normal	<code>&lt;empty&gt;</code>
	replica	2	
	view	0	
	op	0	
	commit	-1	



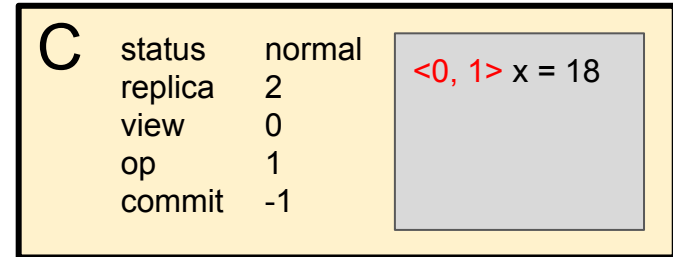
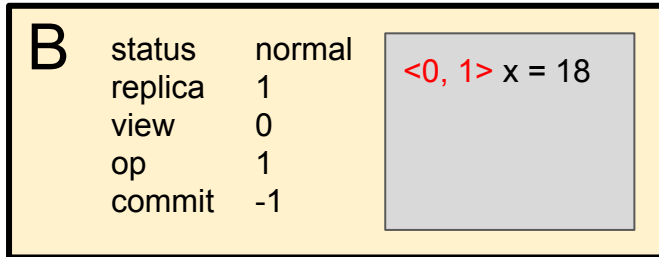
*Primary only needs to wait for  $f = 1$  replies before committing*



*<view, op>*

PrepareOK  
view: 0  
op: 1  
replica: 1

PrepareOK  
view: 0  
op: 1  
replica: 2





Client 136

Reply  
view: 0  
request num: 0  
result: x = 18



<b>A</b>	status	normal	<b>&lt;0, 1&gt; x = 18</b> ✓
	replica	0	
	view	0	
	op	1	
	commit	1	

**<view, op>**  
✓ *committed*

<b>B</b>	status	normal	<b>&lt;0, 1&gt; x = 18</b>
	replica	1	
	view	0	
	op	1	
	commit	-1	

<b>C</b>	status	normal	<b>&lt;0, 1&gt; x = 18</b>
	replica	2	
	view	0	
	op	1	
	commit	-1	

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



<b>A</b>	status	normal	
	replica	0	
	view	0	
	op	1	
	commit	1	

*<view, op>  
✓ committed*

<b>B</b>	status	normal	
	replica	1	
	view	0	
	op	1	
	commit	-1	

<b>C</b>	status	normal	
	replica	2	
	view	0	
	op	1	
	commit	-1	



Client 136

Request  
op: x += 3  
cid: 136  
request num: 1



<b>A</b>	status	normal	<b>&lt;0, 1&gt; x = 18</b> ✓
	replica	0	
	view	0	
	op	1	
	commit	1	

**<view, op>**  
✓ **committed**

<b>B</b>	status	normal	<b>&lt;0, 1&gt; x = 18</b>
	replica	1	
	view	0	
	op	1	
	commit	-1	

<b>C</b>	status	normal	<b>&lt;0, 1&gt; x = 18</b>
	replica	2	
	view	0	
	op	1	
	commit	-1	



<b>A</b>	status	normal	<code>&lt;0, 1&gt; x = 18</code> ✓ <code>&lt;0, 2&gt; x += 3</code>
	replica	0	
	view	0	
	op	2	
	commit	1	

`<view, op>`  
✓ *committed*

Prepare  
view: 0  
op: 2  
commit: 1  
<Request>

<b>B</b>	status	normal	<code>&lt;0, 1&gt; x = 18</code>
	replica	1	
	view	0	
	op	1	
	commit	-1	

<b>C</b>	status	normal	<code>&lt;0, 1&gt; x = 18</code>
	replica	2	
	view	0	
	op	1	
	commit	-1	



<b>A</b>	status	normal	<code>&lt;0, 1&gt; x = 18 ✓</code> <code>&lt;0, 2&gt; x += 3</code>
	replica	0	
	view	0	
	op	2	
	commit	1	

`<view, op>`  
`✓ committed`

PrepareOK  
view: 0  
op: 2  
replica: 1

PrepareOK  
view: 0  
op: 2  
replica: 2

<b>B</b>	status	normal	<code>&lt;0, 1&gt; x = 18 ✓</code> <code>&lt;0, 2&gt; x += 3</code>
	replica	1	
	view	0	
	op	2	
	commit	1	

<b>C</b>	status	normal	<code>&lt;0, 1&gt; x = 18 ✓</code> <code>&lt;0, 2&gt; x += 3</code>
	replica	2	
	view	0	
	op	2	
	commit	1	



Client 136

Reply  
view: 0  
request num: 1  
result: x = 21



<b>A</b>	status	normal	<b>&lt;0, 1&gt;</b> x = 18 ✓ <b>&lt;0, 2&gt;</b> x += 3 ✓
	replica	0	
	view	0	
	op	2	
	commit	2	

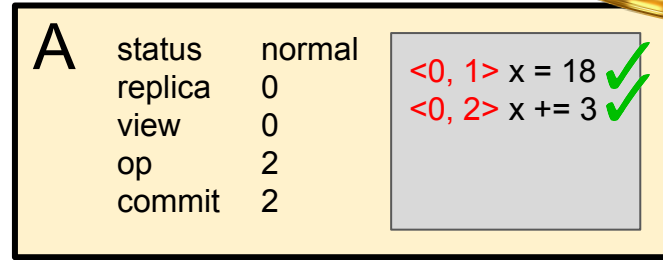
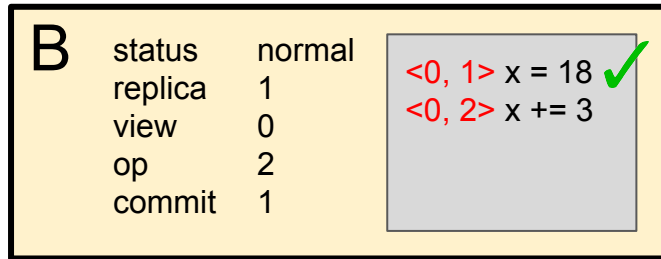
**<view, op>**  
✓ *committed*

<b>B</b>	status	normal	<b>&lt;0, 1&gt;</b> x = 18 ✓ <b>&lt;0, 2&gt;</b> x += 3
	replica	1	
	view	0	
	op	2	
	commit	1	

<b>C</b>	status	normal	<b>&lt;0, 1&gt;</b> x = 18 ✓ <b>&lt;0, 2&gt;</b> x += 3
	replica	2	
	view	0	
	op	2	
	commit	1	

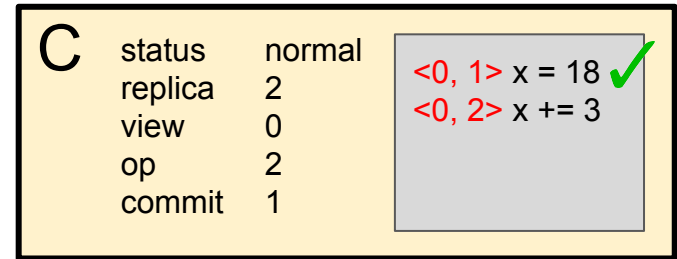
What if the next Prepare never comes?

Primary times out and sends a Commit message to each backup



<view, op>  
✓ committed

Commit  
view: 0  
commit: 2







<b>A</b>	status	normal	
	replica	0	
	view	0	
	op	2	
	commit	2	

<view, op>

✓ committed

<b>B</b>	status	normal	
	replica	1	
	view	0	
	op	2	
	commit	2	

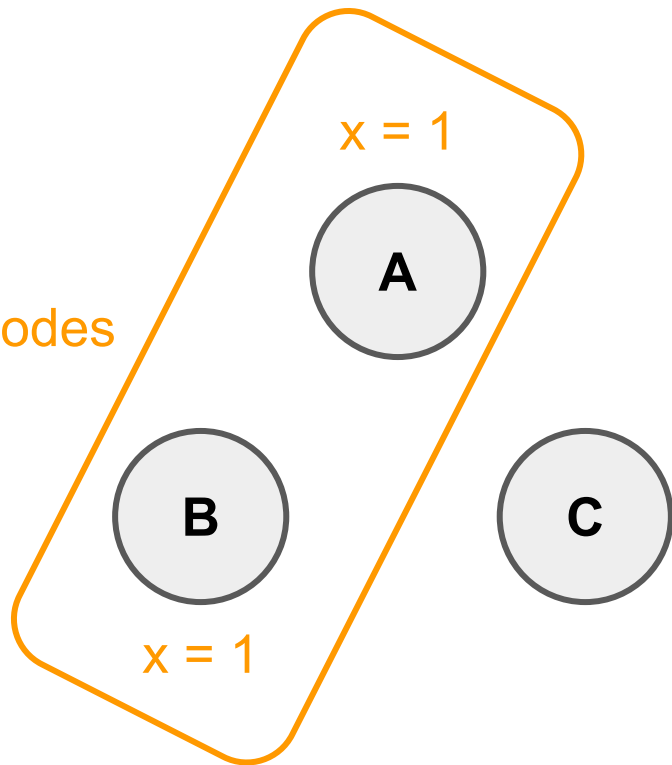
<b>C</b>	status	normal	
	replica	2	
	view	0	
	op	2	
	commit	2	

# Why is waiting for $f$ nodes enough?

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

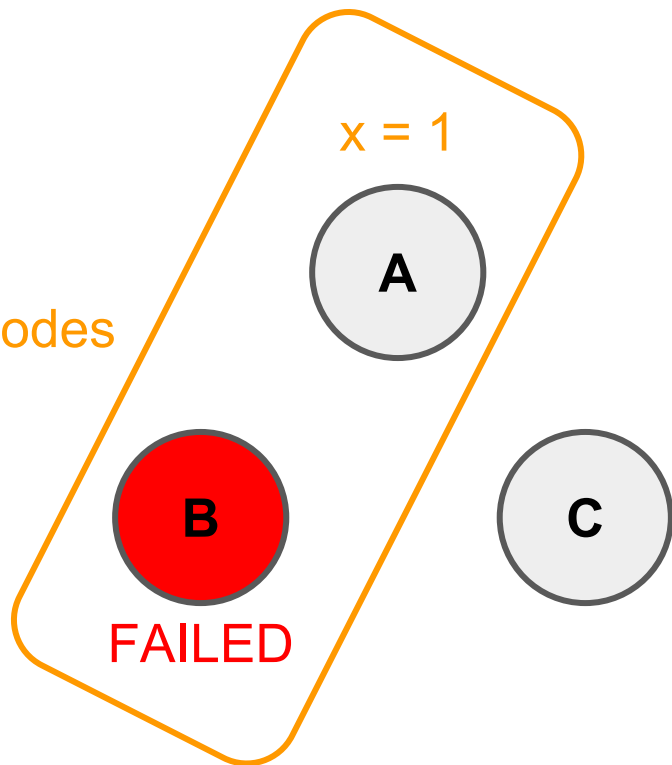
# Overlapping quorums

Write quorum  
contains  $f + 1$  nodes



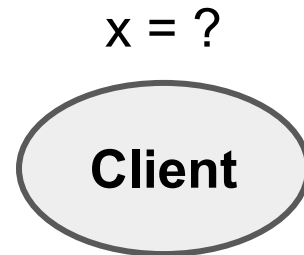
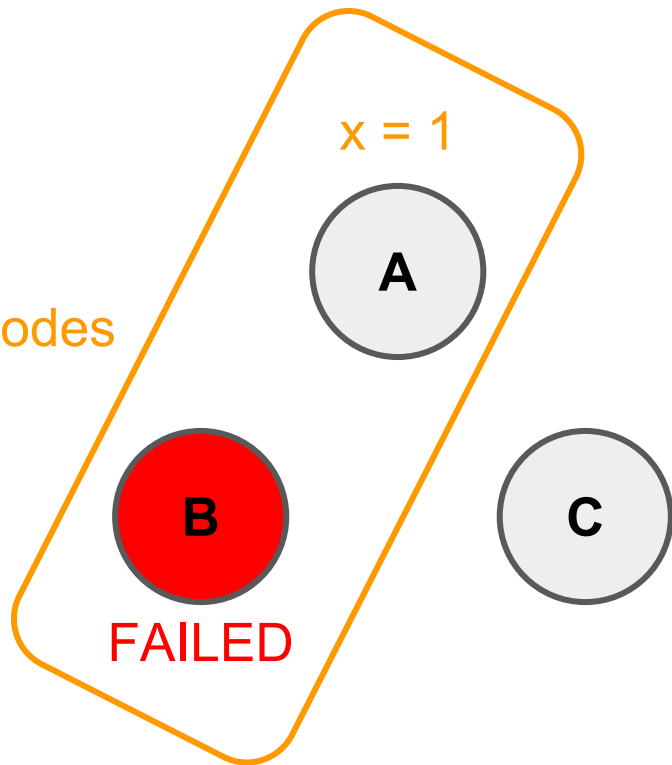
# Overlapping quorums

Write quorum  
contains  $f + 1$  nodes



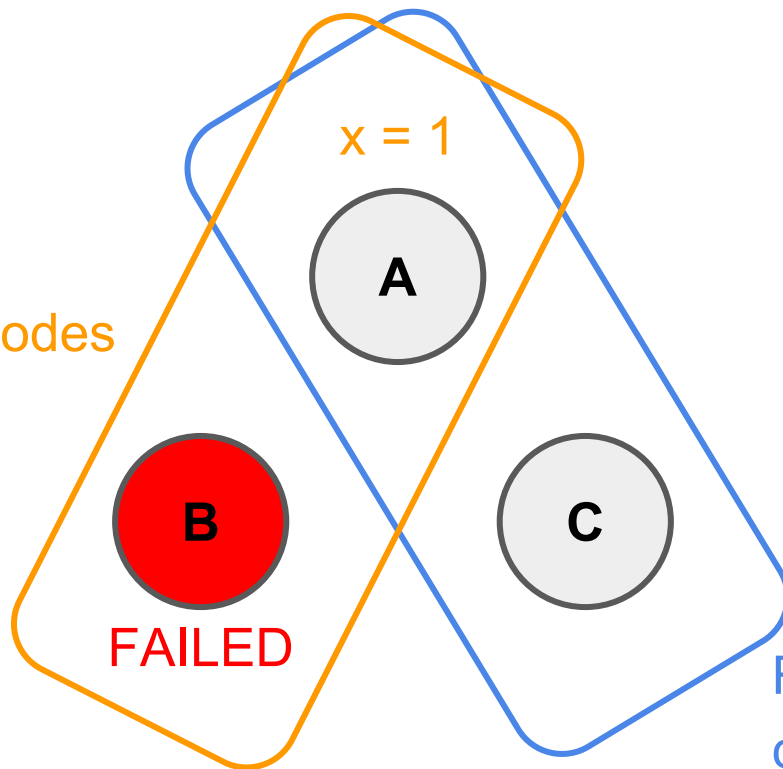
# Overlapping quorums

Write quorum  
contains  $f + 1$  nodes



# Overlapping quorums

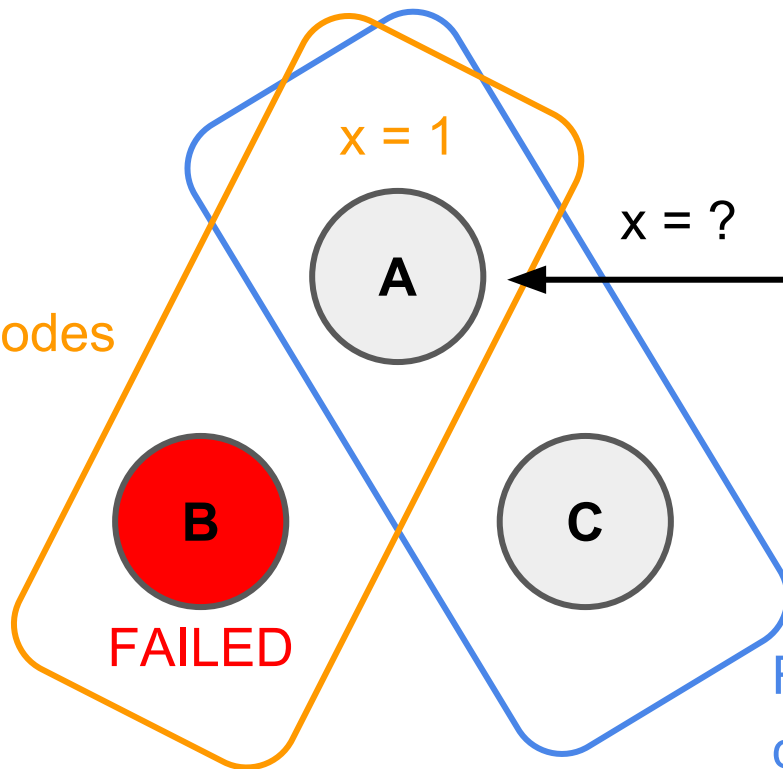
Write quorum  
contains  $f + 1$  nodes



Read quorum  
contains  $f + 1$  nodes

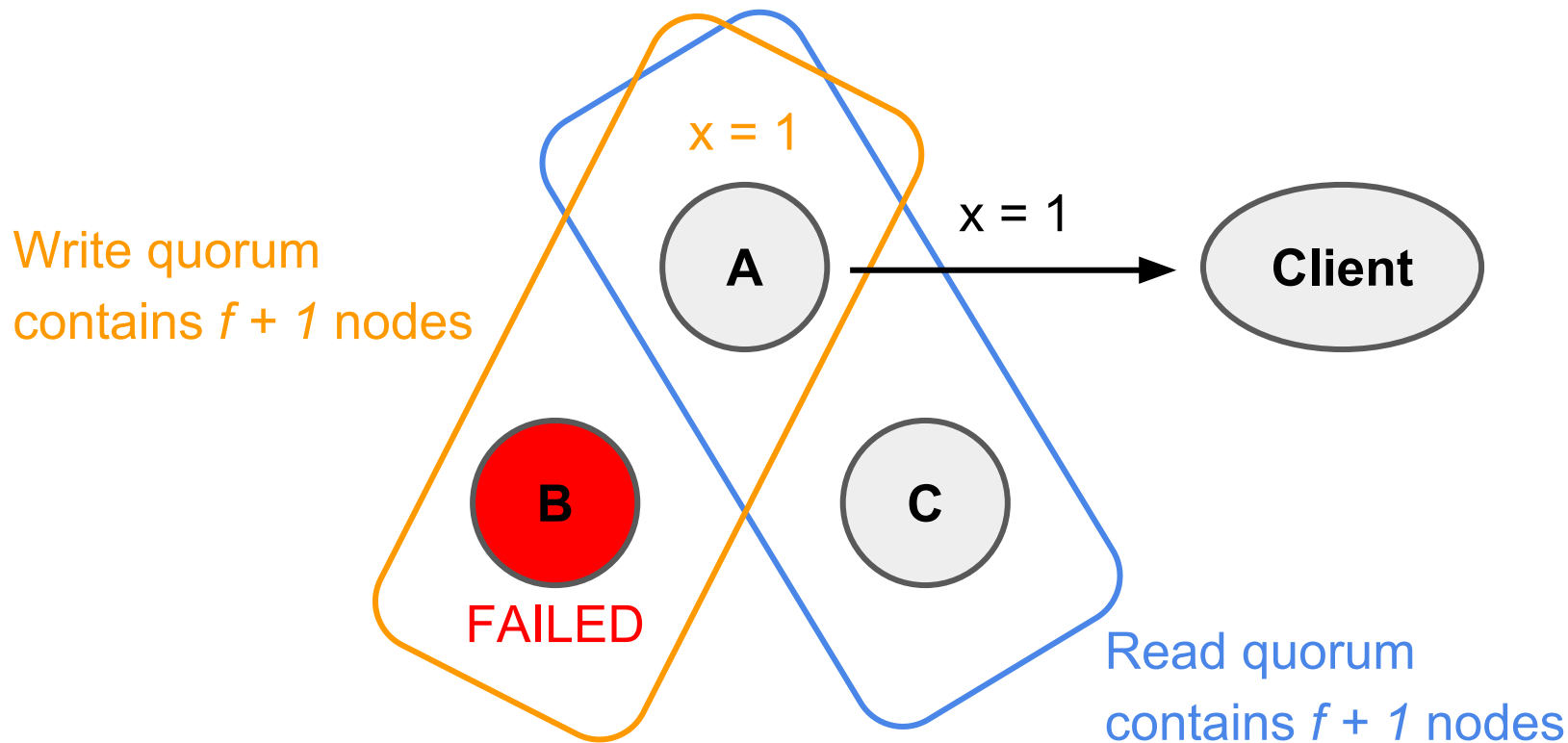
# Overlapping quorums

Write quorum  
contains  $f + 1$  nodes



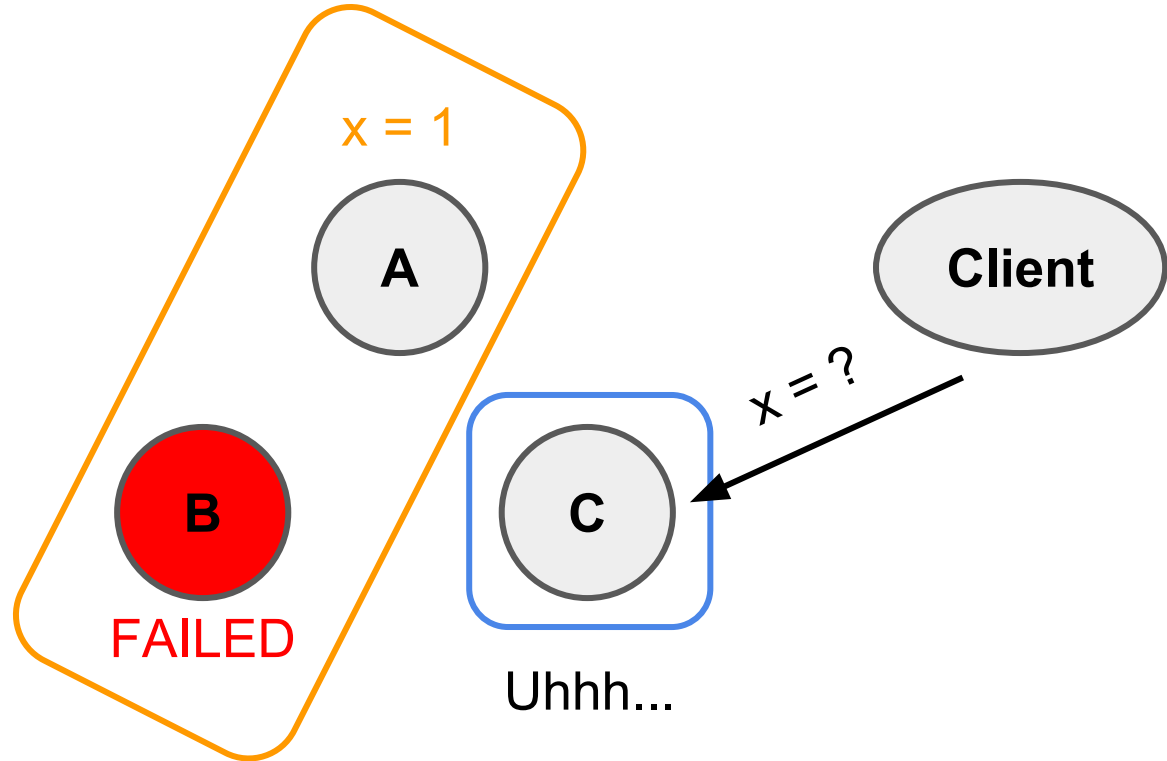
Read quorum  
contains  $f + 1$  nodes

# Overlapping quorums





# Non-overlapping quorums?



# **Viewstamped replication**

View change



Client 25

Request  
op: y = 100  
cid: 25  
request num: 0



<b>A</b>	status	normal	<b>&lt;0, 1&gt;</b> x = 18 ✓ <b>&lt;0, 2&gt;</b> x += 3 ✓
	replica	0	
	view	0	
	op	2	
	commit	2	

**<view, op>**  
✓ *committed*

<b>B</b>	status	normal	<b>&lt;0, 1&gt;</b> x = 18 ✓ <b>&lt;0, 2&gt;</b> x += 3 ✓
	replica	1	
	view	0	
	op	2	
	commit	2	

<b>C</b>	status	normal	<b>&lt;0, 1&gt;</b> x = 18 ✓ <b>&lt;0, 2&gt;</b> x += 3 ✓
	replica	2	
	view	0	
	op	2	
	commit	2	



<b>A</b>	status	normal	<code>&lt;0, 1&gt; x = 18</code> ✓ <code>&lt;0, 2&gt; x += 3</code> ✓ <code>&lt;0, 3&gt; y = 100</code>
	replica	0	
	view	0	
	op	3	
	commit	2	

*<view, op>*

✓ *committed*

<b>B</b>	status	normal	<code>&lt;0, 1&gt; x = 18</code> ✓ <code>&lt;0, 2&gt; x += 3</code> ✓
	replica	1	
	view	0	
	op	2	
	commit	2	

<b>C</b>	status	normal	<code>&lt;0, 1&gt; x = 18</code> ✓ <code>&lt;0, 2&gt; x += 3</code> ✓
	replica	2	
	view	0	
	op	2	
	commit	2	

*Primary fails before sending Prepare to B*



<b>A</b>	status	normal	<pre>&lt;0, 1&gt; x = 18 ✓✓ &lt;0, 2&gt; x += 3 ✓✓ &lt;0, 3&gt; y = 100</pre>
	replica	0	
	view	0	
	op	3	
	commit	2	

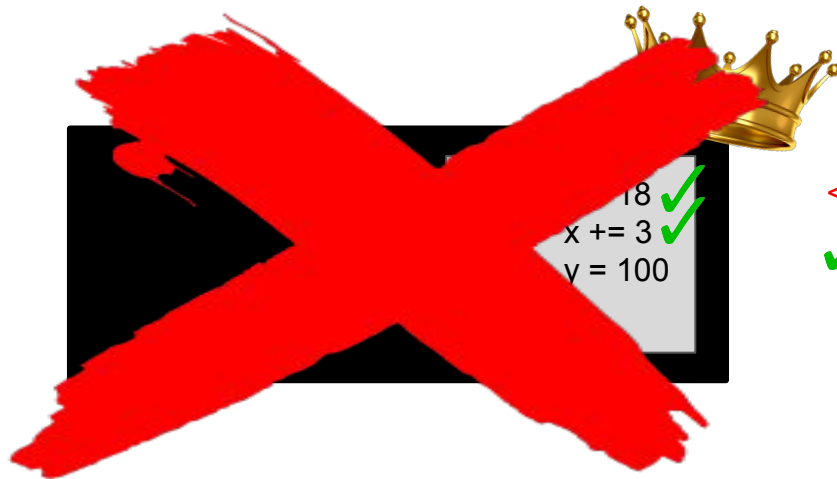
<view, op>  
✓ committed

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

<b>B</b>	status	normal	<pre>&lt;0, 1&gt; x = 18 ✓✓ &lt;0, 2&gt; x += 3 ✓✓</pre>
	replica	1	
	view	0	
	op	2	
	commit	2	

<b>C</b>	status	normal	<pre>&lt;0, 1&gt; x = 18 ✓✓ &lt;0, 2&gt; x += 3 ✓✓ &lt;0, 3&gt; y = 100</pre>
	replica	2	
	view	0	
	op	3	
	commit	2	

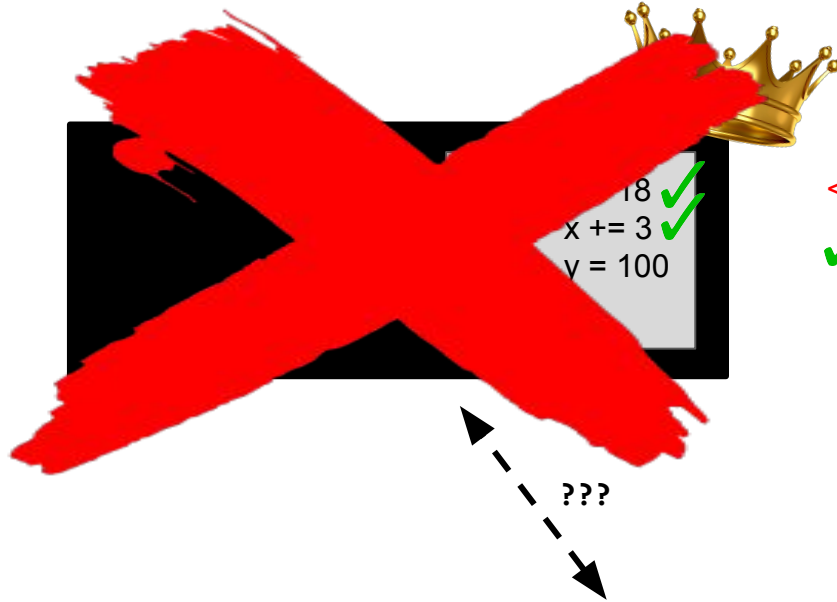
*Logs are out of sync*



<b>B</b>	status	normal	
	replica	1	
	view	0	
	op	2	
	commit	2	

<b>C</b>	status	normal	
	replica	2	
	view	0	
	op	3	
	commit	2	

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

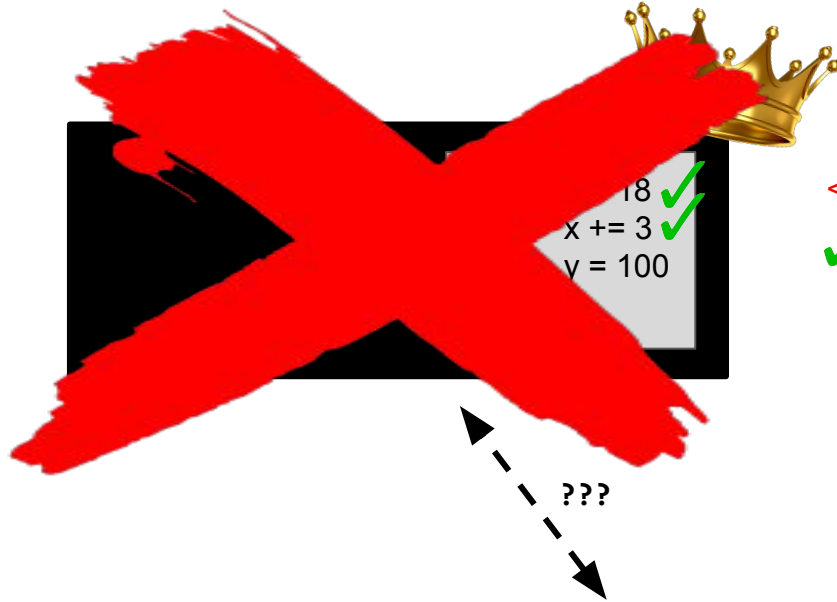


<b>B</b>	status	normal	<pre>&lt;0, 1&gt; x = 18 ✓✓ &lt;0, 2&gt; x += 3 ✓✓</pre>
	replica	1	
	view	0	
	op	2	
	commit	2	

<b>C</b>	status	normal	<pre>&lt;0, 1&gt; x = 18 ✓✓ &lt;0, 2&gt; x += 3 ✓✓ &lt;0, 3&gt; y = 100</pre>
	replica	2	
	view	0	
	op	3	
	commit	2	

Who is the new primary?

Go through the list of sorted IP addresses and find the next one (i.e. B)



<view, op>

✓ committed

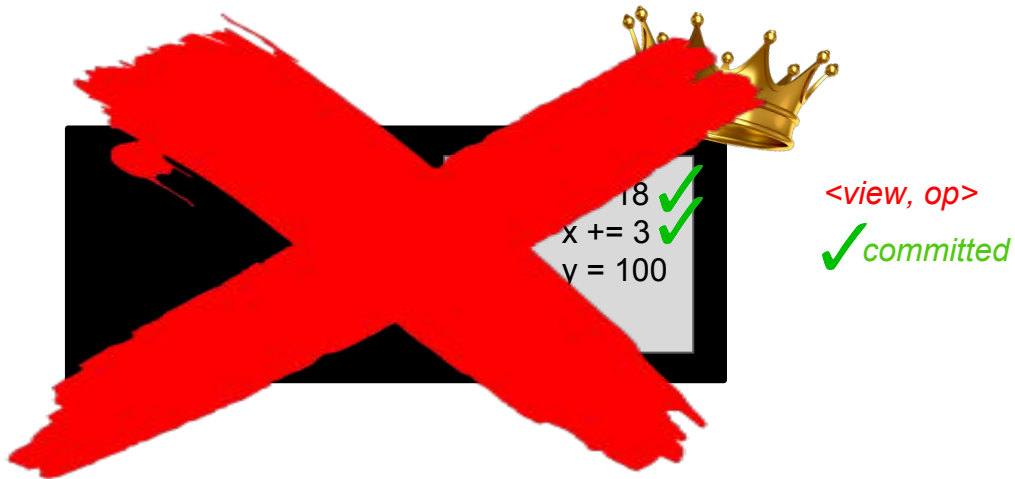
<b>B</b>	status	normal	<pre>&lt;0, 1&gt; x = 18 ✓ &lt;0, 2&gt; x += 3 ✓</pre>
	replica	1	
	view	0	
	op	2	
	commit	2	

<b>C</b>	status	normal	<pre>&lt;0, 1&gt; x = 18 ✓ &lt;0, 2&gt; x += 3 ✓ &lt;0, 3&gt; y = 100</pre>
	replica	2	
	view	0	
	op	3	
	commit	2	



## Start view change:

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

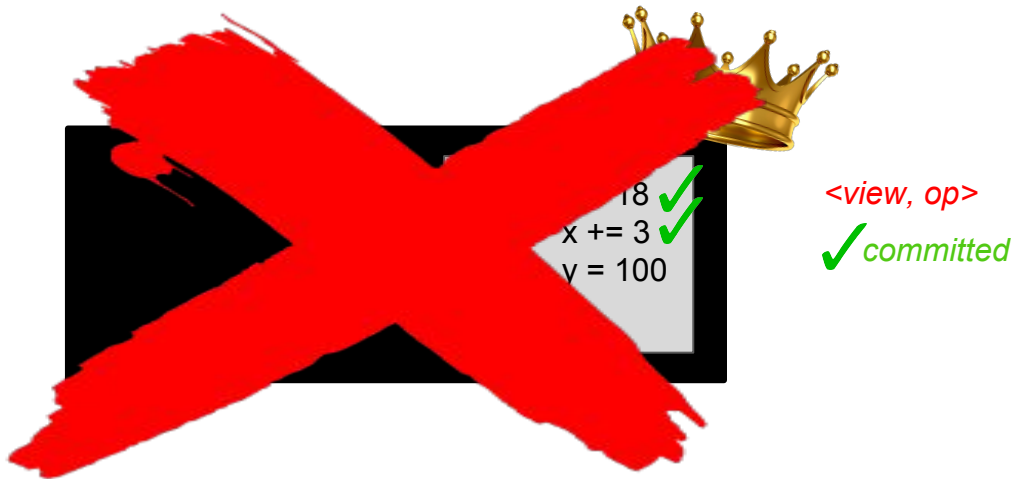


<b>B</b>	status	normal	<i>&lt;0, 1&gt; x = 18 ✓ &lt;0, 2&gt; x += 3 ✓</i>
	replica	1	
	view	0	
	op	2	
	commit	2	

<b>C</b>	status	normal	<i>&lt;0, 1&gt; x = 18 ✓ &lt;0, 2&gt; x += 3 ✓ &lt;0, 3&gt; y = 100</i>
	replica	2	
	view	0	
	op	3	
	commit	2	

## Start view change:

Status = change  
Increment local view  
Send SVC to all nodes



<b>B</b>	status	normal	
	replica	1	
	view	0	
	op	2	
	commit	2	

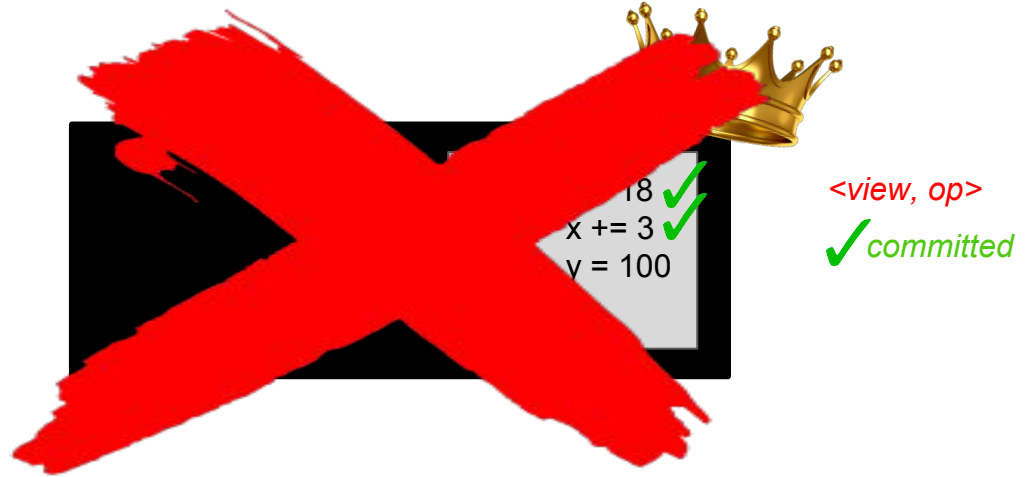
StartViewChange  
view: 1  
replica: 2



<b>C</b>	status	change	
	replica	2	
	view	1	
	op	3	
	commit	2	

## Receive SVC where:

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



<b>B</b>	status	normal	
	replica	1	
	view	0	
	op	2	
	commit	2	

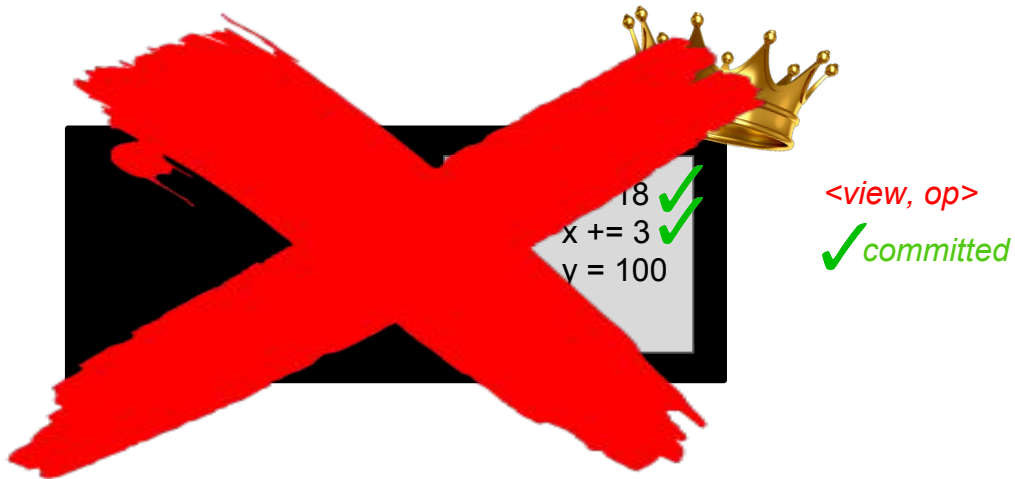
StartViewChange  
view: 1  
replica: 2



<b>C</b>	status	change	
	replica	2	
	view	1	
	op	3	
	commit	2	

## Receive SVC where:

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



<b>B</b>	status	change	
	replica	1	
	view	1	
	op	2	
	commit	2	

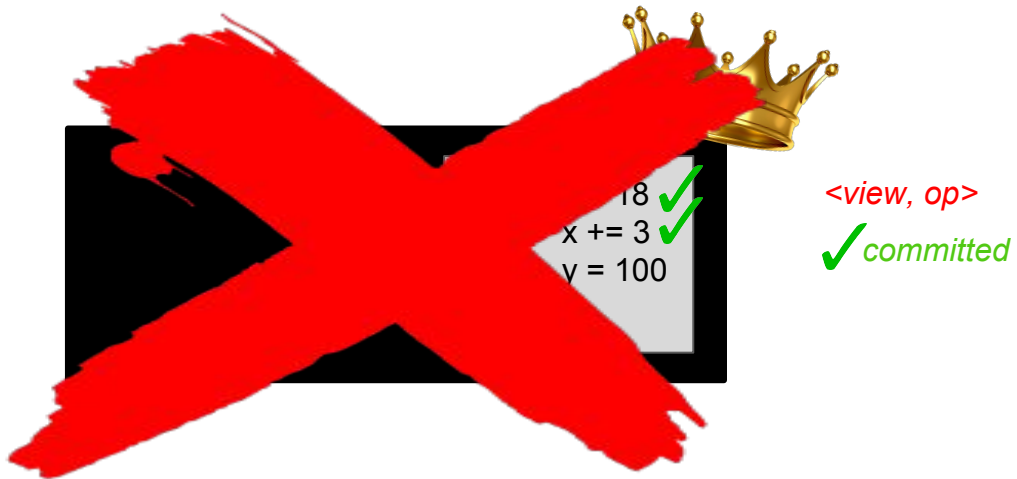
StartViewChange  
view: 1  
replica: 1



<b>C</b>	status	change	
	replica	2	
	view	1	
	op	3	
	commit	2	

**Receive f SVCs where:**

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



<b>B</b>	status	change	
	replica	1	
	view	1	
	op	2	
	commit	2	

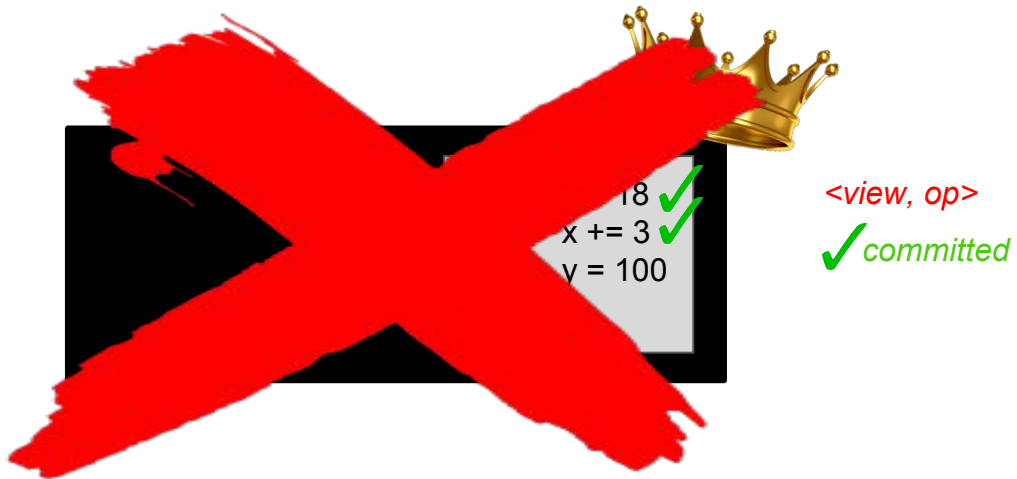
**StartViewChange**  
view: 1  
replica: 1



<b>C</b>	status	change	
	replica	2	
	view	1	
	op	3	
	commit	2	

Receive f SVCs where:

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



DoViewChange

replica: 2

view: 1

op: 3

commit: 2

<log>



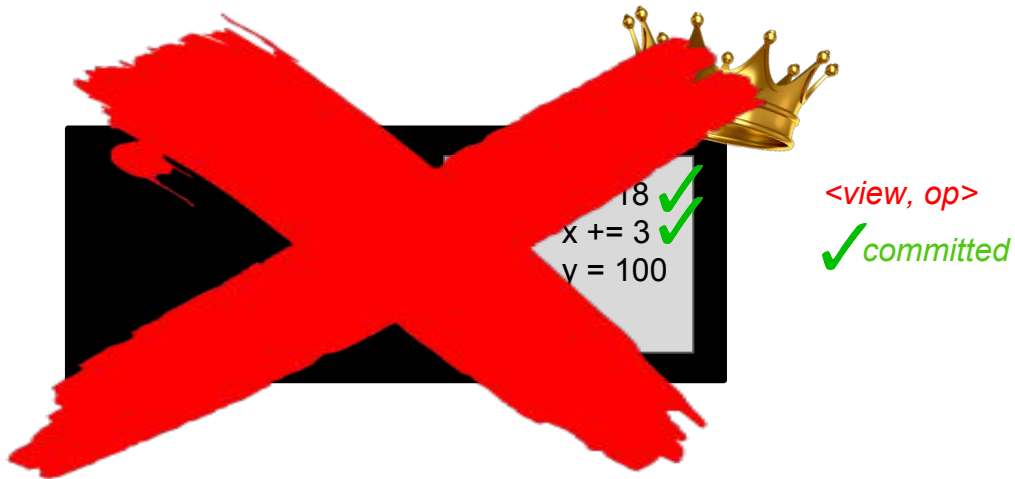
<b>B</b>	status	change	
	replica	1	
	view	1	
	op	2	
	commit	2	

<b>C</b>	status	change	
	replica	2	
	view	1	
	op	3	
	commit	2	

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



<b>B</b>	status	change	<pre>&lt;0, 1&gt; x = 18 ✓ &lt;0, 2&gt; x += 3 ✓ &lt;0, 3&gt; y = 100</pre>
	replica	1	
	view	1	
	op	3	
	commit	2	

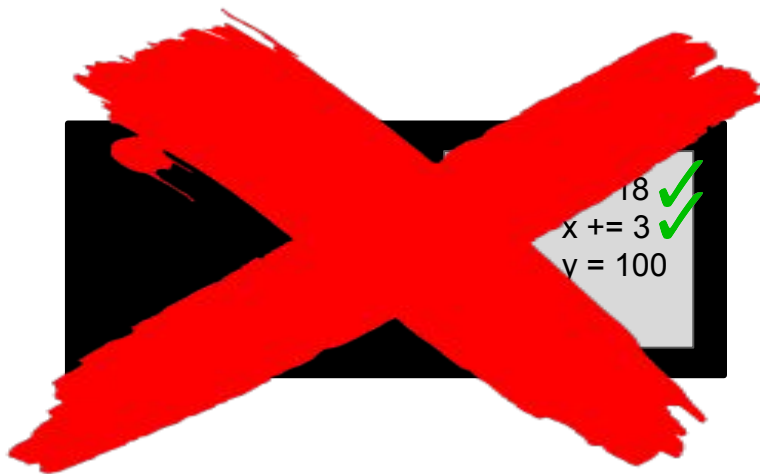
<b>C</b>	status	change	<pre>&lt;0, 1&gt; x = 18 ✓ &lt;0, 2&gt; x += 3 ✓ &lt;0, 3&gt; y = 100</pre>
	replica	2	
	view	1	
	op	3	
	commit	2	

## Receive f DVCs:

*Become new primary*

*Send StartView to others*

Why do we send the log here?



<view, op>

✓ committed



<b>B</b>	status	normal	<pre>&lt;0, 1&gt; x = 18 ✓ &lt;0, 2&gt; x += 3 ✓ &lt;0, 3&gt; y = 100</pre>
	replica	1	
	view	1	
	op	3	
	commit	2	

### StartView

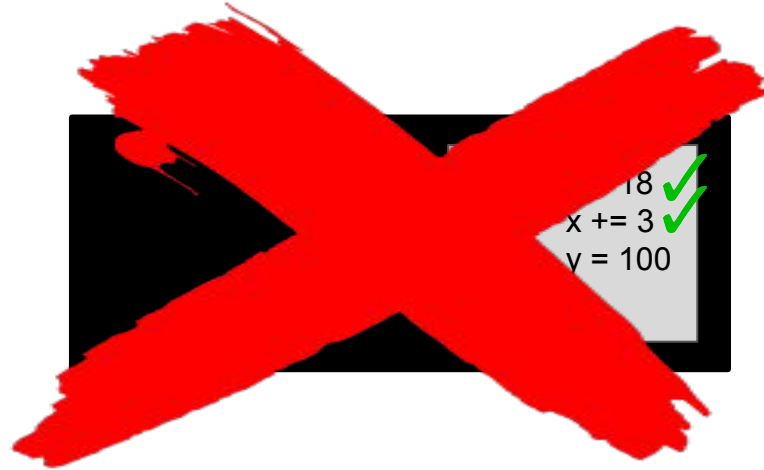
view: 1  
replica: 1  
op: 3  
commit: 2  
<log>

<b>C</b>	status	change	<pre>&lt;0, 1&gt; x = 18 ✓ &lt;0, 2&gt; x += 3 ✓ &lt;0, 3&gt; y = 100</pre>
	replica	2	
	view	1	
	op	3	
	commit	2	



Notice  $\langle 0, 3 \rangle$  is uncommitted and from an old view...

Do we commit it?



$\langle \text{view}, \text{op} \rangle$   
✓ committed



<b>B</b>	status	normal	$\langle 0, 1 \rangle$ x = 18 ✓ $\langle 0, 2 \rangle$ x += 3 ✓ $\langle 0, 3 \rangle$ y = 100
	replica	1	
	view	1	
	op	3	
	commit	2	

PrepareOK  
view: 0  
op: 3  
replica: 2



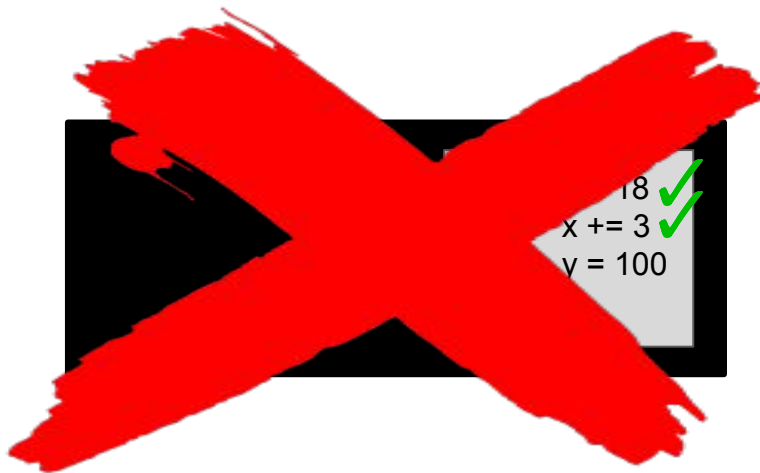
<b>C</b>	status	normal	$\langle 0, 1 \rangle$ x = 18 ✓ $\langle 0, 2 \rangle$ x += 3 ✓ $\langle 0, 3 \rangle$ y = 100
	replica	2	
	view	1	
	op	3	
	commit	2	

Are uncommitted ops like  $\langle 0, 3 \rangle$  guaranteed to survive into the new view?

What about committed ops? (e.g.  $\langle 0, 1 \rangle$  and  $\langle 0, 2 \rangle$ )



<b>B</b>	status	normal	$\langle 0, 1 \rangle$ x = 18 ✓ $\langle 0, 2 \rangle$ x += 3 ✓ $\langle 0, 3 \rangle$ y = 100 ✓
	replica	1	
	view	1	
	op	3	
	commit	3	



$\langle \text{view}, \text{op} \rangle$   
✓ committed

<b>C</b>	status	normal	$\langle 0, 1 \rangle$ x = 18 ✓ $\langle 0, 2 \rangle$ x += 3 ✓ $\langle 0, 3 \rangle$ y = 100
	replica	2	
	view	1	
	op	3	
	commit	2	

# 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?

# Failure detection

# Two kinds of failures

Server failures

Network partitions

These two are indistinguishable from a single machine!

# Failure detection goals

**Completeness:** Each failure is detected

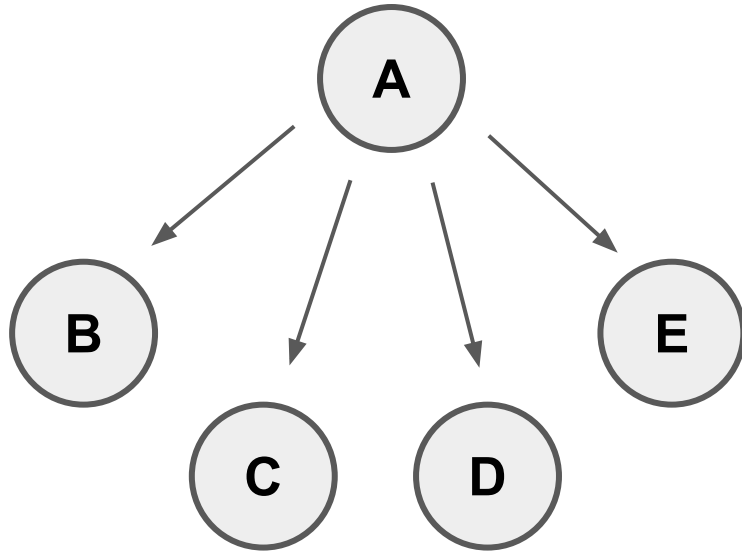
**Accuracy:** There is no mistaken detection

**Speed:** Time to first detection of a failure

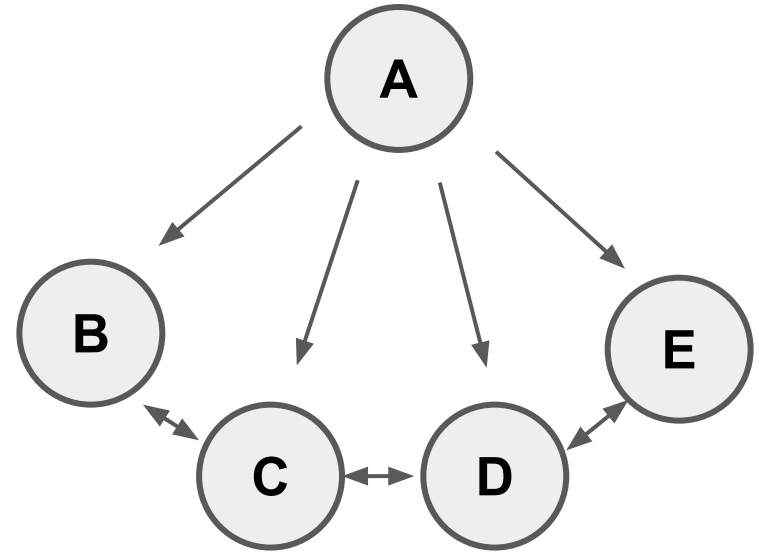
**Scale:** Equal load on each node

... in terms of CPU and network bandwidth

Centralized detection

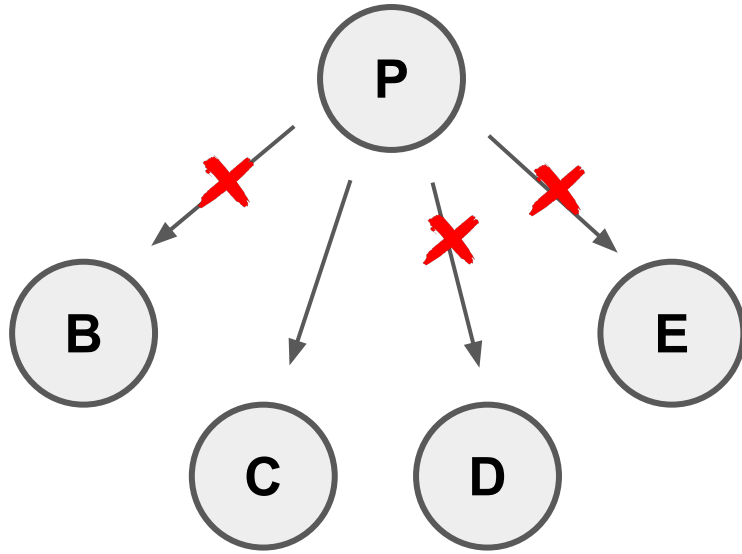


Gossip detection

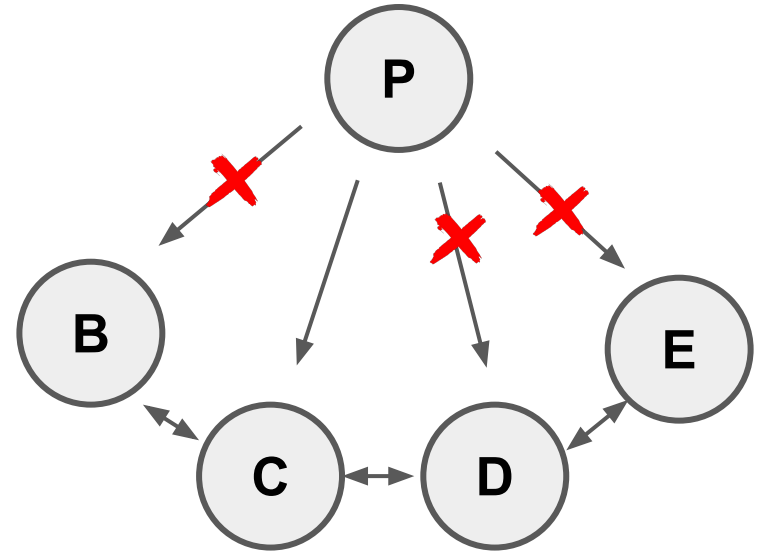


**Completeness, accuracy, speed, load?**

Centralized detection



Gossip detection



If we're running the view change protocol, what happens in each case?



# What is gossip detection good for?

Certainly not viewstamped replication!

May cause *liveness* issues; primary cannot reach  $f$  nodes

Dynamo uses gossip for membership and failure detection

More suitable for completely decentralized environments

# 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