

# Spanner pt. 2

Doug Woos

# Logistics notes

Problem Set 4 out tonight

# Spanner pt 2

1. Distributed transactions, in detail
  - On one Paxos group
  - Between Paxos groups
2. Fast read-only transactions with TrueTime
3. Discussion

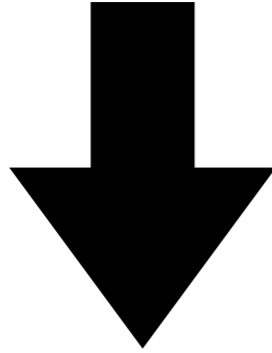
Leader



Client



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

- 1.
- 2.
- 3.
- 4.
- 5.
- 6.
- 7.

Leader



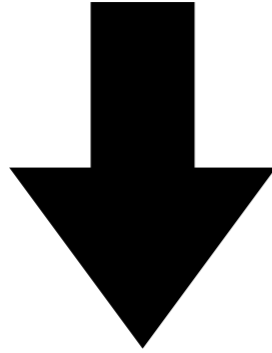
Client



read(checking\_bal)



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

- 1.
- 2.
- 3.
- 4.
- 5.
- 6.
- 7.

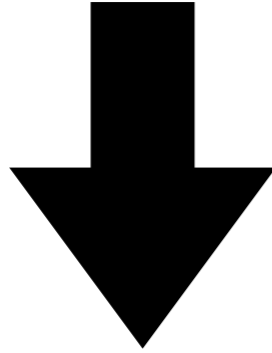
Leader



Client



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
- 2.
- 3.
- 4.
- 5.
- 6.
- 7.

Leader



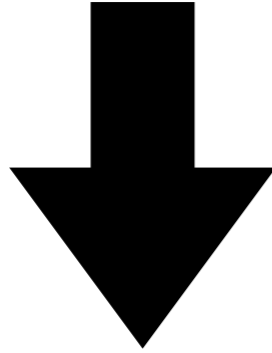
200



Client



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
- 2.
- 3.
- 4.
- 5.
- 6.
- 7.

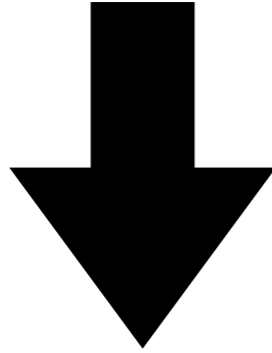
Leader



Client



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
- 2.
- 3.
- 4.
- 5.
- 6.
- 7.



Leader



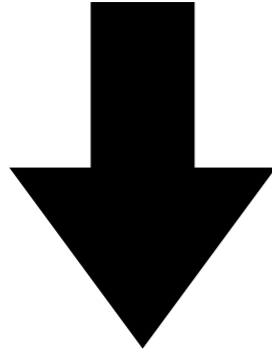
Client



`read(savings_bal)`



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
- 2.
- 3.
- 4.
- 5.
- 6.
- 7.

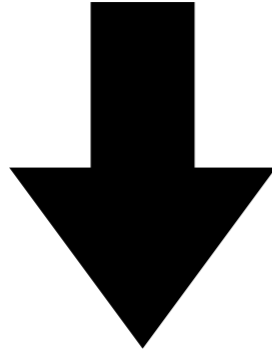
Leader



Client



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(R, savings\_bal)
- 3.
- 4.
- 5.
- 6.
- 7.

Leader

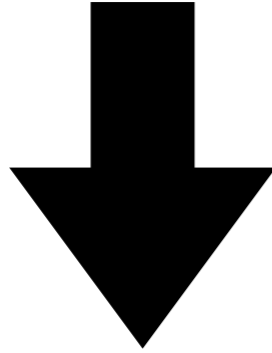


0

Client



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(R, savings\_bal)
- 3.
- 4.
- 5.
- 6.
- 7.

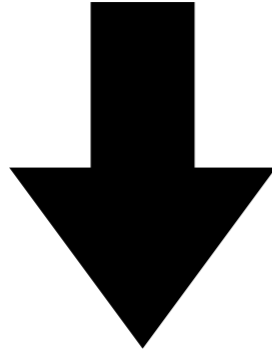
Leader



Client



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(R, savings\_bal)
- 3.
- 4.
- 5.
- 6.
- 7.

Leader



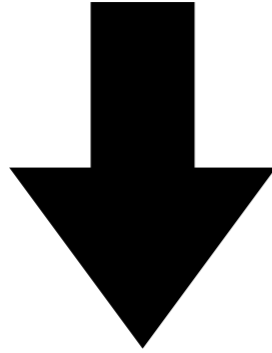
Client



`write(checking_bal, 100)`



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(R, savings\_bal)
- 3.
- 4.
- 5.
- 6.
- 7.

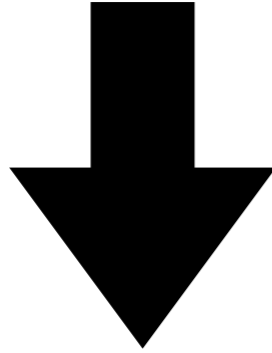
Leader



Client



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(R, savings\_bal)
3. A: Lock(W, checking\_bal)
4. A: Write(checking\_bal, 100)
- 5.
- 6.
- 7.

Leader

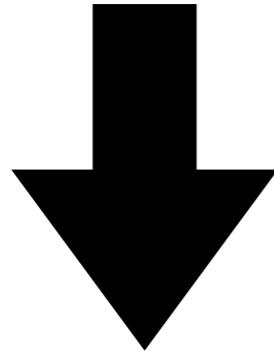


OK

Client



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(R, savings\_bal)
3. A: Lock(W, checking\_bal)
4. A: Write(checking\_bal, 100)
- 5.
- 6.
- 7.

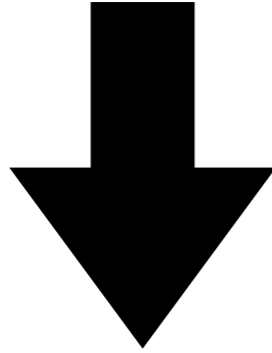
Leader



Client



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(R, savings\_bal)
3. A: Lock(W, checking\_bal)
4. A: Write(checking\_bal, 100)
- 5.
- 6.
- 7.



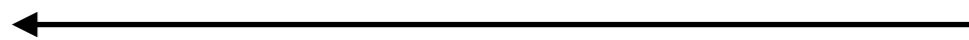
Leader



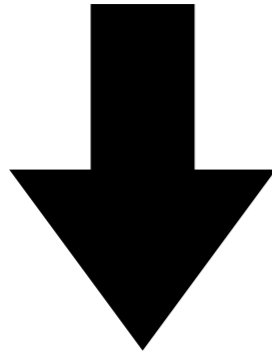
Client



`write(savings_bal, 100)`



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(R, savings\_bal)
3. A: Lock(W, checking\_bal)
4. A: Write(checking\_bal, 100)
- 5.
- 6.
- 7.

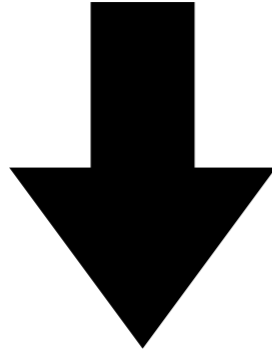
Leader



Client



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(R, savings\_bal)
3. A: Lock(W, checking\_bal)
4. A: Write(checking\_bal, 100)
5. A: Lock(W, savings\_bal)
6. A: Write(savings\_bal, 100)
- 7.

Leader



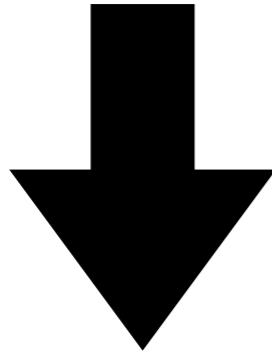
OK



Client



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(R, savings\_bal)
3. A: Lock(W, checking\_bal)
4. A: Write(checking\_bal, 100)
5. A: Lock(W, savings\_bal)
6. A: Write(savings\_bal, 100)
- 7.

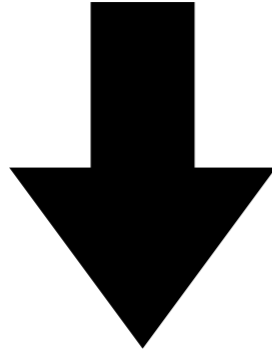
Leader



Client



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(R, savings\_bal)
3. A: Lock(W, checking\_bal)
4. A: Write(checking\_bal, 100)
5. A: Lock(W, savings\_bal)
6. A: Write(savings\_bal, 100)
- 7.

Leader

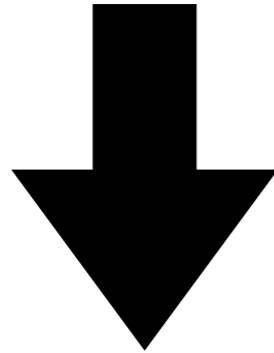


commit

Client



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(R, savings\_bal)
3. A: Lock(W, checking\_bal)
4. A: Write(checking\_bal, 100)
5. A: Lock(W, savings\_bal)
6. A: Write(savings\_bal, 100)
- 7.

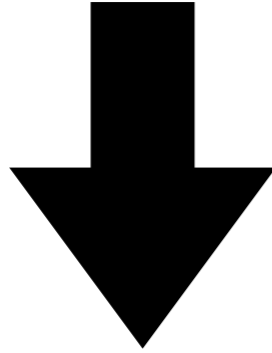
Leader



Client



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(R, savings\_bal)
3. A: Lock(W, checking\_bal)
4. A: Write(checking\_bal, 100)
5. A: Lock(W, savings\_bal)
6. A: Write(savings\_bal, 100)
7. A: Commit

Leader



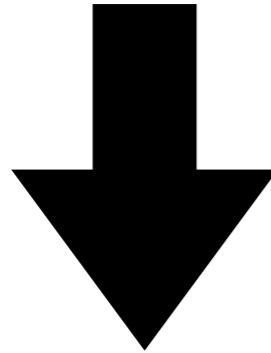
OK



Client



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(R, savings\_bal)
3. A: Lock(W, checking\_bal)
4. A: Write(checking\_bal, 100)
5. A: Lock(W, savings\_bal)
6. A: Write(savings\_bal, 100)
7. A: Commit

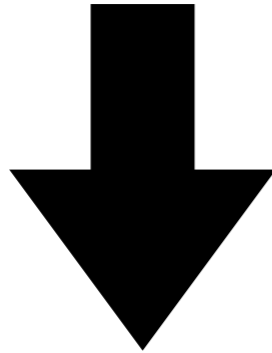
Leader



Client



Follower



Follower



```
x = read(checking_bal)
if (x > 100) {
    y = read(savings_bal)
    write(checking_bal, x - 100)
    write(savings_bal, y + 100)
}
```

LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(R, savings\_bal)
3. A: Lock(W, checking\_bal)
4. A: Write(checking\_bal, 100)
5. A: Lock(W, savings\_bal)
6. A: Write(savings\_bal, 100)
7. A: Commit
8. A: Unlock(checking\_bal)
9. A: Unlock(savings\_bal)



Leader 1



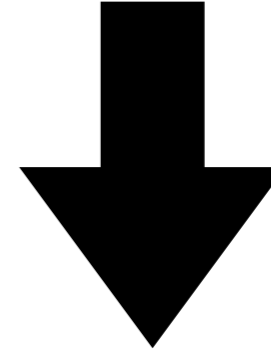
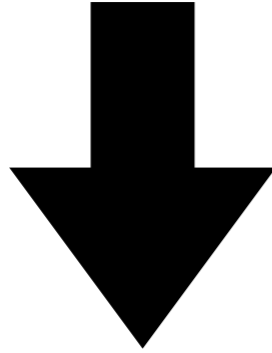
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

LOG:

Leader 1



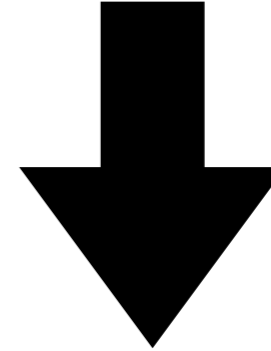
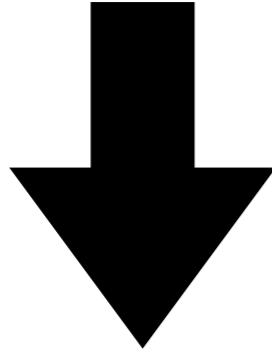
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

LOG:

Leader 1



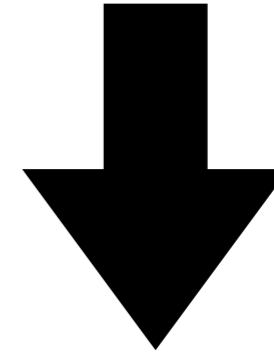
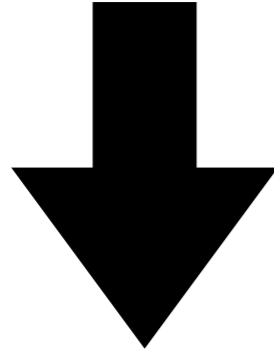
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)

LOG:

Leader 1



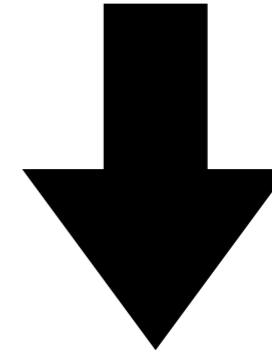
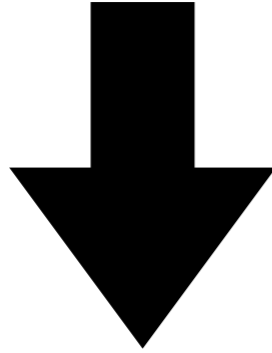
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)

LOG:

Leader 1



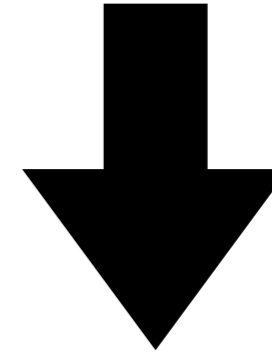
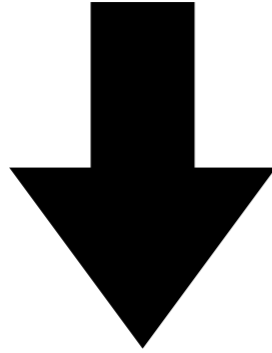
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)

LOG:

1. A: Lock(R, savings\_bal)

Leader 1



Client

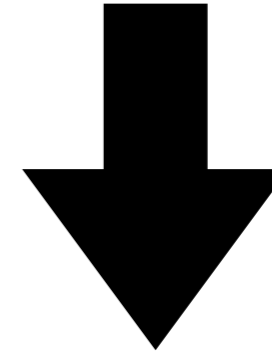
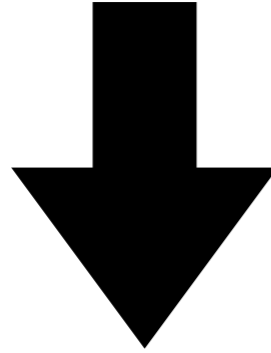


Leader 2



←

```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)

LOG:

1. A: Lock(R, savings\_bal)

Leader 1



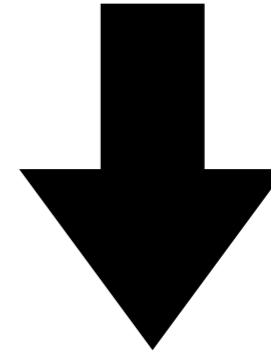
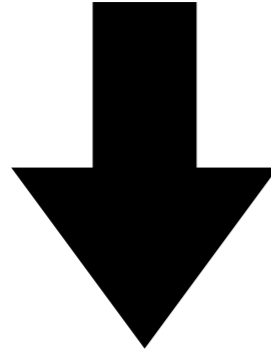
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)

LOG:

1. A: Lock(R, savings\_bal)

Leader 1



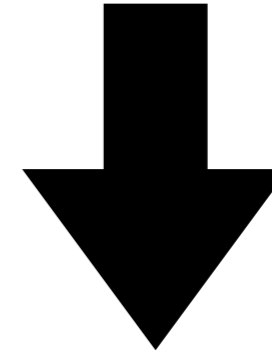
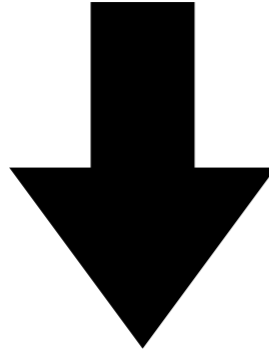
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)

LOG:

1. A: Lock(R, savings\_bal)



Leader 1



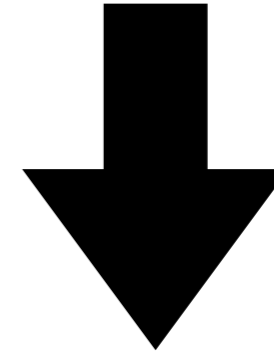
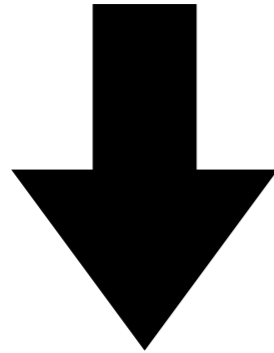
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)

Leader 1



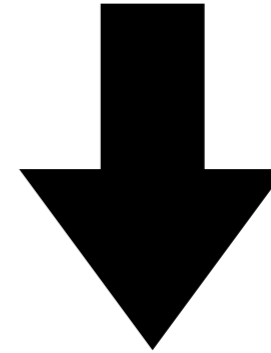
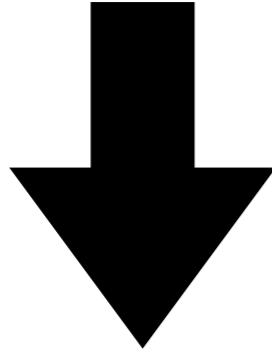
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)

Leader 1



commit(1)



Client



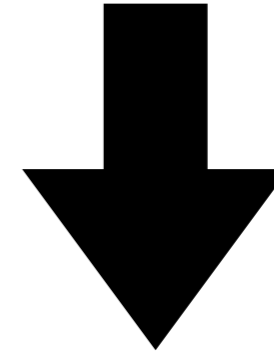
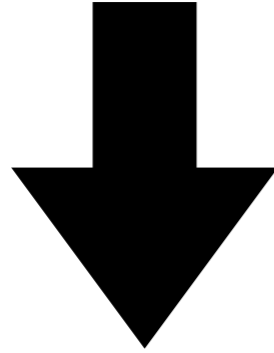
commit(1)



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)

Leader 1



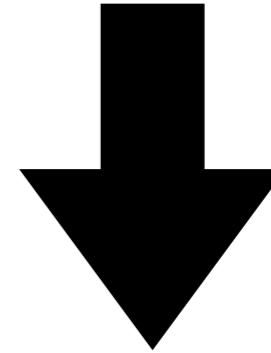
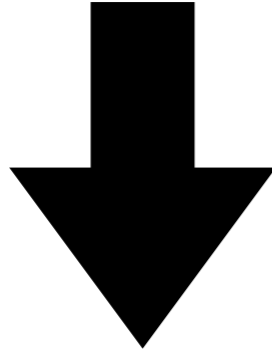
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)
4. A: Prepare

Leader 1



Client

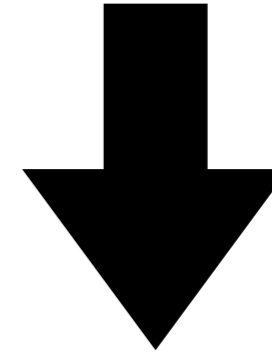
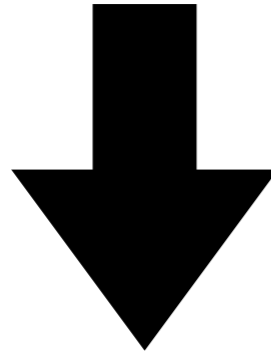


OK

Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)
4. A: Prepare

Leader 1



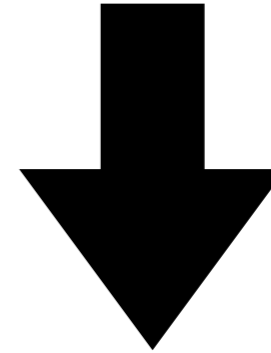
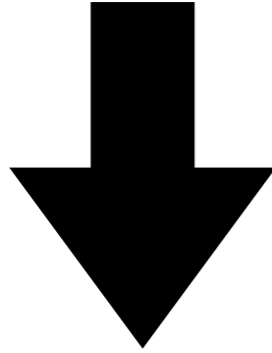
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)
4. A: Commit

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)
4. A: Prepare

Leader 1



OK



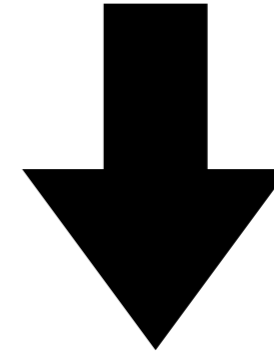
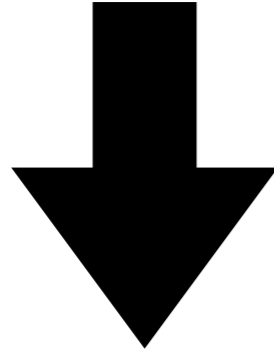
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)
4. A: Commit

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)
4. A: Prepare

Leader 1



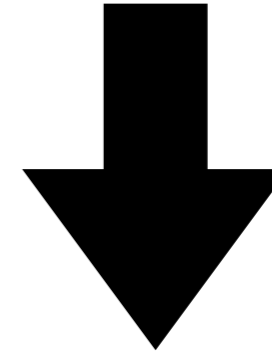
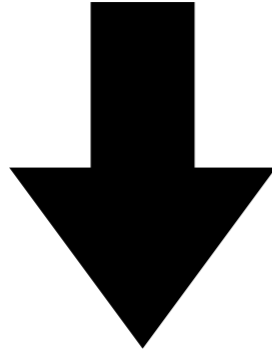
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)
4. A: Commit
5. A: Unlock(checking\_bal)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)
4. A: Prepare



Leader 1



commit

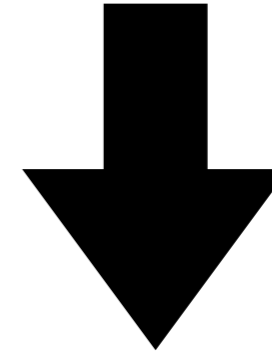
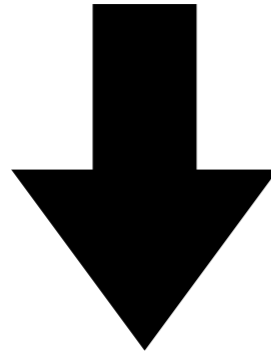
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)
4. A: Commit
5. A: Unlock(checking\_bal)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)
4. A: Prepare

Leader 1



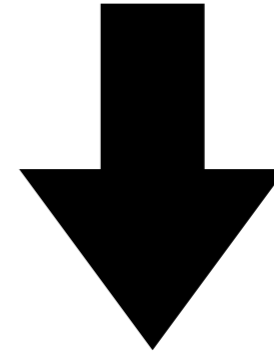
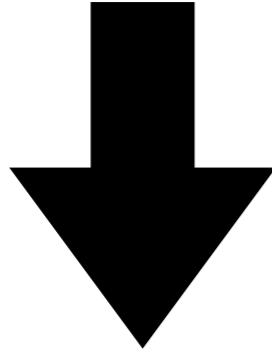
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)
4. A: Commit
4. A: Unlock(checking\_bal)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)
4. A: Prepare
5. A: Commit

Leader 1



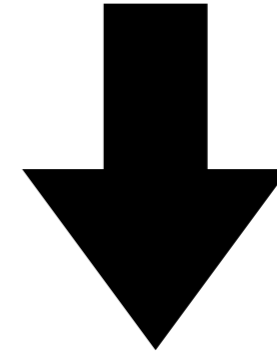
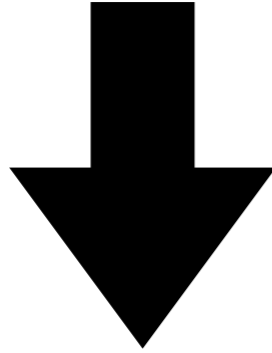
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)
4. A: Commit
4. A: Unlock(checking\_bal)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)
4. A: Prepare
5. A: Commit
6. A: Unlock(savings\_bal)

# How can we get fast reads?

R/W transactions are complicated!

- And slow

Can we do fast, lock-free reads?

- Real time to the rescue

# TrueTime

API that exposes real time, with uncertainty

```
{earliest: e, latest: l} = TT.now()
```

“Real time” is between `earliest` and `latest`

Time is an illusion; lunchtime, doubly so

If I call `TT.now()` on two nodes simultaneously, intervals *guaranteed* to overlap!

And: if intervals don't overlap, the later one really happened later!

# TrueTime usage

Assign a timestamp to each transaction

- At each Paxos group, timestamp increases monotonically
- Globally, if T1 returns before T2 starts,  
 $\text{timestamp}(T1) < \text{timestamp}(T2)$

# TrueTime usage

Timestamp for an RW transaction chosen by coordinator leader

Timestamps for R/W transactions is max of:

- Local time
- Prepare timestamps at every participant
- Timestamp of any previous local transaction

# Commit wait

Need to ensure that all future transactions will get a higher timestamp

Therefore, need to wait until

$TT.now() > \text{transaction timestamp}$



Leader 1



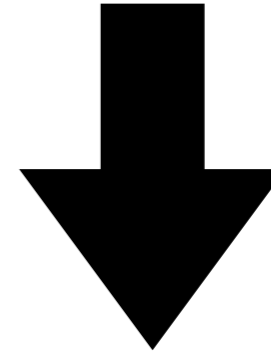
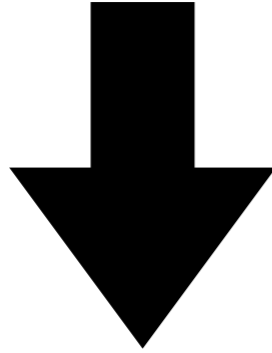
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

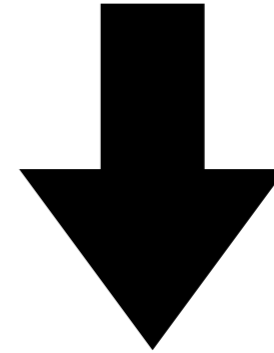
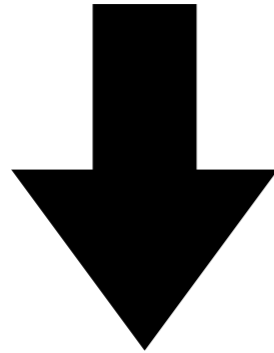
1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

- 1. A: Lock(R, checking\_bal)
- 2. A: Lock(W, checking\_bal)
- 3. A: Write(checking\_bal, 100)

LOG:

- 1. A: Lock(R, savings\_bal)
- 2. A: Lock(W, savings\_bal)
- 3. A: Write(savings\_bal, 100)

Leader 1



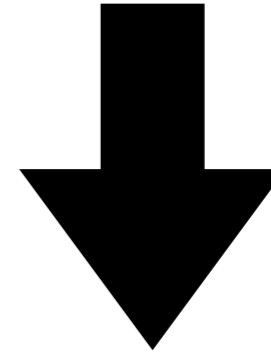
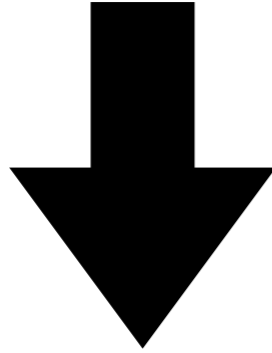
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)
4. A: Prepare@t1

Leader 1



Client

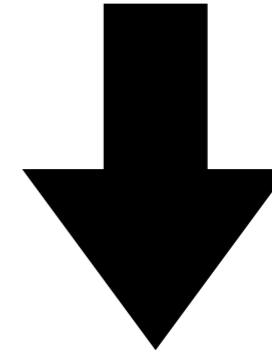
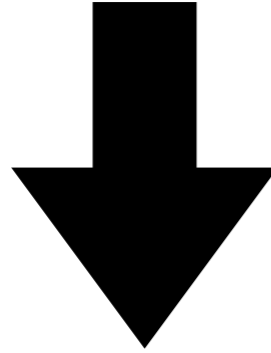


Leader 2



OK@t1

```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)
4. A: Prepare@t1

Leader 1



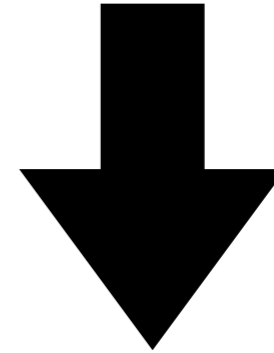
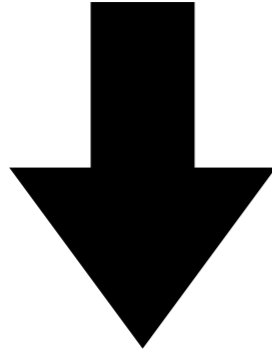
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)
4. A: Commit@max(t1, tleader)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)
4. A: Prepare@t1

Leader 1



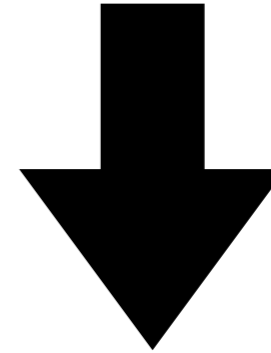
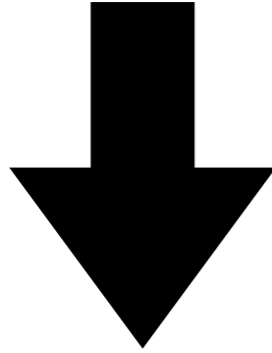
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)
4. A: Commit@max(t1, tleader)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)
4. A: Prepare@t1

Waiting...

Leader 1



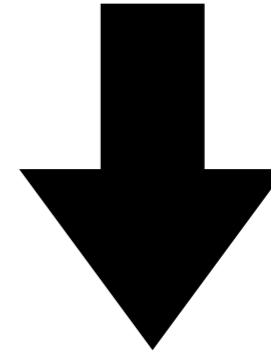
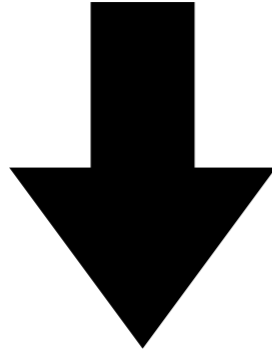
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)
4. A: Commit@max(t1, tleader)
5. A: Unlock(checking\_bal)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)
4. A: Prepare@t1

Leader 1



OK



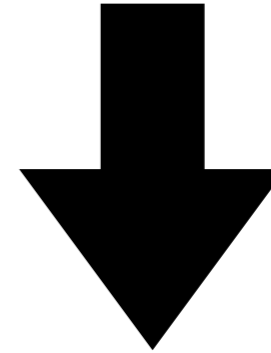
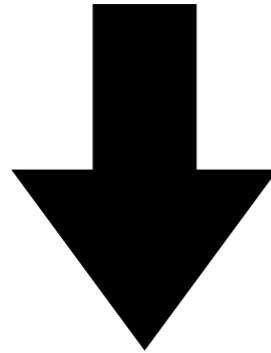
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)
4. A: Commit@max(t1, tleader)
5. A: Unlock(checking\_bal)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)
4. A: Prepare@t1



Leader 1



commit@T

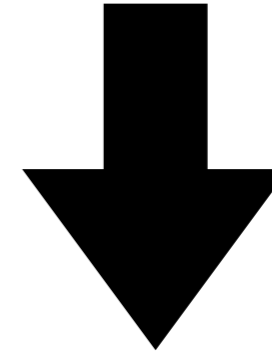
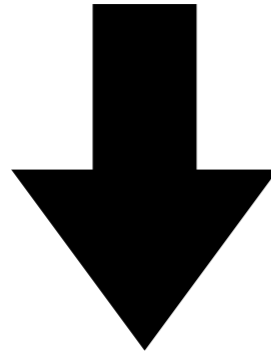
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)
4. A: Commit@max(t1, tleader)
5. A: Unlock(checking\_bal)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)
4. A: Prepare@t1

Leader 1



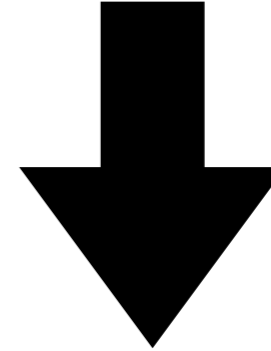
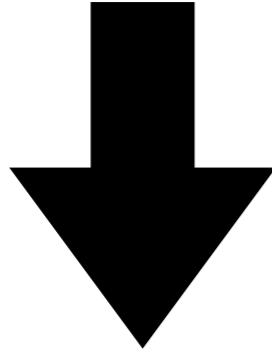
Client



Leader 2



```
x = read(checking_bal)
if (x > 100) {
  y = read(savings_bal)
  write(checking_bal, x - 100)
  write(savings_bal, y + 100)
}
```



LOG:

1. A: Lock(R, checking\_bal)
2. A: Lock(W, checking\_bal)
3. A: Write(checking\_bal, 100)
4. A: Commit@max(t1, tleader)
5. A: Unlock(checking\_bal)

LOG:

1. A: Lock(R, savings\_bal)
2. A: Lock(W, savings\_bal)
3. A: Write(savings\_bal, 100)
4. A: Prepare@t1
5. A: Commit@max(t1, tleader)
6. A: Unlock(checking\_bal)

# TrueTime usage

Timestamps for RO transactions

- Client can choose a timestamp
- For external consistency, always safe to choose

`TT.now().latest`

A response from a replica is a promise: will never have a new transaction that commits before timestamp

Can read from any replica that has seen a Paxos write after timestamp

# TrueTime implementation

GPS, atomic clocks

Armageddon masters and timeslave daemons

All local clocks synced with masters, and expose uncertainty to local apps

Assumptions made about local clock drift

# Conclusions

What's cool about Spanner?

- Distributed transactions with decent performance
- What makes that possible?
- Read-only transactions with great performance
- What makes that possible?

Clocks are a form of communication!

# Discussion

CPU errors and bad clocks

- When does Google discover errors?

Disaster preparedness

Is Spanner's complexity scary?

What happens if there is high clock skew?

- If we know about it
- If we don't?

