

The Paxos Register

Harry C. Li, Allen Clement, Amitanand Aiyer, Lorenzo Alvisi

The University of Texas at Austin

1

Consensus & Paxos

- Validity
- Integrity
- Agreement
- Termination

2

Consensus & Paxos

Part-time Parliament

[Lamport '98]

- Validity
- Integrity
- Agreement
- Termination

2

Consensus & Paxos

Part-time Parliament

[Lamport '98]

Safety

- Validity
- Integrity
- Agreement
- Termination

Liveness

2

Consensus & Paxos

Safety

- Validity
- Integrity
- Agreement
- Termination

Liveness

Part-time Parliament
Frangipani

[Lamport '98]
[Thekkath et al. '97]

2

Consensus & Paxos

Safety

- Validity
- Integrity
- Agreement
- Termination

Liveness

Part-time Parliament
Frangipani
Byzantine Paxos

[Lamport '98]
[Thekkath et al. '97]
[Castro & Liskov '99]

2

Consensus & Paxos

Safety

- Validity
- Integrity
- Agreement
- Termination

Liveness

Part-time Parliament
Frangipani
Byzantine Paxos
Disk Paxos

[Lamport '98]
[Thekkath et al. '97]
[Castro & Liskov '99]
[Gafni & Lamport '00]

2

Consensus & Paxos

Safety

- Validity
- Integrity
- Agreement
- Termination

Liveness

Part-time Parliament
Frangipani
Byzantine Paxos
Disk Paxos
Deconstructing Paxos
Reconstructing Paxos

[Lamport '98]
[Thekkath et al. '97]
[Castro & Liskov '99]
[Gafni & Lamport '00]
[Boichat et al. '01]
[Boichat et al. '01]

2

Consensus & Paxos

Safety

- Validity
- Integrity
- Agreement
- Termination

Liveness

Part-time Parliament
Frangipani
Byzantine Paxos
Disk Paxos
Deconstructing Paxos
Reconstructing Paxos
Active Disk Paxos

[Lamport '98]
[Thekkath et al. '97]
[Castro & Liskov '99]
[Gafni & Lamport '00]
[Boichat et al. '01]
[Boichat et al. '01]
[Chockler & Malkhi '02]

2

Consensus & Paxos

Safety

- Validity
- Integrity
- Agreement
- Termination

Liveness

Part-time Parliament
Frangipani
Byzantine Paxos
Disk Paxos
Deconstructing Paxos
Reconstructing Paxos
Active Disk Paxos
Agreement & Execution

[Lamport '98]
[Thekkath et al. '97]
[Castro & Liskov '99]
[Gafni & Lamport '00]
[Boichat et al. '01]
[Boichat et al. '01]
[Chockler & Malkhi '02]
[Yin et al. '03]

2

Consensus & Paxos

Safety

- Validity
- Integrity
- Agreement
- Termination

Liveness

Part-time Parliament
Frangipani
Byzantine Paxos
Disk Paxos
Deconstructing Paxos
Reconstructing Paxos
Active Disk Paxos
Agreement & Execution
Byzantine Disk Paxos

[Lamport '98]
[Thekkath et al. '97]
[Castro & Liskov '99]
[Gafni & Lamport '00]
[Boichat et al. '01]
[Boichat et al. '01]
[Chockler & Malkhi '02]
[Yin et al. '03]
[Abraham et al. '04]

2

Consensus & Paxos

Safety

- Validity
- Integrity
- Agreement
- Termination

Liveness

Part-time Parliament
Frangipani
Byzantine Paxos
Disk Paxos
Deconstructing Paxos
Reconstructing Paxos
Active Disk Paxos
Agreement & Execution
Byzantine Disk Paxos
Fast Byzantine Paxos
Fast Paxos

[Lamport '98]
[Thekkath et al. '97]
[Castro & Liskov '99]
[Gafni & Lamport '00]
[Boichat et al. '01]
[Boichat et al. '01]
[Chockler & Malkhi '02]
[Yin et al. '03]
[Abraham et al. '04]
[Martin & Alvisi '05]
[Lamport '05]

2

Consensus & Paxos

Safety

- Validity
- Integrity
- Agreement
- Termination

Liveness

Part-time Parliament	[Lamport '98]
Frangipani	[Thekkath et al. '97]
Byzantine Paxos	[Castro & Liskov '99]
Disk Paxos	[Gafni & Lamport '00]
Deconstructing Paxos	[Boichat et al. '01]
Reconstructing Paxos	[Boichat et al. '01]
Active Disk Paxos	[Chockler & Malkhi '02]
Agreement & Execution	[Yin et al. '03]
Byzantine Disk Paxos	[Abraham et al. '04]
Fast Byzantine Paxos	[Martin & Alvisi '06]
Fast Paxos	[Lamport '05]
Hybrid Quorums	[Cowling et al. '06]
Chubby	[Burrows '06]
Zyzyva	[Kotla et al. '07]

2

What Makes Paxos Paxos?

3

What Makes Paxos Paxos?

- Guarantee safety
- Provide liveness during periods of synchrony
- Use leaders to coordinate quorums

3

What Makes Paxos Paxos?

- Guarantee safety
- Provide liveness during periods of synchrony
- Use leaders to coordinate quorums

Is that it?

3

Understanding Paxos

Abstract Paxos	[Lamport '01]
◇ Register & ◇ Leader	[Boichat et al. '01]
Ranked Register	[Chockler & Malkhi '05]
Alpha of Indulgent Consensus	[Guerraoui & Raynal '07]

4

Understanding Paxos

Crash and Byzantine

Abstract Paxos	[Lamport '01]
◇ Register & ◇ Leader	[Boichat et al. '01]
Ranked Register	[Chockler & Malkhi '05]
Alpha of Indulgent Consensus	[Guerraoui & Raynal '07]

4

Understanding Paxos

Crash and Byzantine

Abstract Paxos	[Lampson '01]
◇ Register & ◇ Leader	[Boichat et al. '01]
Ranked Register	[Chockler & Malkhi '05]
Alpha of Indulgent Consensus	[Guerraoui & Raynal '07]

Register-based

4

Understanding Paxos

Crash and Byzantine

Abstract Paxos	[Lampson '01]
◇ Register & ◇ Leader	[Boichat et al. '01]
Ranked Register	[Chockler & Malkhi '05]
Alpha of Indulgent Consensus	[Guerraoui & Raynal '07]

Register-based

No register-based description encompassing
both crash and Byzantine variants

4

Paxos Register

5

Paxos Register

- Captures crash and Byzantine Paxos variants

5

Paxos Register

- Captures crash and Byzantine Paxos variants
- Focuses on most difficult part of Paxos

5

Paxos Register

- Captures crash and Byzantine Paxos variants
- Focuses on most difficult part of Paxos
- Implements a write-once register

5

Paxos Register

- Captures crash and Byzantine Paxos variants
- Focuses on most difficult part of Paxos
- Implements a write-once register

read \longrightarrow write

5

Paxos Register

- Captures crash and Byzantine Paxos variants
- Focuses on most difficult part of Paxos
- Implements a write-once register

read \longrightarrow write
●
token

5

Paxos Register

- Captures crash and Byzantine Paxos variants
- Focuses on most difficult part of Paxos
- Implements a write-once register

read \longrightarrow write
●
token

- Describes Classic, Byzantine, and FaB Paxos
- Leads to Byzantine Paxos using secret sharing

6

Outline

- Motivation
- Paxos Register specification
- Crash implementation
- Byzantine implementation

7

Consensus & Paxos Register

world time \longrightarrow

8

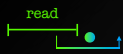
Consensus & Paxos Register

read
└───┘

world time \longrightarrow

8

Consensus & Paxos Register



world time

8

Consensus & Paxos Register



world time

8

Consensus & Paxos Register



world time

8

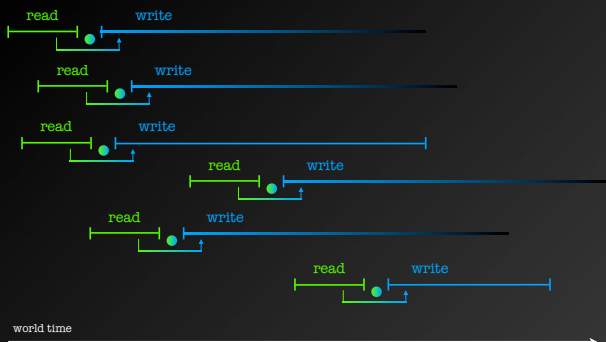
Consensus & Paxos Register



world time

8

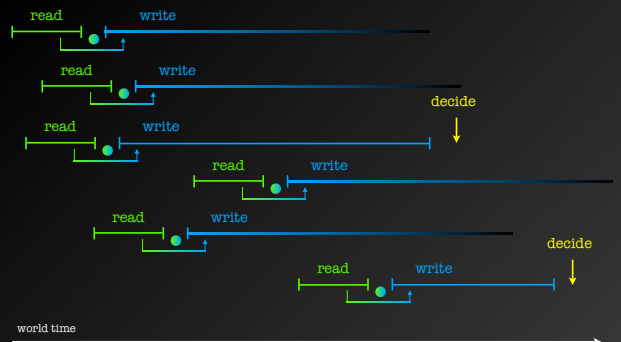
Consensus & Paxos Register



world time

8

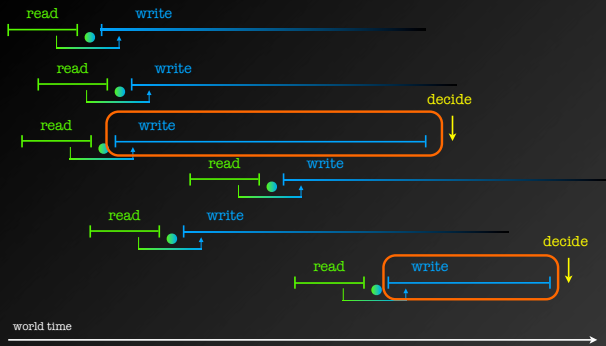
Consensus & Paxos Register



world time

8

Consensus & Paxos Register



8

Details

world time

9

Details

read
token ($v=1, ts=0$)

world time

9

Details

read
token ($v=1, ts=0$)

write
 $v=a, token = (1,0)$

write starts as **partial**

world time

9

Details

read
token ($v=1, ts=0$)

write
 $v=a, token = (1,0)$

write starts as **partial**
partial write may become **visible**

world time

9

Details

read
token ($v=1, ts=0$)

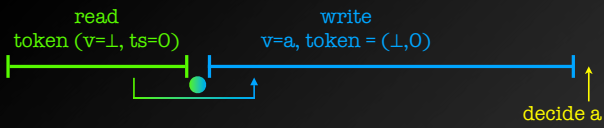
write
 $v=a, token = (1,0)$

write starts as **partial**
partial write may become **visible**
visible write may become **total**

world time

9

Details

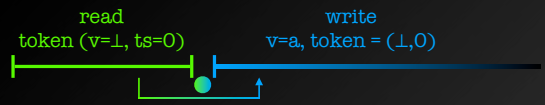


write starts as **partial**
 partial write may become **visible**
 visible write may become **total**

world time

9

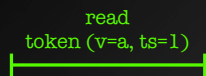
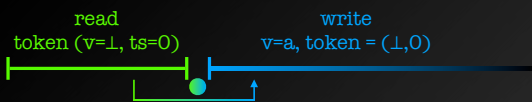
Details



world time

10

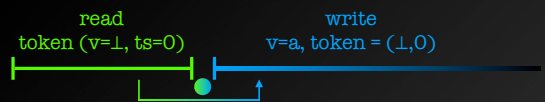
Details



world time

10

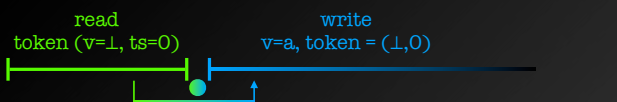
Details



world time

10

Details



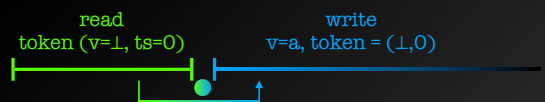
decide a



world time

10

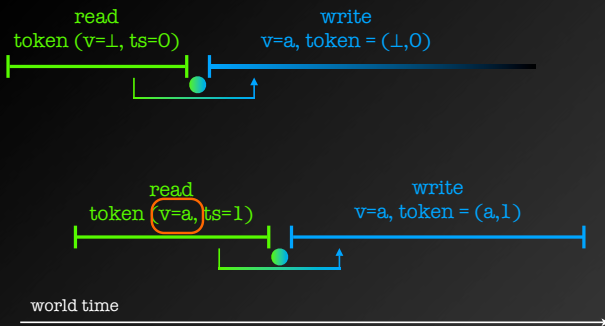
Details



world time

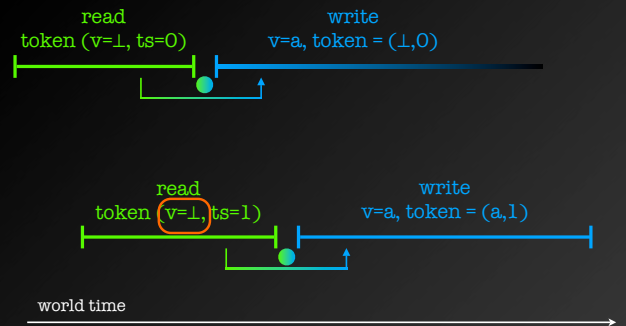
11

Details



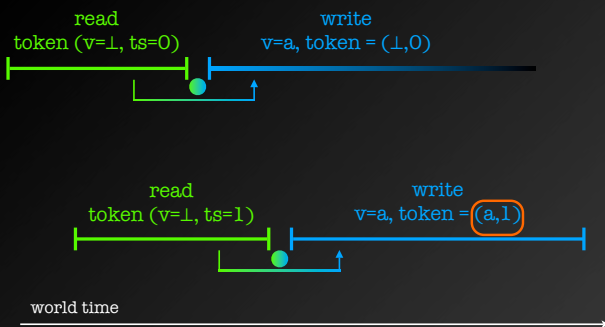
11

Details



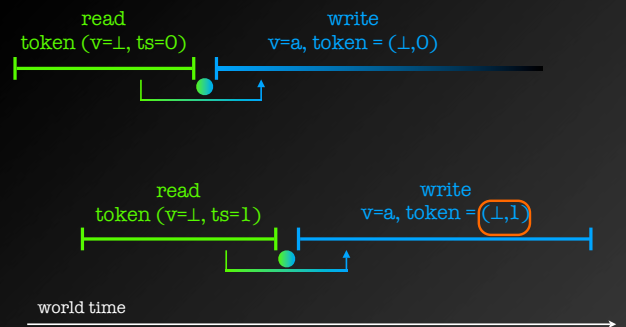
11

Details



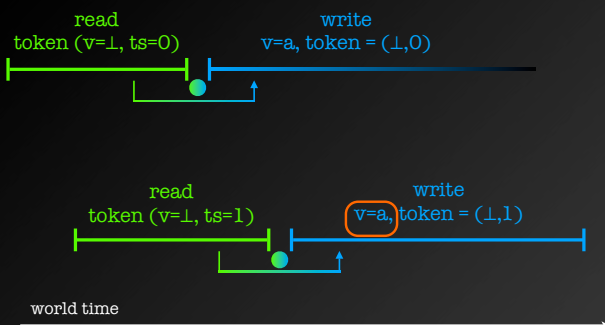
11

Details



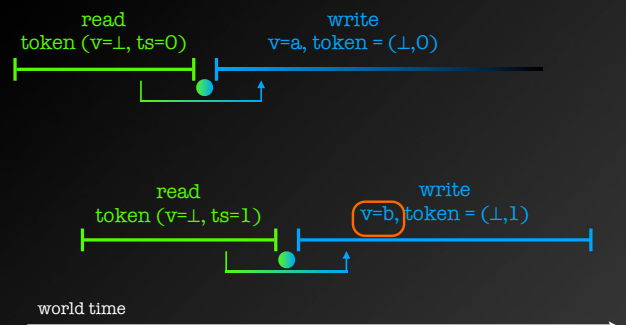
11

Details



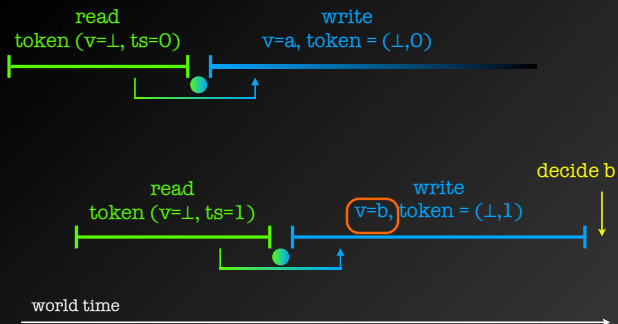
11

Details



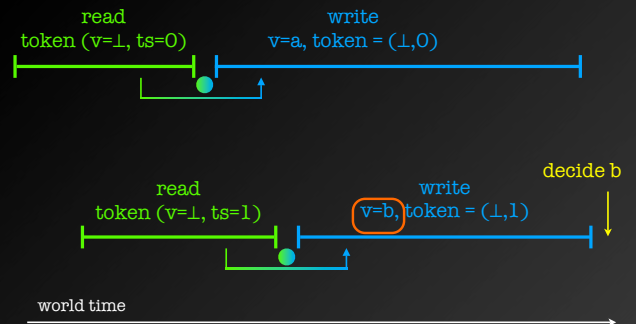
11

Details



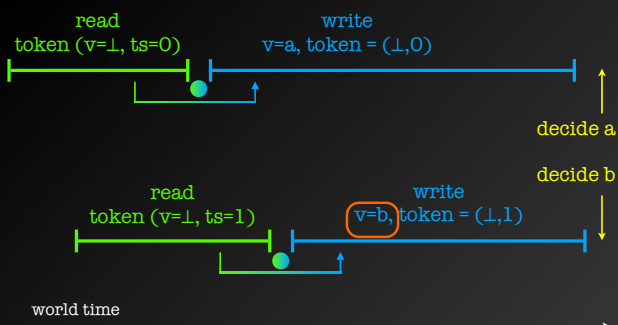
11

Details



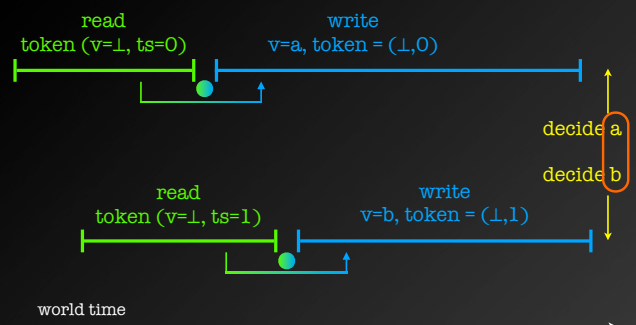
11

Details



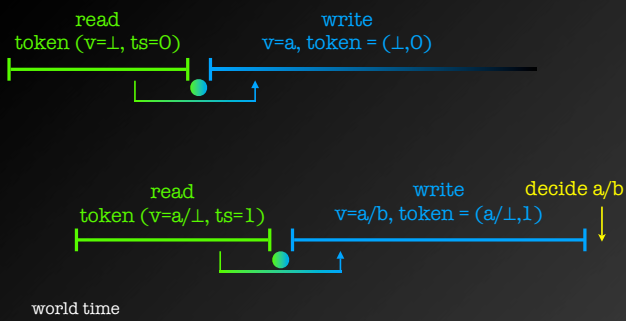
11

Details



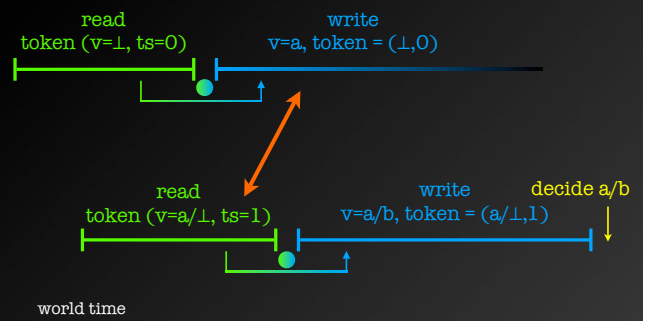
11

Concurrent Operations



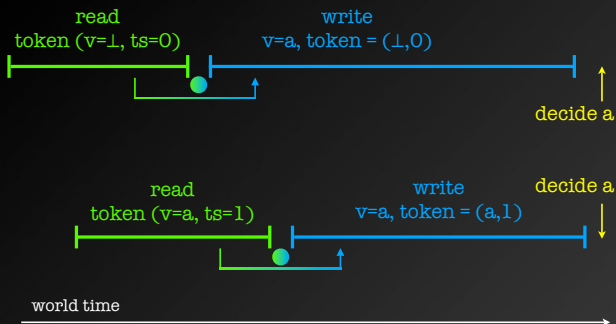
12

Concurrent Operations



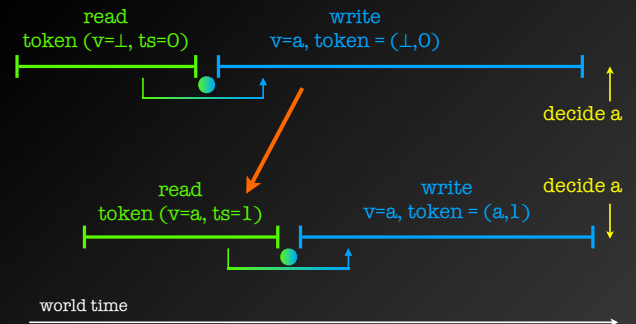
12

Almost Regular



13

Almost Regular



13

Consistency Semantics

read returns last written value or any value being concurrently written

14

Consistency Semantics

read returns last written value or any value being concurrently written

- Paxos Register **redefines** 'concurrent'

14

Consistency Semantics

read returns last written value or any value being concurrently written

- Paxos Register **redefines** 'concurrent'
- Regular w.r.t. **modified partial order**
 1. Writes ordered by increasing timestamp
 2. Total write precedes higher timestamped read
 3. Read precedes higher or same timestamped write

14

Consistency Semantics

read returns last written value or any value being concurrently written

- Paxos Register **redefines** 'concurrent'
- Regular w.r.t. **modified partial order**
 1. Writes ordered by increasing timestamp
 2. Total write precedes higher timestamped read
 3. Read precedes higher or same timestamped write

14

Modified Partial Order

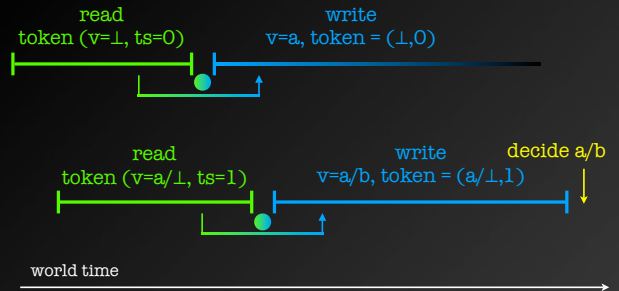
2. Total write precedes higher timestamped read



15

Modified Partial Order

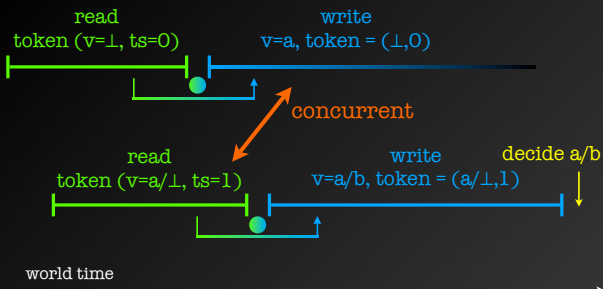
2. Total write precedes higher timestamped read



15

Modified Partial Order

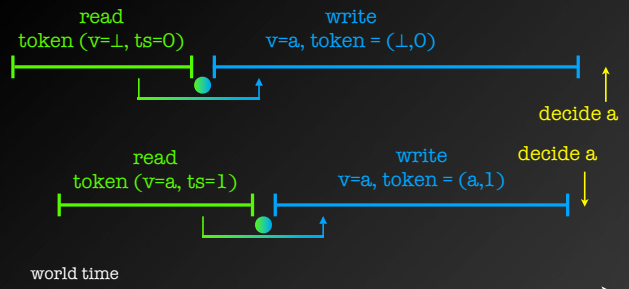
2. Total write precedes higher timestamped read



15

Modified Partial Order

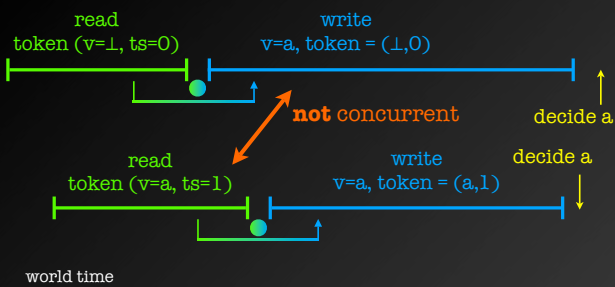
2. Total write precedes higher timestamped read



16

Modified Partial Order

2. Total write precedes higher timestamped read



16

Write-Once Register

17

Write-Once Register

Theorem: All total writes write the same value.

17

Write-Once Register

Theorem: All total writes write the same value.

- Paxos variants share high-level protocol

17

Write-Once Register

Theorem: All total writes write the same value.

- Paxos variants share high-level protocol
- Paxos variants use common abstractions
 - Timestamps, tokens
 - Partial, total, and visible writes

17

Write-Once Register

Theorem: All total writes write the same value.

- Paxos variants share high-level protocol
- Paxos variants use common abstractions
 - Timestamps, tokens
 - Partial, total, and visible writes
 - Implementations differ

18

Outline

- Motivation
- Paxos Register specification
- Crash Implementation
- Byzantine Implementation

19

Crash Model

- Asynchronous system
- Proposers, acceptors, learners
- At least one proposer correct
- Majority of acceptors correct

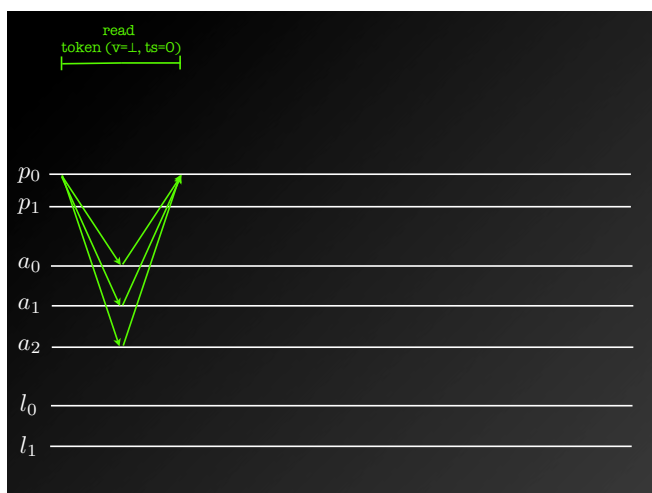
20



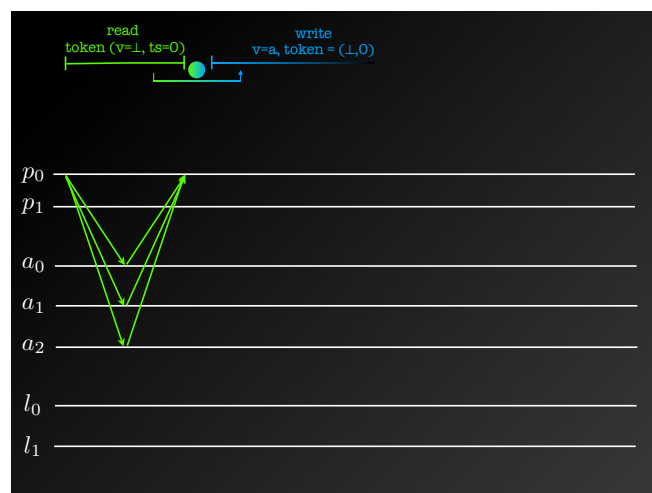
21



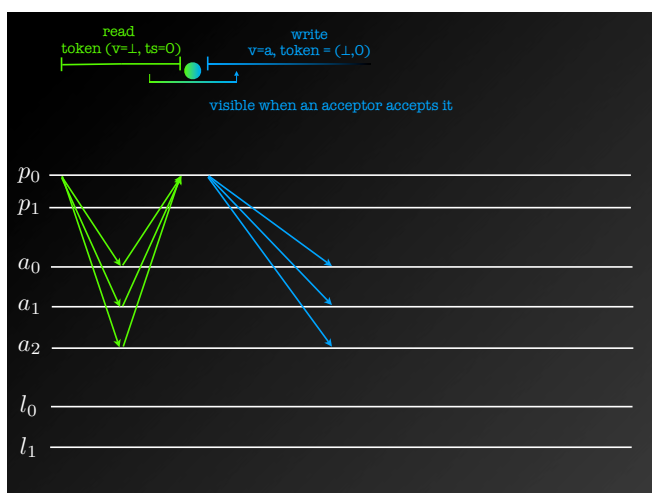
21



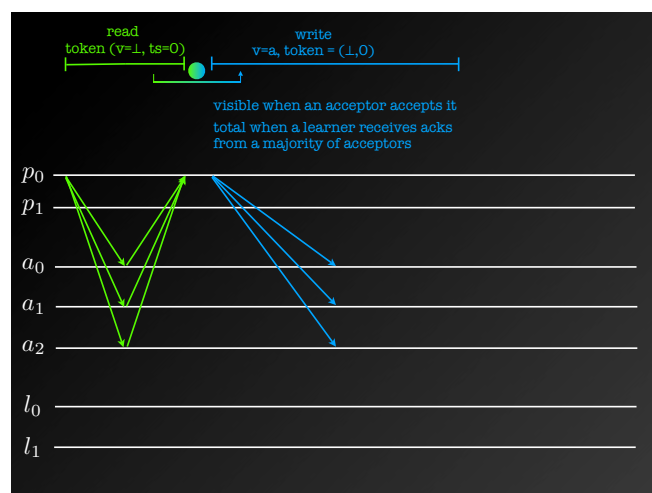
21



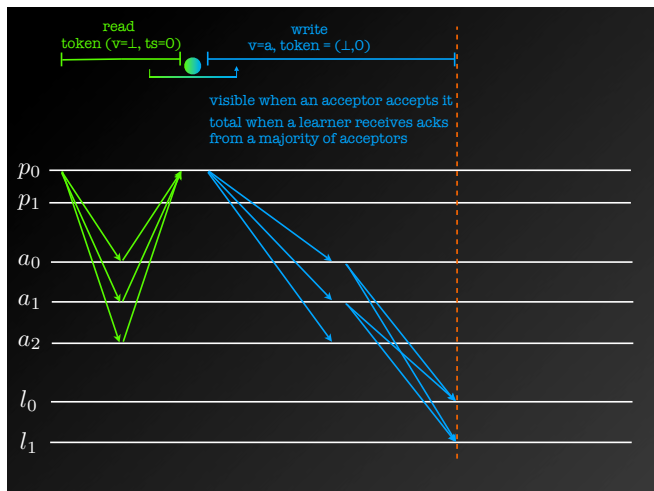
21



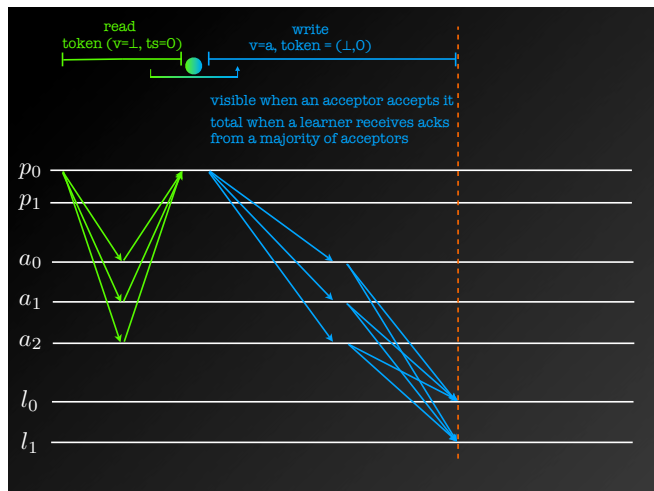
21



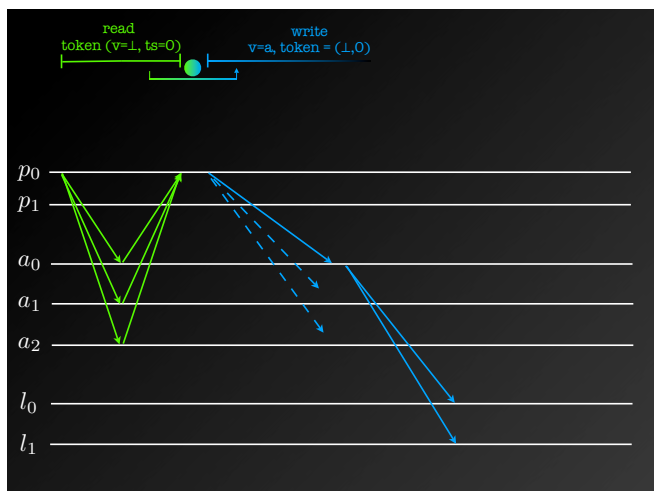
21



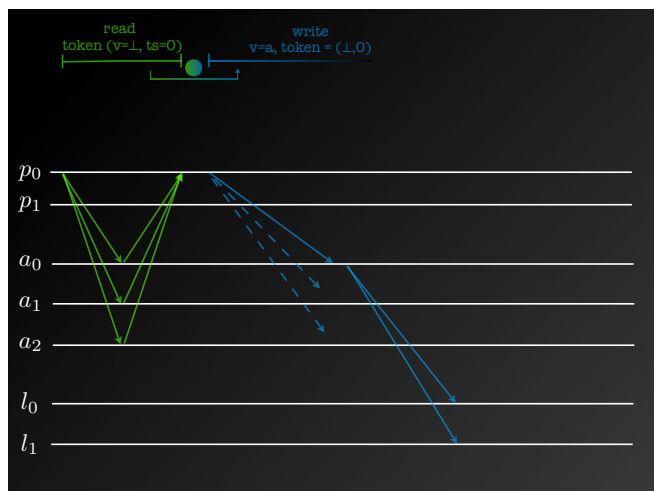
21



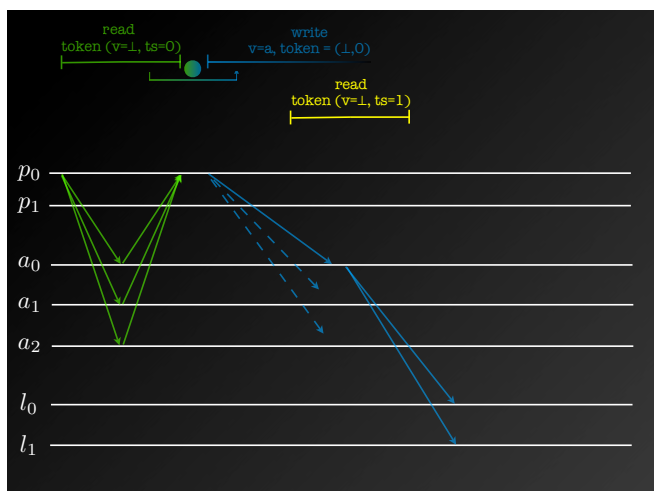
21



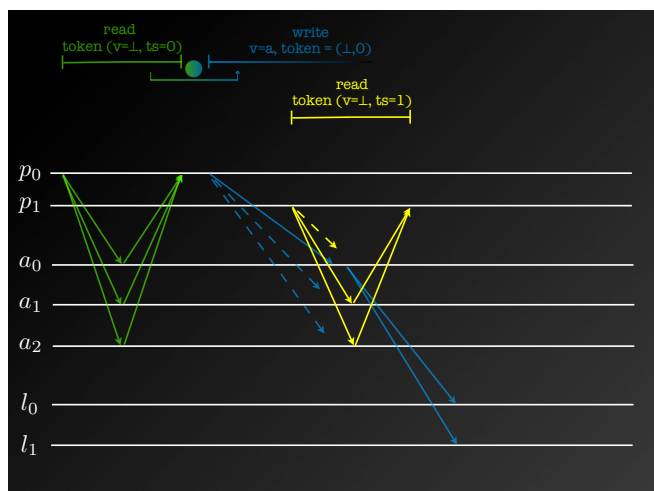
22



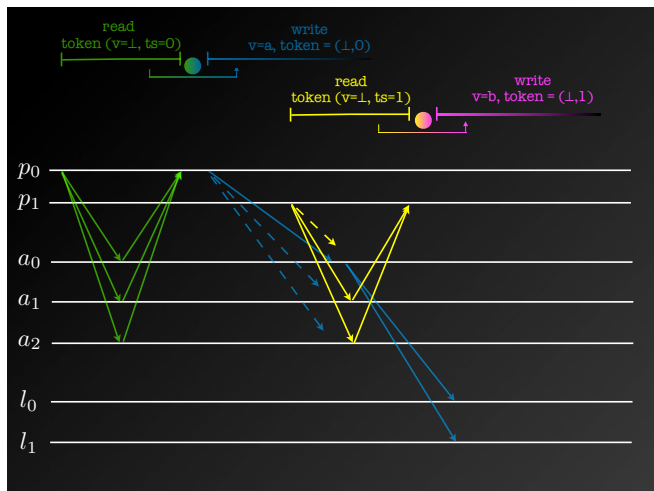
23



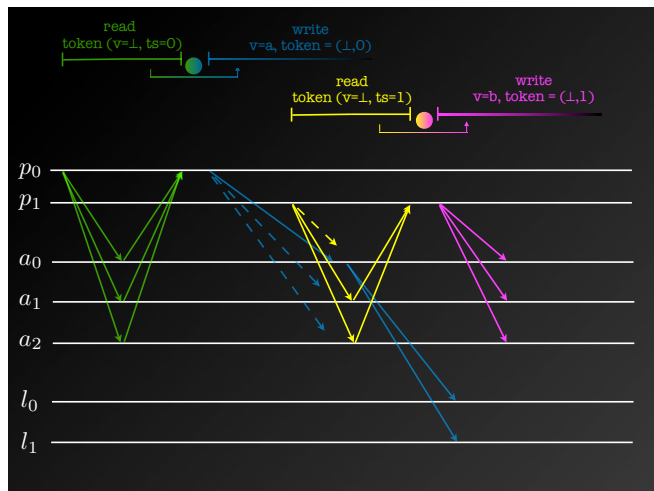
23



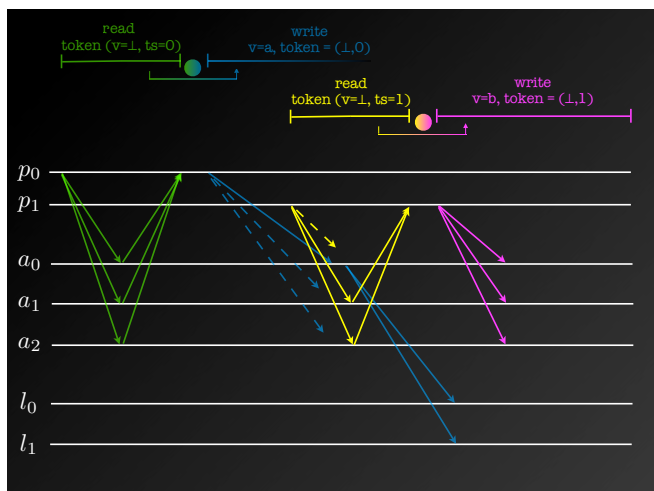
23



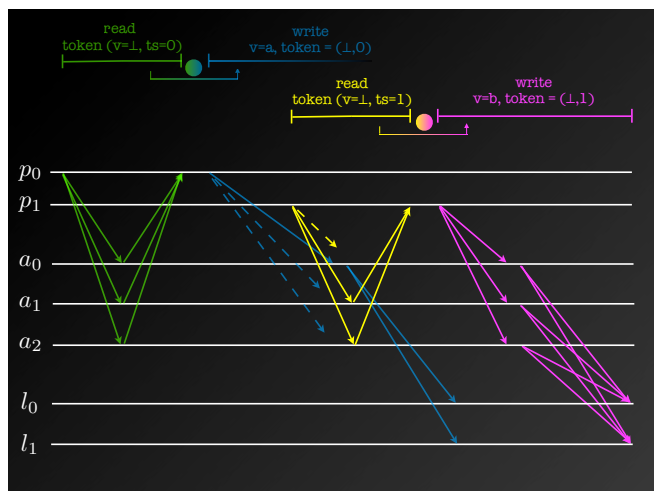
23



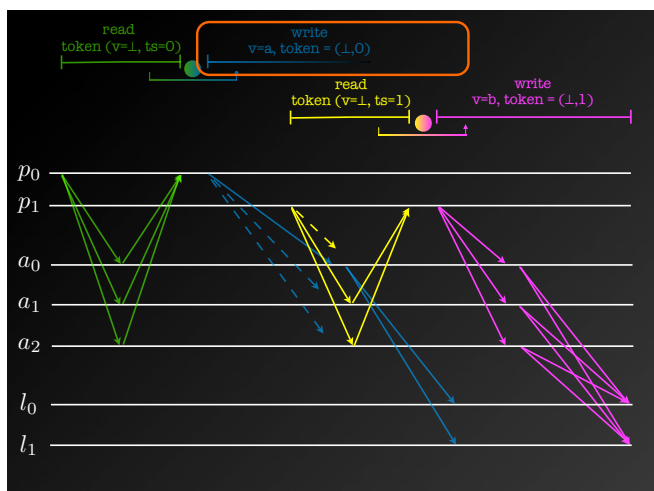
23



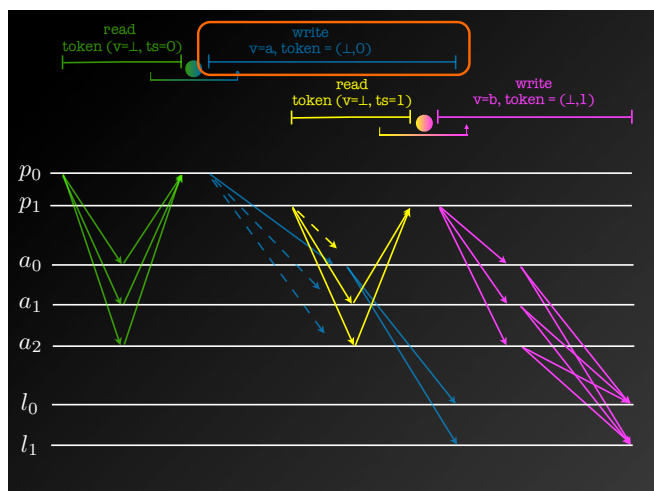
23



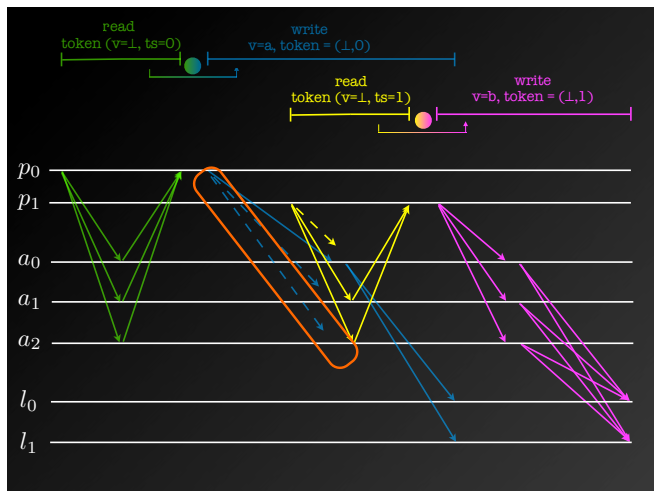
23



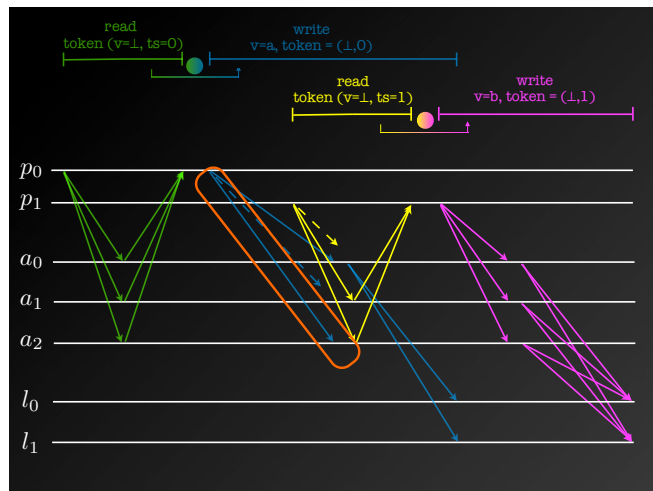
23



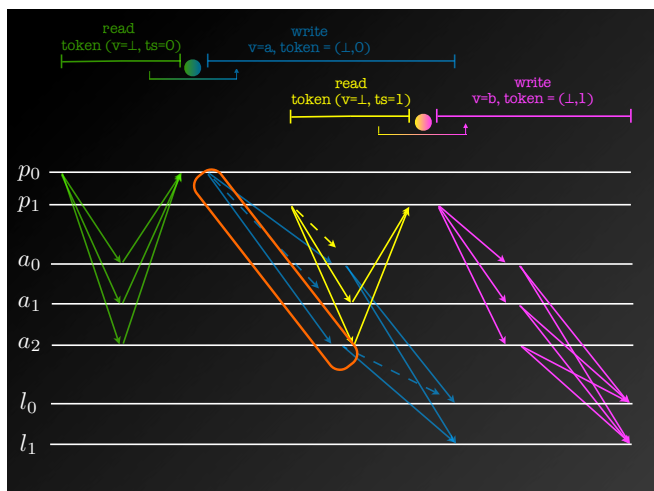
23



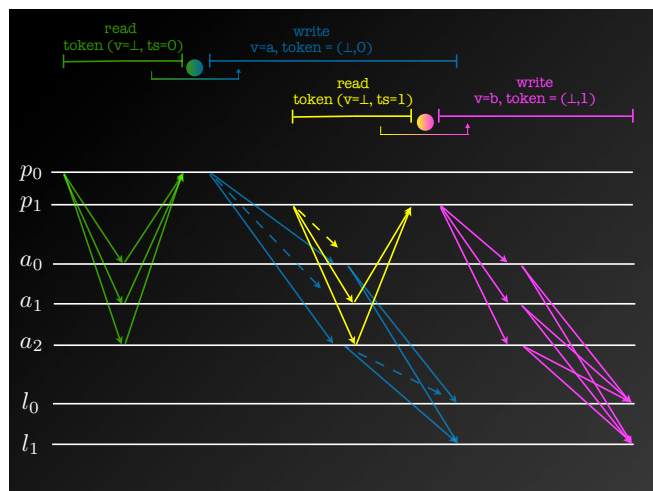
23



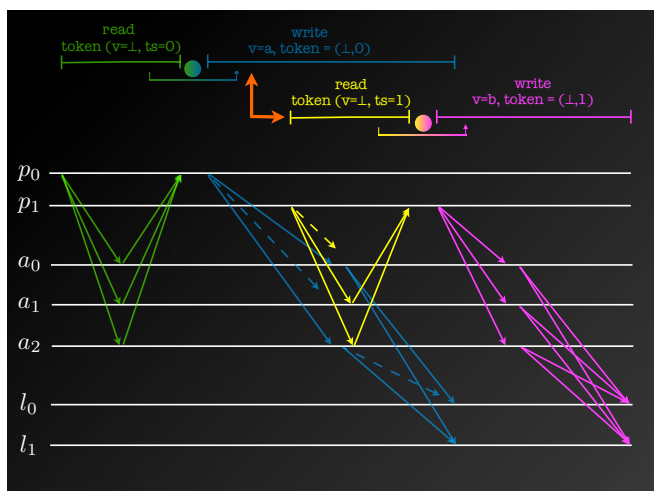
23



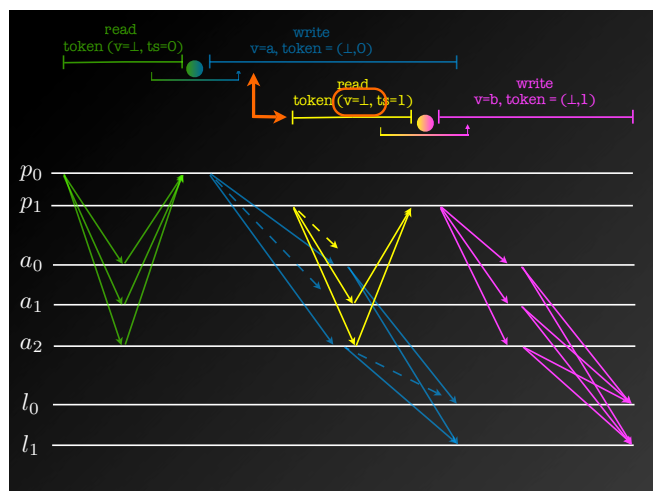
23



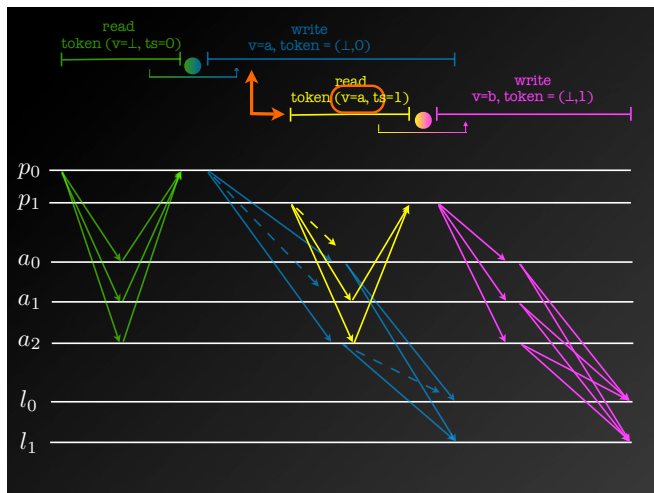
24



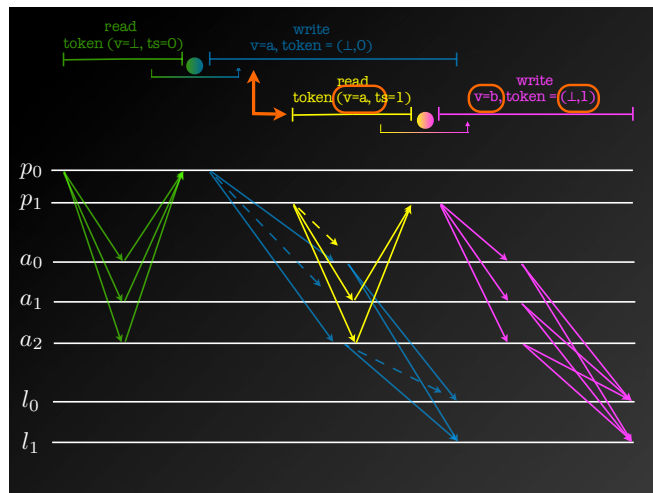
24



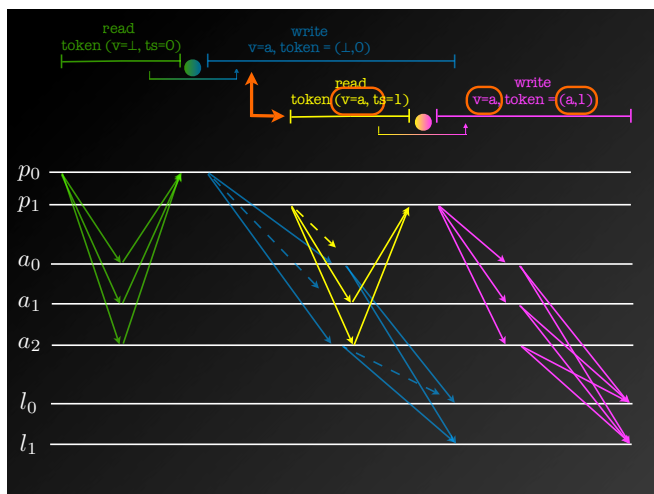
24



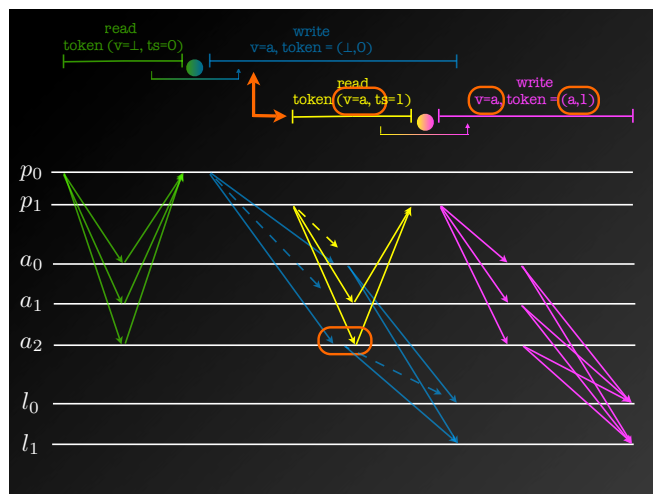
24



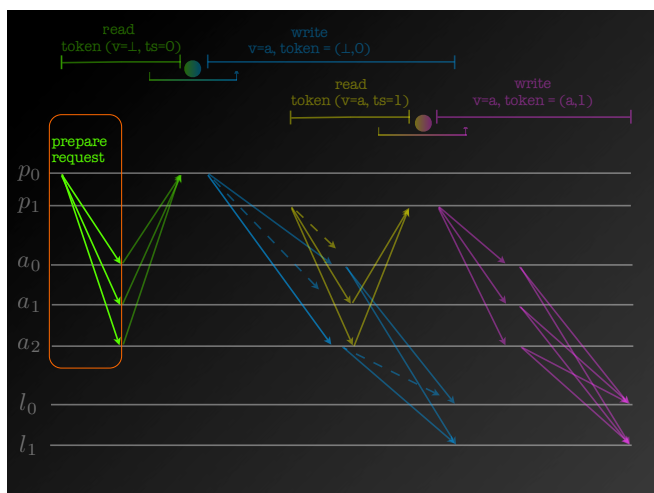
24



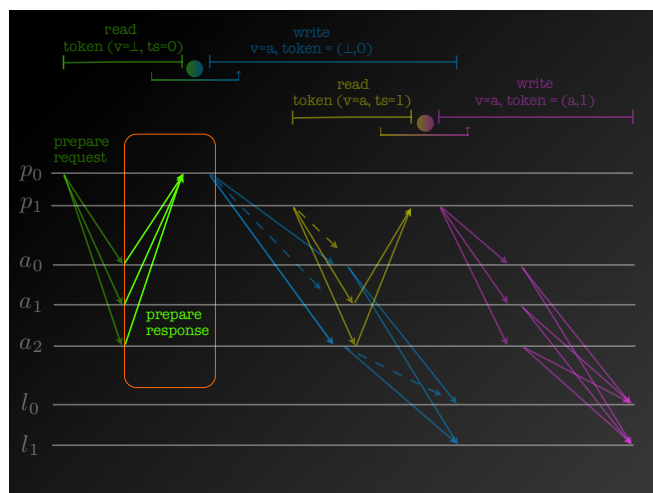
24



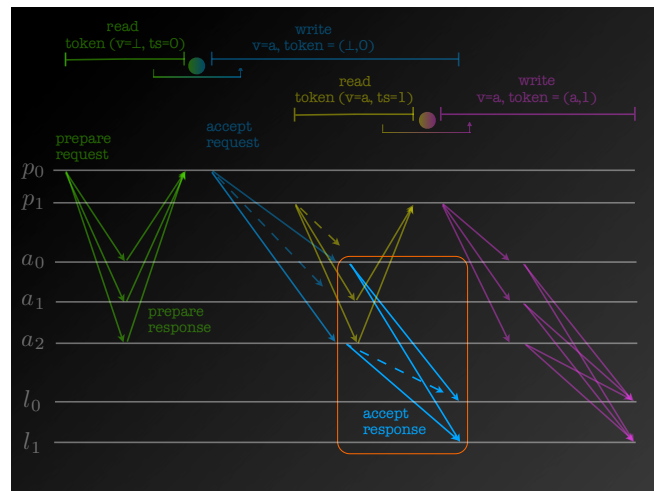
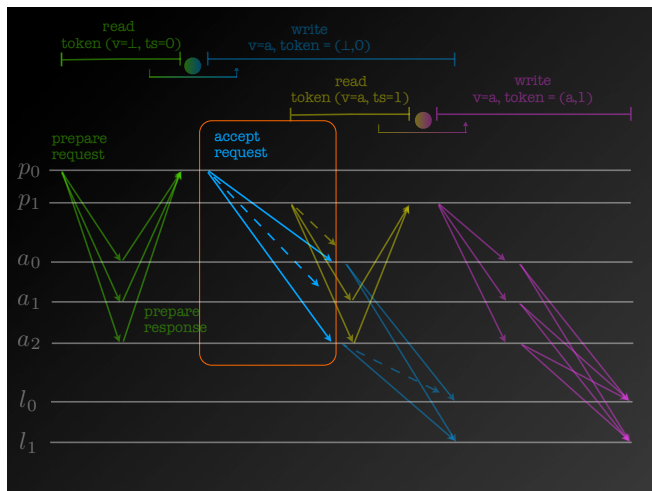
24



25



26

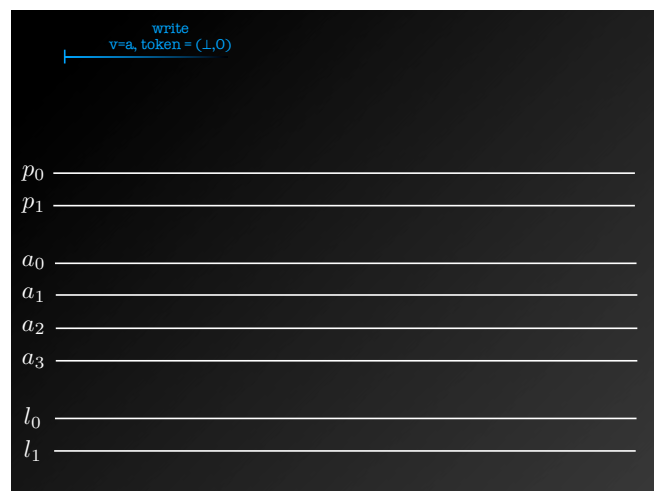


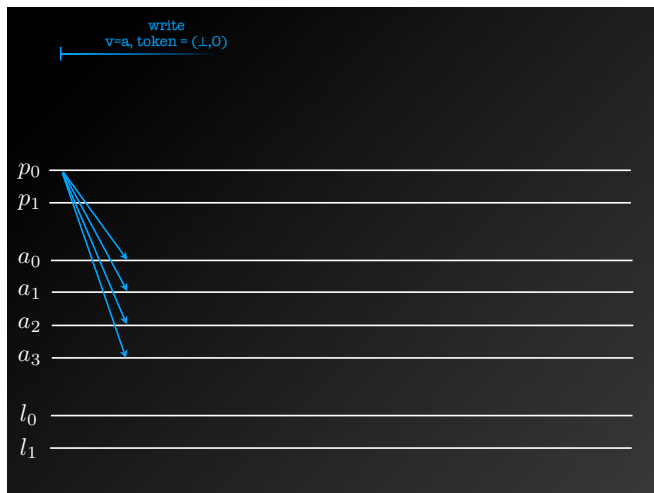
Outline

- Motivation
- Paxos Register specification
- Crash implementation
- Byzantine implementation

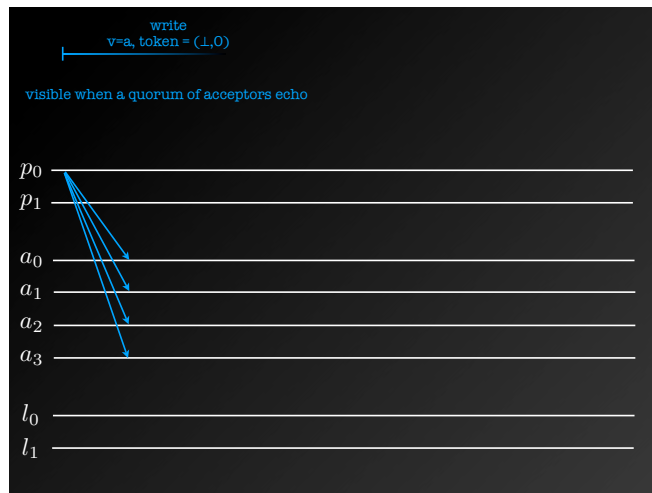
Byzantine Model

- Asynchronous system
- Digital signatures
- Proposers, acceptors, learners
- At least one proposer correct
- More than two thirds acceptors correct

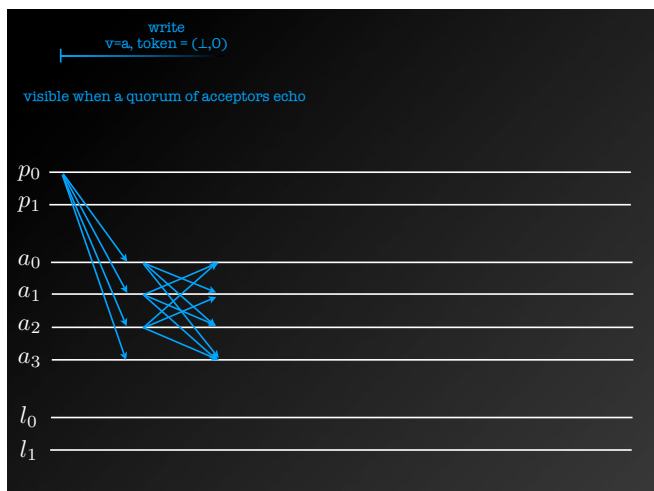




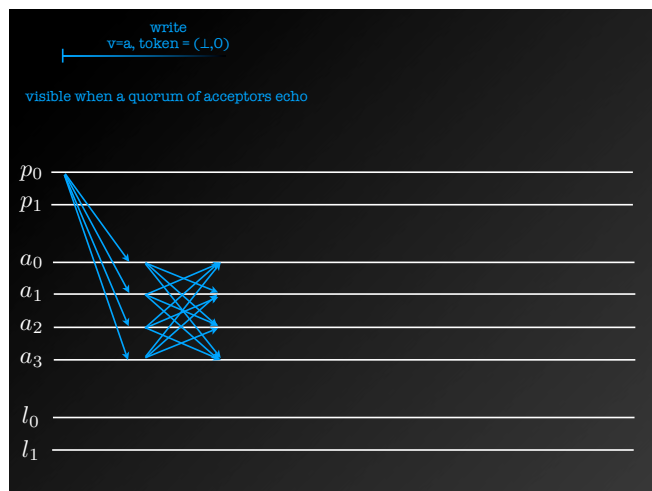
31



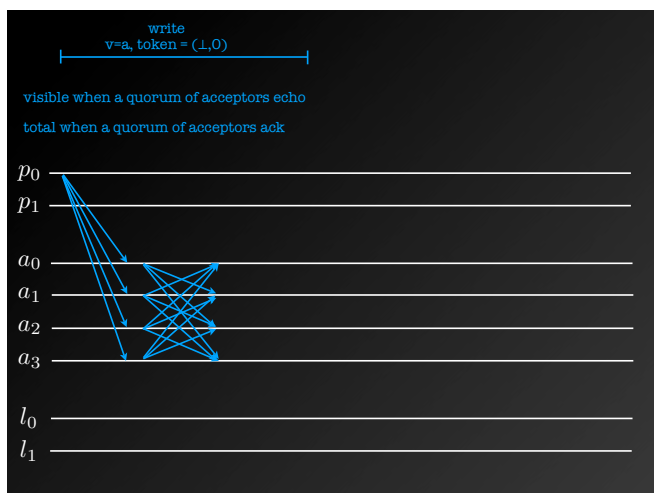
31



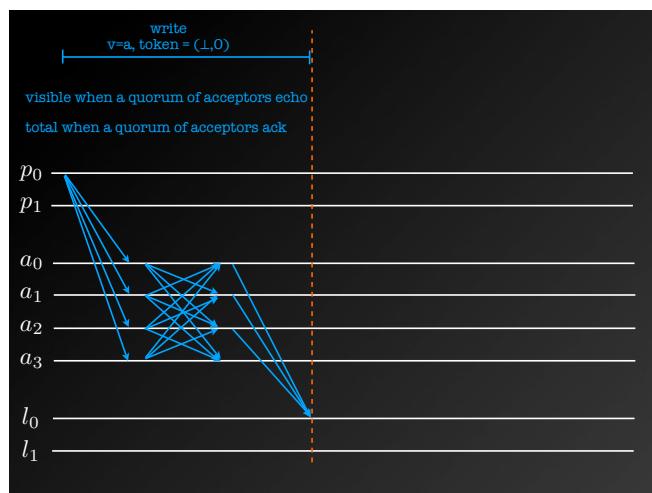
31



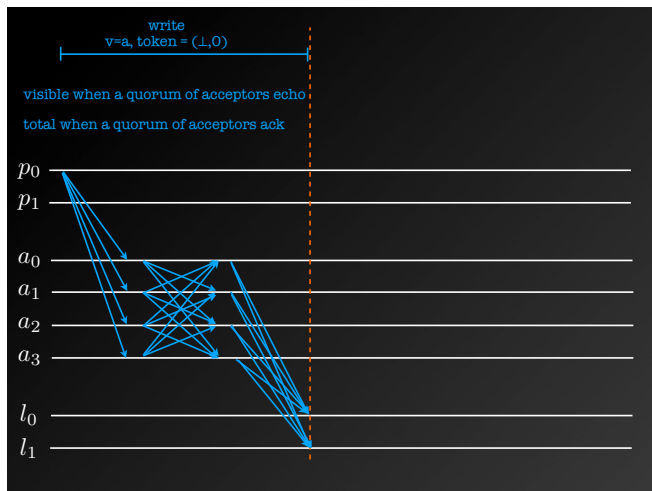
31



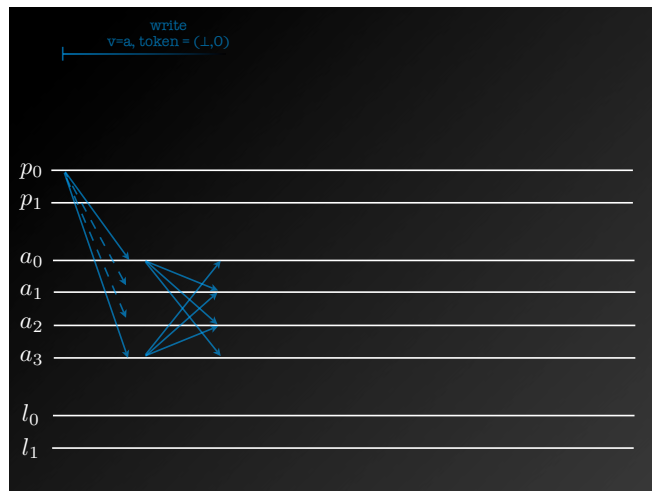
31



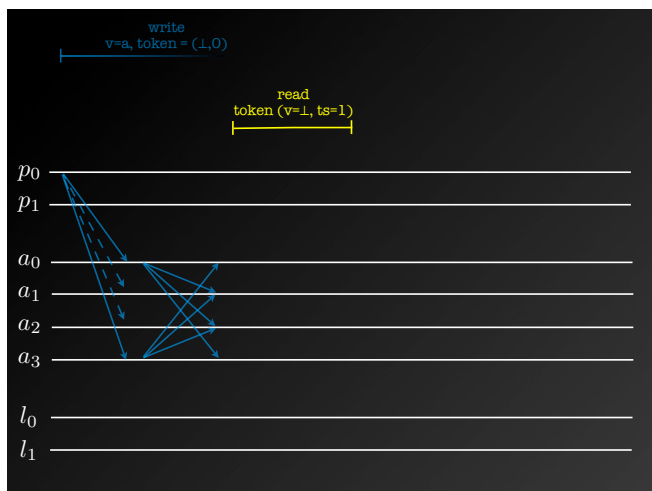
31



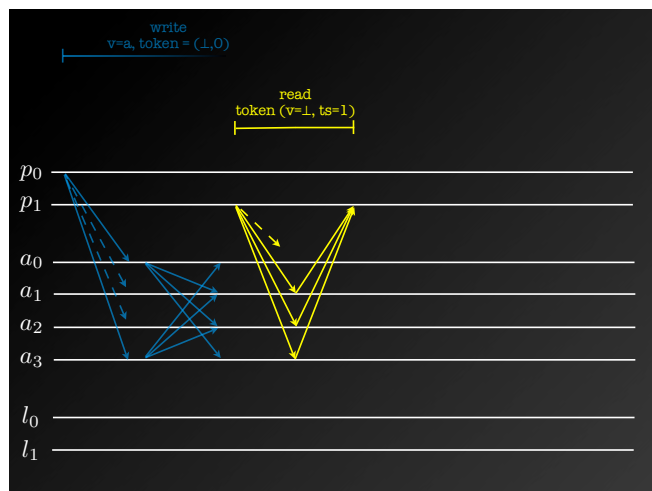
31



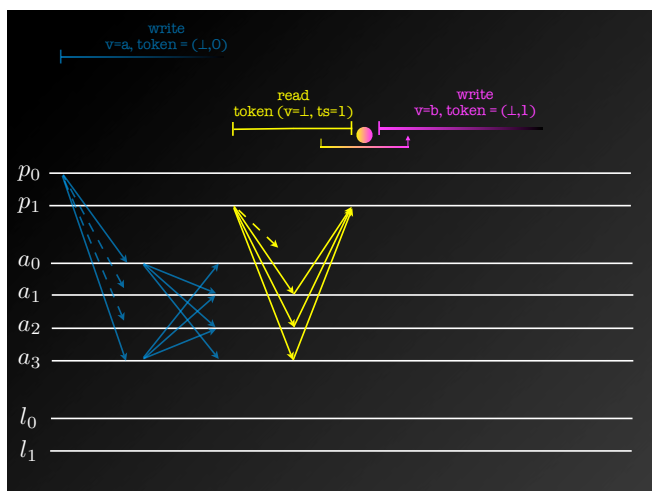
32



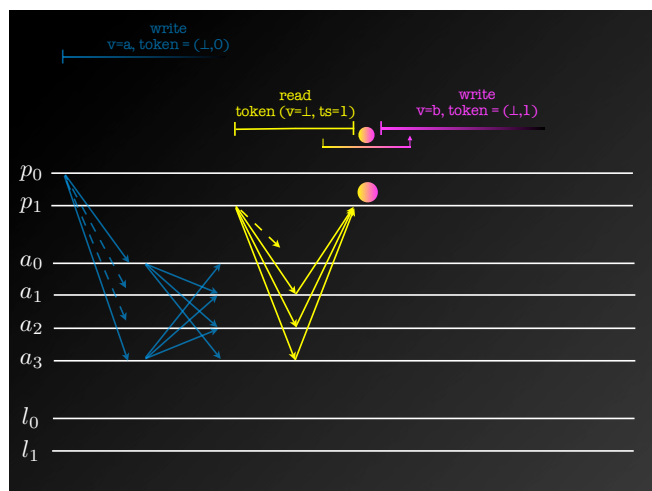
32



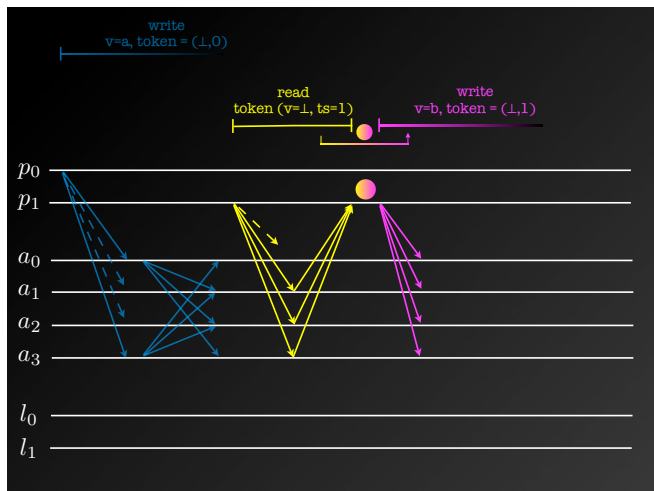
32



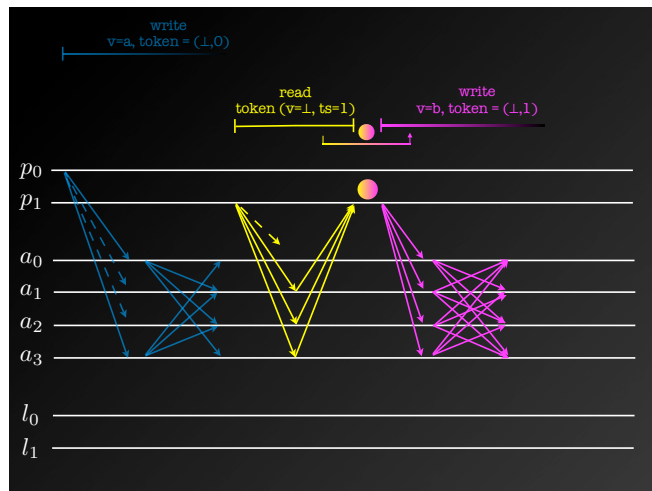
32



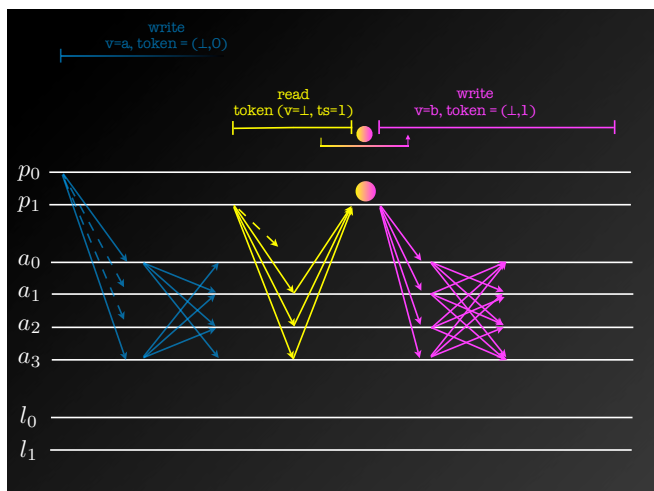
32



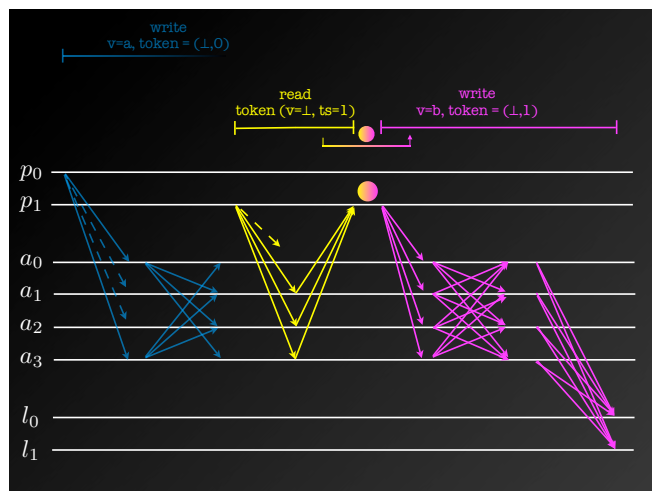
32



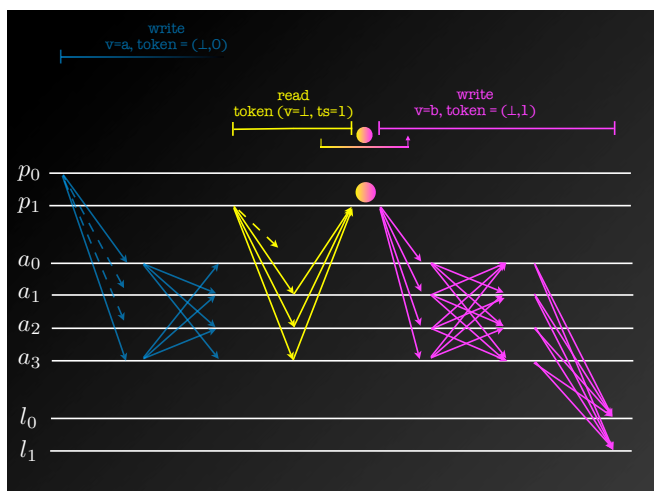
32



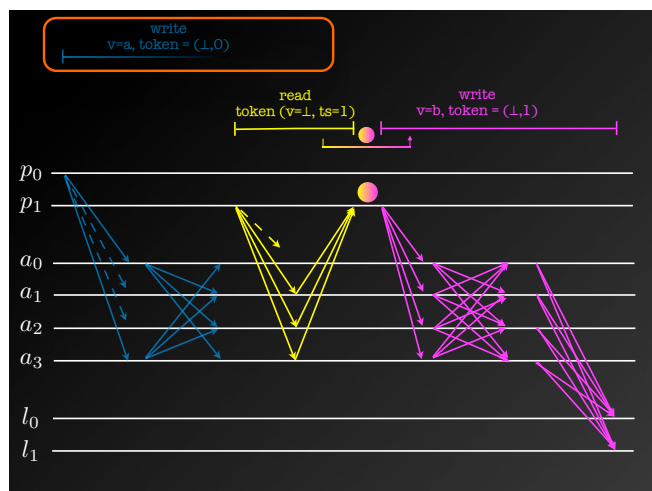
32



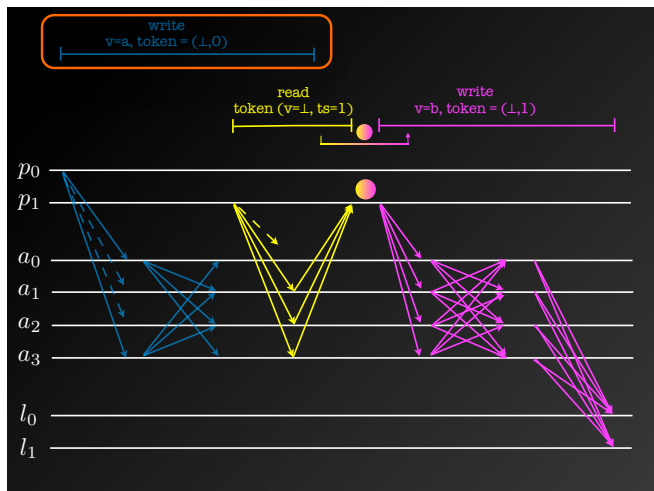
32



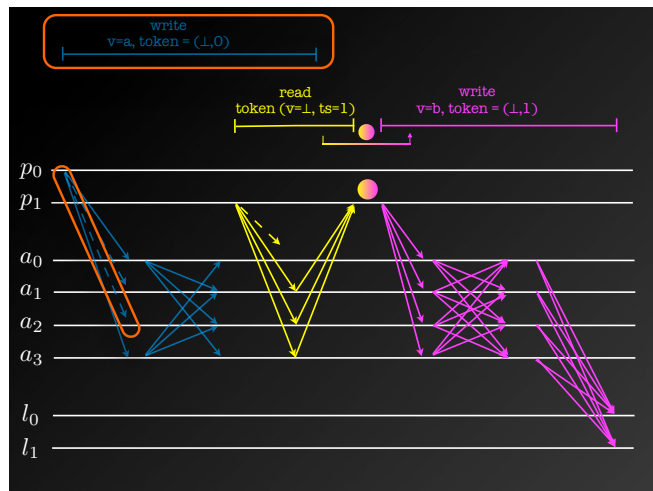
33



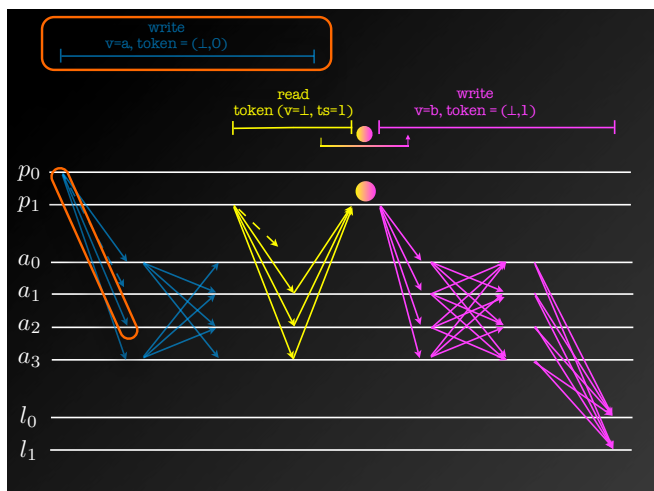
33



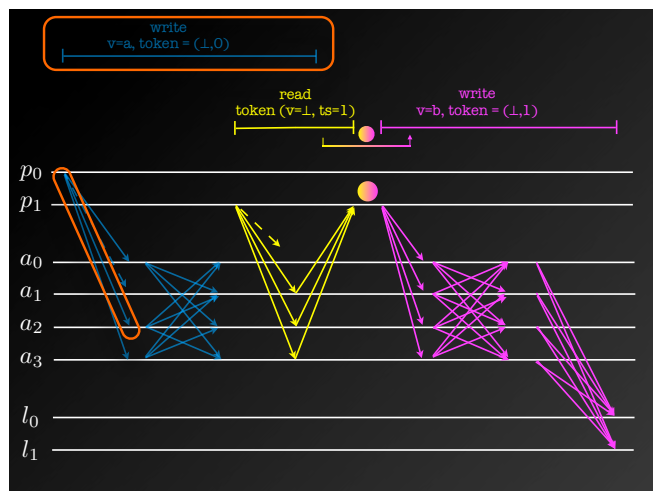
33



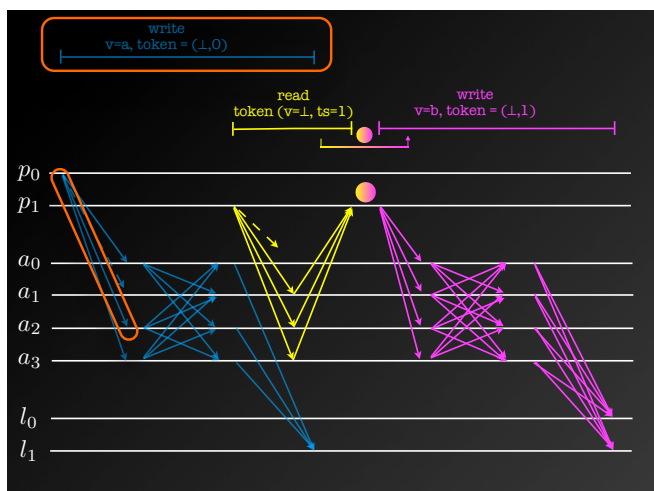
33



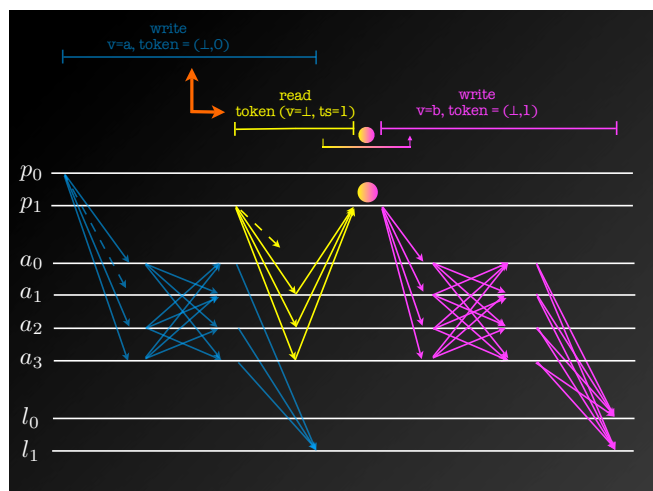
33



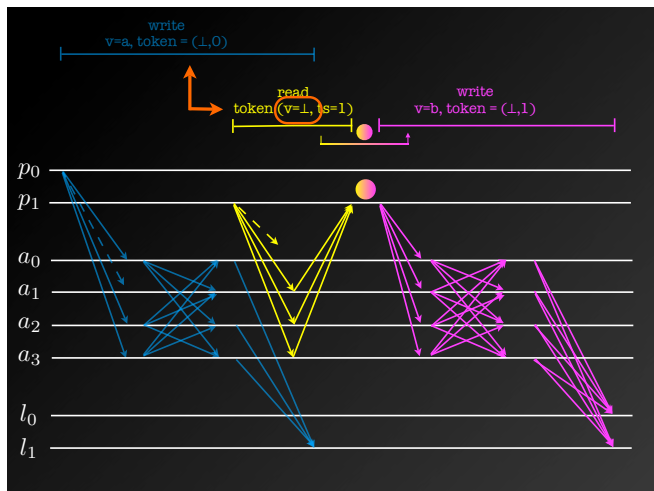
33



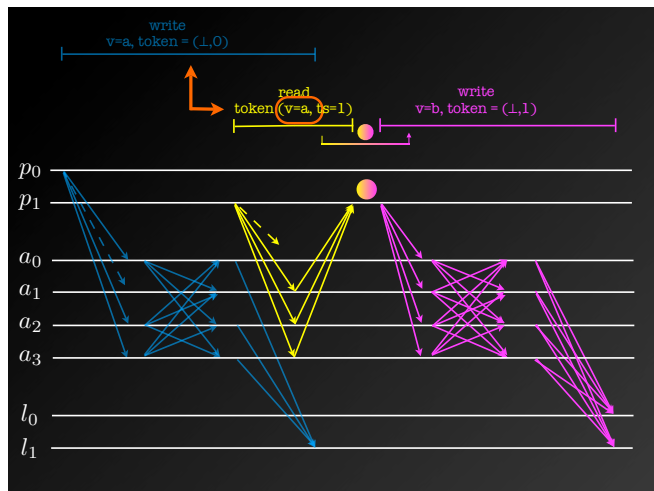
33



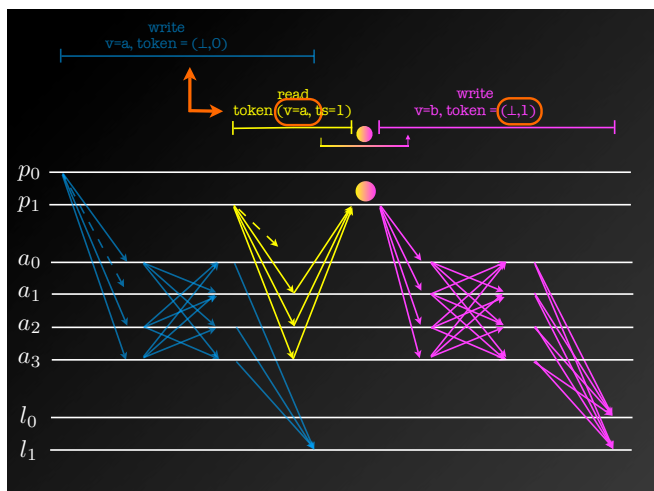
34



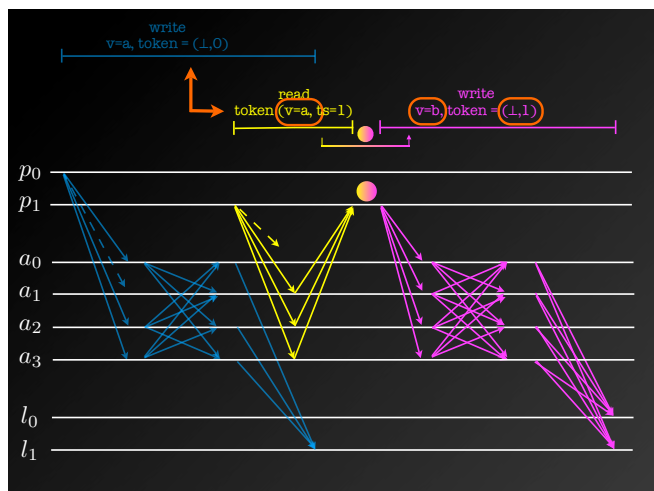
34



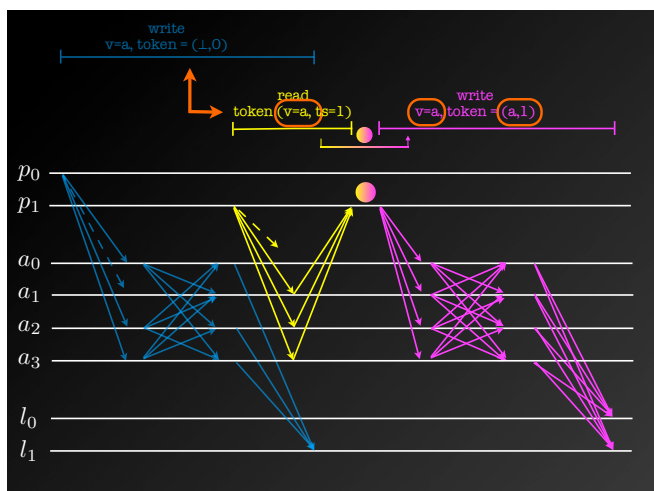
34



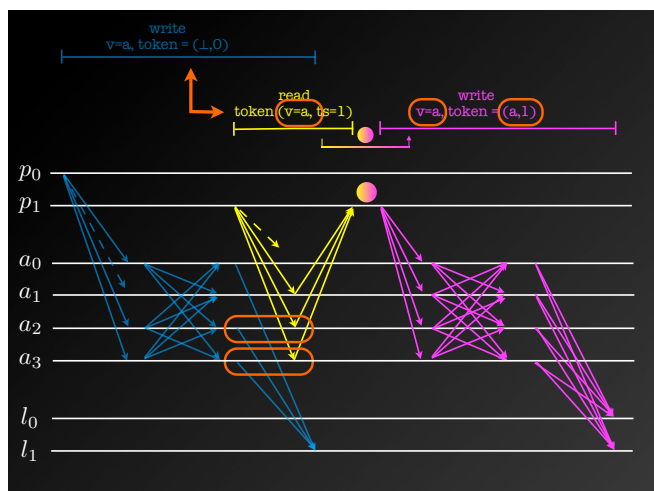
34



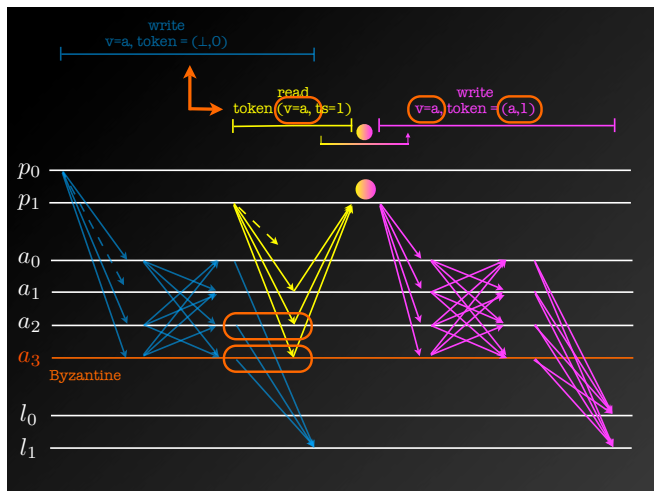
34



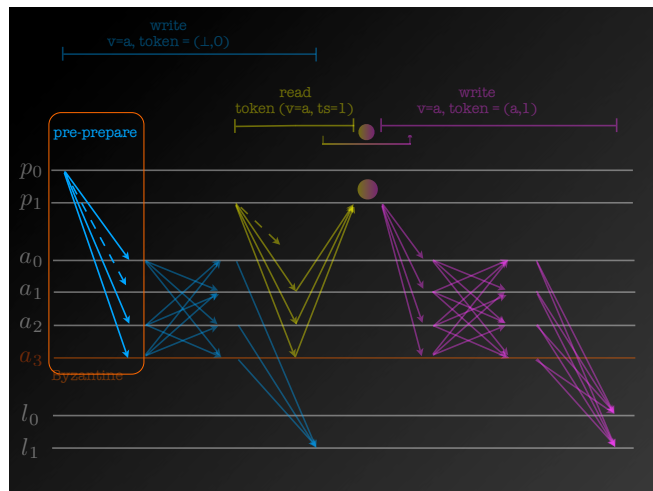
34



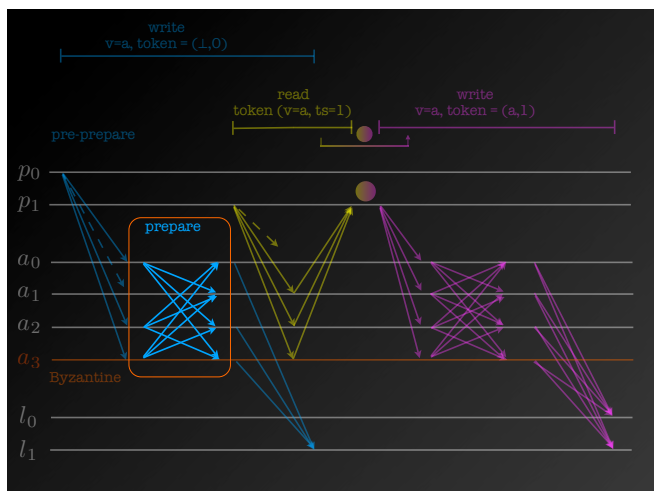
34



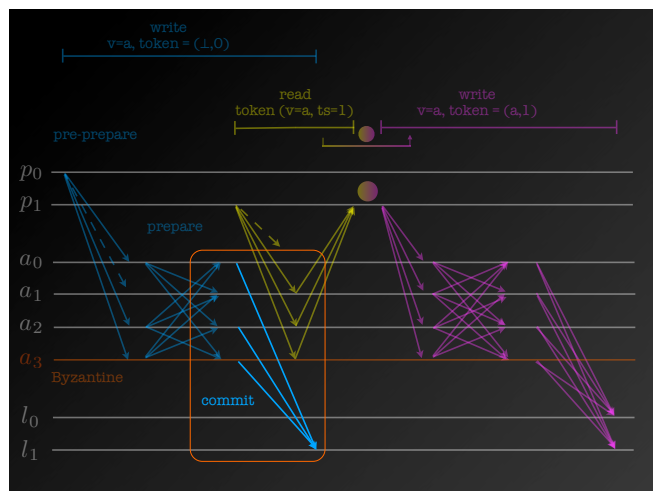
34



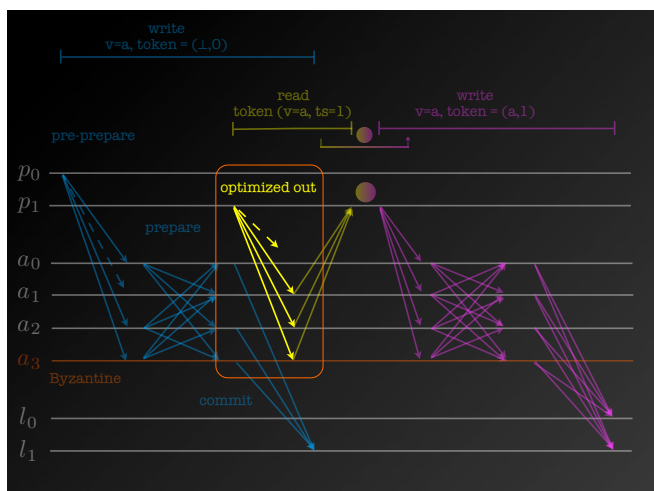
35



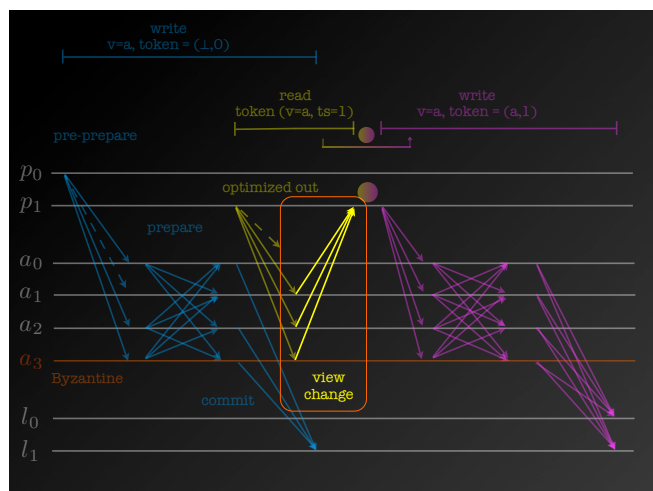
36



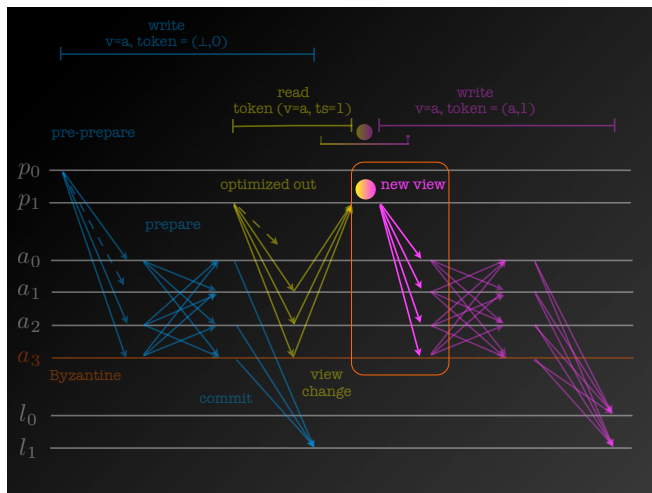
37



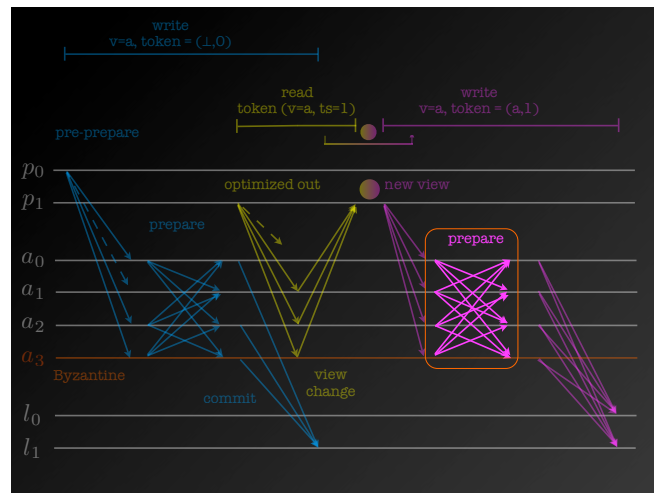
38



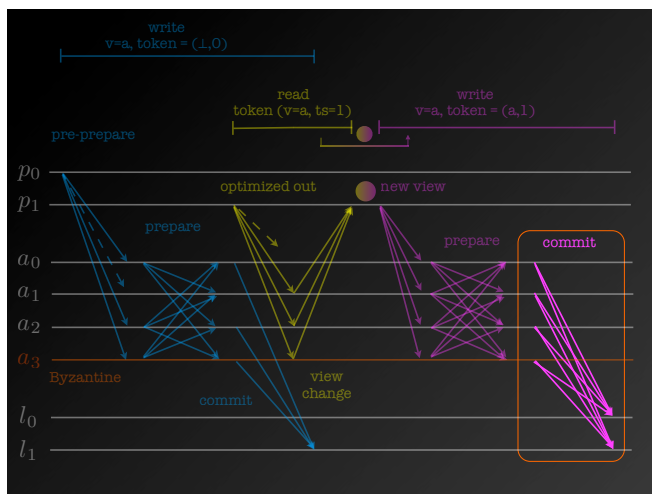
39



40



41



42

Take Away Points

Paxos Register

- Captures similarities in Paxos protocols
- Hides implementation details
- Important abstractions:
tokens, visible writes, total writes

43



44