



Corfu: A Cloud Scale Consistency Platform

Michael Wei
VMware Research Group
Code Mesh 2016

How a Typical Application Becomes Distributed

The screenshot shows a search results page from Stack Overflow. At the top right, there is a grey bar with the word "Questions". On the left, the Stack Overflow logo is visible. Below the logo, the word "Search" is displayed. A search bar contains the query "how to scale an application". To the right of the search bar is a blue button labeled "search". Below the search bar, the text "9,110 results" is shown. To the right of this, there are four filter buttons: "relevance" (which is highlighted with an orange border), "newest", "votes", and "active".

How a Typical Application Becomes Distributed



How a Typical Application Becomes Distributed



How a Typical Application Becomes Distributed



Questions

Search

build a distributed system from scratch

search

31 results

relevance

newest

votes

active

Q: What would a clean slate design of a distributed system look like?

Why would you ever want to do such a thing?

[distributed-computing](#)

[distributed-system](#)

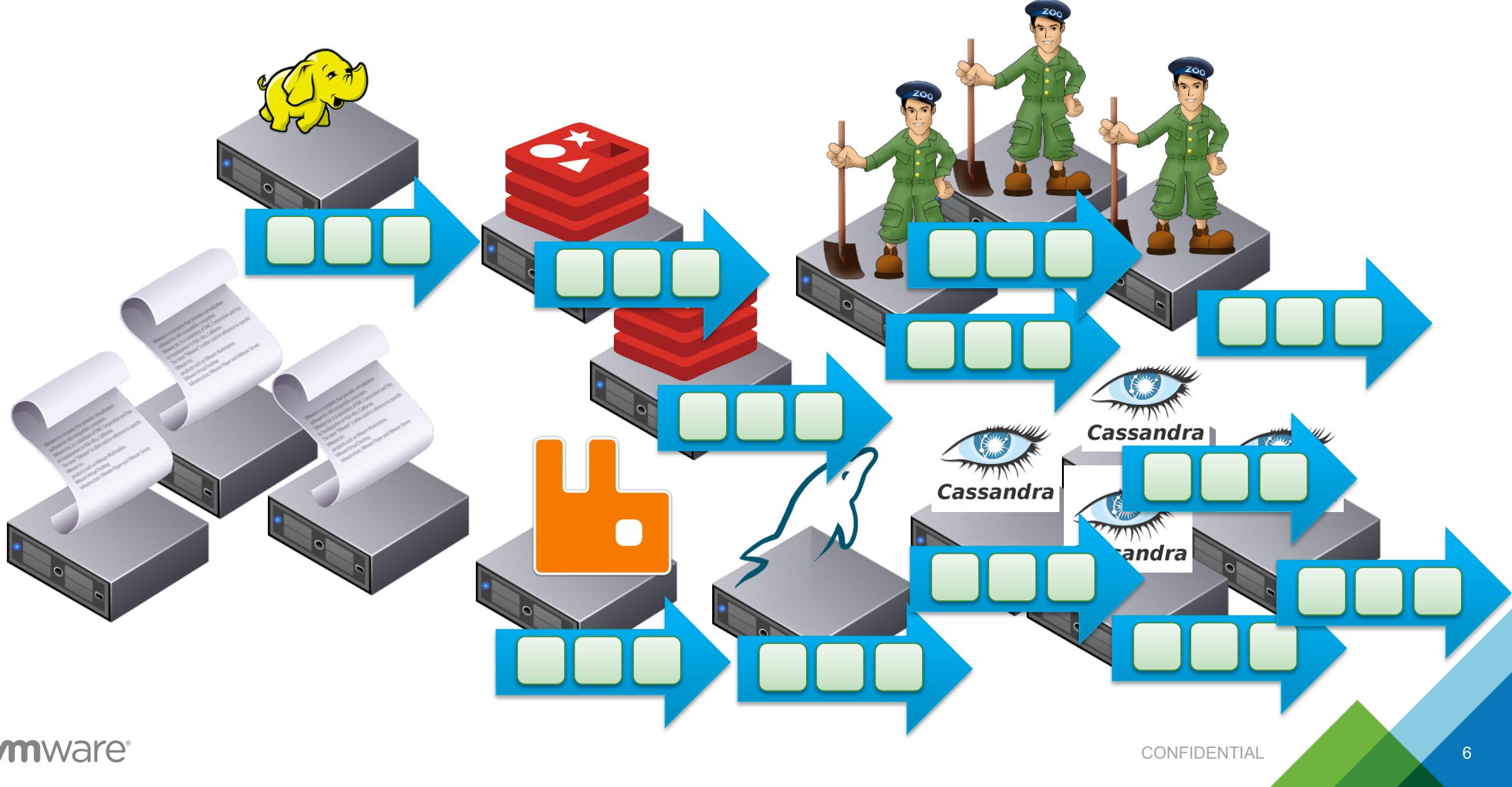
1

vote

0

answers

How a Typical Application Becomes Distributed



What is Corfu?

Corfu



What is Corfu?

Corfu



Paxos

Paxos – A Family of Consensus Protocols

A = 0



A = 0;

Paxos – A Family of Consensus Protocols

A = 1



**A = 0;
A = 1;**

Paxos – A Family of Consensus Protocols

A = 1



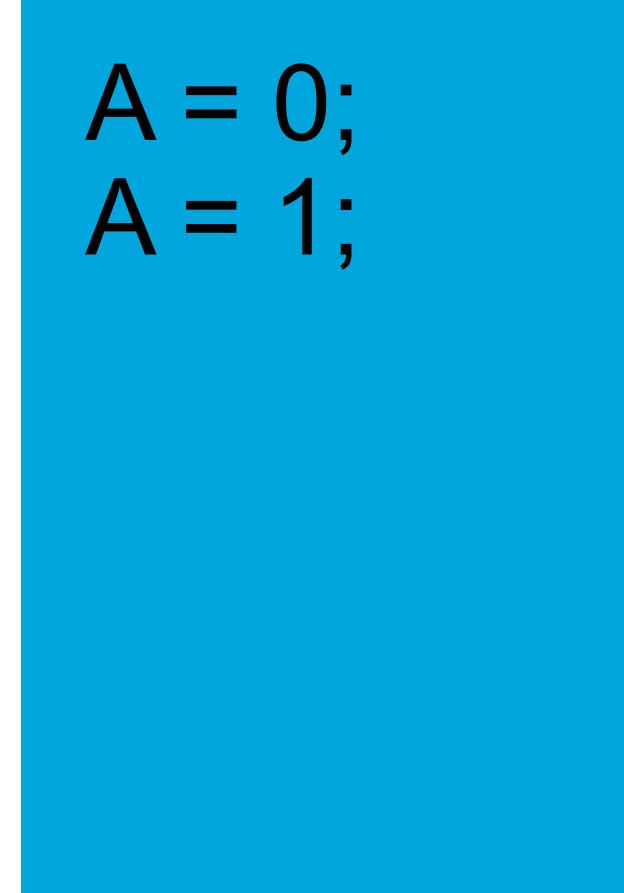
A = 1



A = 1



**A = 0;
A = 1;**



Paxos – A Family of Consensus Protocols

A = 2



A = 1

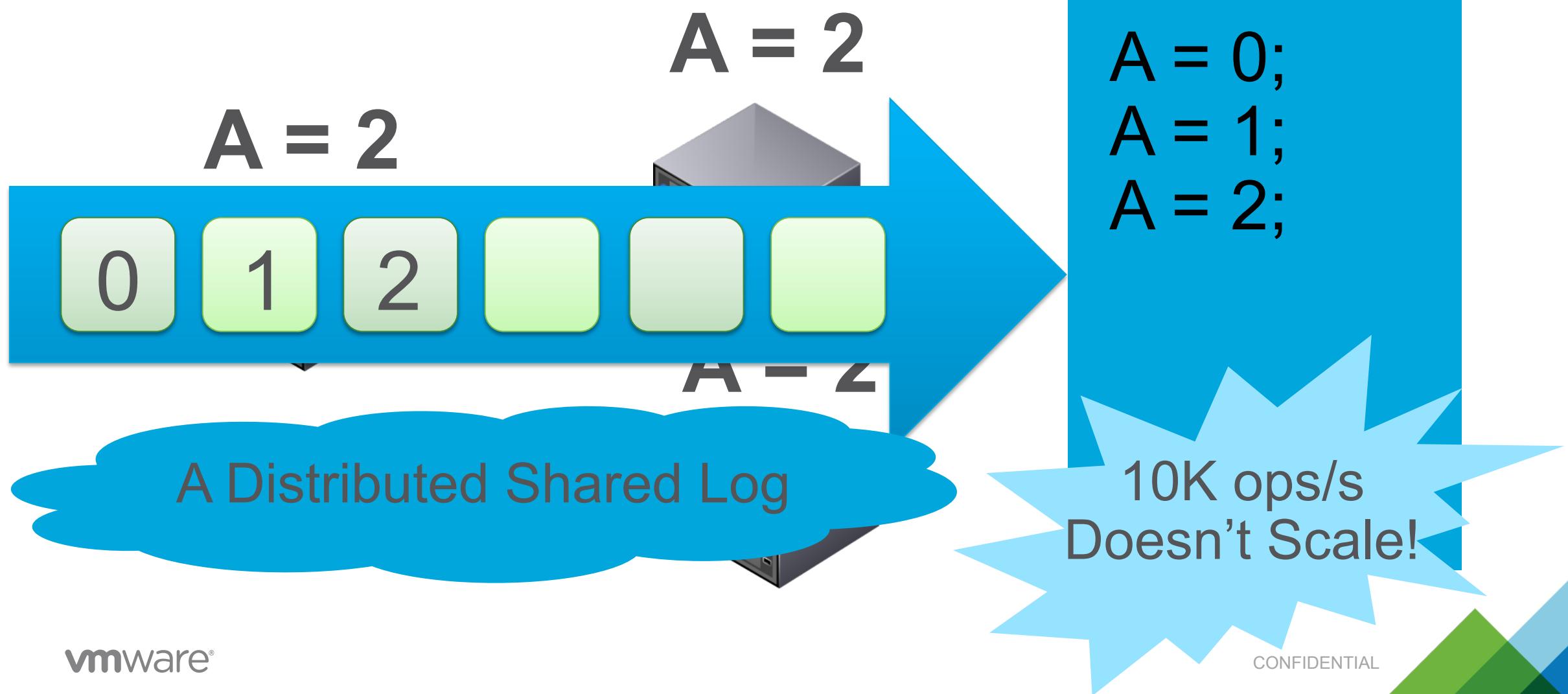


A = 1

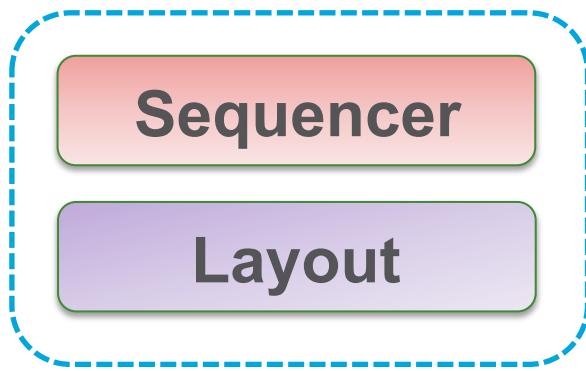


**A = 0;
A = 1;
A = 2;**

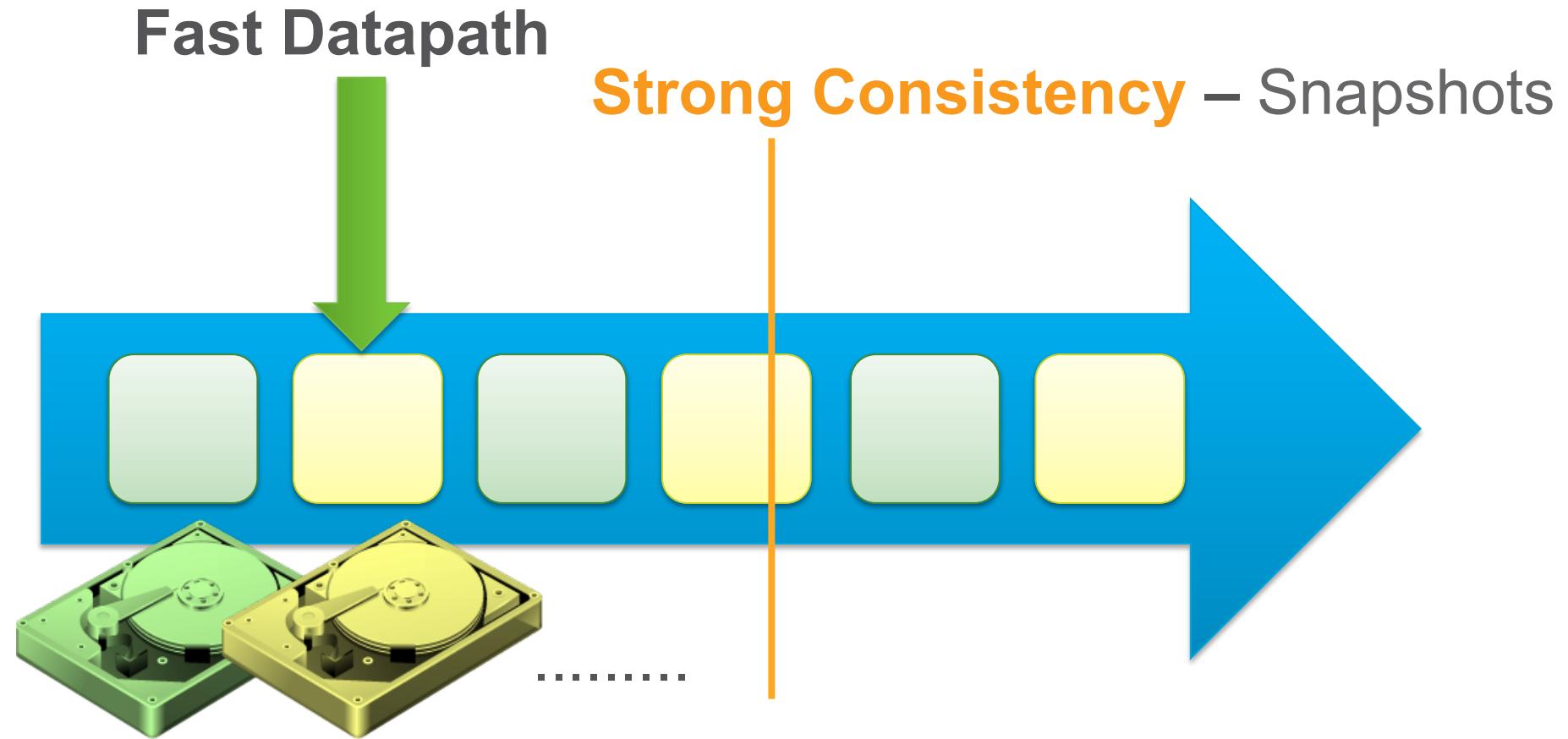
Paxos – A Family of Consensus Protocols



Corfu: A Distributed Shared Log



“Control Plane”
Traditional I/O
bottleneck



Write Scaling – Fast Updates

Read Scaling – Fast Random Reads

Corfu: A Distributed Shared Log

Sequencer

Layout

Sequencer is fast

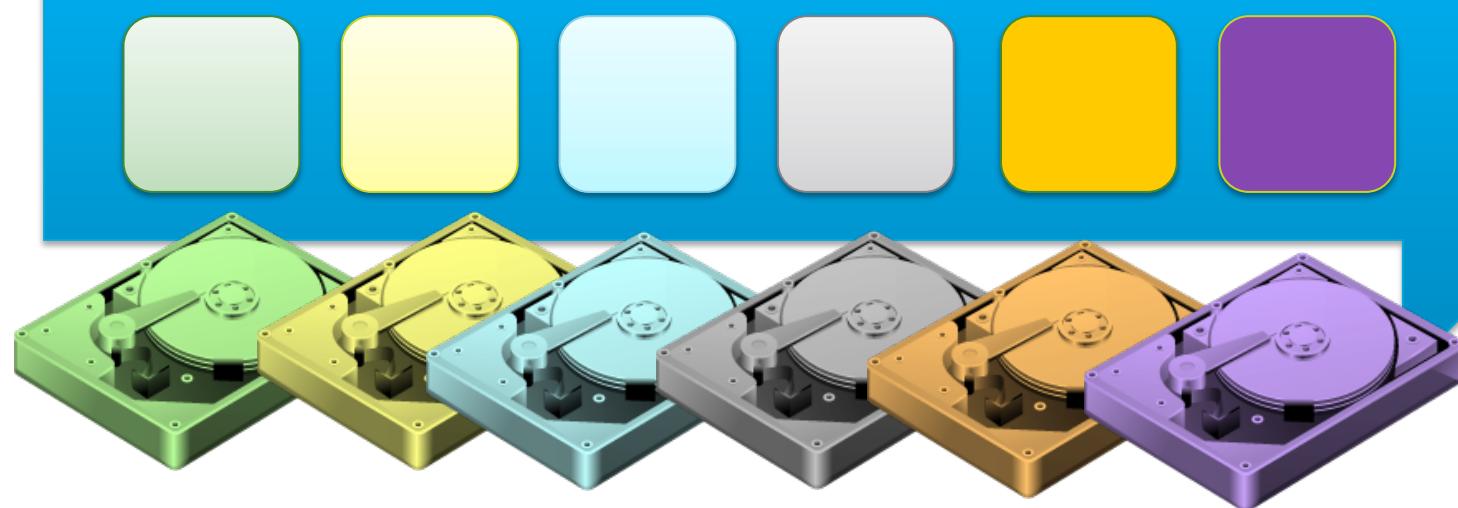
600K op/s – not part of I/O path

Corfu: A Distributed Shared Log

Sequencer

Layout

Scalable
Add more stripes...

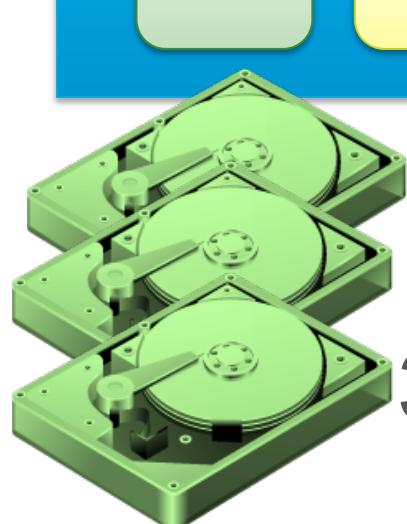


Corfu: A Distributed Shared Log

Sequencer

Layout

Fault Tolerant
Via chain replication...



F+1 Replicas

3 Replicas = 2 failures “per chain”

Corfu: A Distributed Shared Log

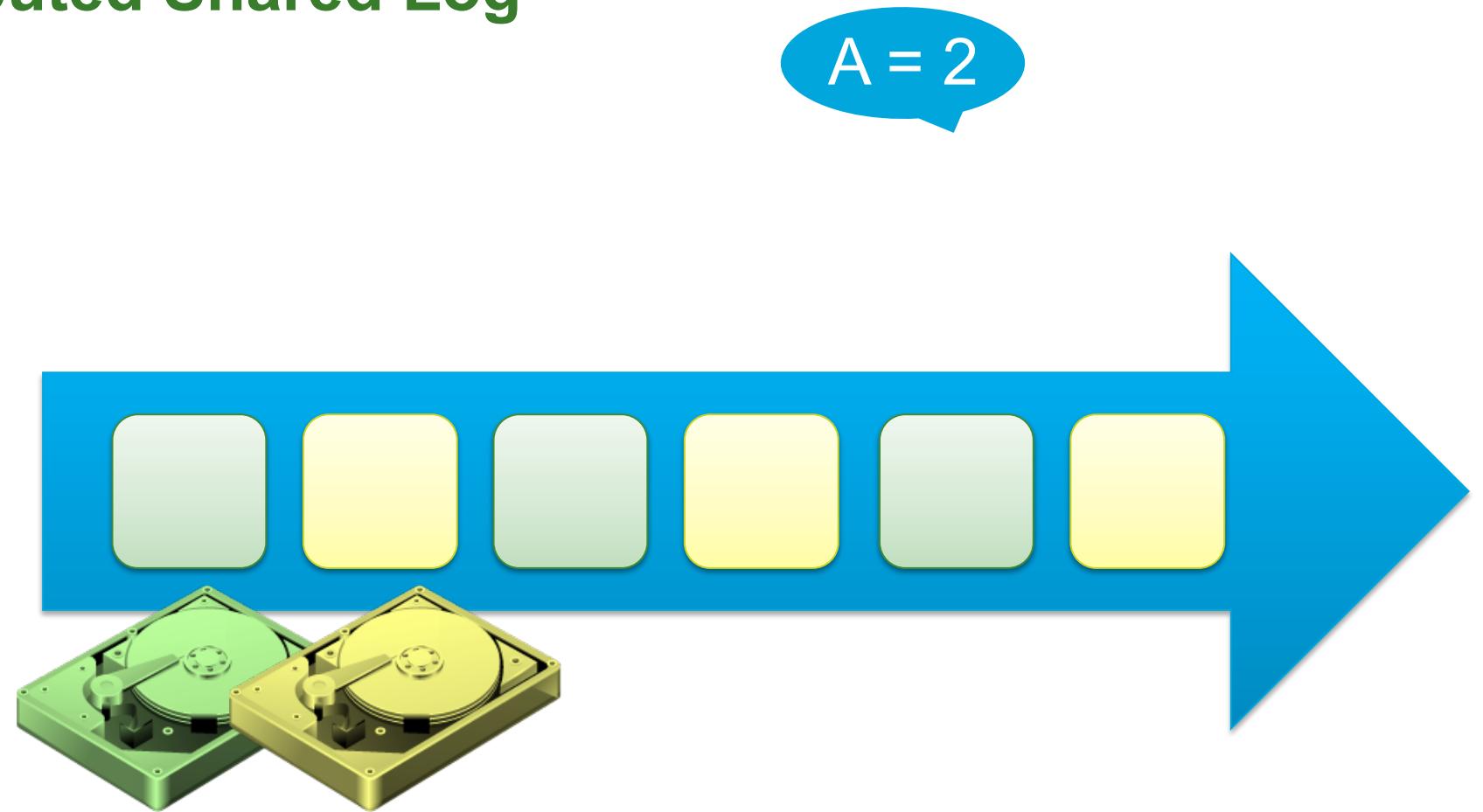
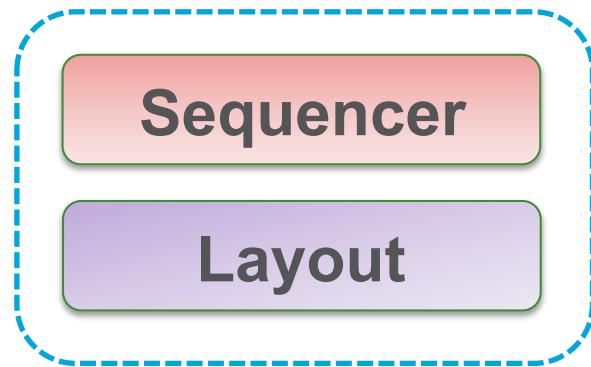
Sequencer

Layout

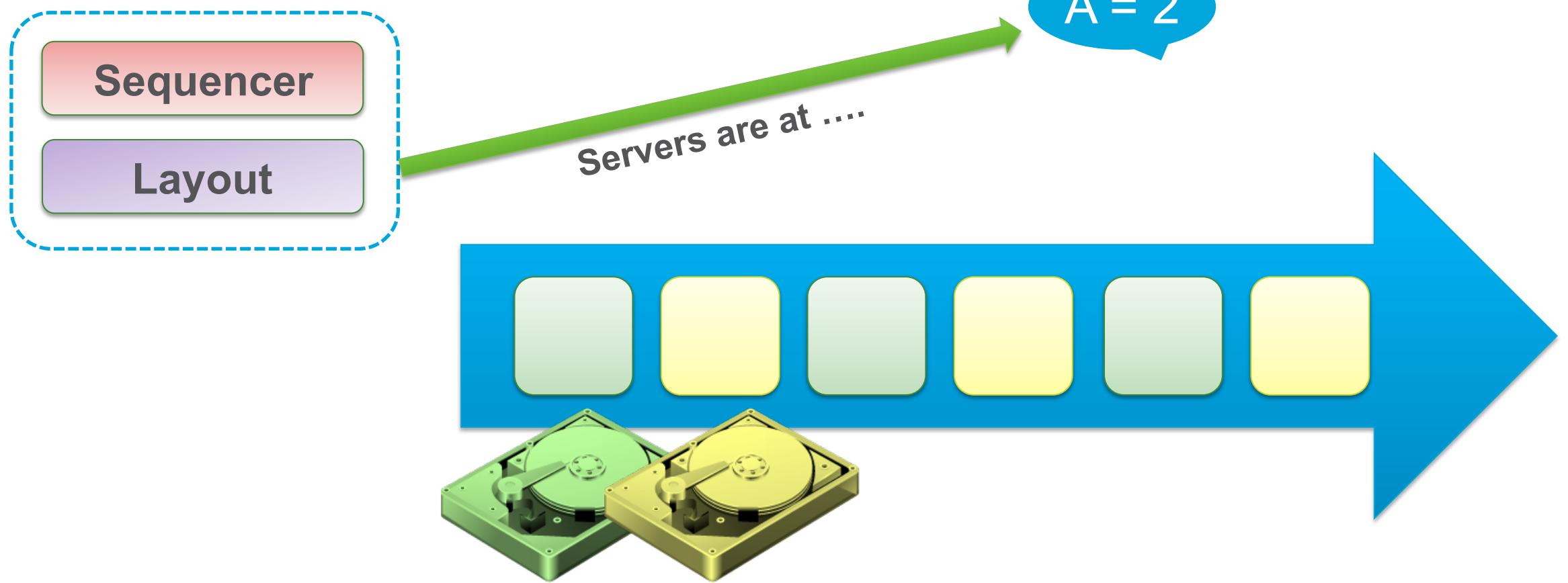
Strong Consistency
The log ... is the truth!



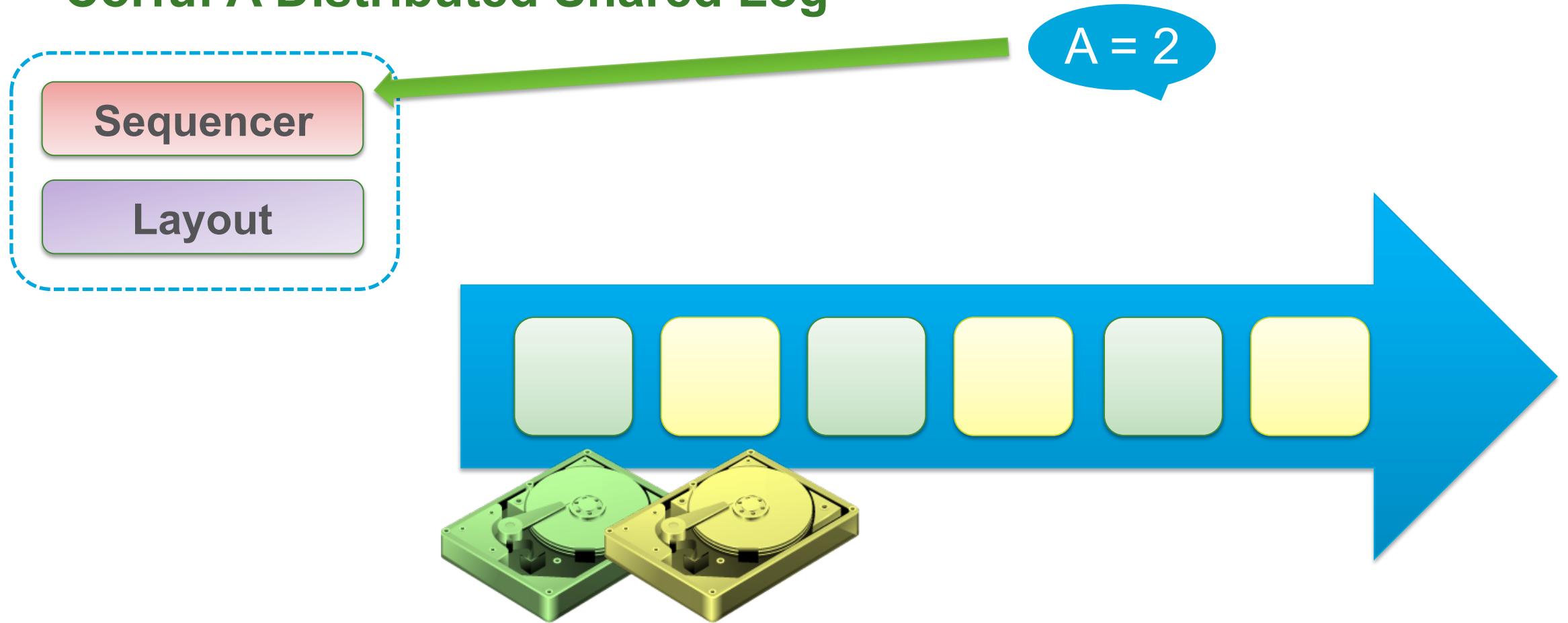
Corfu: A Distributed Shared Log



Corfu: A Distributed Shared Log



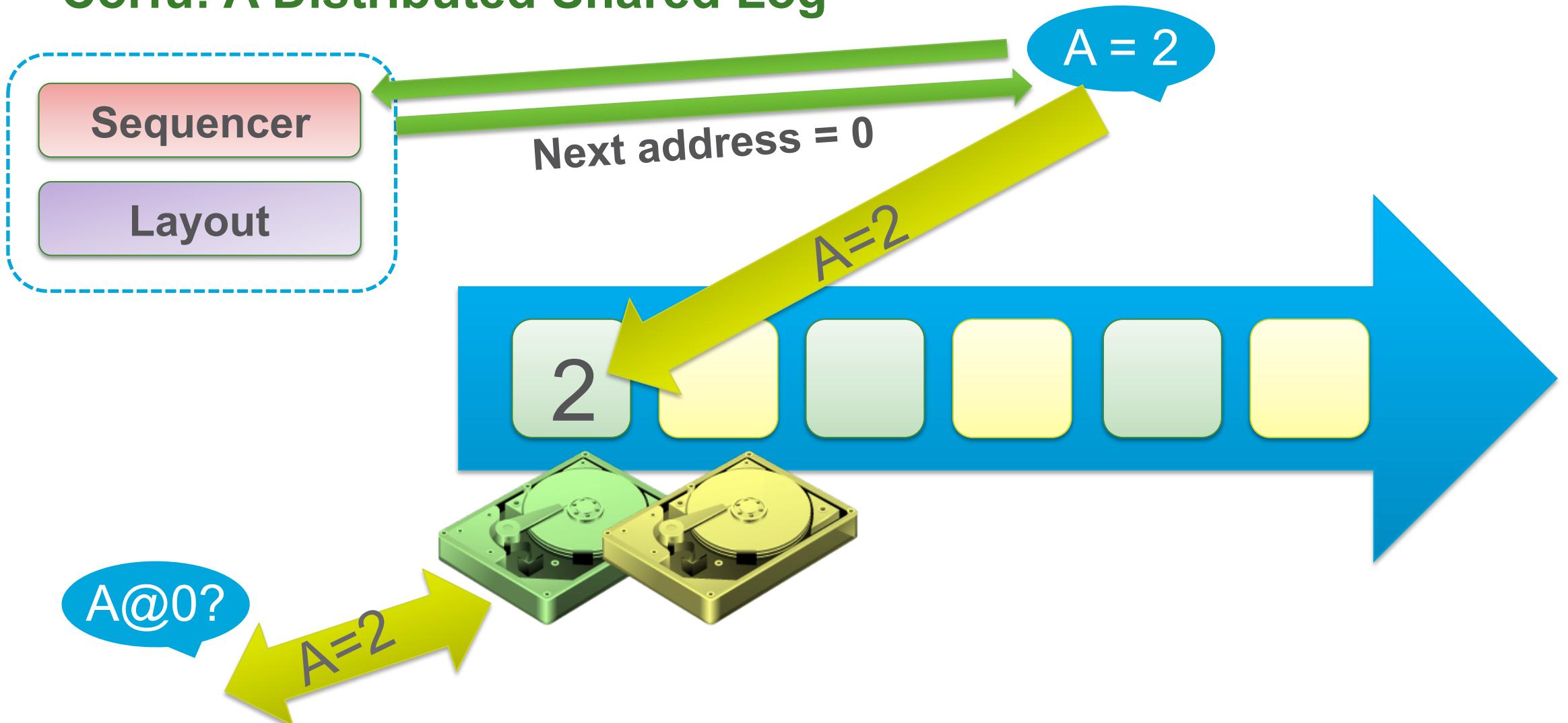
Corfu: A Distributed Shared Log



Corfu: A Distributed Shared Log



Corfu: A Distributed Shared Log



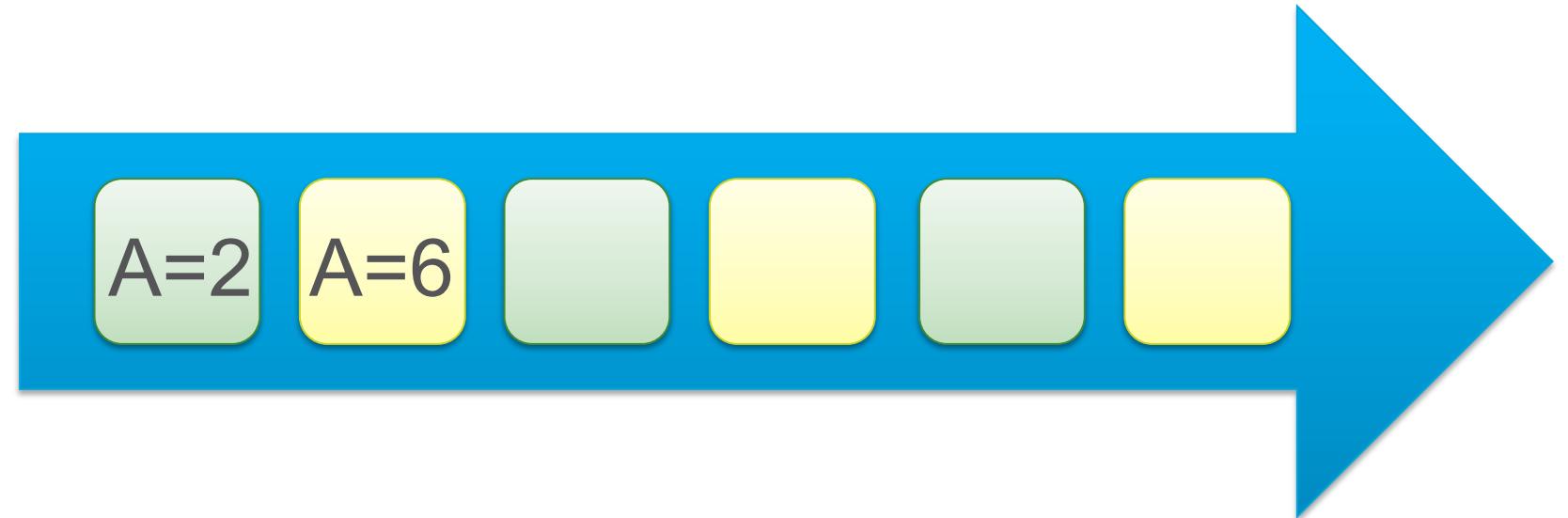
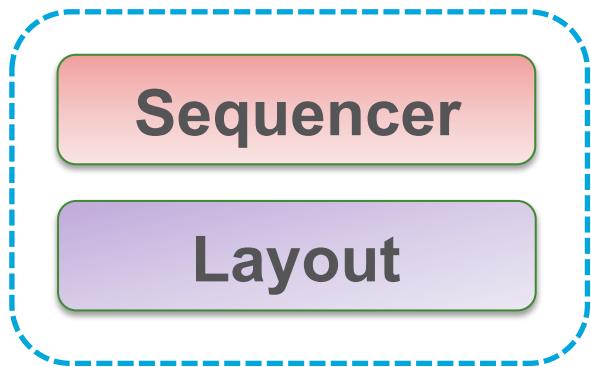
Corfu: A Distributed Shared Log



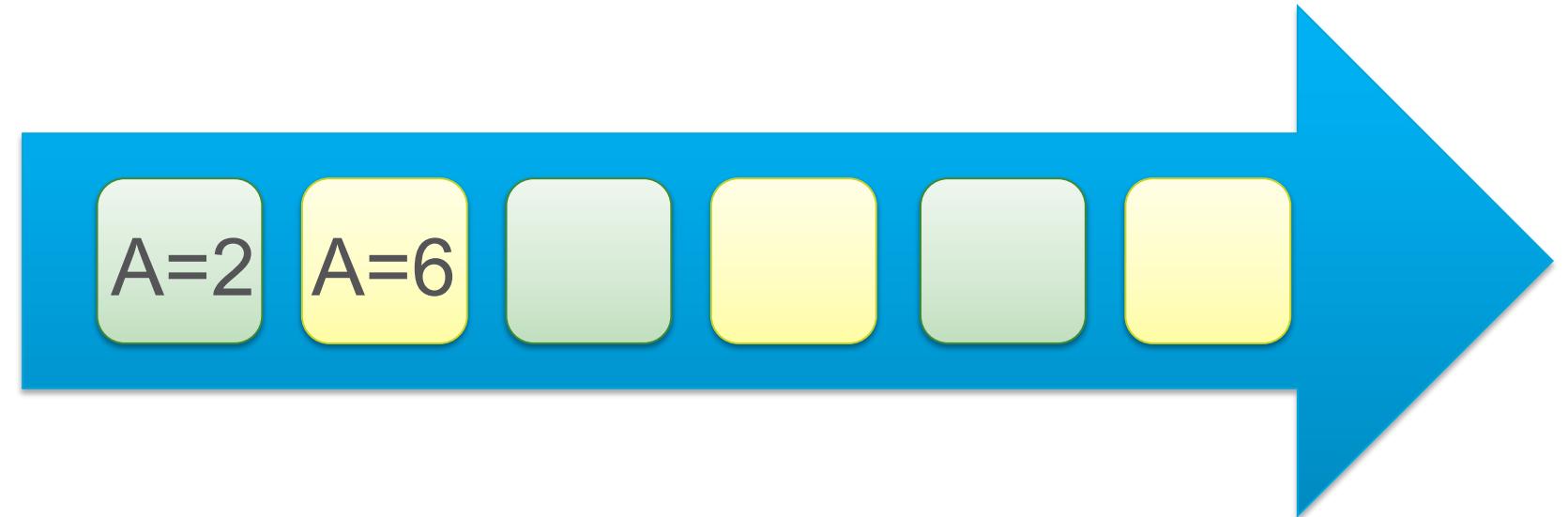
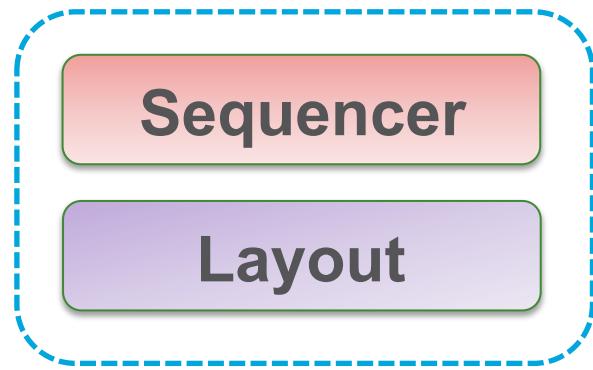
Corfu: A Distributed Shared Log



Corfu: A Distributed Shared Log



Corfu: A Distributed Shared Log

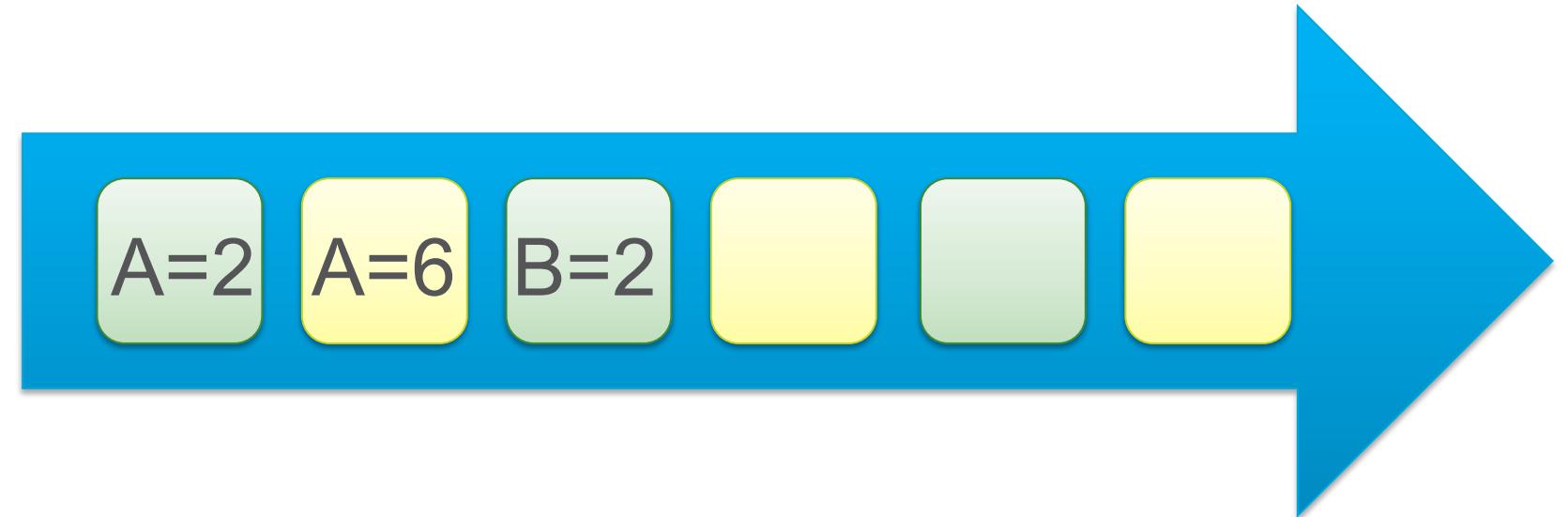


Corfu: A Distributed Shared Log

A@1 B@2

Sequencer

Layout



Corfu: A Distributed Shared Log

A@1 B@2

Sequencer

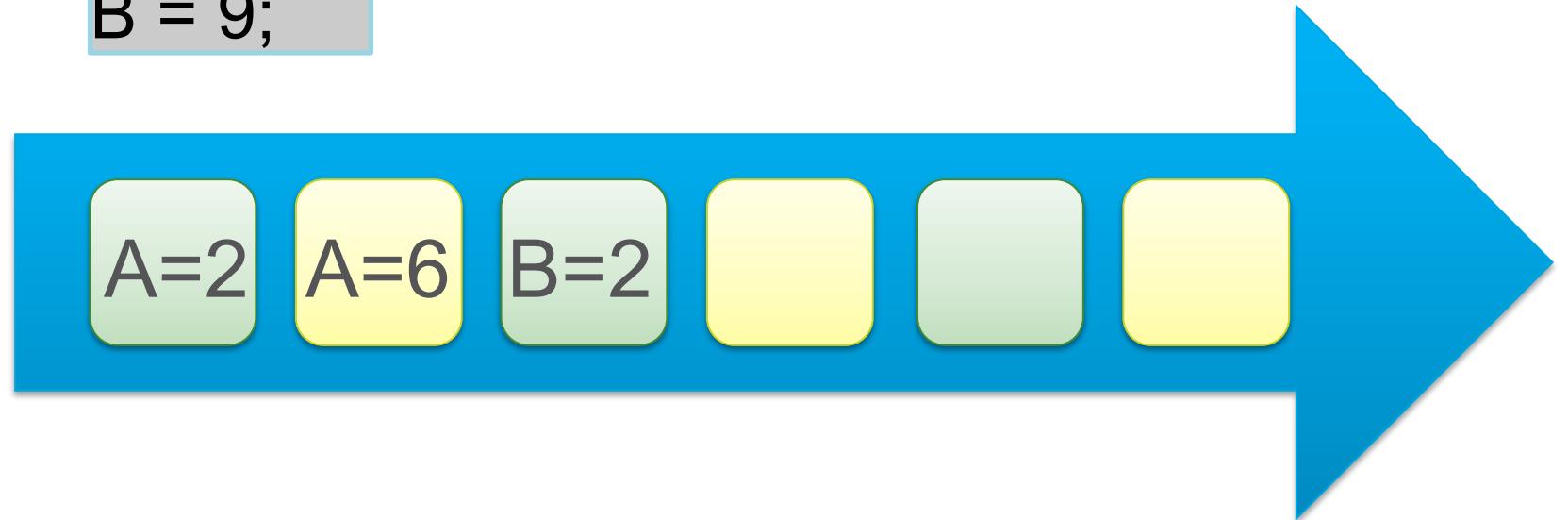
Layout

Transaction

A = 1;
B = 9;

A
6

B
2



Corfu: A Distributed Shared Log

A@3 B@3

Sequencer

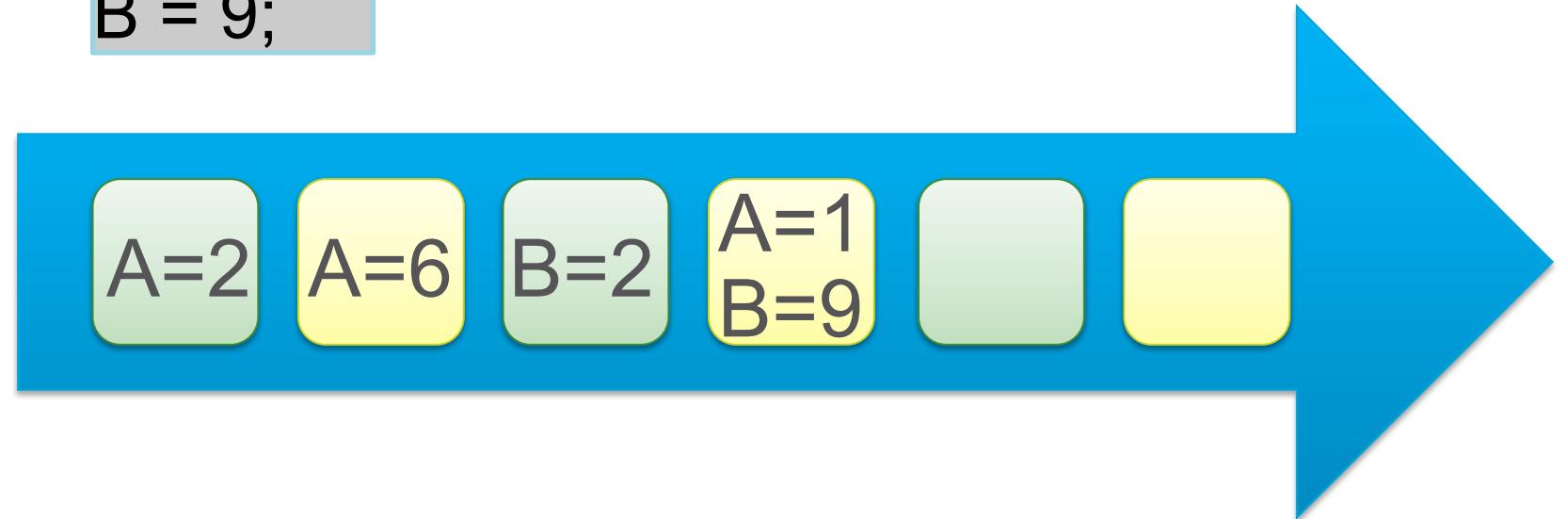
Layout

Transaction

A = 1;
B = 9;

A
1

B
9



Corfu: A Distributed Shared Log

A@3 B@3

Sequencer

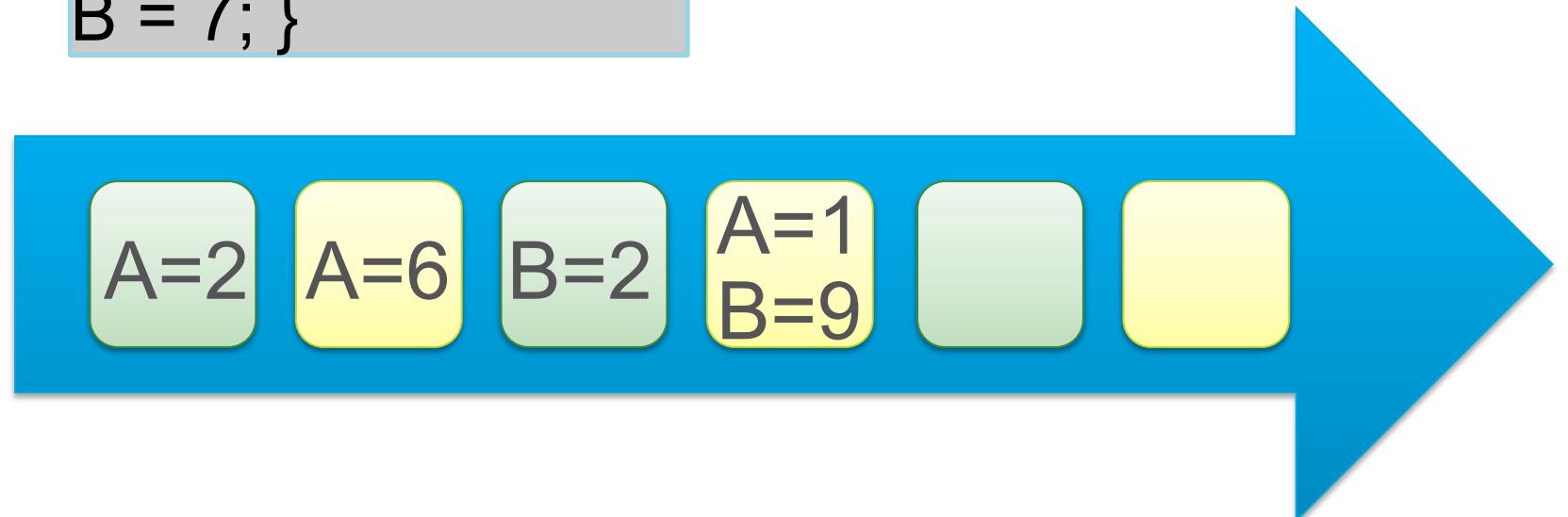
Layout

Transaction

```
if (A == 1) {  
    B = 7; }
```

A
1

B
9



Corfu: A Distributed Shared Log

A@3 B@3

Sequencer

Layout

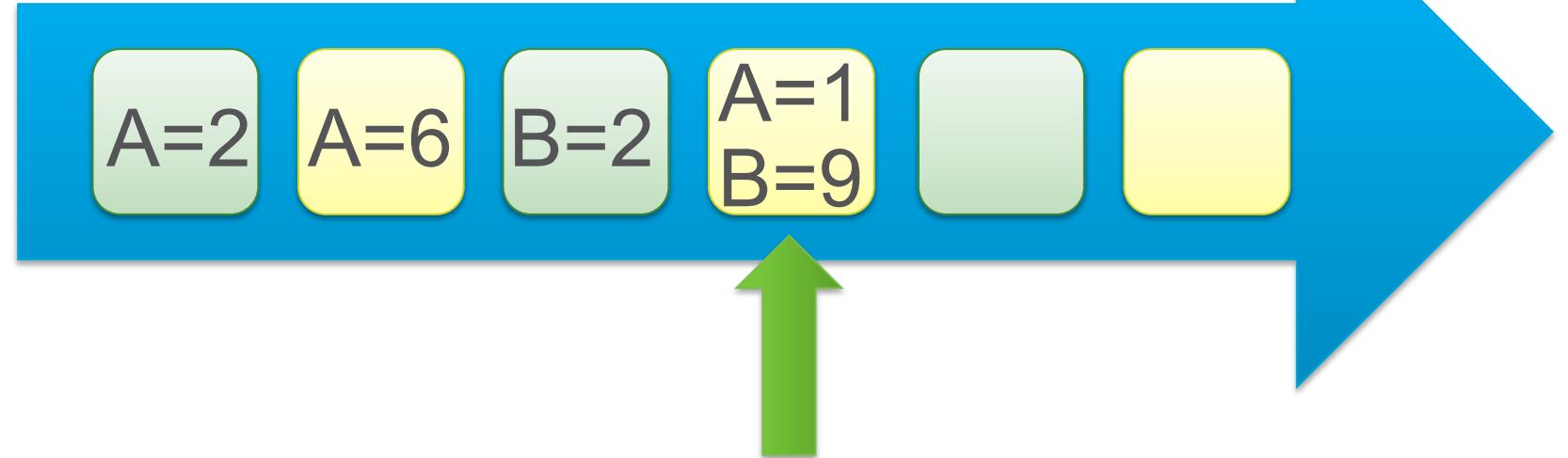
Transaction

```
if (A == 1) {  
    B = 7; }
```

A
1

B
9

RS = A @ 3
WS = { B = 7 }



Corfu: A Distributed Shared Log

A@3 B@4

Sequencer

Layout

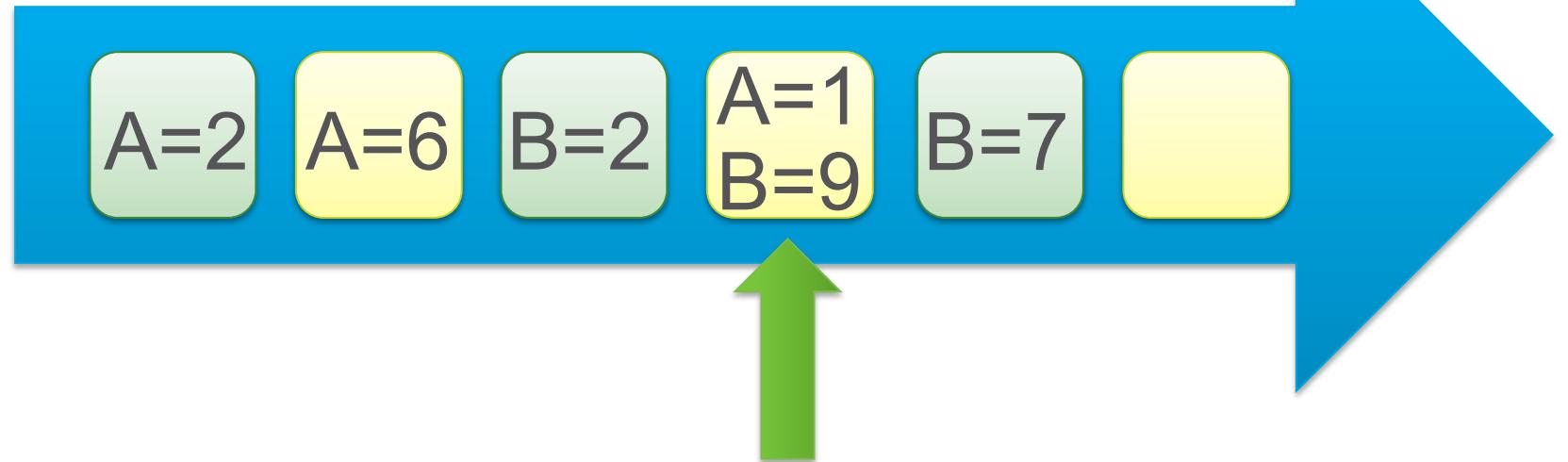
Transaction

```
if (A == 1) {  
    B = 7; }
```

A
1

B
7

RS = A @ 3
WS = { B = 7 }



Corfu: A Distributed Shared Log

A@4 B@3

Sequencer

Layout

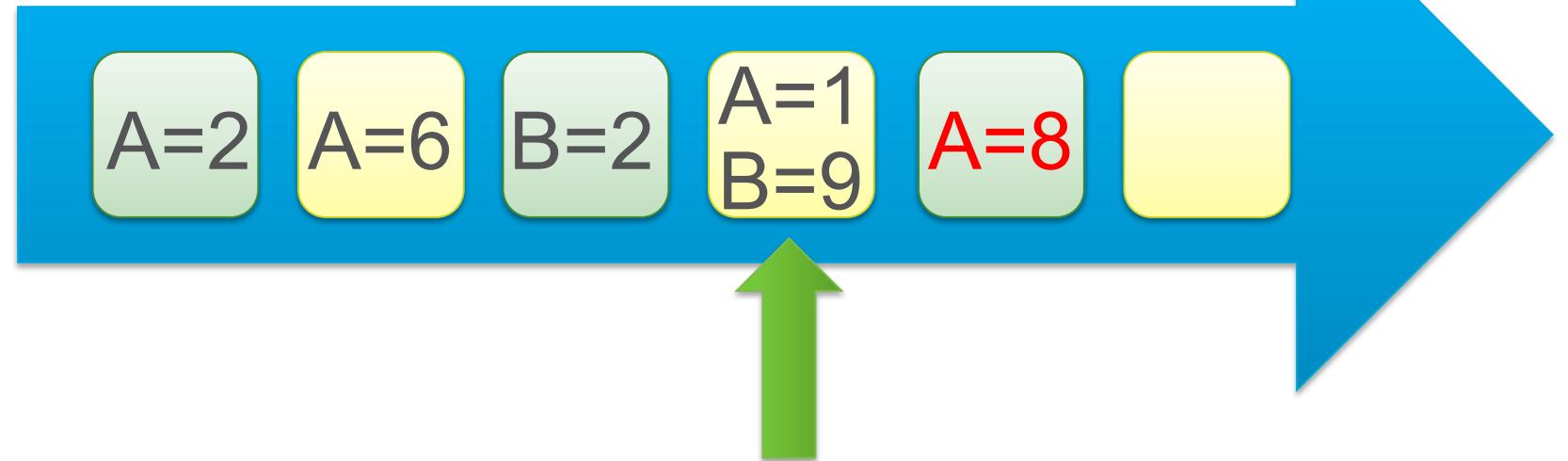
Transaction

```
if (A == 1) {  
    B = 7; }
```

A
1

B
7

RS = A @ 3
WS = { B = 7 }



Address = 3

Corfu: A Distributed Shared Log

A@4 B@3

Sequencer

Layout

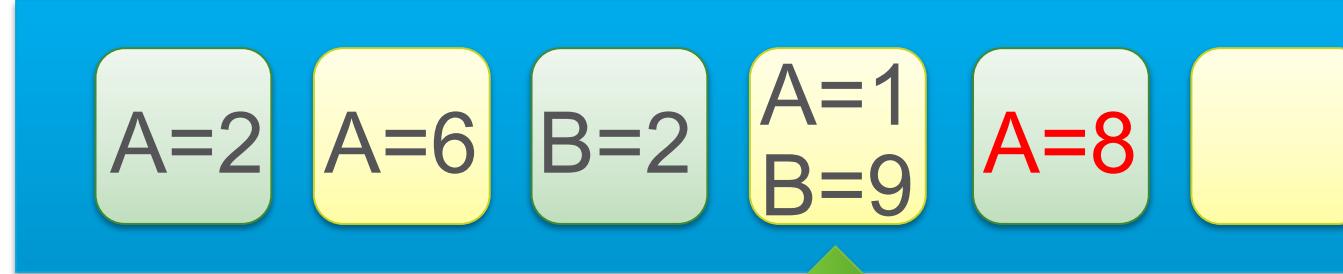
Transaction

```
if (A == 1) {  
    B = 7; }
```

A
1

B
7

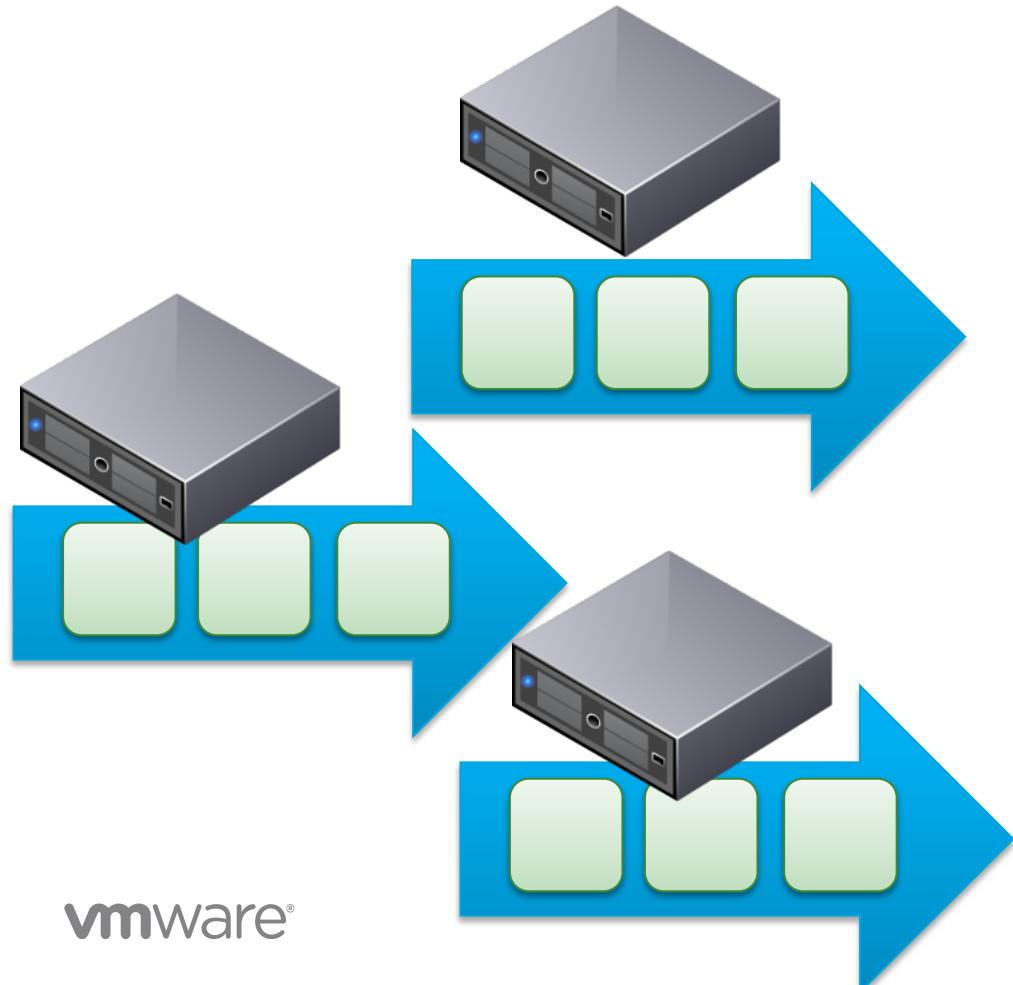
R - A @ 3
W = { B = 7 }



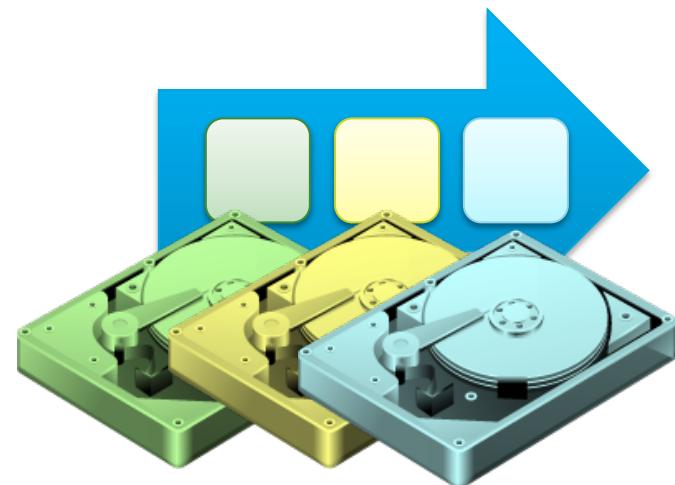
Address = 3

Corfu vs other “Distributed Logs”

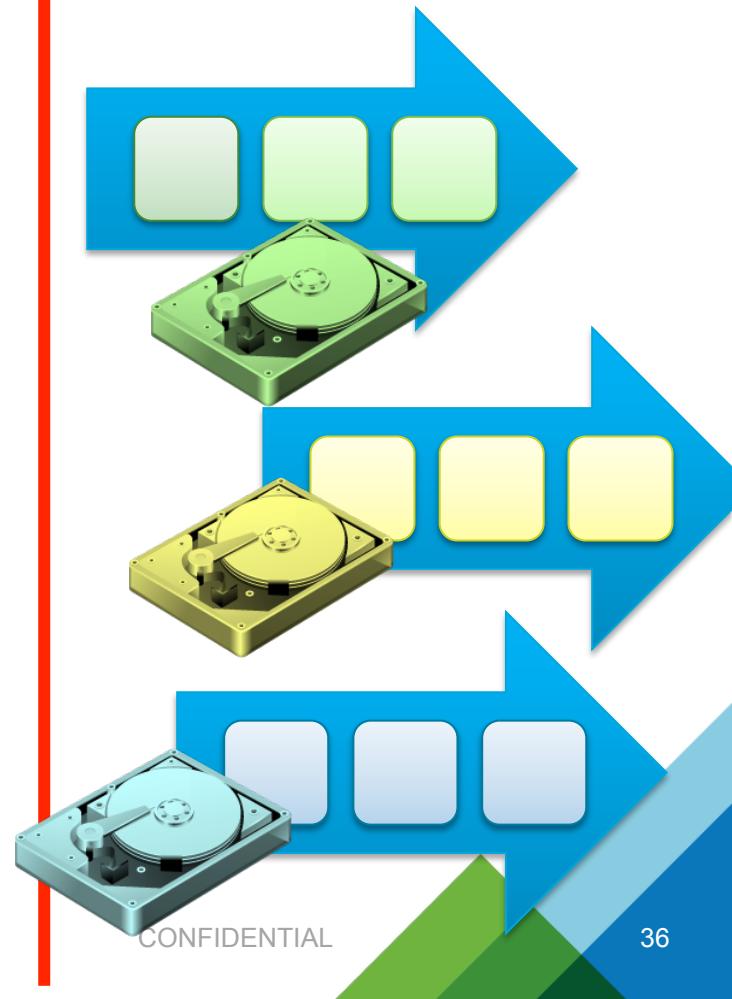
Paxos



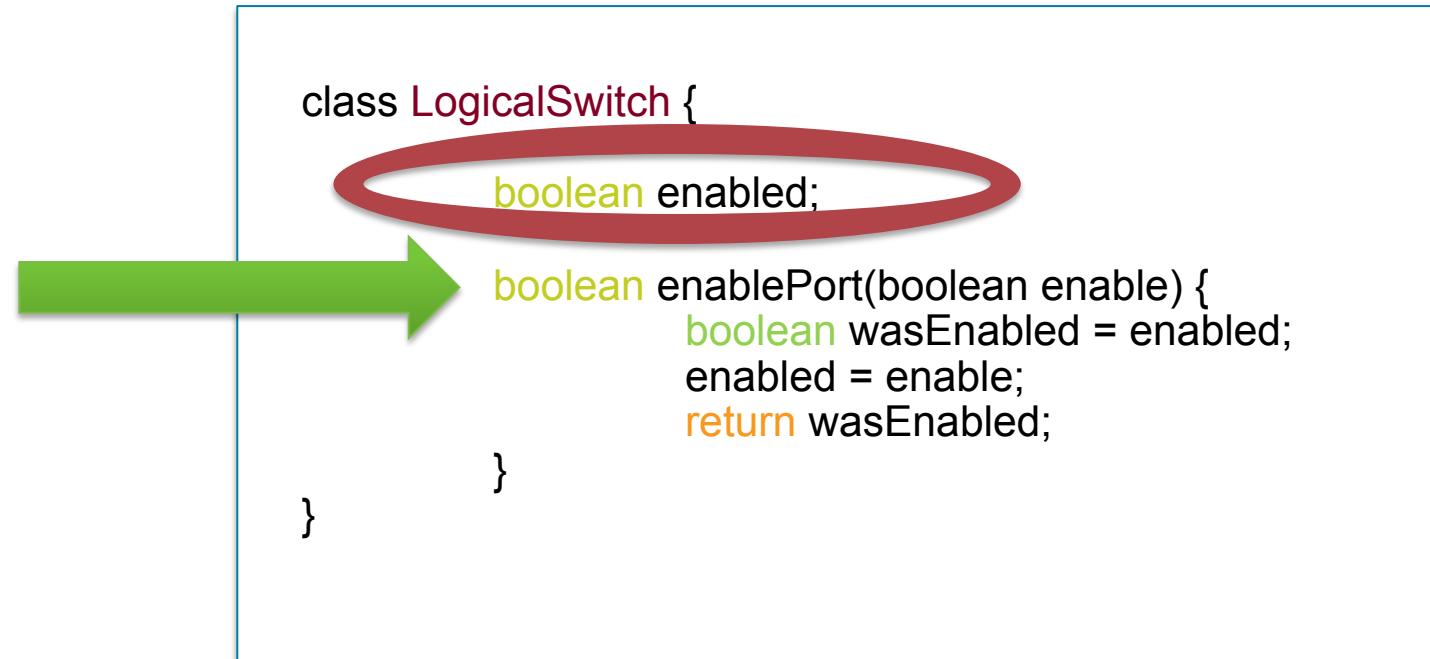
Corfu



Kafka



New Programming Model



```
class LogicalSwitch {  
    boolean enabled;  
    boolean enablePort(boolean enable) {  
        boolean wasEnabled = enabled;  
        enabled = enable;  
        return wasEnabled;  
    }  
}
```

New Programm

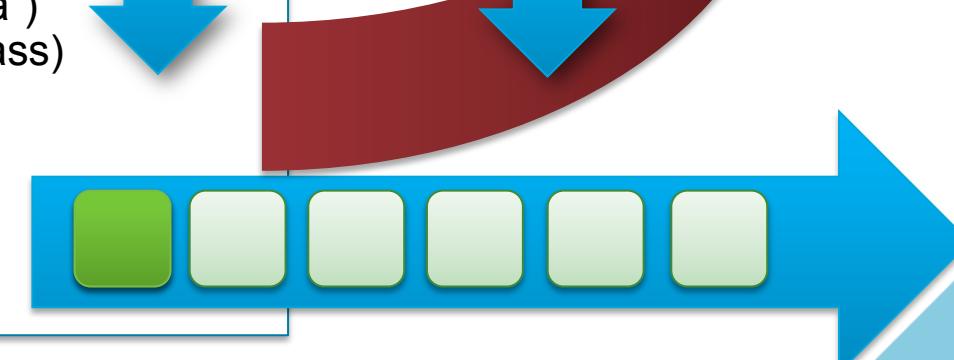
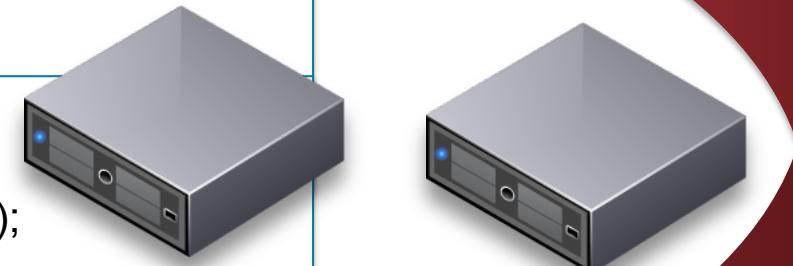
```
class LogicalSwitch {  
    boolean enabled;  
  
    boolean enablePort(boolean enable)  
        boolean wasEnabled = enabled;  
        enabled = enable;  
        return wasEnabled;  
}
```

```
CorfuRuntime cr =  
    new CorfuRuntime("10.0.0.1:9000");  
  
LogicalSwitch mySwitch =  
    cr.build()  
        .setStreamName("switch-a")  
        .setType(LogicalSwitch.class)  
        .open();  
  
mySwitch.enablePort(true);
```

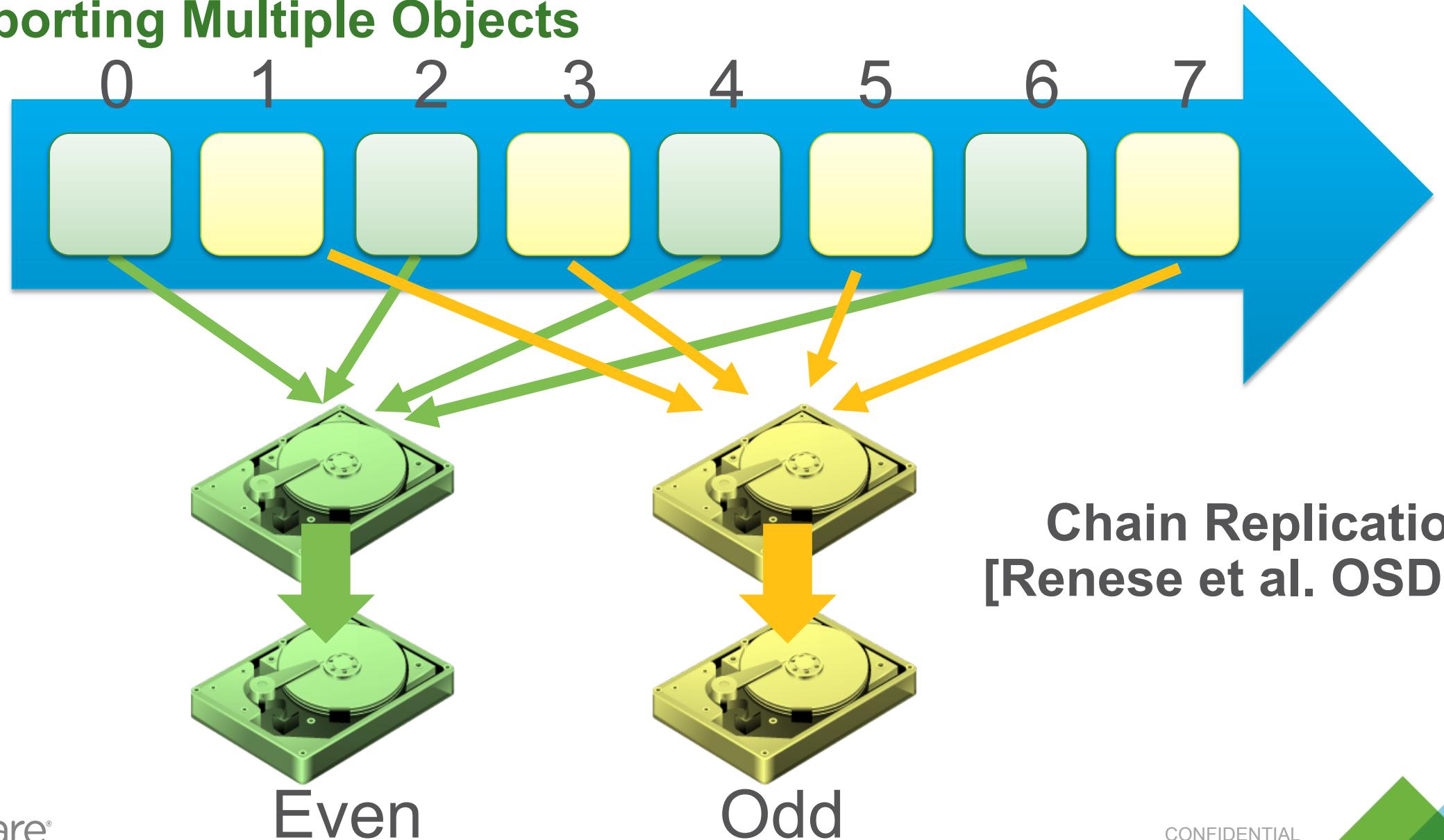
New Programm

```
class LogicalSwitch {  
    boolean enabled;  
  
    boolean enablePort(boolean enable)  
        boolean wasEnabled = enabled;  
        enabled = enable;  
        return wasEnabled;  
}
```

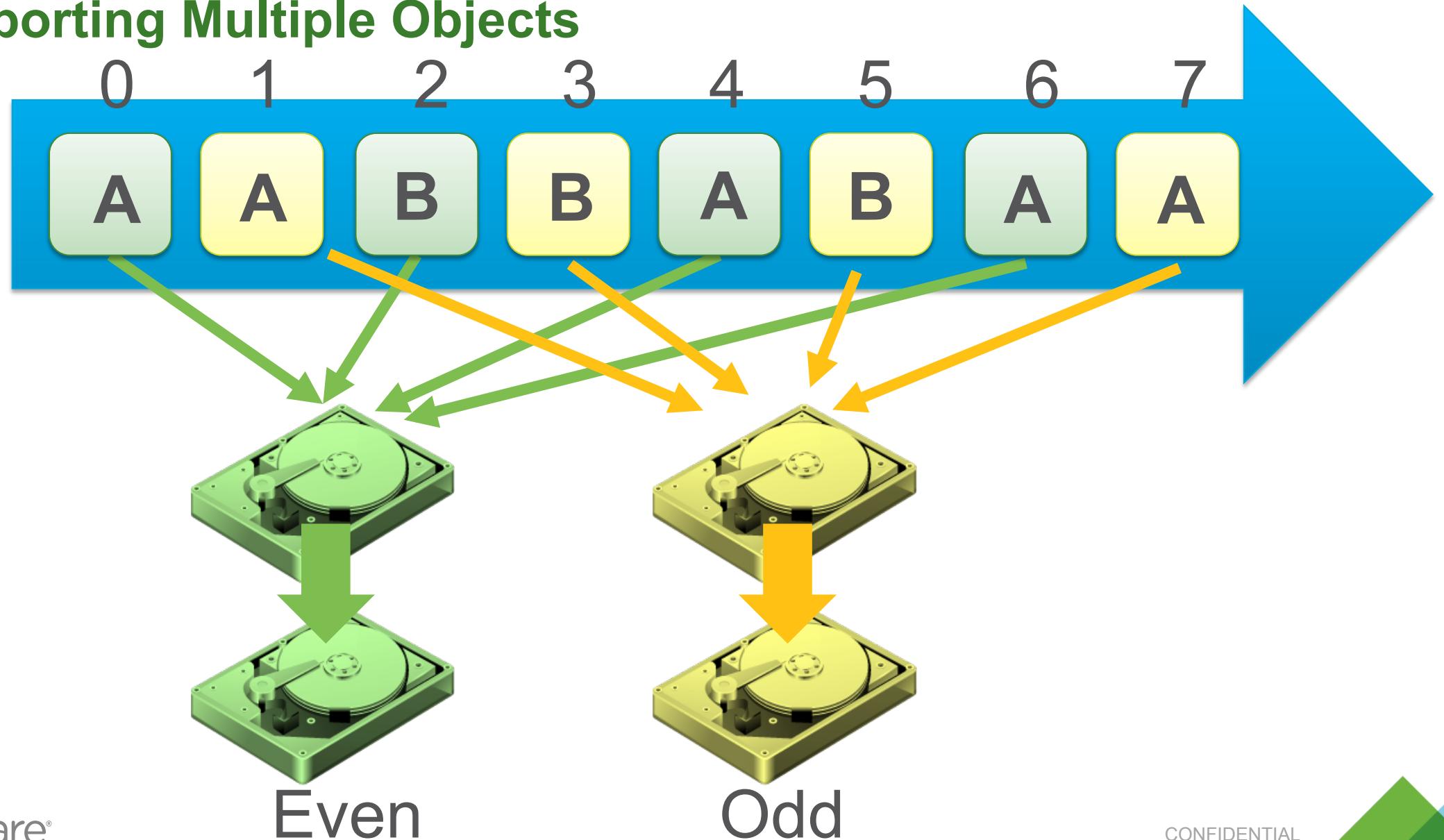
```
CorfuRuntime cr =  
    new CorfuRuntime("127.0.0.1:9000");  
  
LogicalSwitch mySwitch  
    cr.build()  
        .setStreamName("switch-a")  
        .setType(LogicalSwitch.class)  
        .open();  
  
mySwitch.enablePort(true);
```



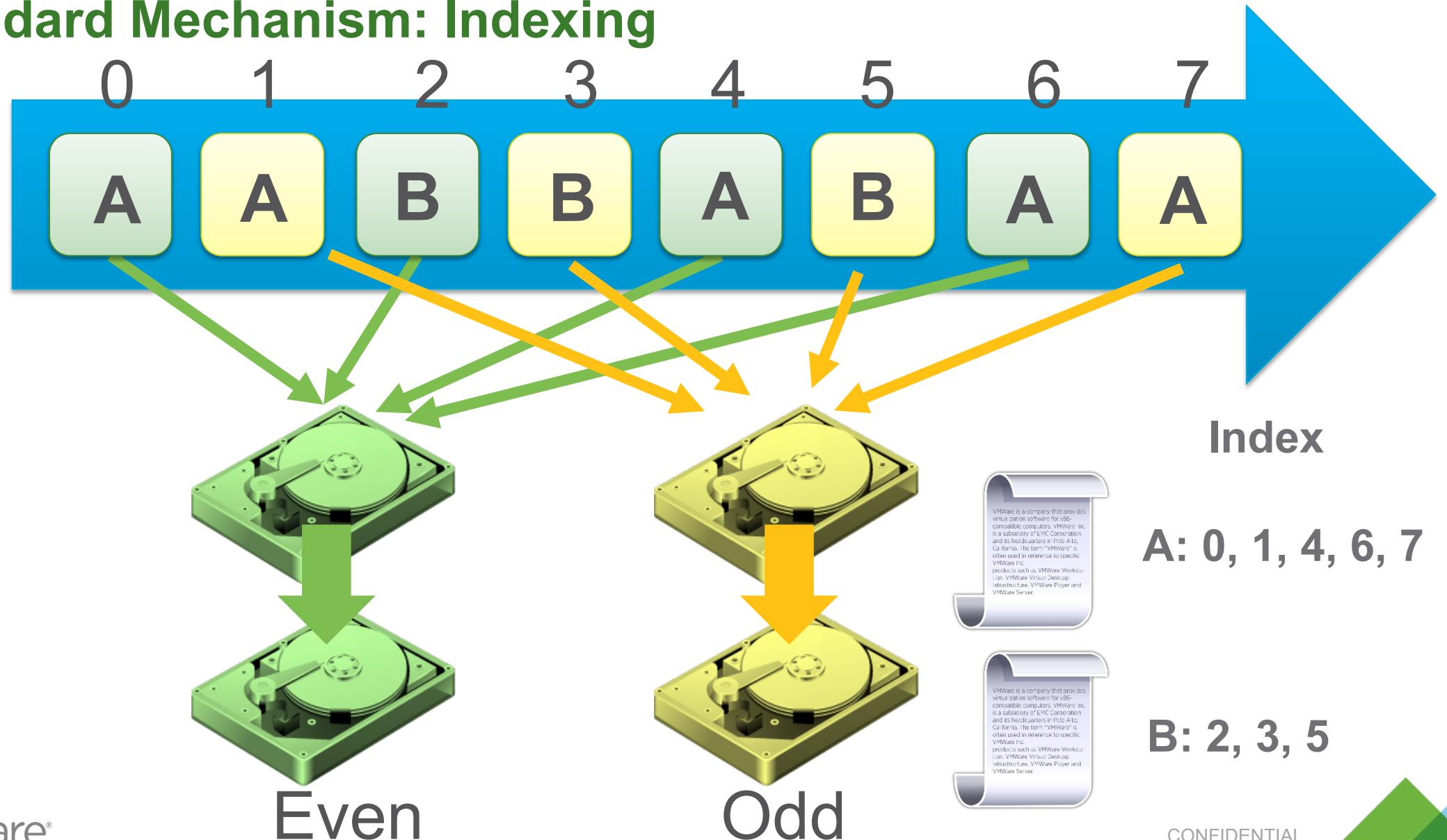
Supporting Multiple Objects



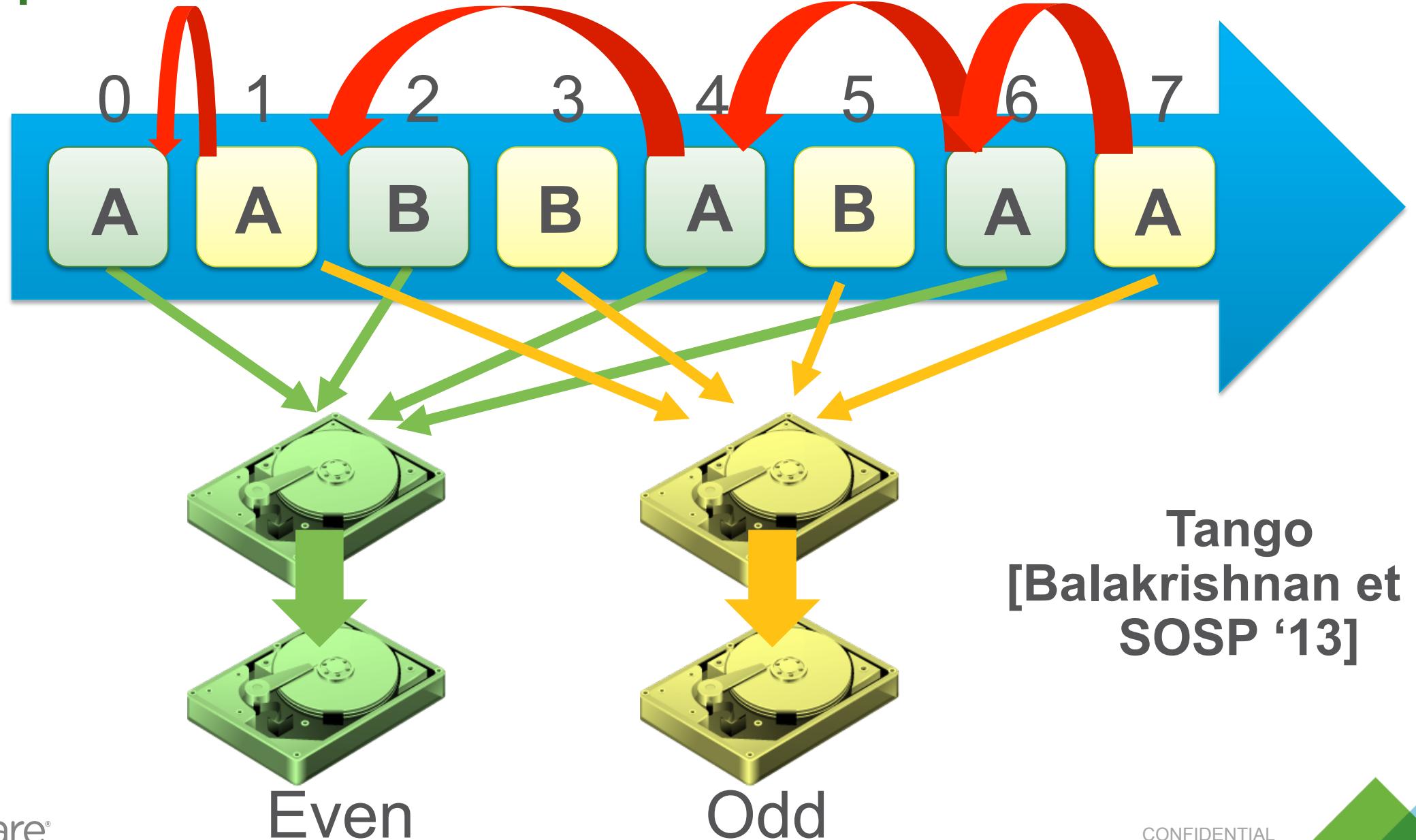
Supporting Multiple Objects



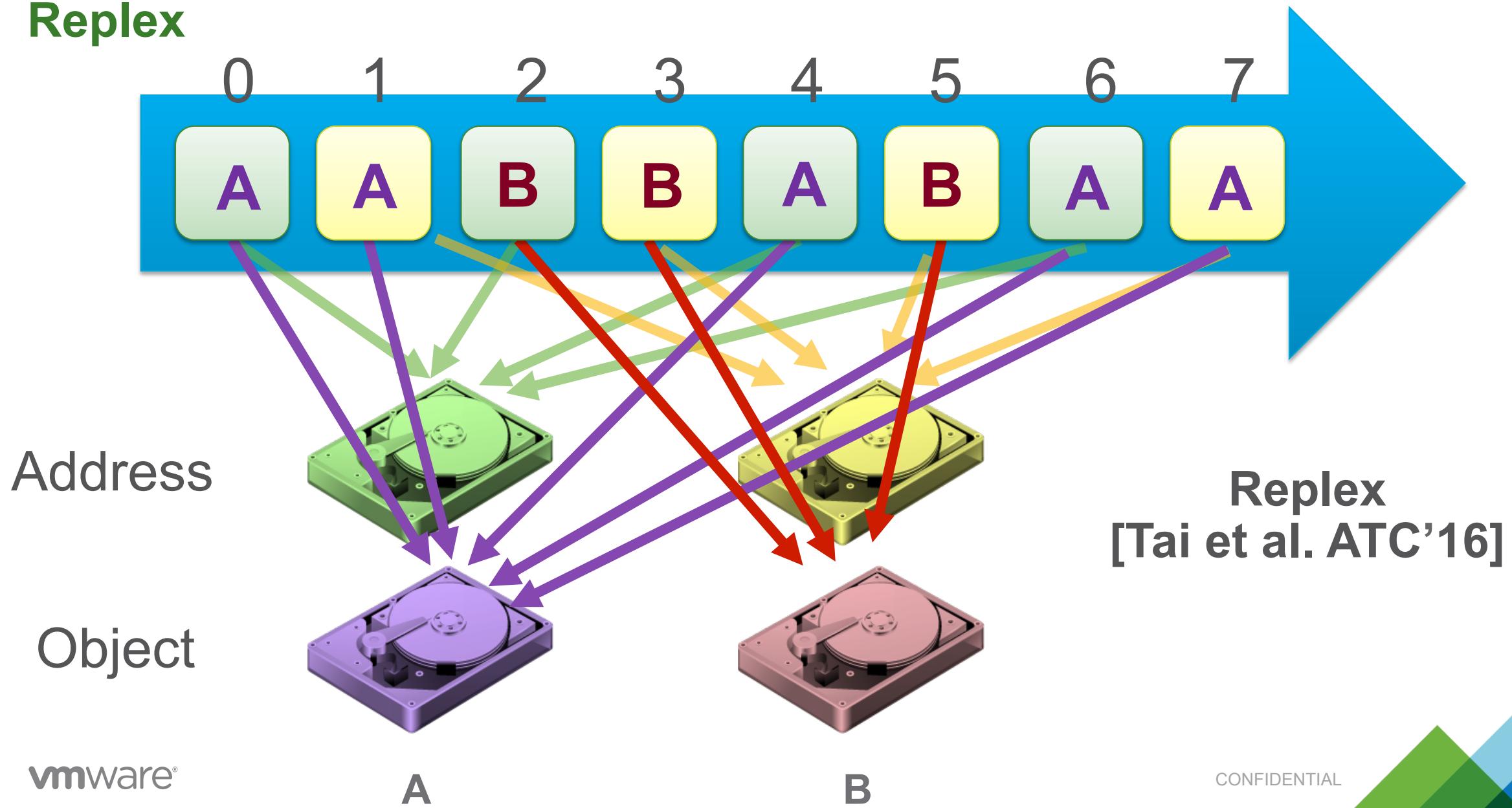
Standard Mechanism: Indexing



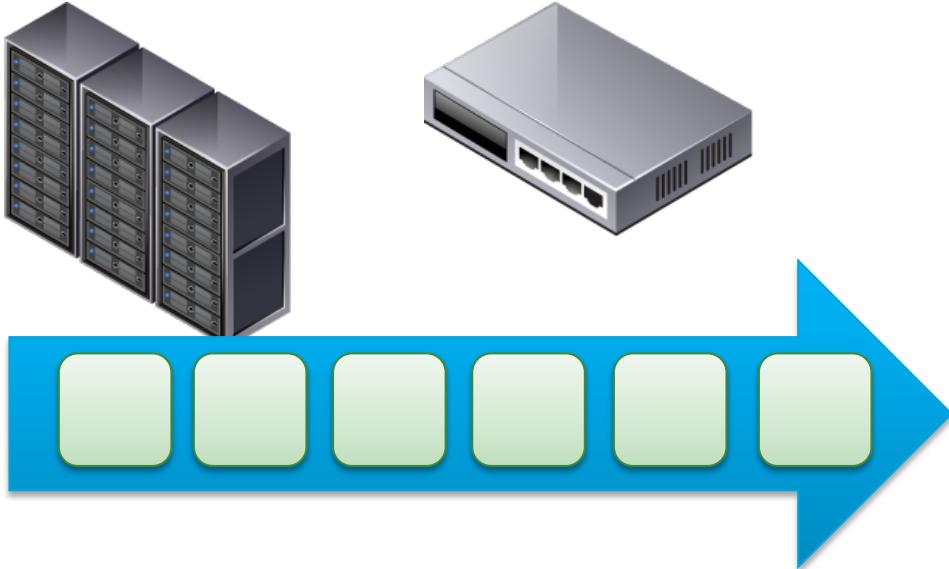
Backpointers



Replex



CMPLAT: Cluster Management Platform for NSX



Left to Right:
Scott Fritchie
Michael Wei(PoC),
Chris Rossbach,
Dahlia Malkhi (PoC),
Udi Wieder



Jim Stabile, Ittai Abraham,
Sandeep Uttamchandani
Medhavi Dawan,
Maithem Munshed



- Re-engineer NSX cluster management to use Corfu “under the covers”
- Avoid the need for every NSX installation to include Zookeeper, etc.
- Corfu evolved significantly as a result of this joint effort

CMPLAT: Cluster Management Platform for NSX

Left to Right:
Scott Fritchie
Michael Wei(PoC),
Chris Rossbach,
Dahlia Malkhi (PoC),
Udi Wieder



Jim Stabile, Ittai Abraham,
Sandeep Uttamchandani
Medhavi Dawan,
Maithem Munshed



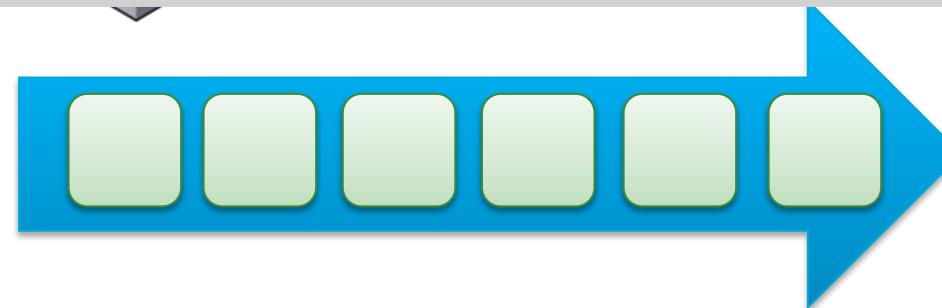
API Requests



NSX Scale-Out Management Plane

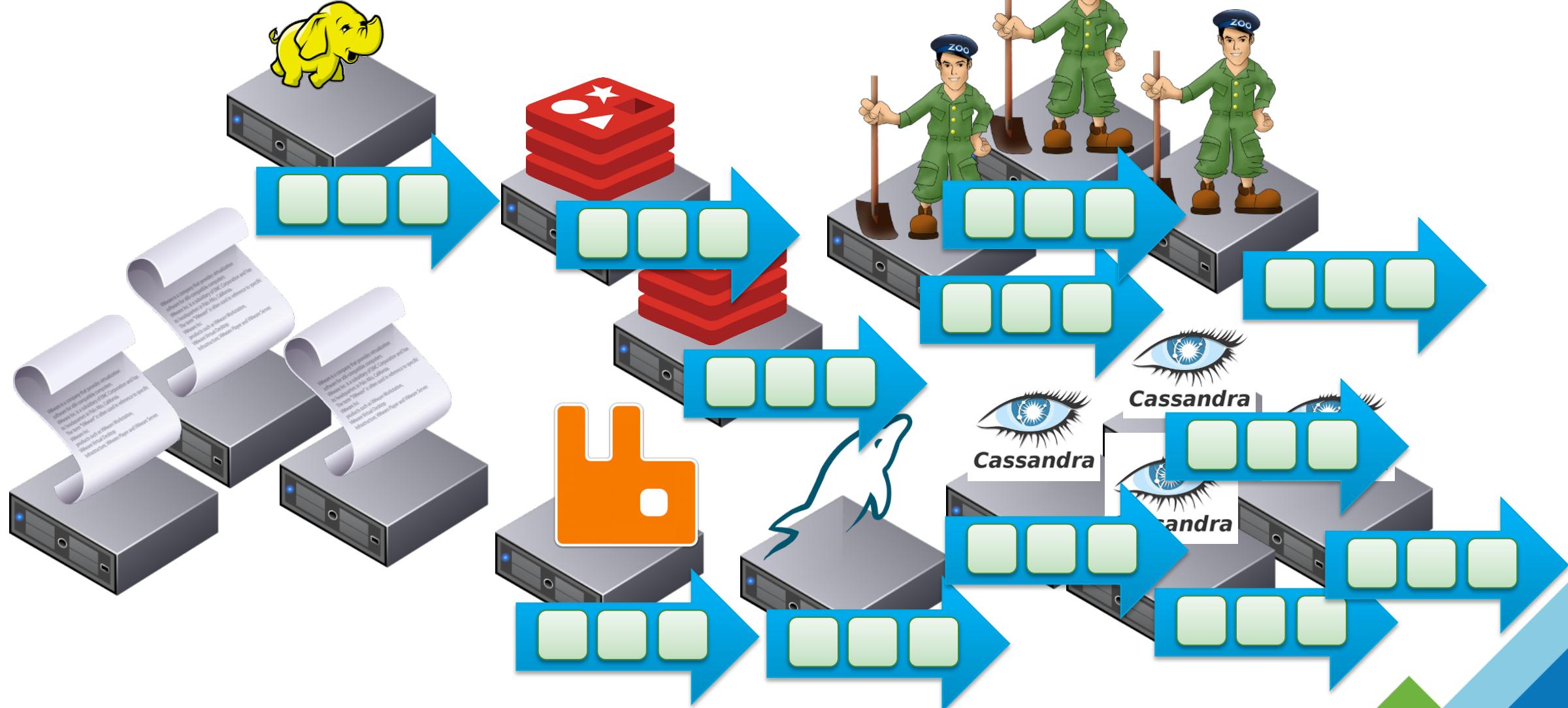


NSX Java Object Model



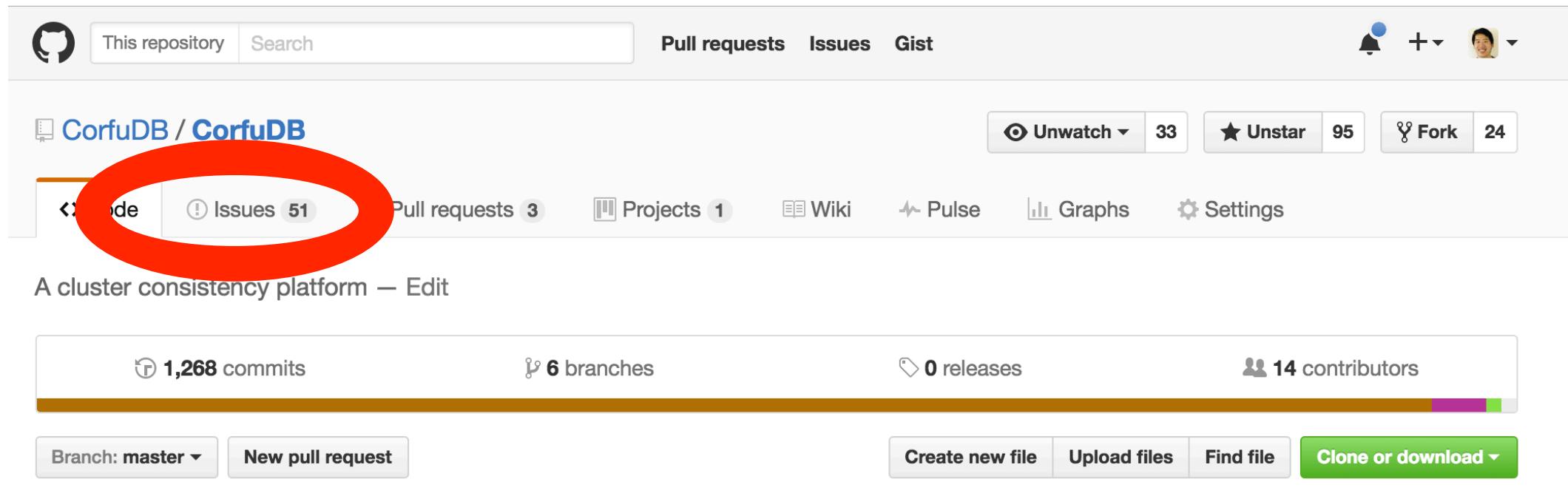
Corfu Distributed Log

OneData



Open Source Efforts

Yes, we're on Github...
github.com/CorfuDB/CorfuDB



No, this is not a source code dump...