

# Exo: Atomic Broadcast for the Rack-Scale Computer

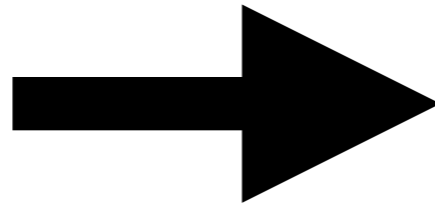
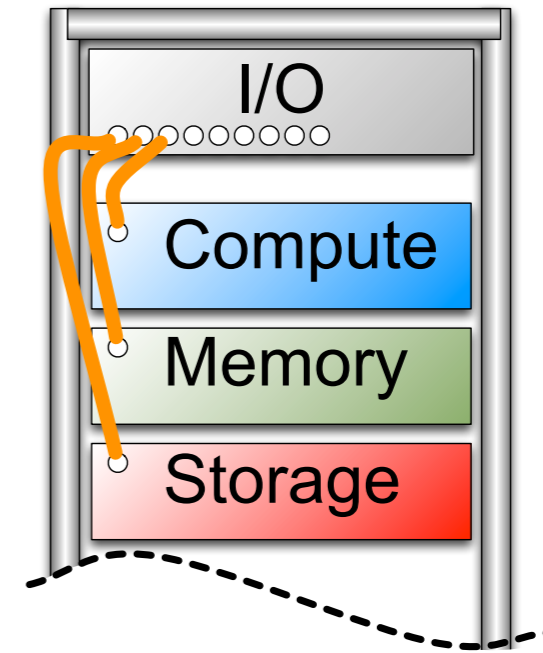
**Matthew P. Grosvenor** Marwan Fayed Andrew W. Moore



## Today



## Tomorrow

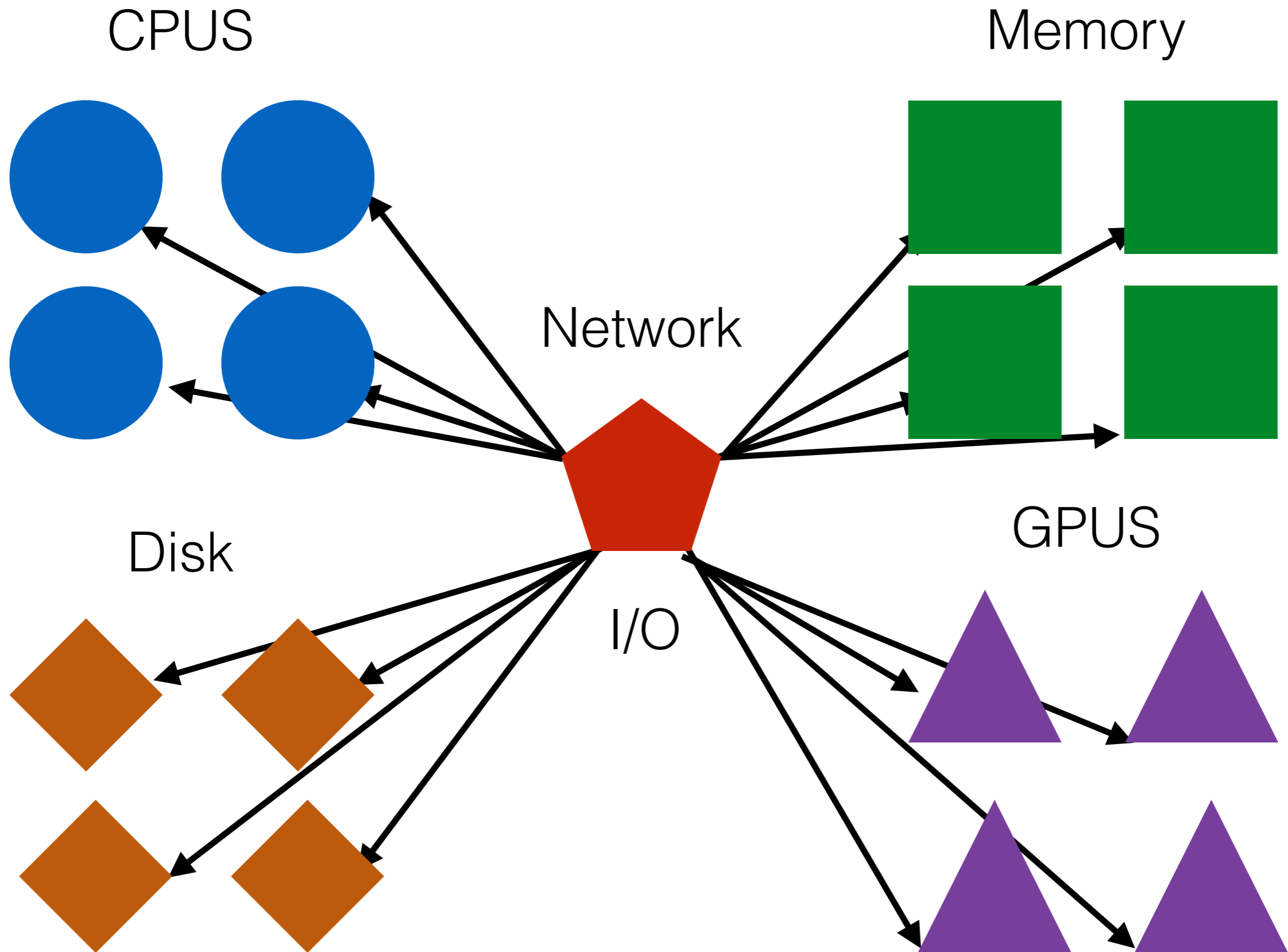


- 50 - 200 machines
- Commodity hardware
- Commodity network

- 500-2000 “nodes”
- Disaggregated hardware
- Custom (photonic?) interconnect(s)



# The problem with disaggregation





## Single Machines

- Coordination using simple MESI/MOSI protocols
- Specialised hardware over a reliable, low latency interconnect.
- ✓ High performance
- x Not failure tolerant
- x Doesn't scale well



## Single Machines

- Coordination using simple MESI/MOSI protocols
- Specialised hardware over a reliable, low latency interconnect.
- ✓ High performance
- x Not failure tolerant
- x Doesn't scale well

## Cluster Systems

- Coordination using software protocols
- General purpose, high latency, error prone network
- x Low performance
- ✓ Fault Tolerant
- x Doesn't scale well



## Single Machines

- Coordination using simple MESI/MOSI protocols
- Specialised hardware over a reliable interconnect

## Cluster Systems

- Coordination using software protocols
- General purpose, high latency, error

**Something in the middle ???**ork

✓ High performance

x Not failure tolerant

x Doesn't scale well



x Low performance

✓ Fault Tolerant

x Doesn't scale well



## Single Machines

- Coordination using simple MESI/MOSI protocols
- Specialised hardware over a reliable interconnect

## Cluster Systems

- Coordination using software protocols
- General purpose, high latency, error

**Something in the middle ???**ork

✓ High performance

x Not failure tolerant



x Low performance

✓ Fault Tolerant

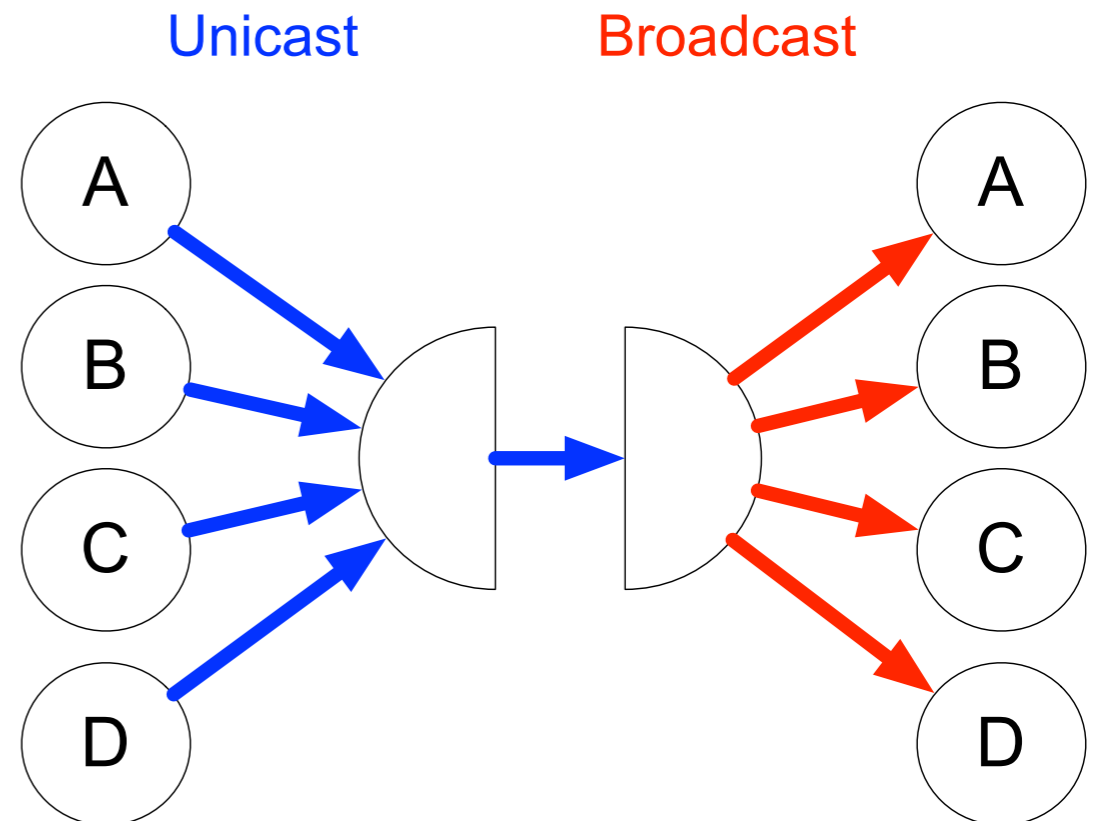
x Doesn't scale well

x Doesn't scale well

**That scales?**



- ✓ Silicon Photonic Interfaces on chip
- ✗ High radix optical switches?
- ✓ Burst mode transceivers
- ✓ Passive all optical interconnect (PON)
- ✓ Network broadcast as a primitive for building Atomic Broadcast



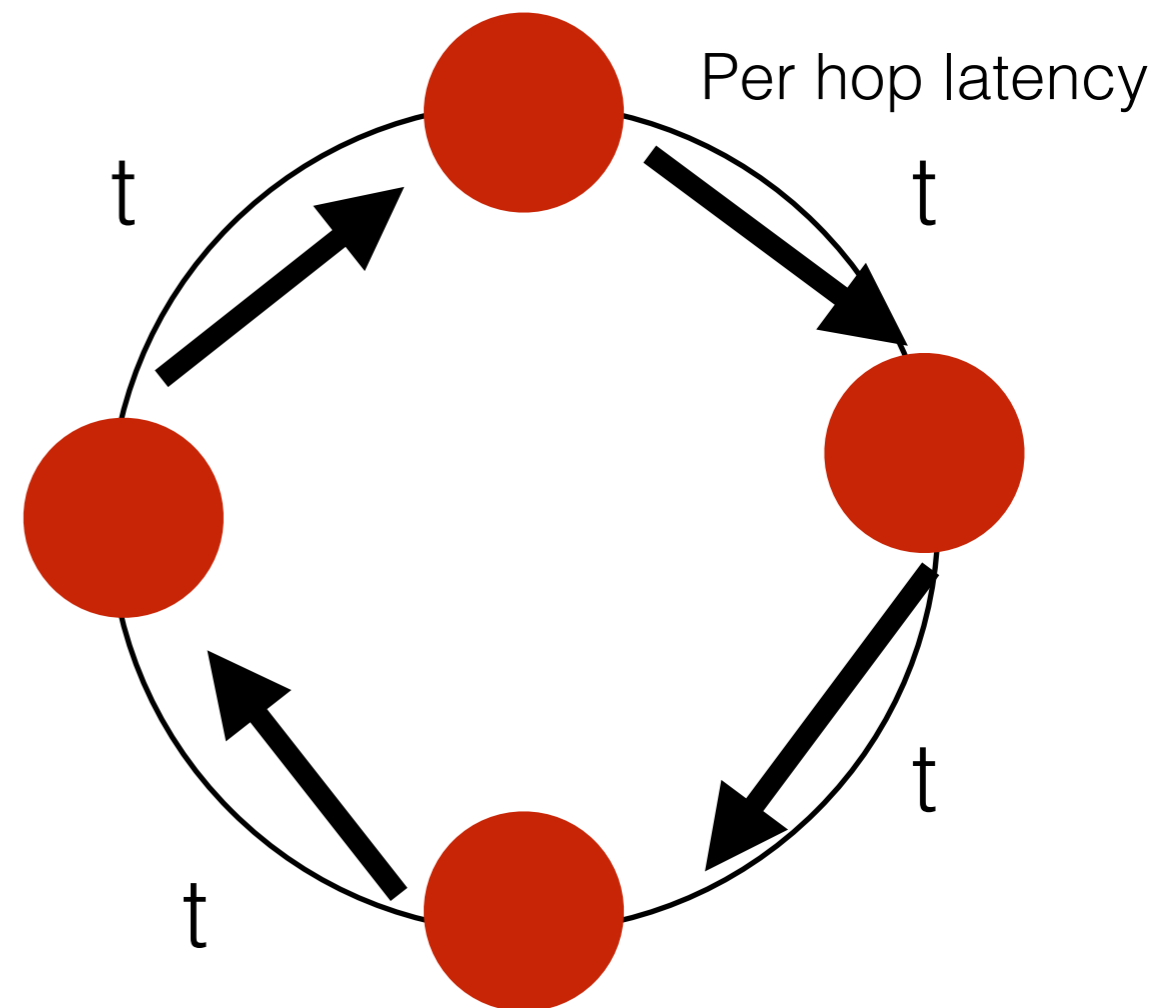




- How do we mediate access to the network?
- What agreement protocol do we use?
- Use token rings
  - Two birds with one stone
    - ✓ Well established MAC protocol
    - ✓ Well established agreement protocol. Very high (optimal) throughput performance.



- Token passing is latency bound
- Each host must wait at least  $n-1$  hops for agreement =  $n \times$  per-hop latency ( $t$ )
- Per hop latency must be minimised



$$\text{Total latency} = t + t + t + t = 4t$$

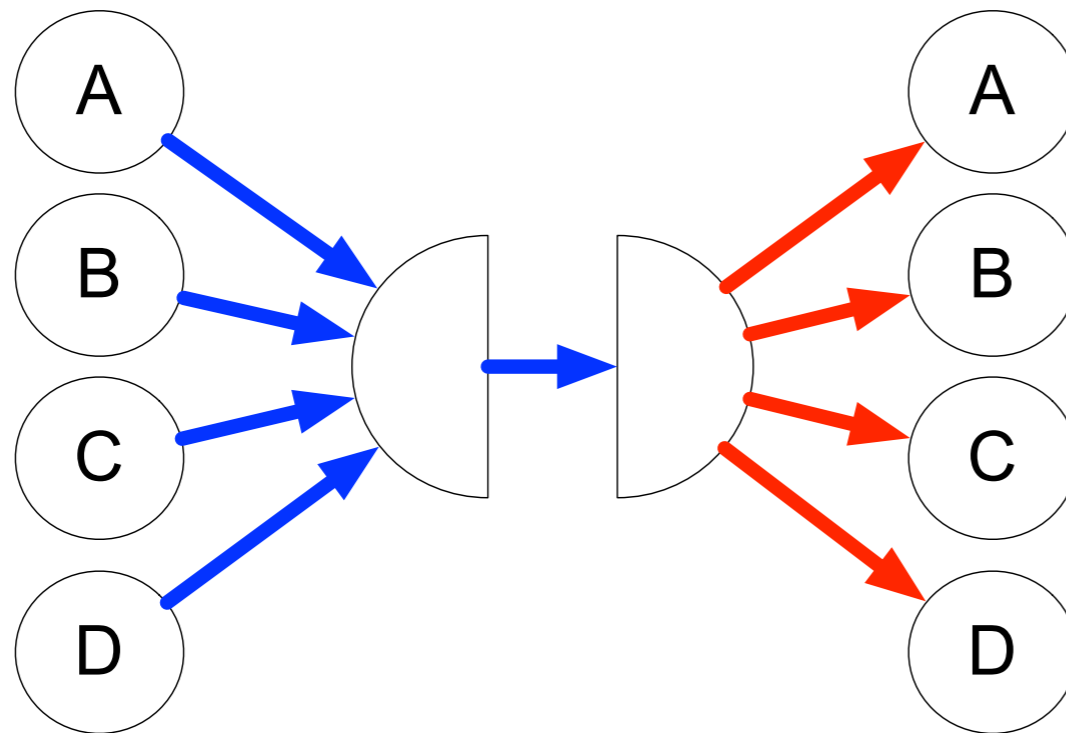
$$\text{Throughput (per host)} \approx 1/4t$$



# What's the latency?

Assuming 3m of fibre @ 0.75C

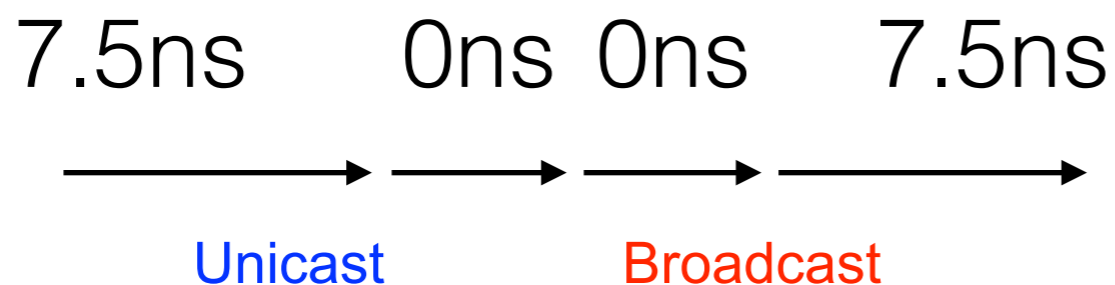
7.5ns    0ns   0ns    7.5ns    = 15ns per hop  
→      →      →      →  
Unicast      Broadcast





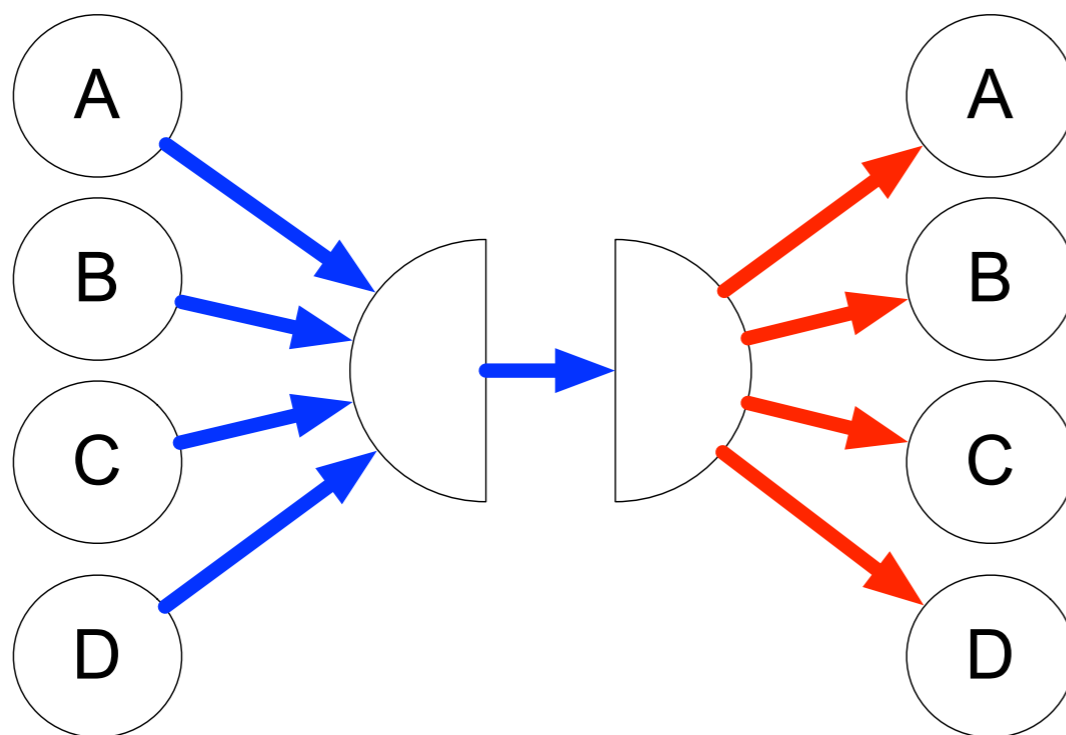
# What's the latency?

Assuming 3m of fibre @ 0.75C



= 15ns per hop

@10G 64B frames  
= 14.4M msg/sec

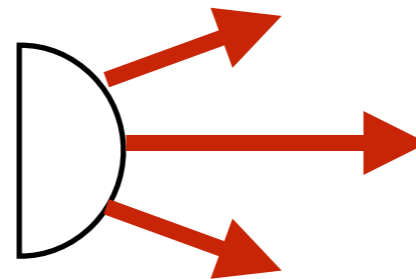




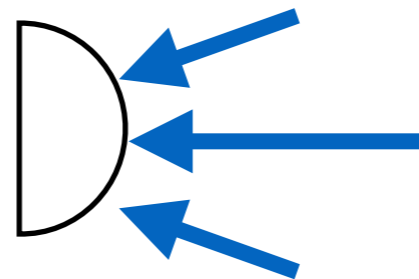
Use Exalink Fusion Switch to emulate all optical network



- Broadcast - 5ns



- Aggregation - 95ns

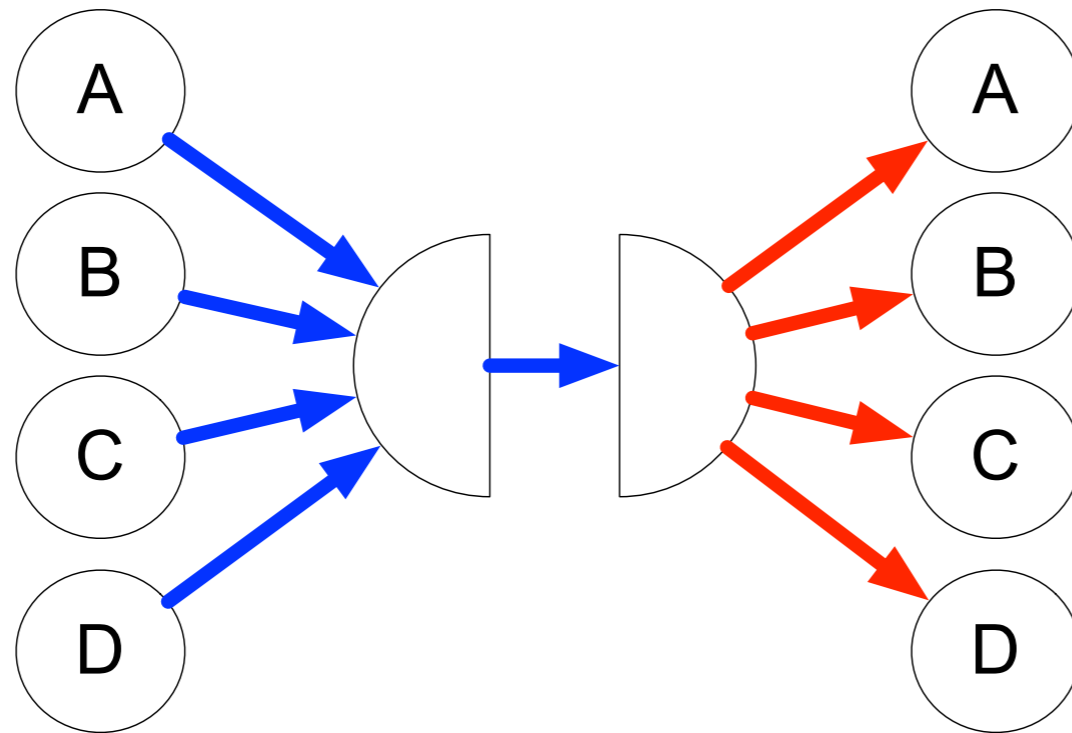




7.5ns    95ns    5ns    7.5ns

Unicast                      Broadcast

= 115ns per hop  
= 8.6M msg/sec



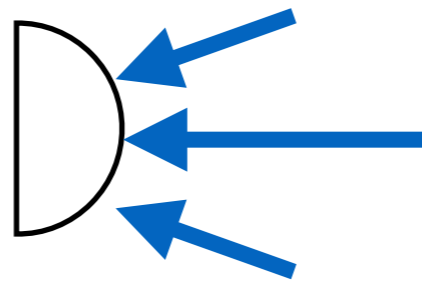
Fusion Aggregate

Fusion Broadcast



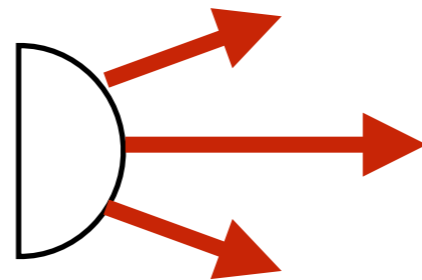
- Use ethernet switch and matrix switch to emulate Exablaze Fusion

- Arista 7124FX



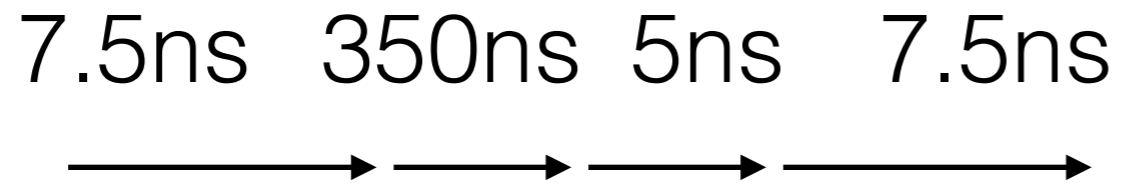
- Aggregation/switching 350ns

- Exablaze ExaLink 50



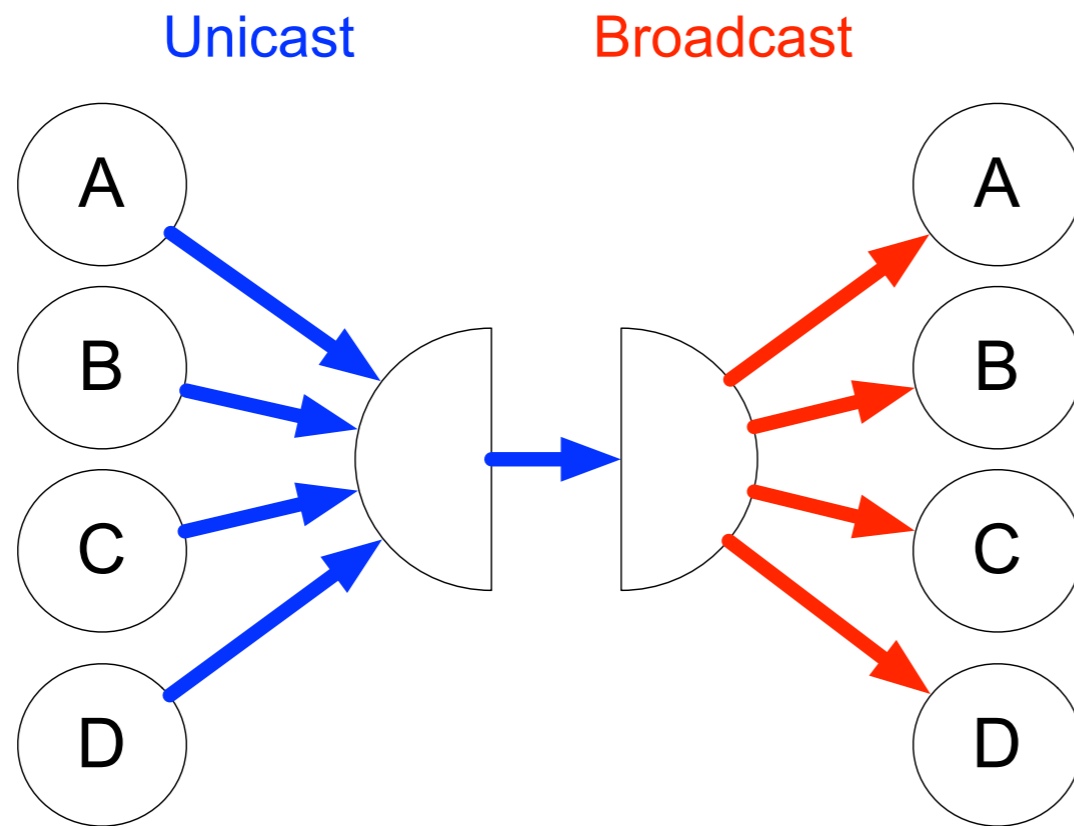
- Broadcast <5ns





= 370ns per hop

= 2.7M msg/sec



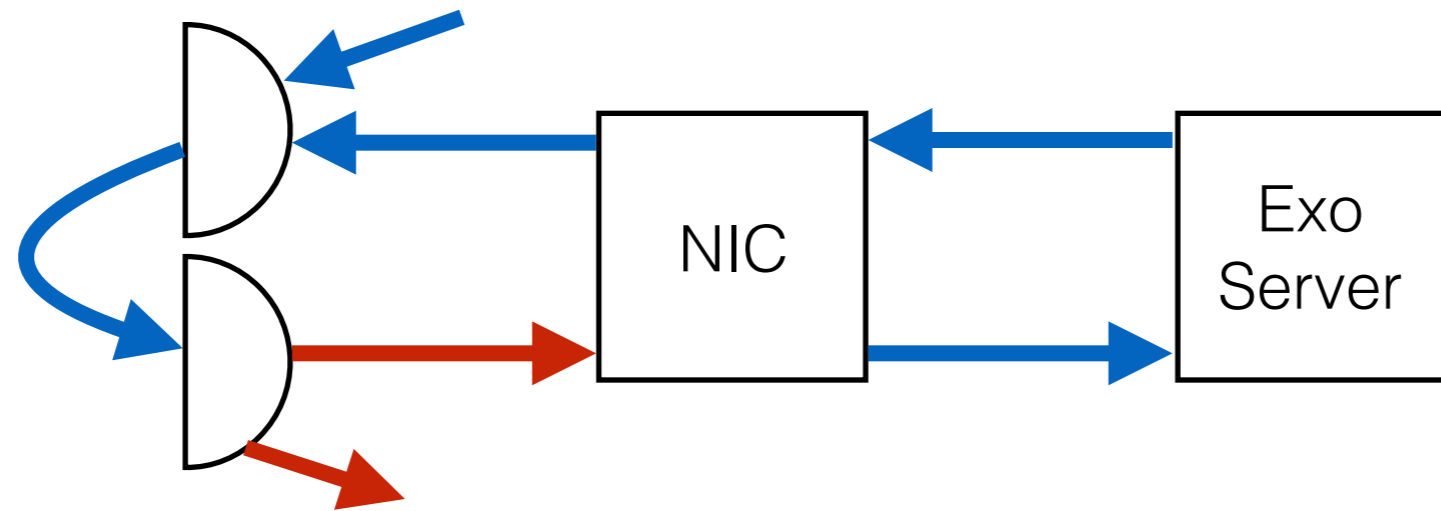
Artist 7124FX

Exablaxe ExaLink50





# But wait, there's more... (latency)

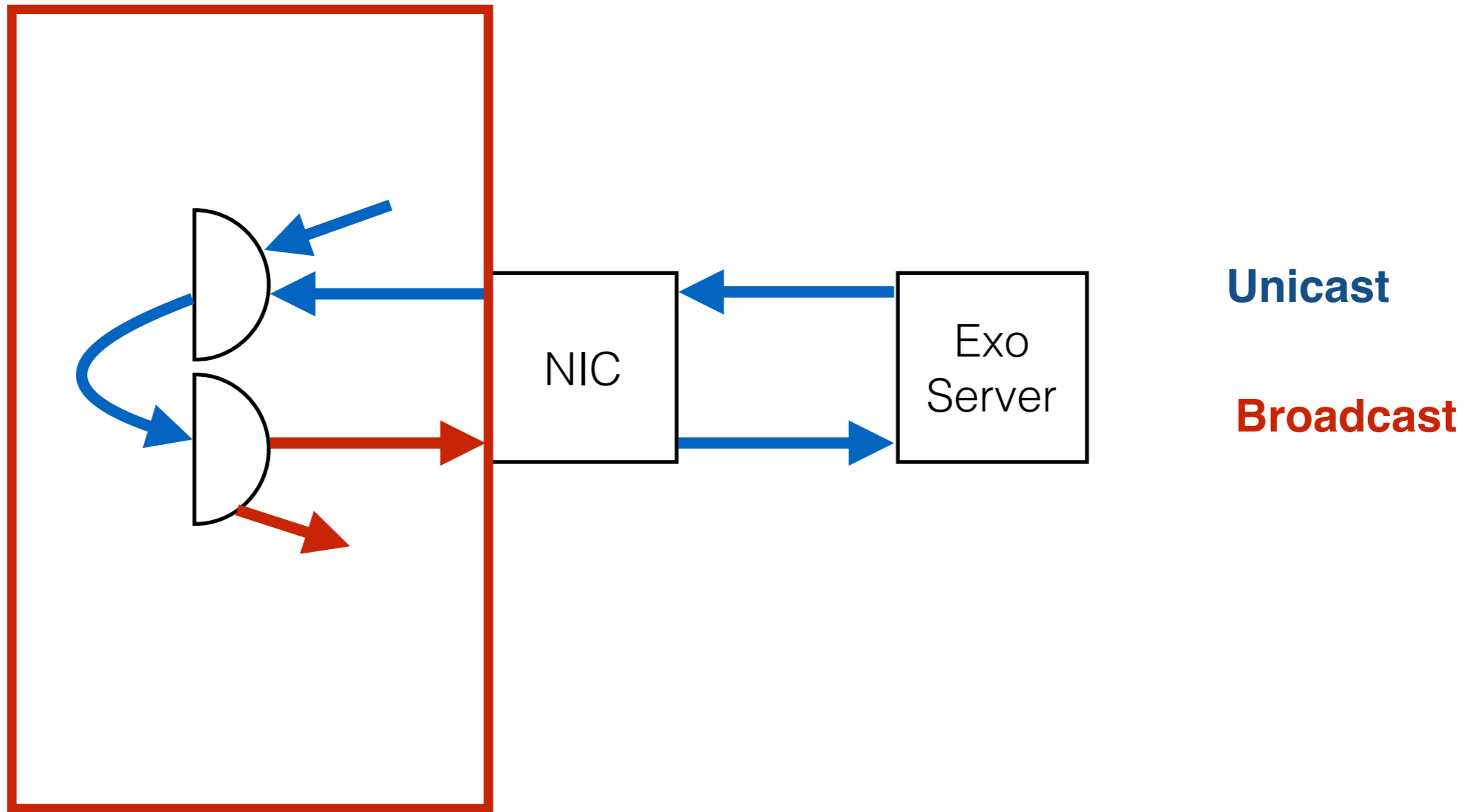


**Unicast**

**Broadcast**



# But wait, there's more... (latency)

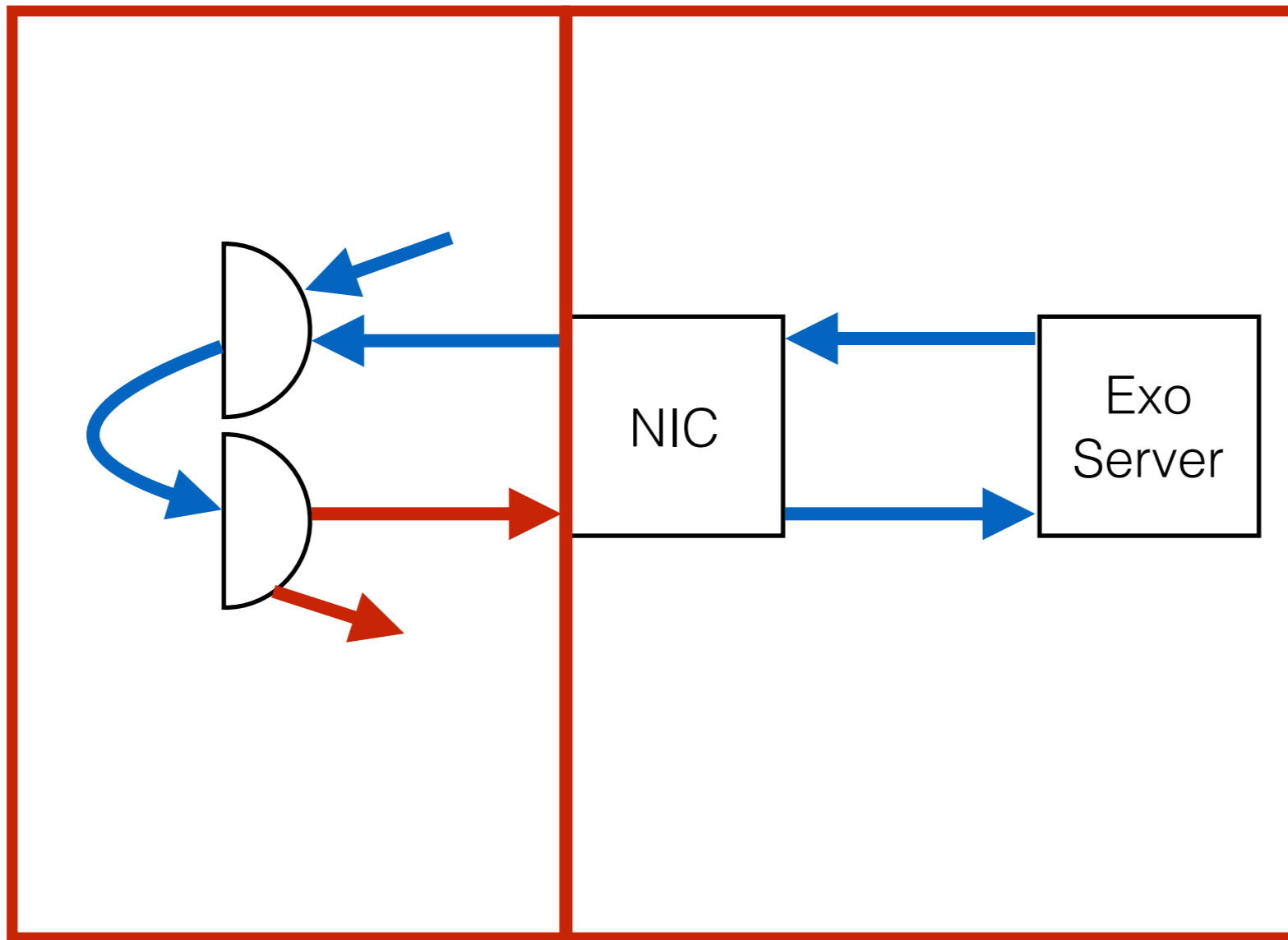


Network Latency

370ns



# But wait, there's more (latency)



**Unicast**

**Broadcast**

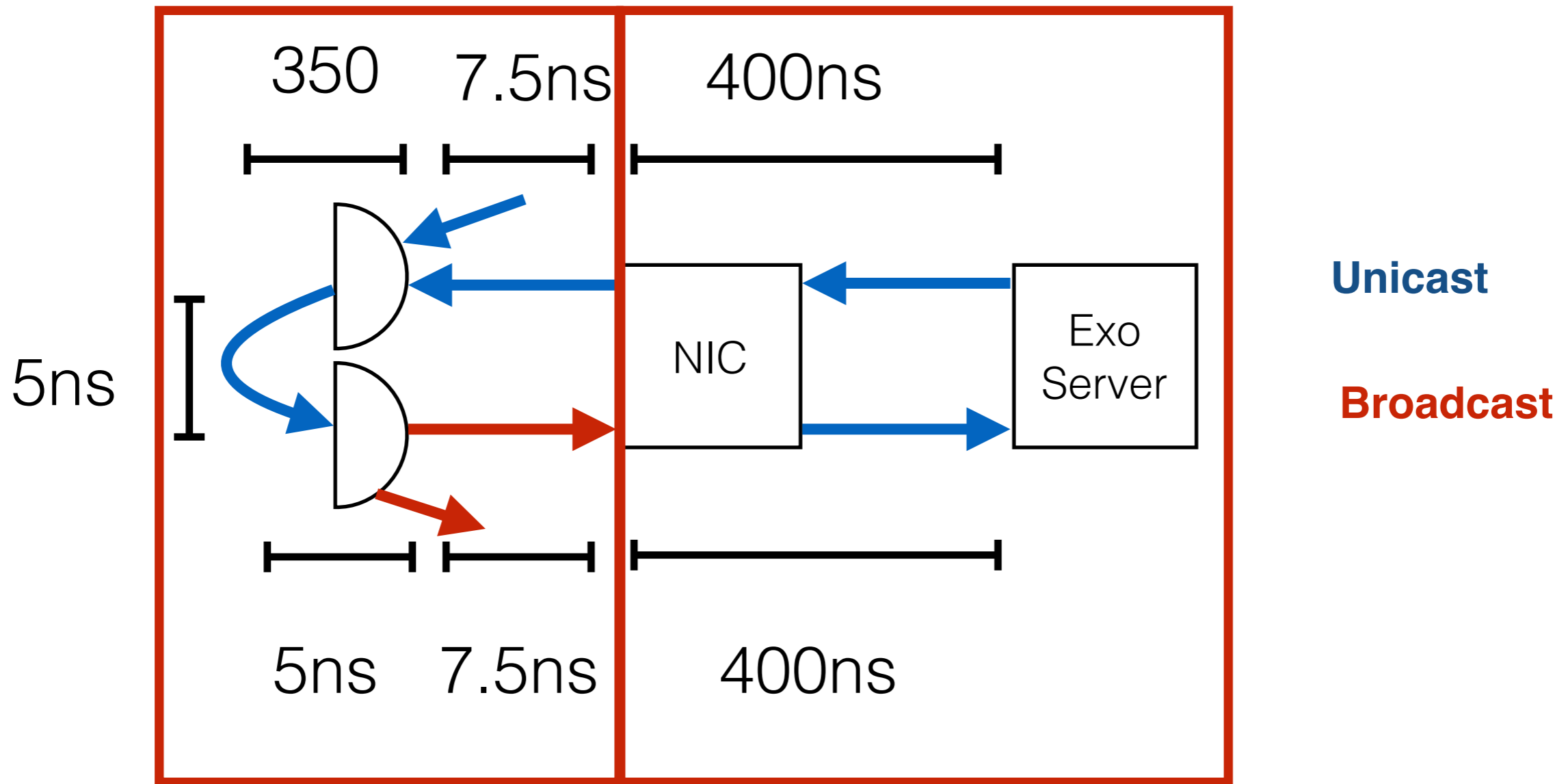
Network Latency

PCIe Latency

370ns



# But wait, there's more... (latency)

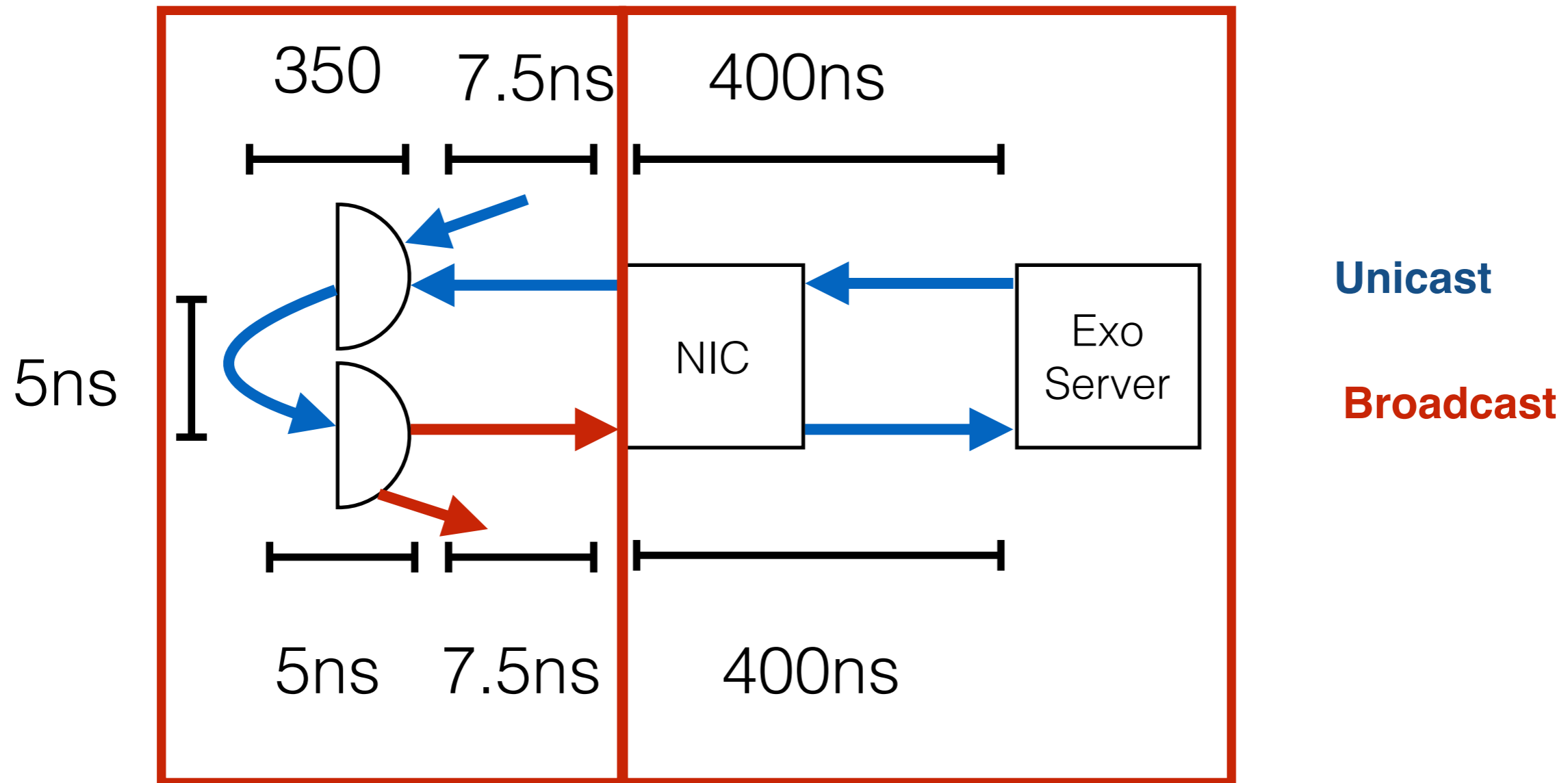


Network Latency

PCIe Latency

370ns

>800ns



- total per-hop-latency  $\approx 1\mu\text{s}$
- $200 \times 1\mu\text{s} \approx 200\mu\text{s}$   
 $\approx 5000$  msgs per host per second.



## 1. Make it go away

- e.g use a Fusion
- e.g use all optical

## 2. Hide it

- Introduce pipelining

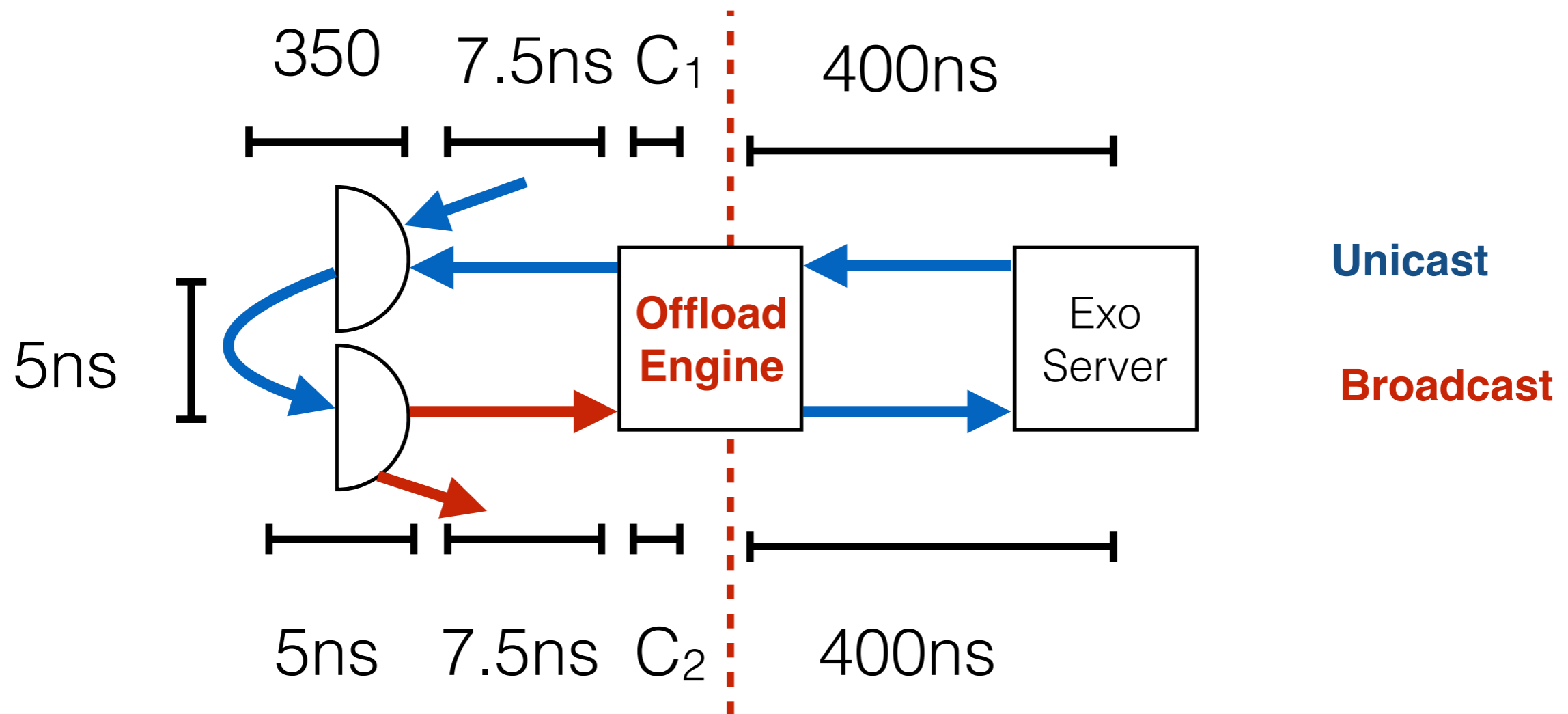


## 1. Make it go away

- e.g use a Fusion
- e.g use all optical

## 2. Hide it

- **Introduce pipelining**

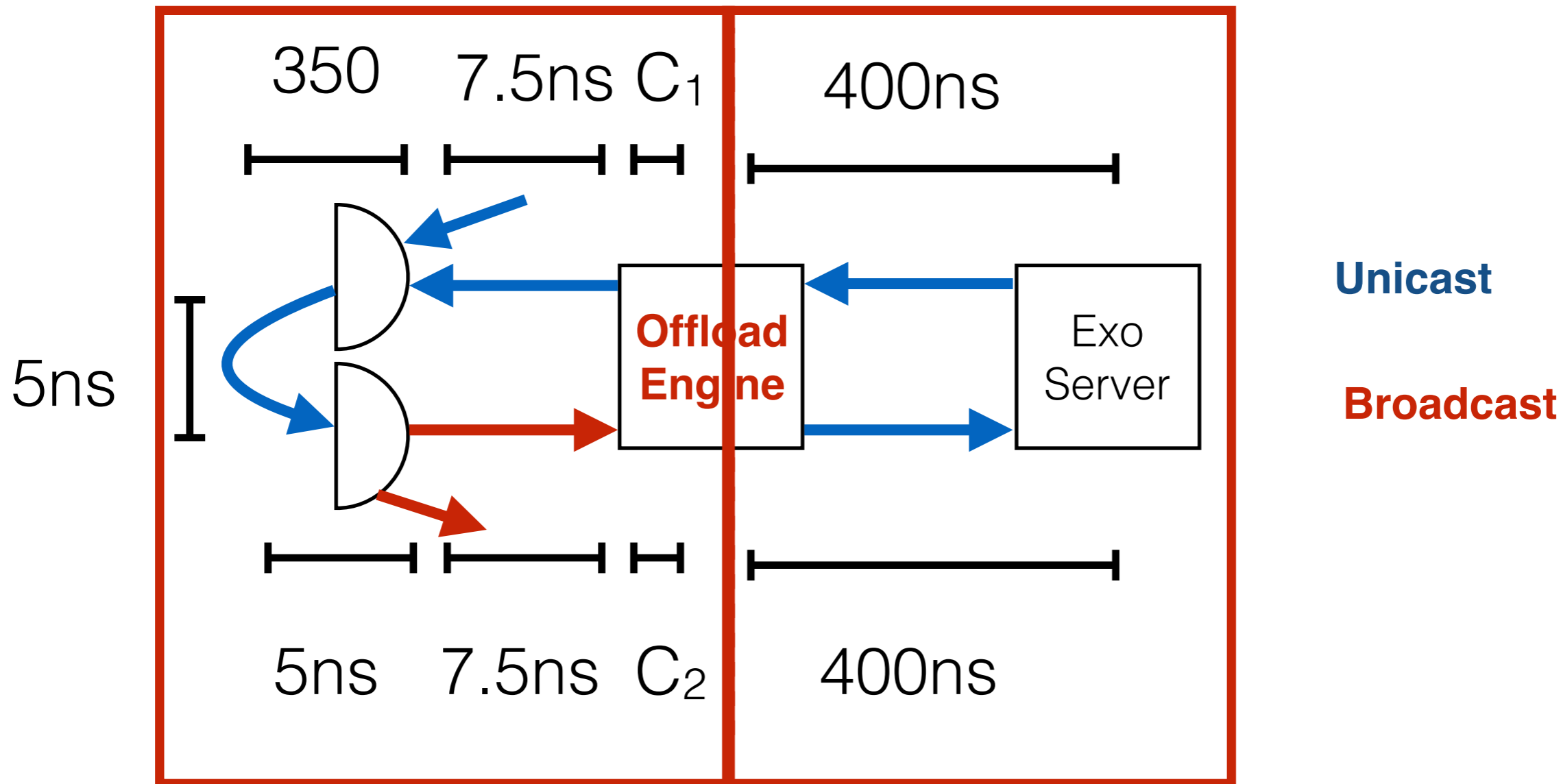


- Introduce pipelining to hide latency
- Host to NIC and NIC to network transmissions run in parallel





# Hiding Latency



- In network latency  $\approx 400\text{ns}$
- Per-hop latency  $\approx 1\mu\text{s}$  (pipelined for 1 message per 400ns)
- $200 \times 400\text{ns} \approx 80\mu\text{s}$   
 $\approx 12,500$  messages per second / per host.



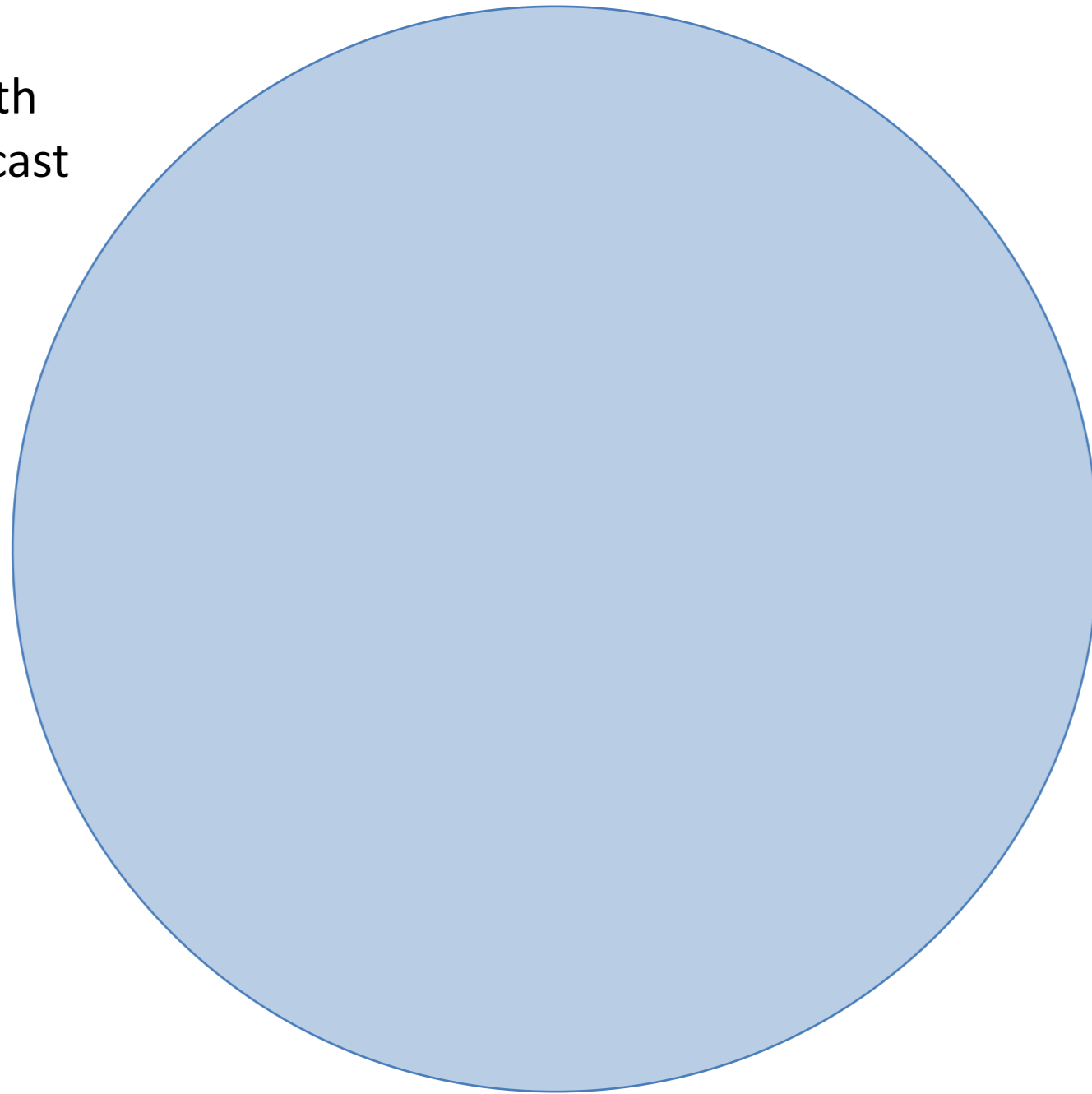
- Failure is very unlikely.
  - No packet loss due to congestion
  - Bit errors on the wire  $1/10^{12}$  ...  $1/10^{14}$
  - Partitioning - Single chip? Can it half fail?  
Byzantine?
    - $1/10^{20}$  ??
  - Partitioning - requires n-way x m - correlated failure
    - $(1/10^{12})^n * m$
- Therefore, optimise for the common case



- Protocol is implemented correctly on all hosts
- Hosts may stop, crash, restart, loose packets or become partitioned
- There is a fixed upper number of hosts ( $n$ ) in the network. Host may leave and come back, but more hosts may not be added.

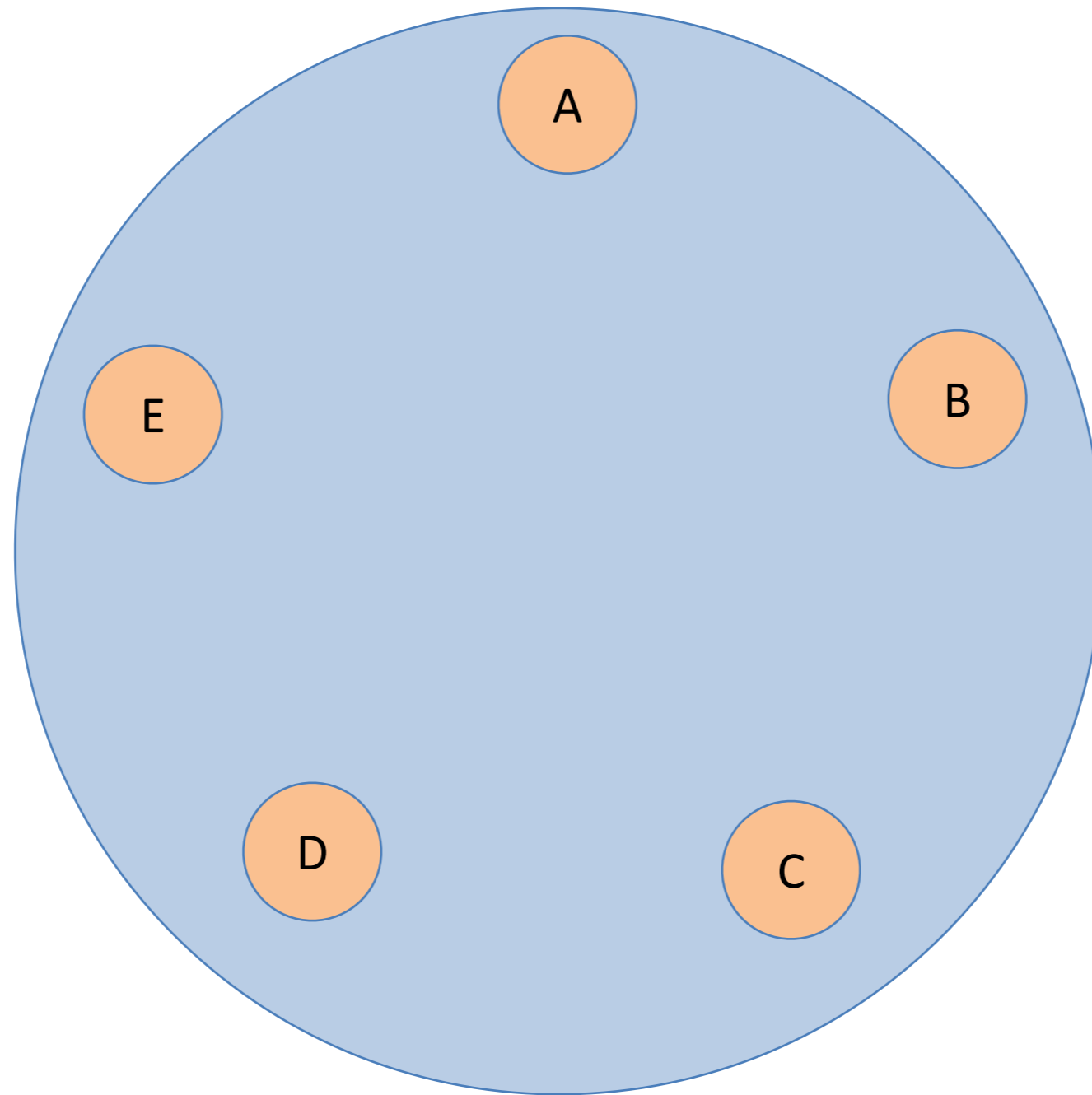


Assume a network with  
fast and cheap broadcast  
messaging



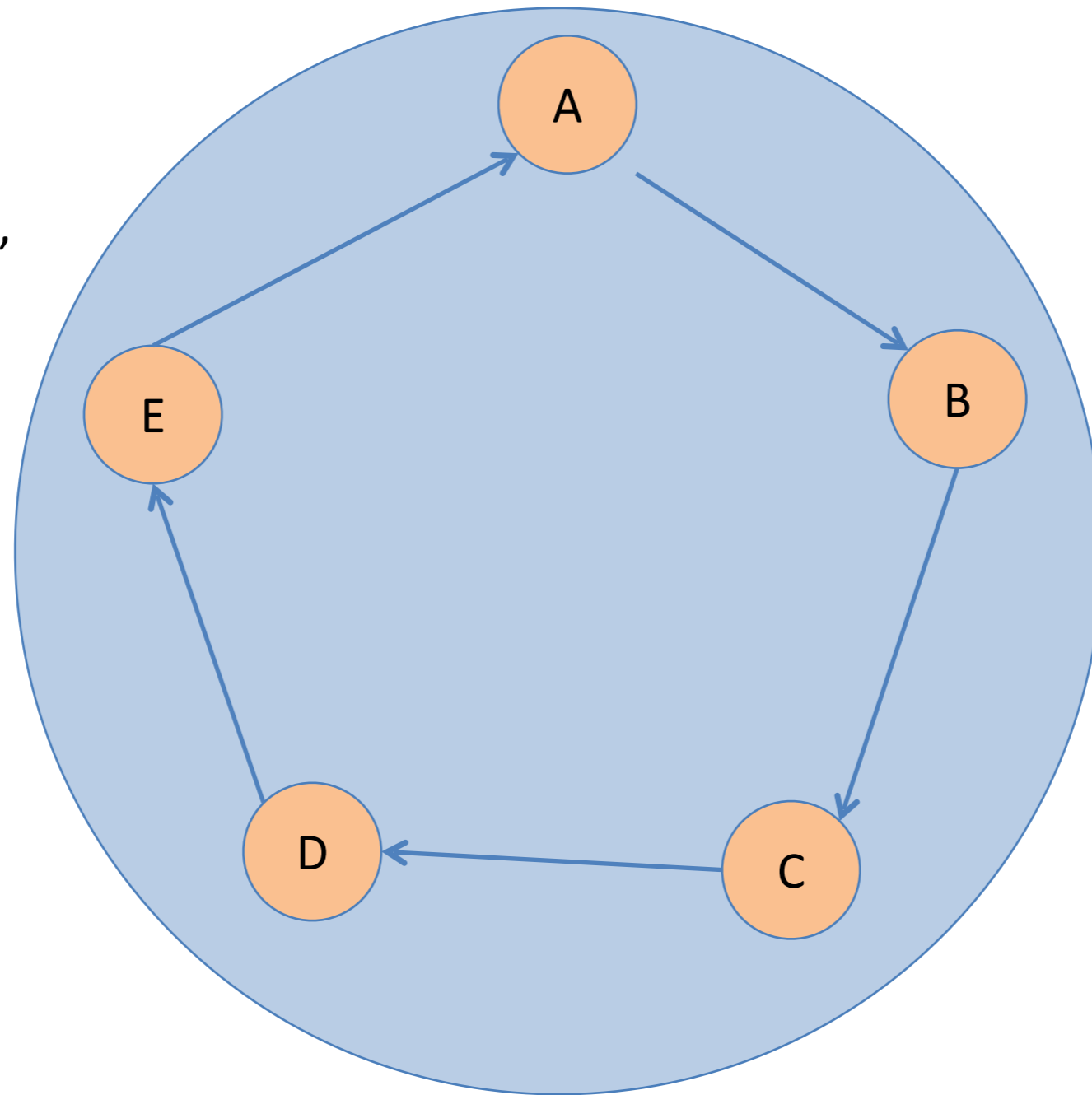


And some maximum  
number of hosts  $n$  (5)  
labeled **A-E**





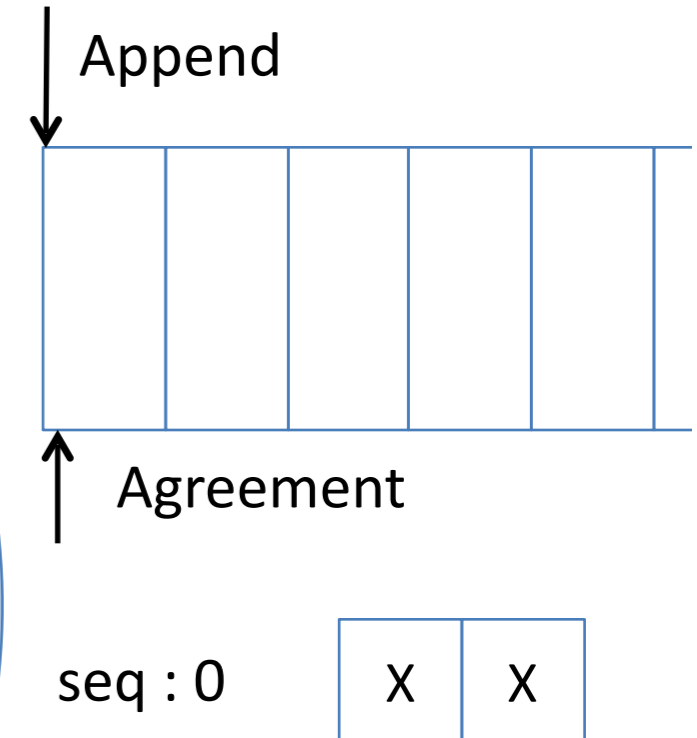
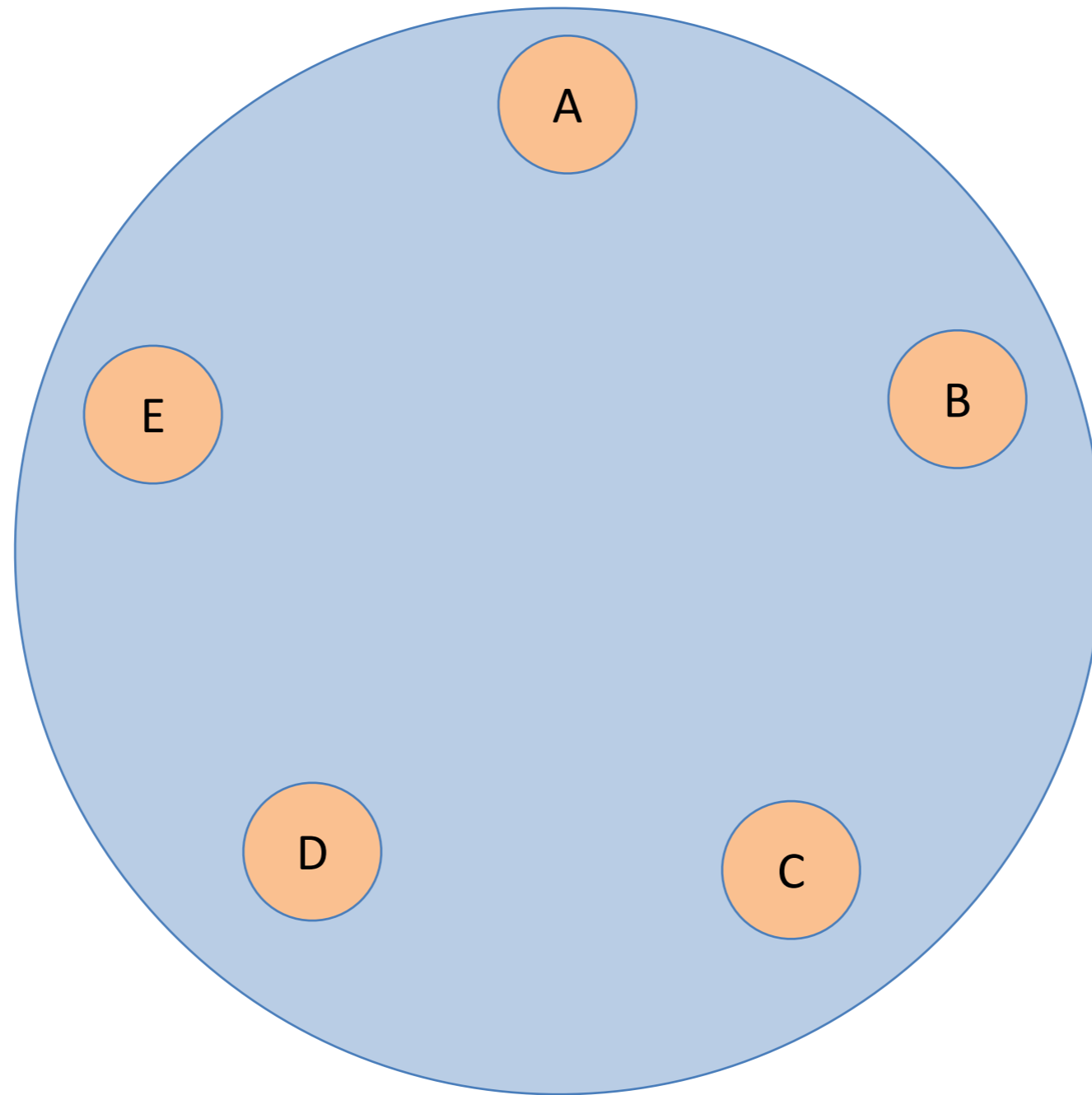
And some maximum number of hosts  $n$  (5) labeled **A-E**, in a fixed, predetermined order





# Anatomy of the Exo Protocol

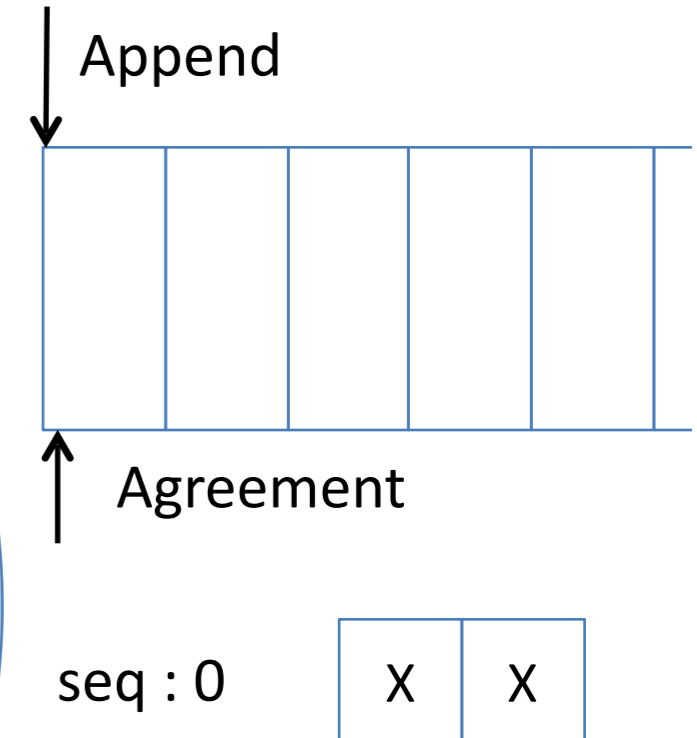
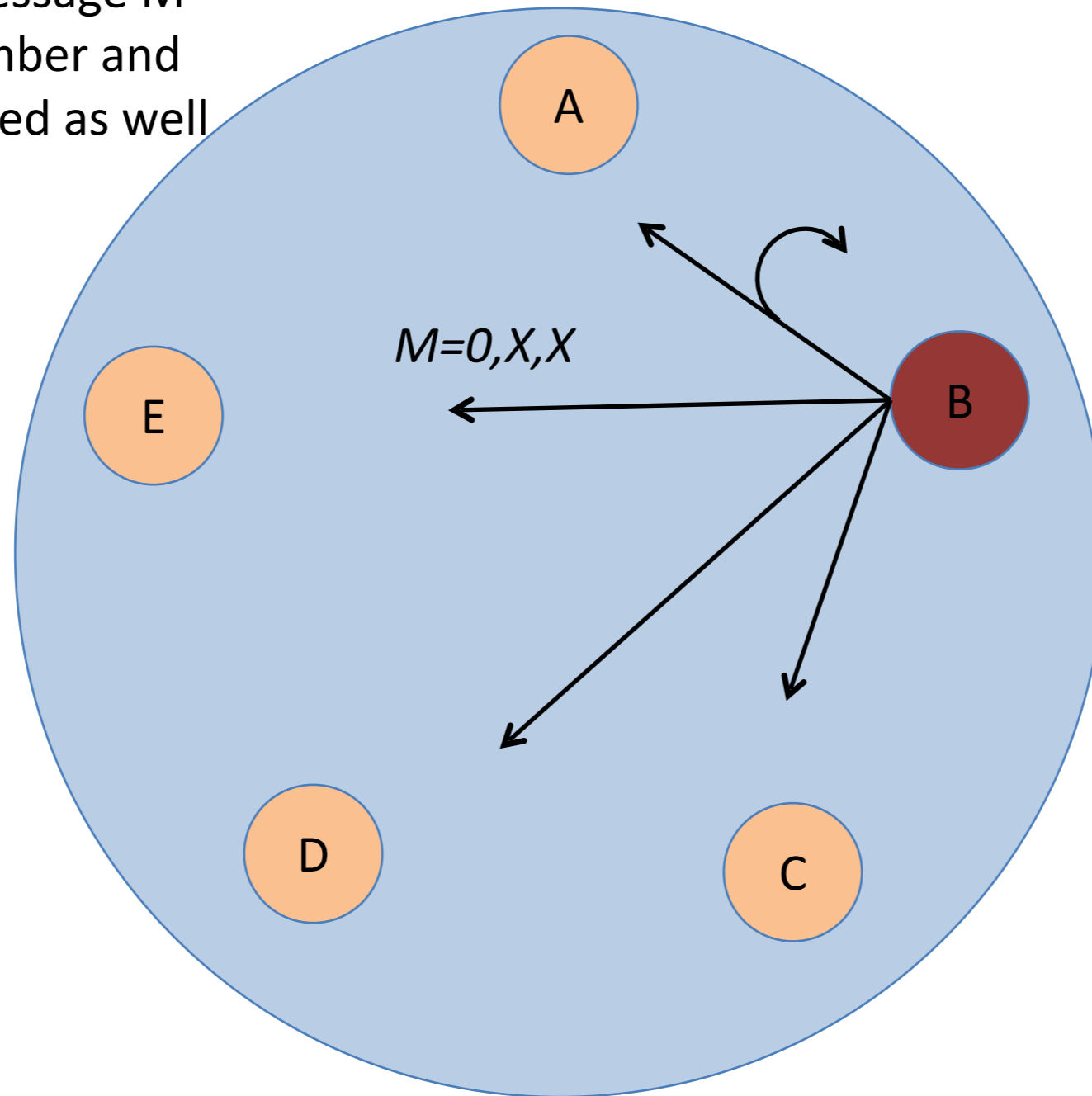
- Each host keeps the following state:
- append only log
  - sequence number
  - $n/2$  message history





# Anatomy of the Exo Protocol

Host *B* broadcasts a message *M* with the sequence number and message history included as well as any data (optional)

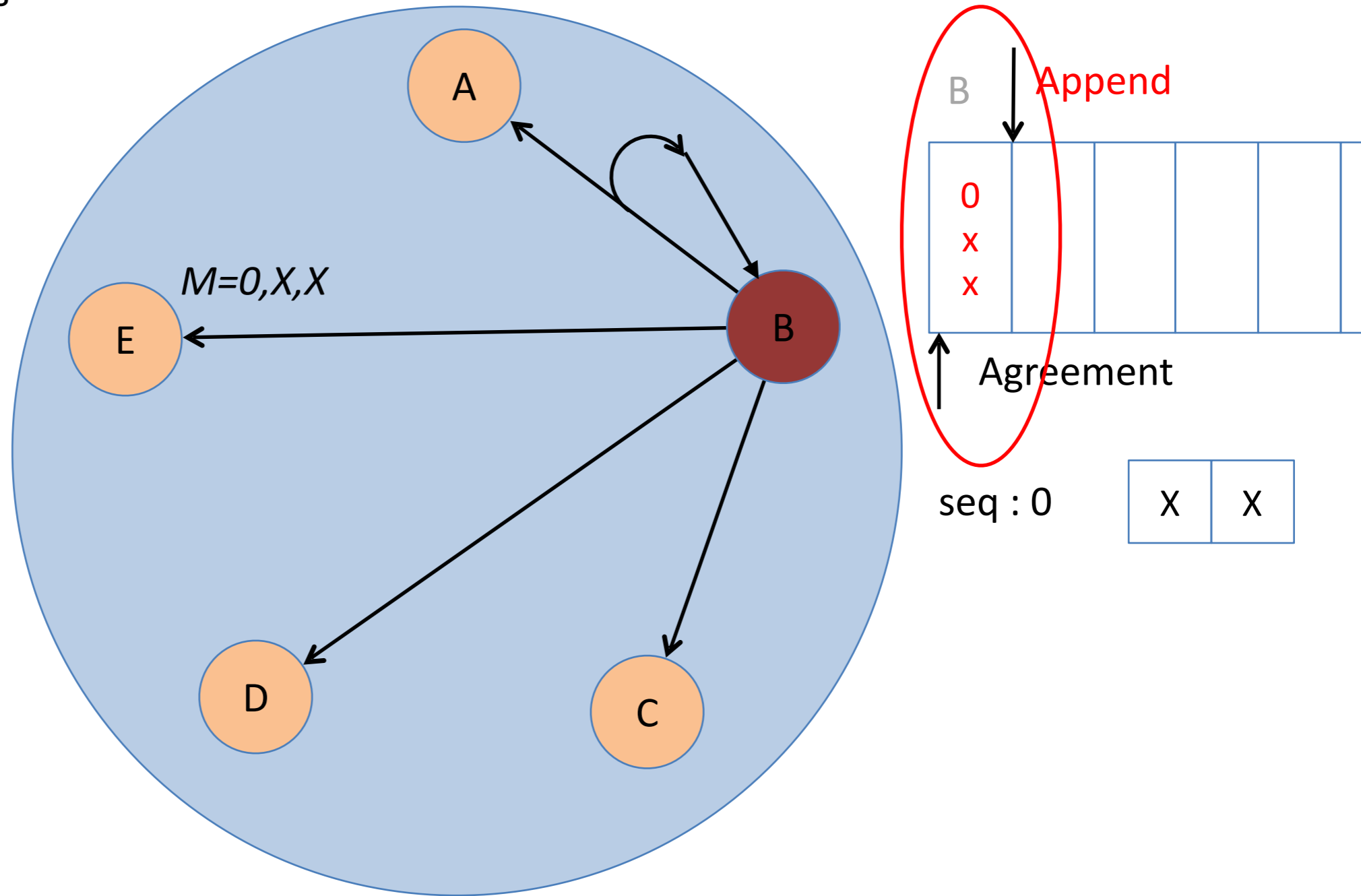






# Anatomy of the Exo Protocol

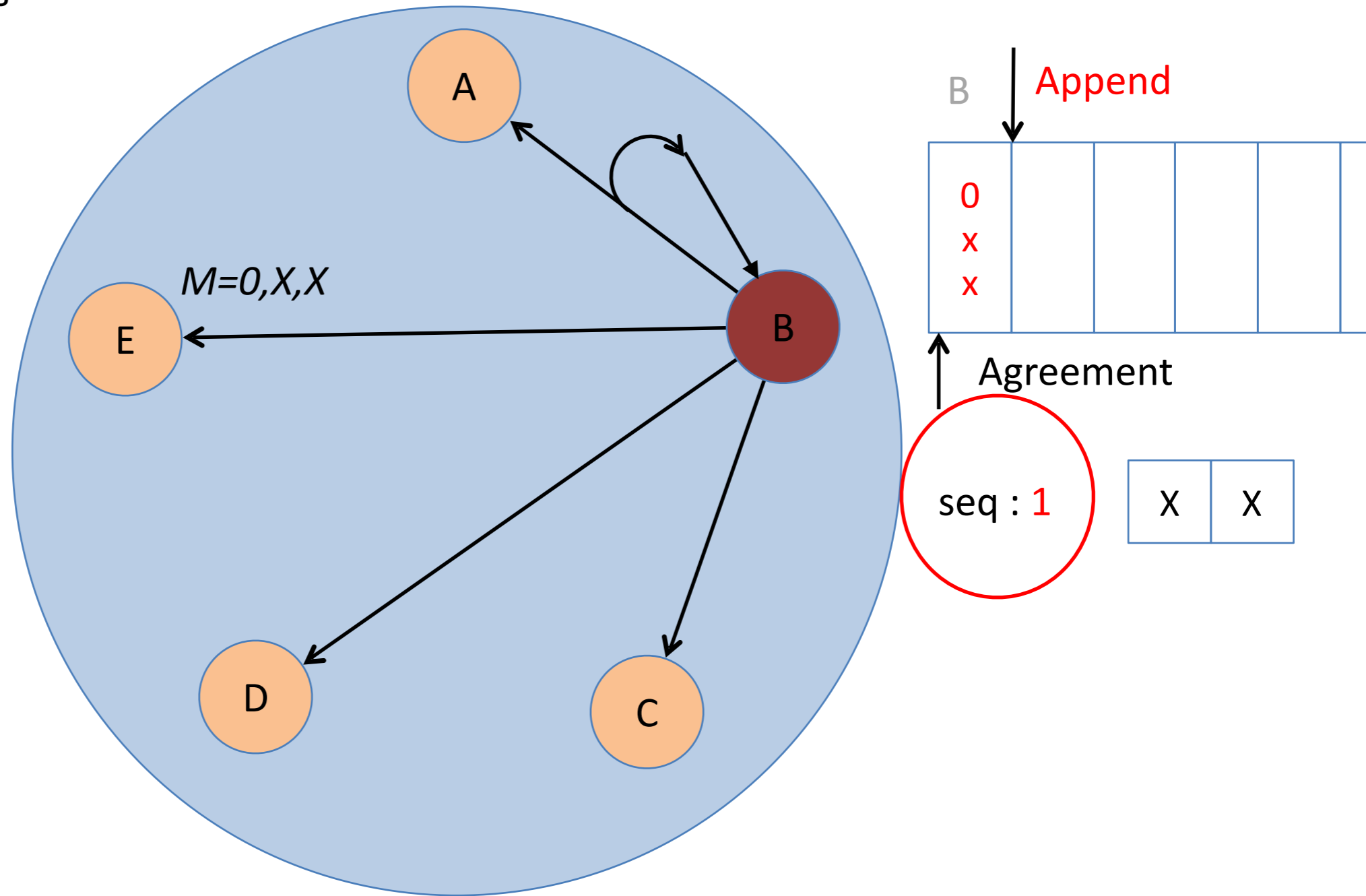
Message from B arrives at all hosts, causes a:  
**-log append**





Message from B arrives at all hosts, causes a:

- log append
- sequence number update

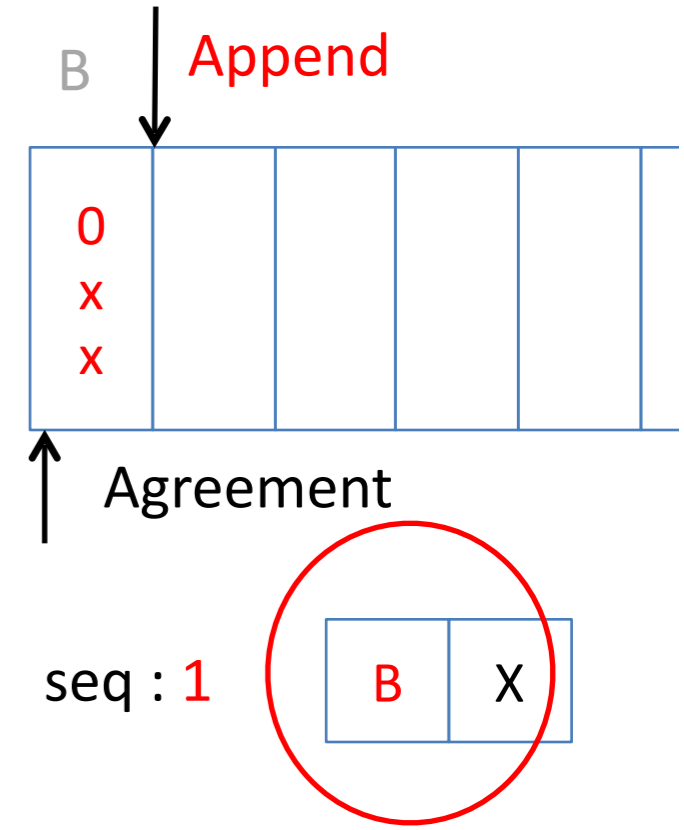
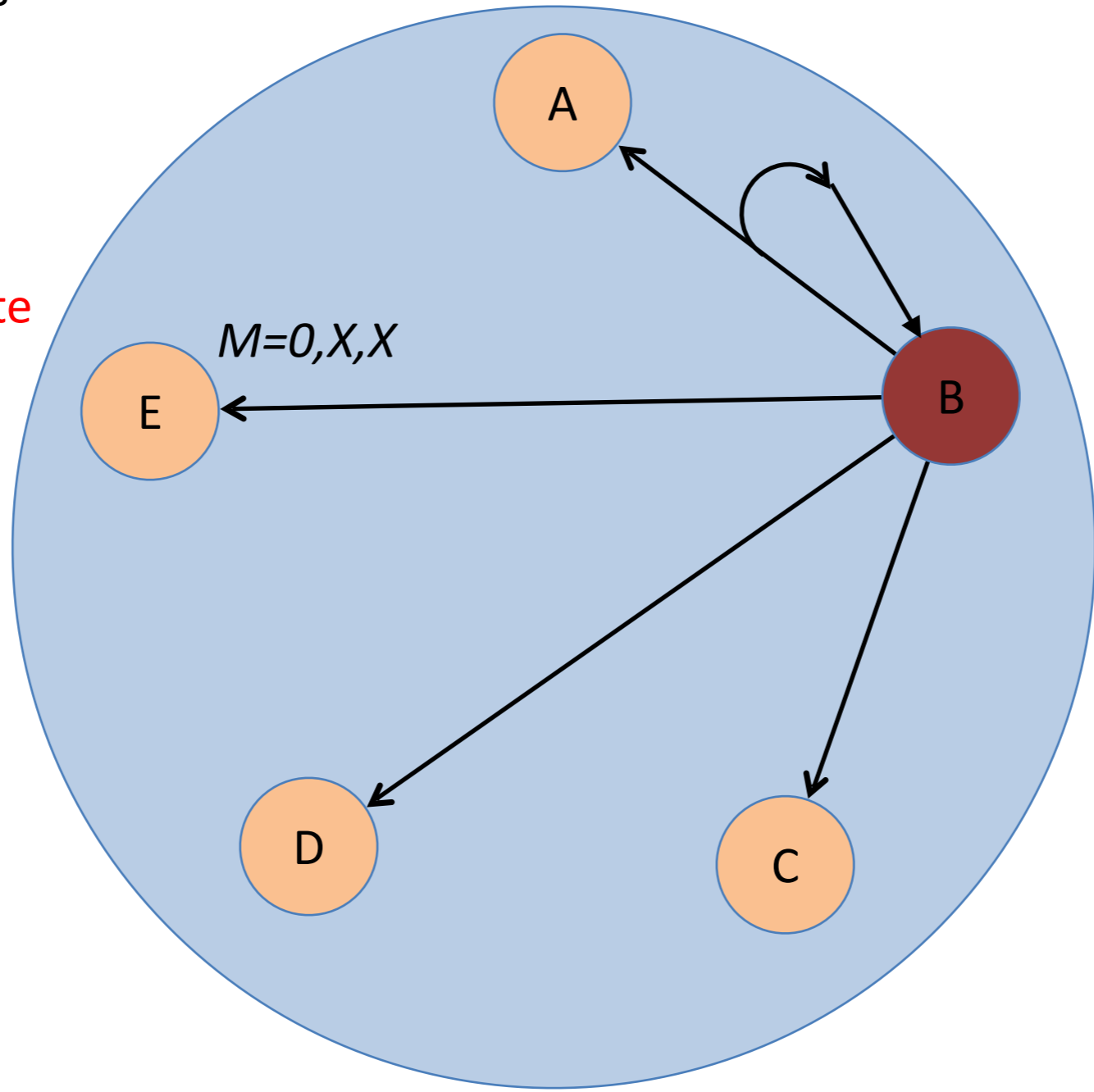




# Anatomy of the Exo Protocol

Message from B arrives at all hosts, causes a:

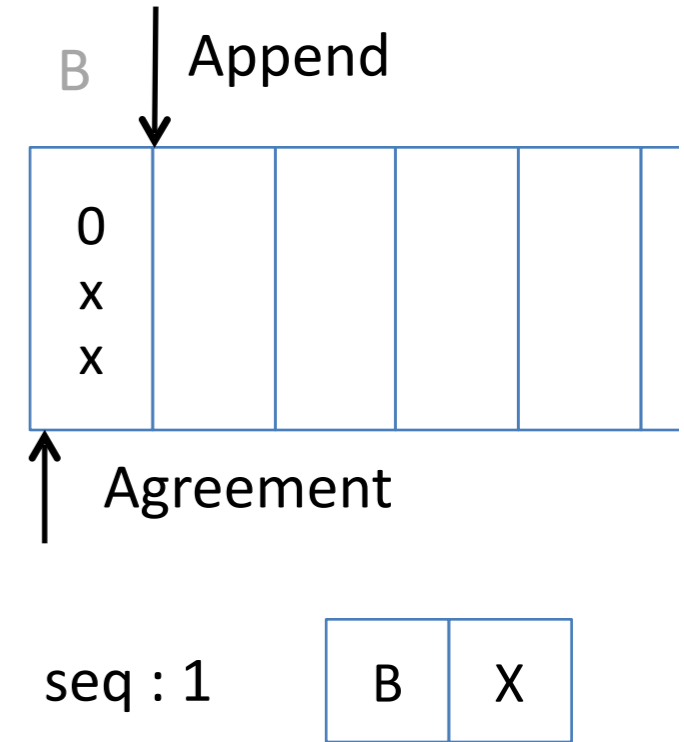
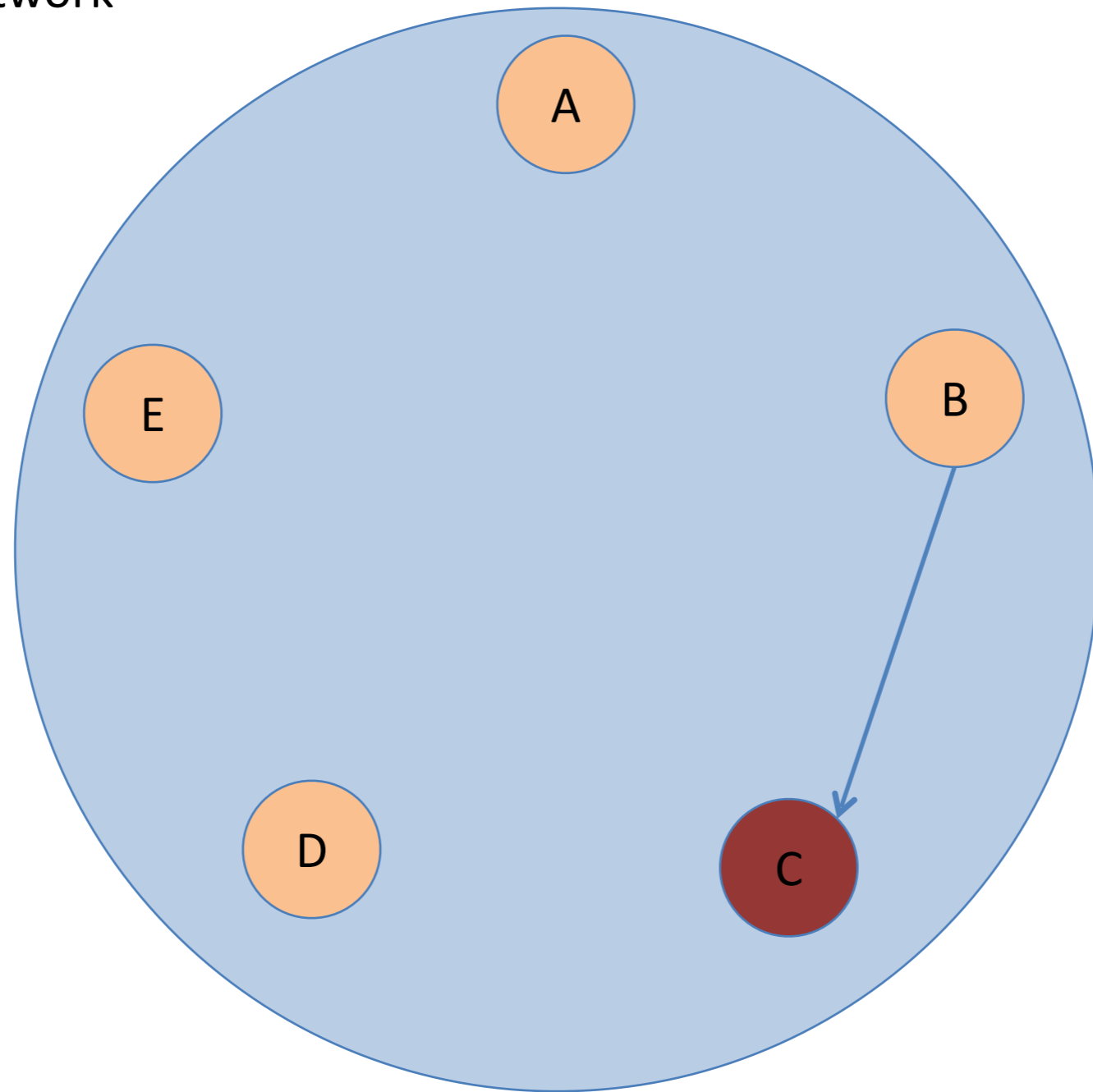
- log append
- sequence number update
- message history update





# Anatomy of the Exo Protocol

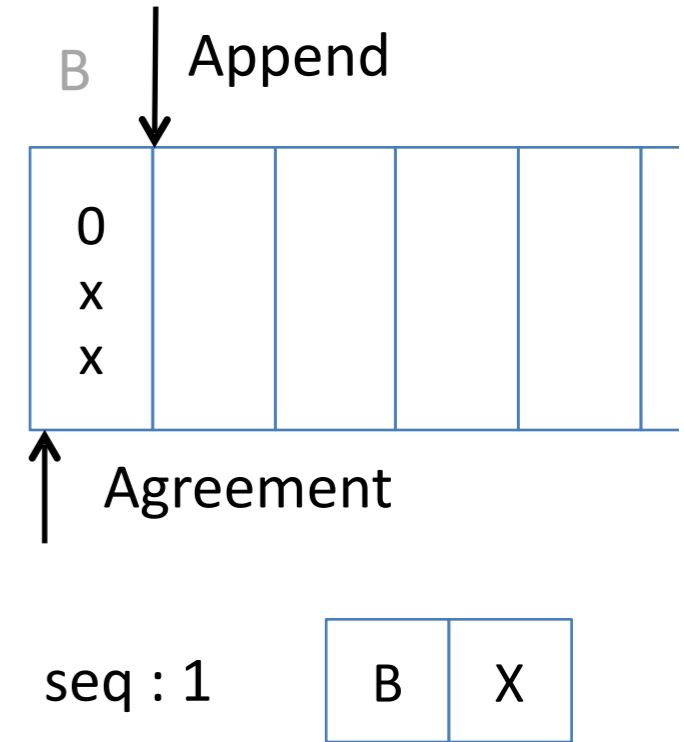
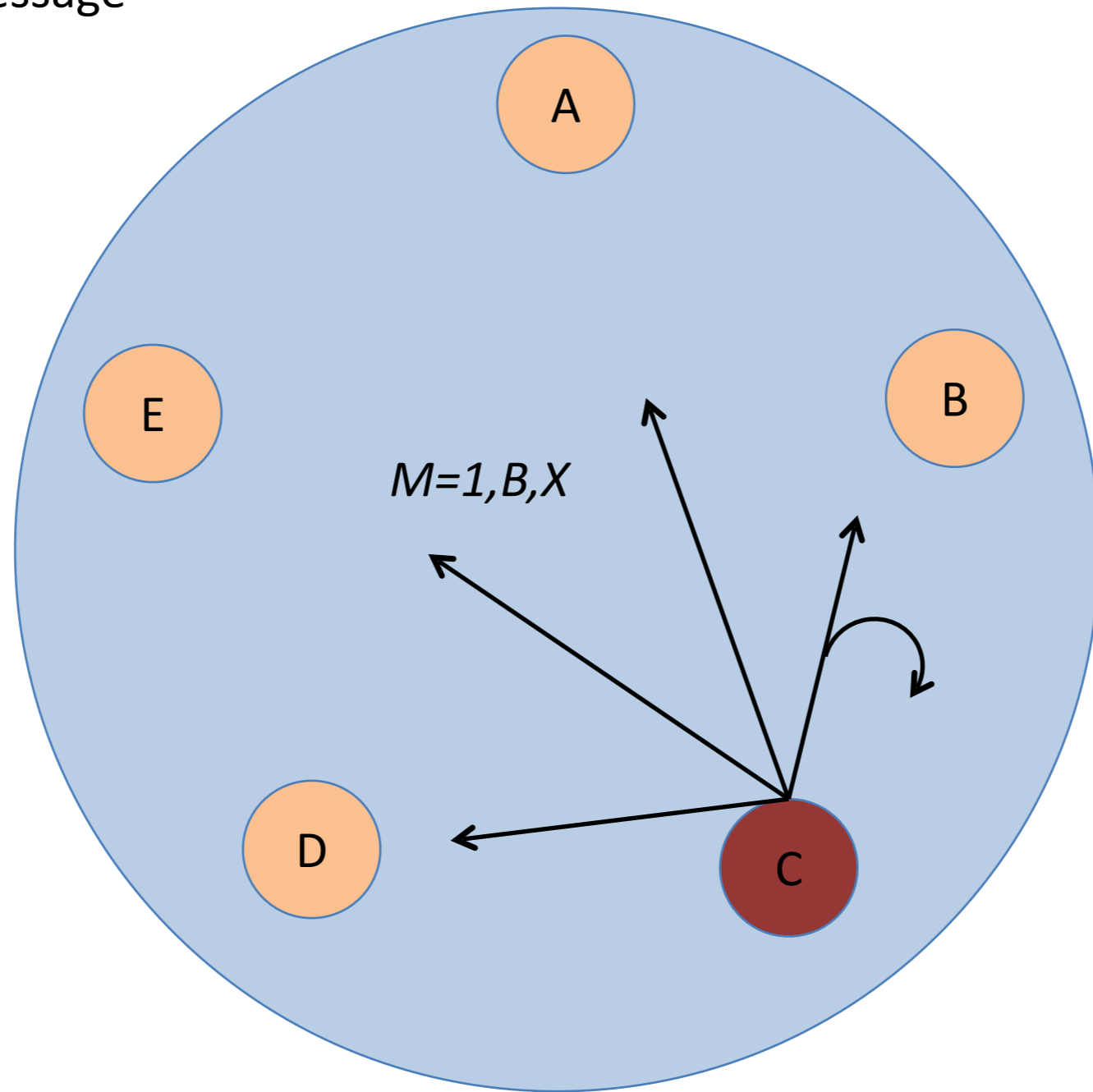
Host **C** now has the network





# Anatomy of the Exo Protocol

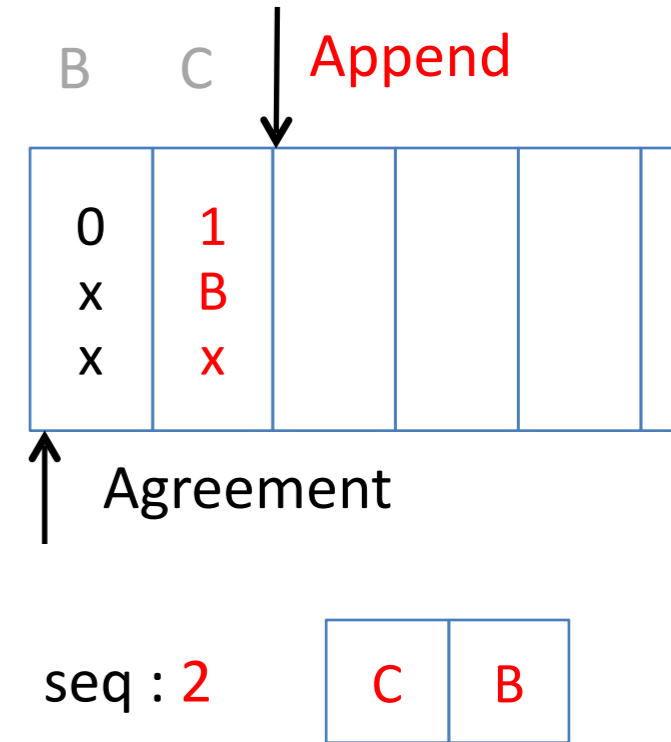
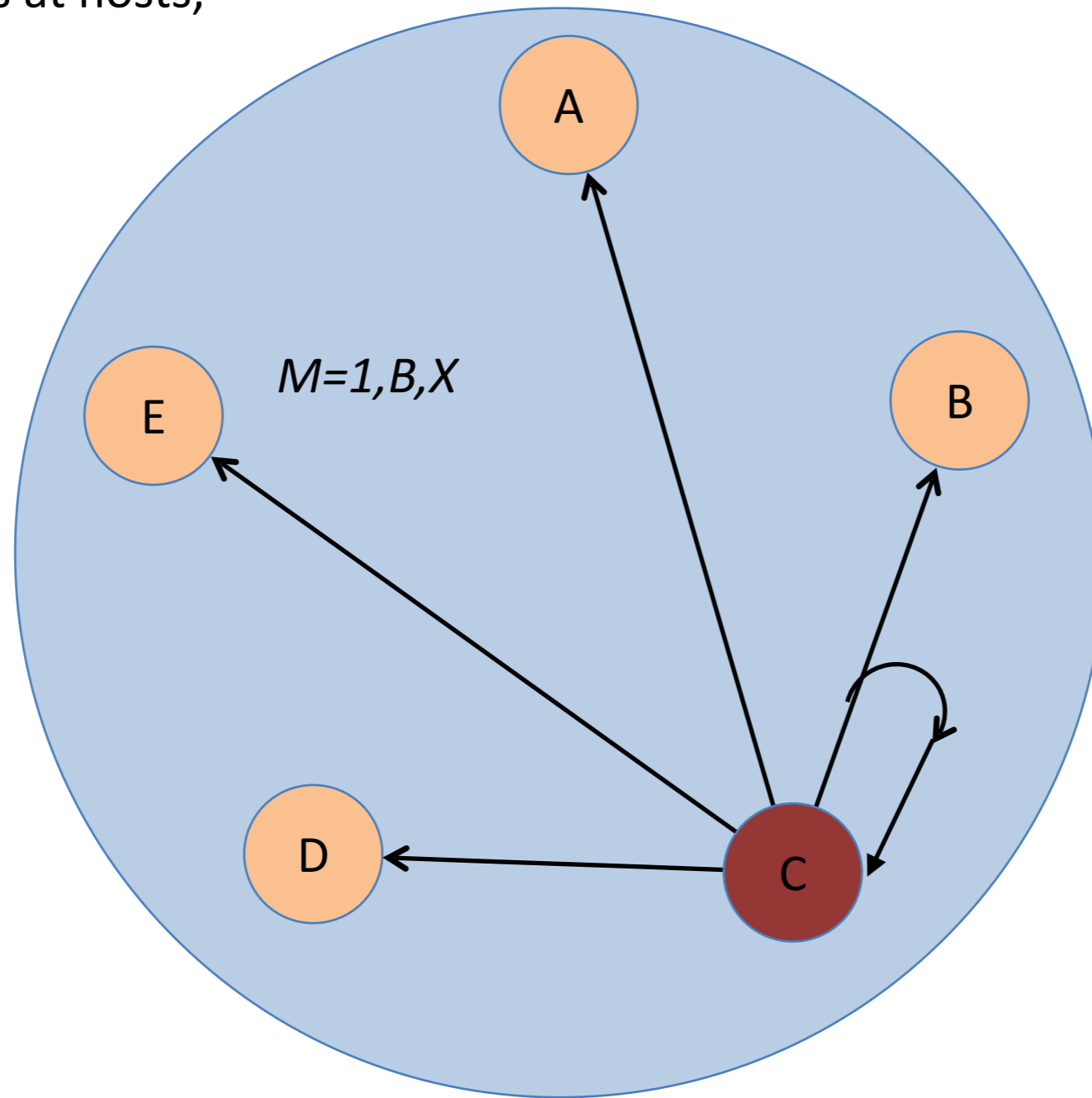
Host C broadcasts a message





# Anatomy of the Exo Protocol

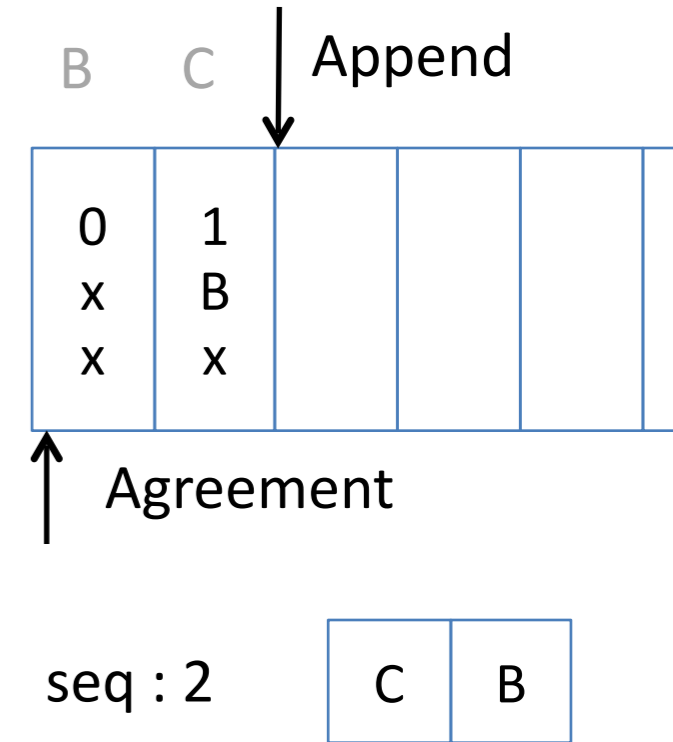
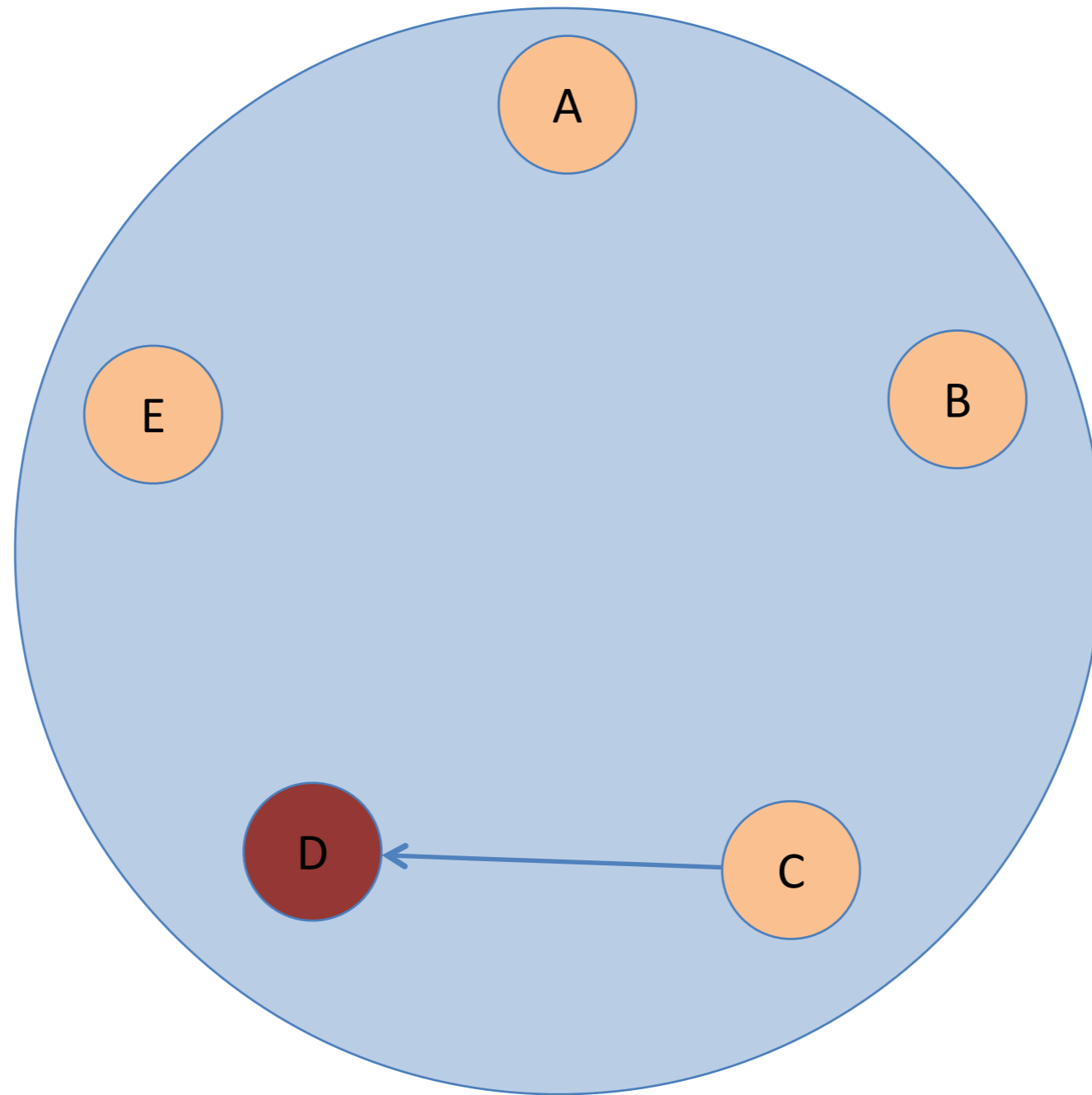
Message from C arrives at hosts, causing **update**





# Anatomy of the Exo Protocol

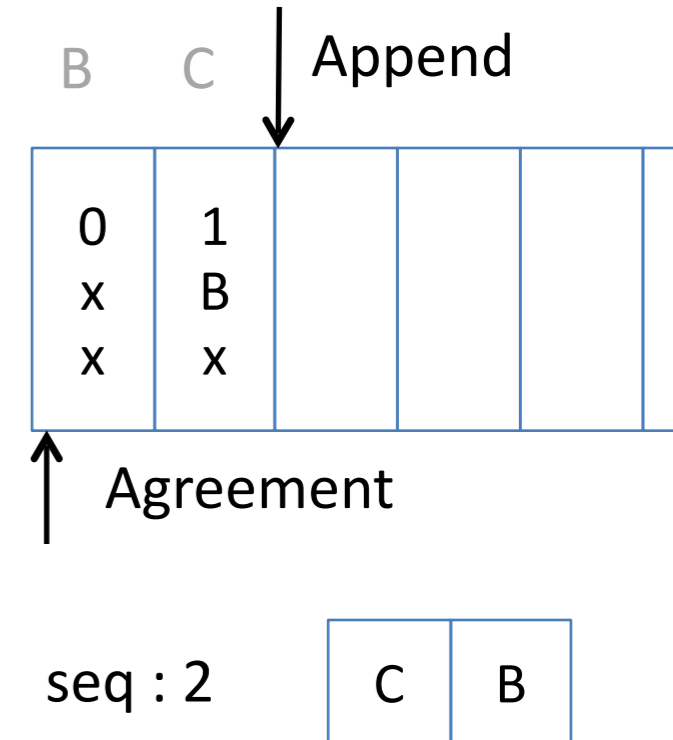
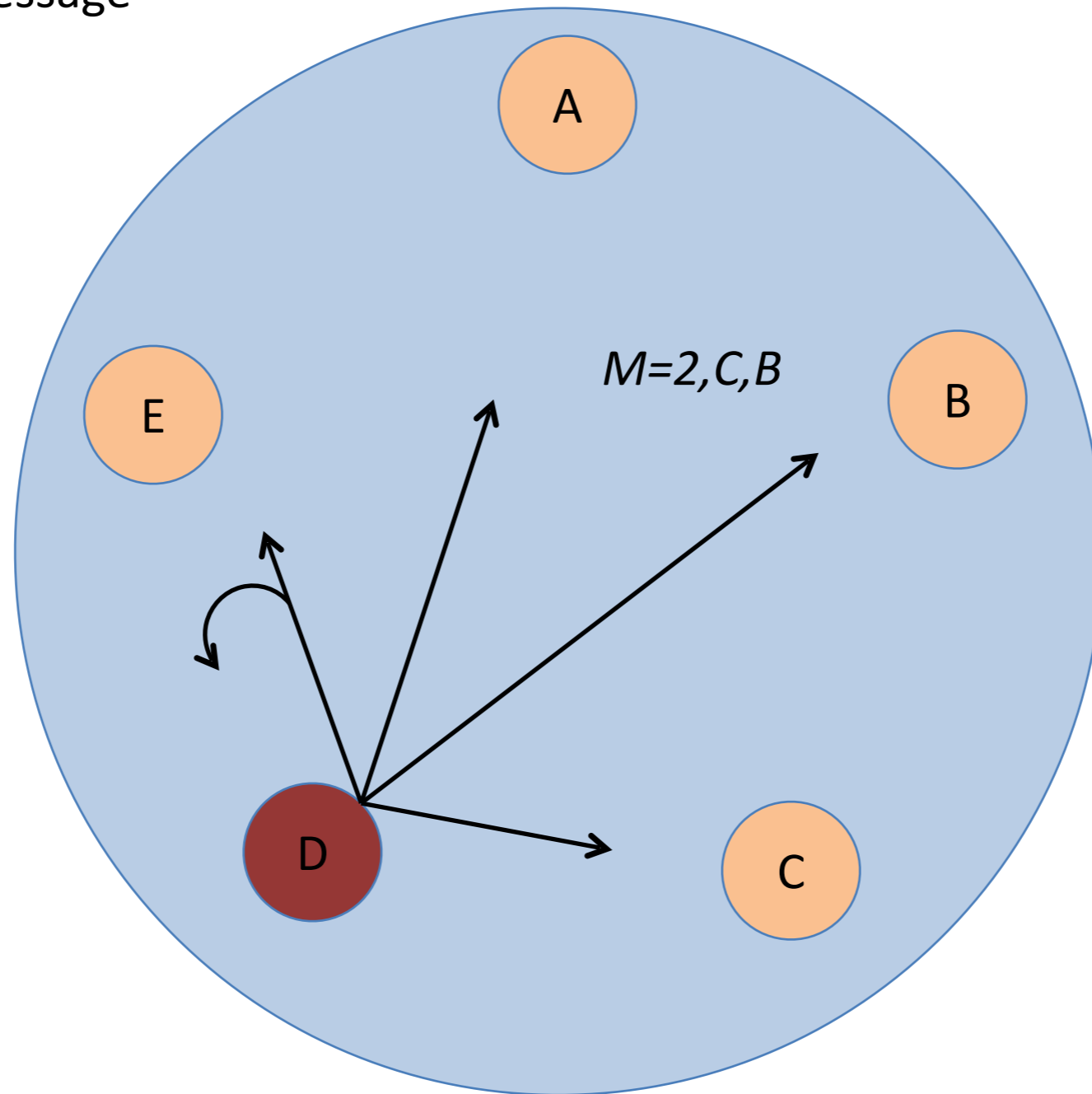
D now has the token





# Anatomy of the Exo Protocol

Host *D* broadcasts a message

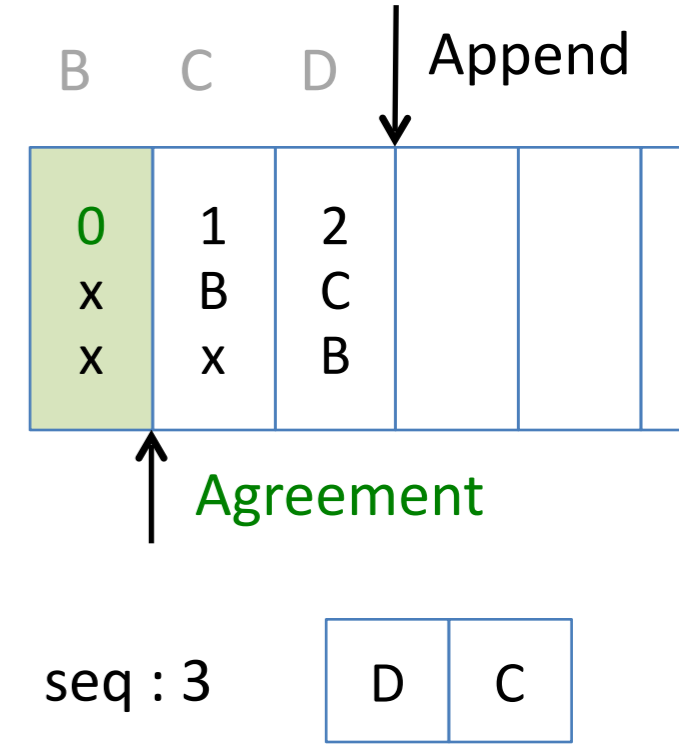
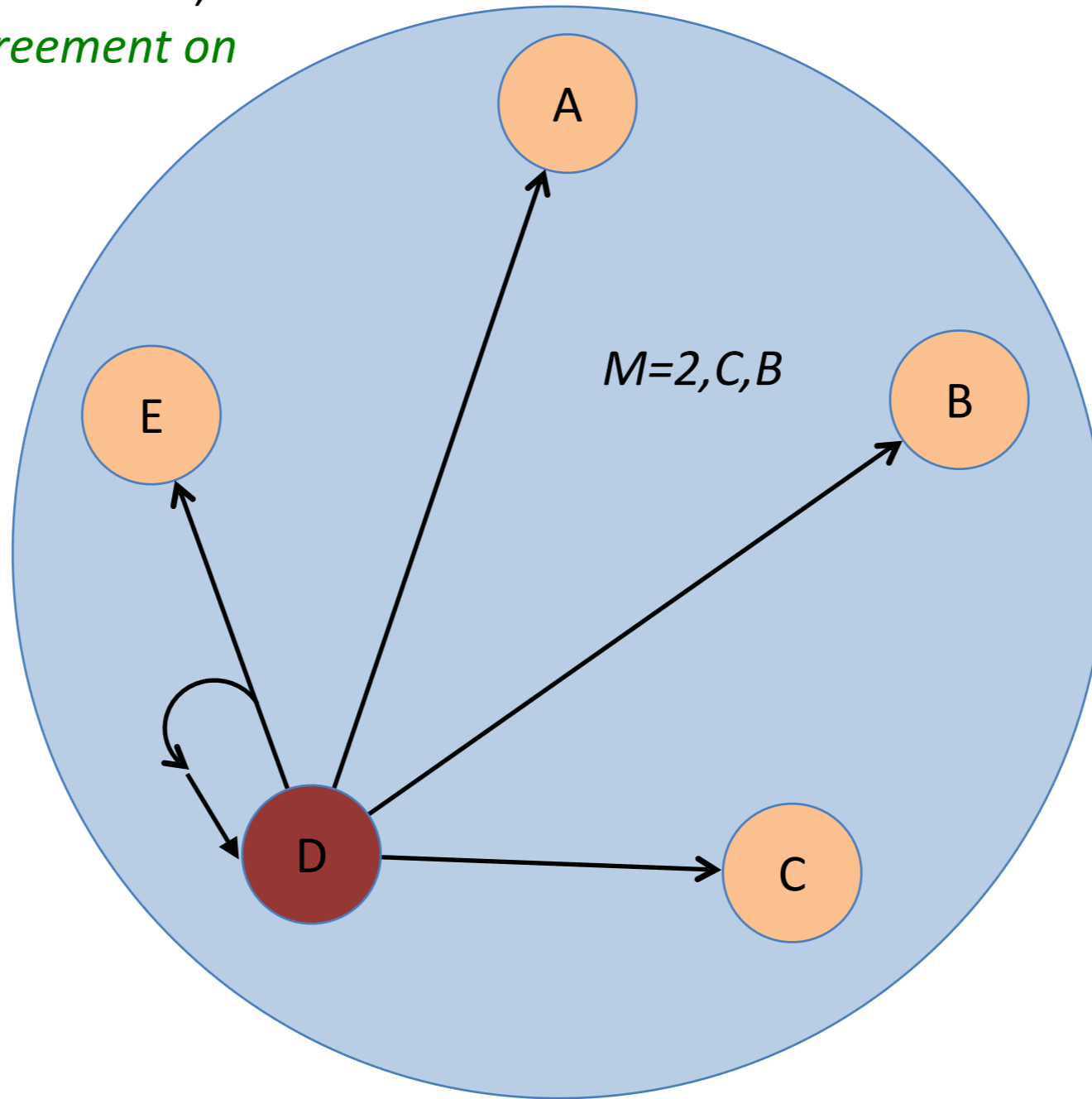






# Anatomy of the Exo Protocol

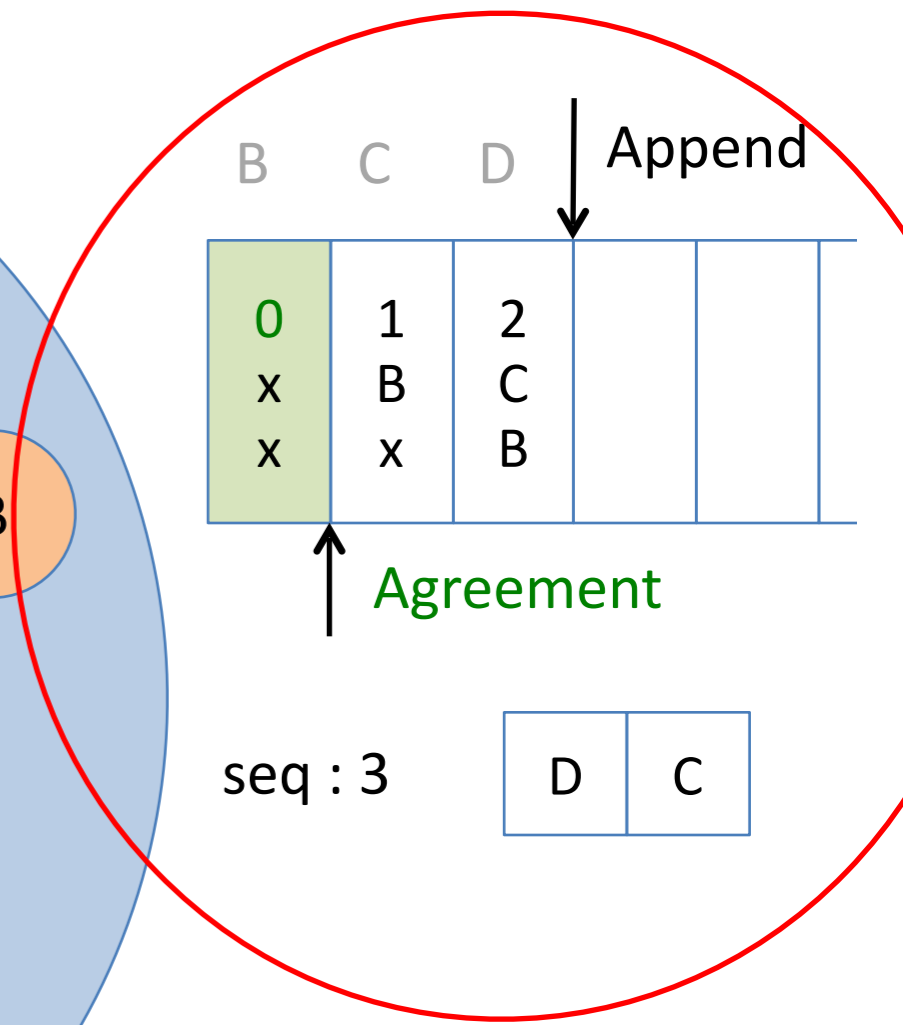
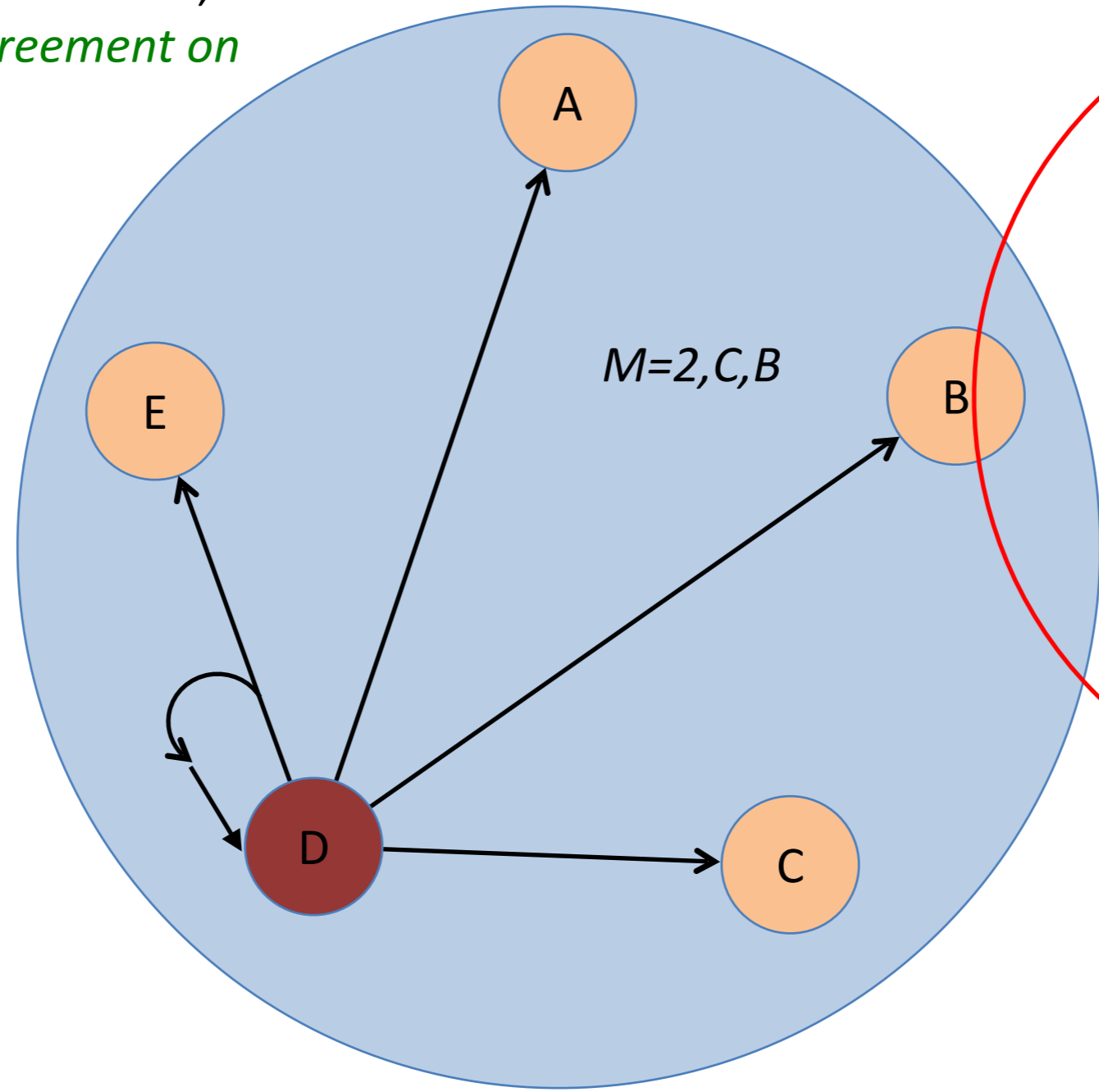
Message from D arrives at hosts, causing *update and agreement on sequence number 0*

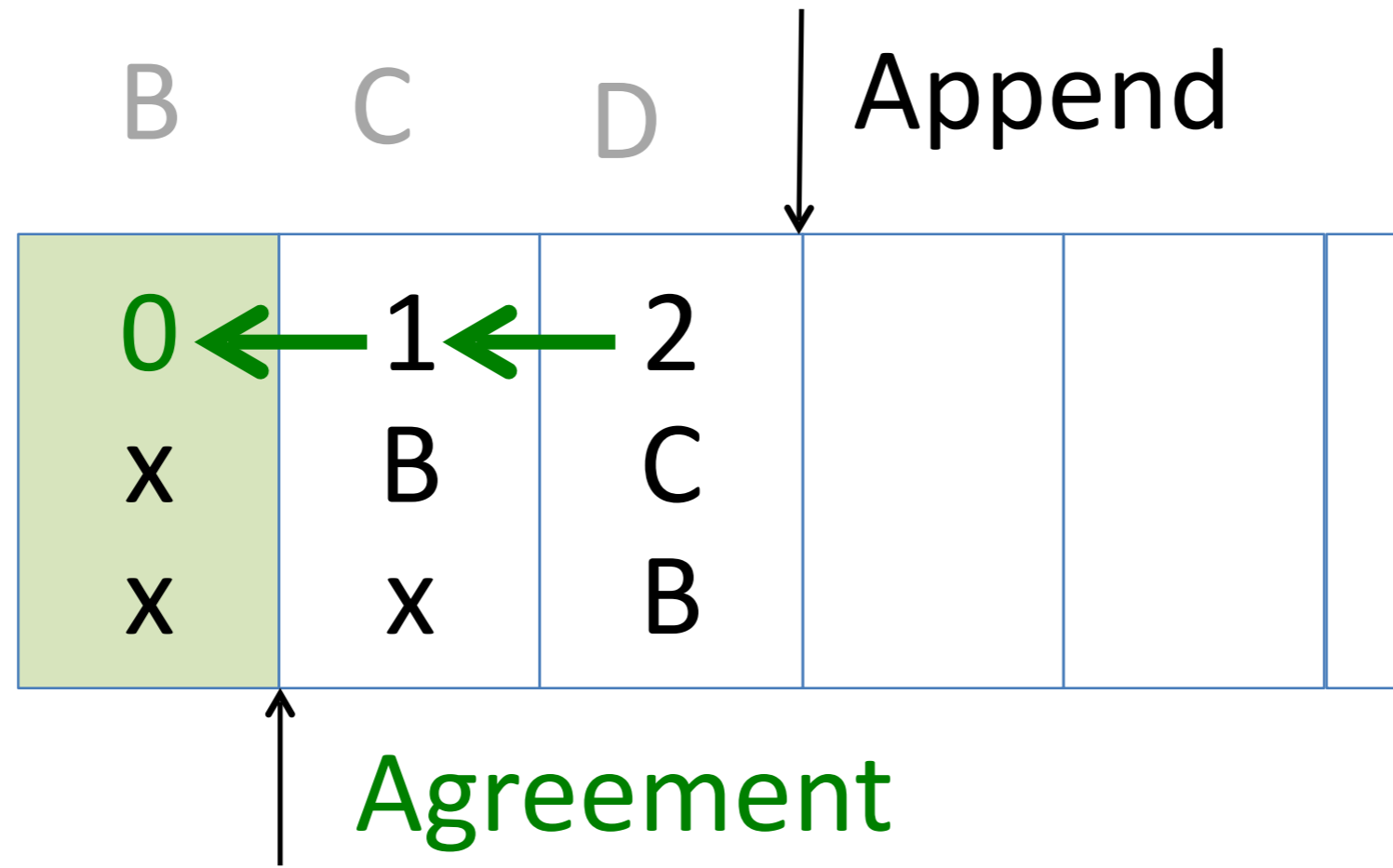




# Zooming in...

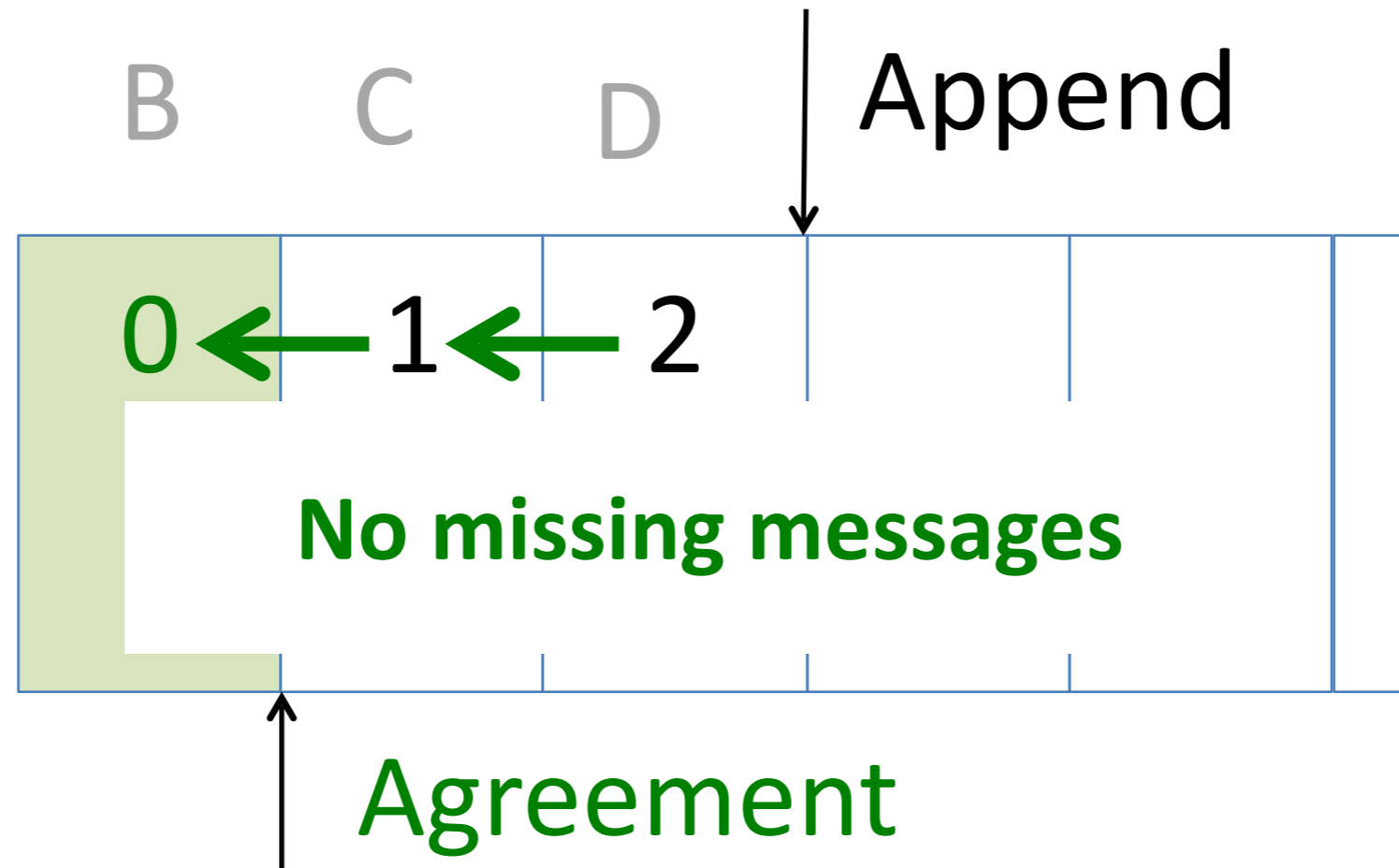
Message from D arrives at hosts, causing *update and agreement on sequence number 0*





seq : 3

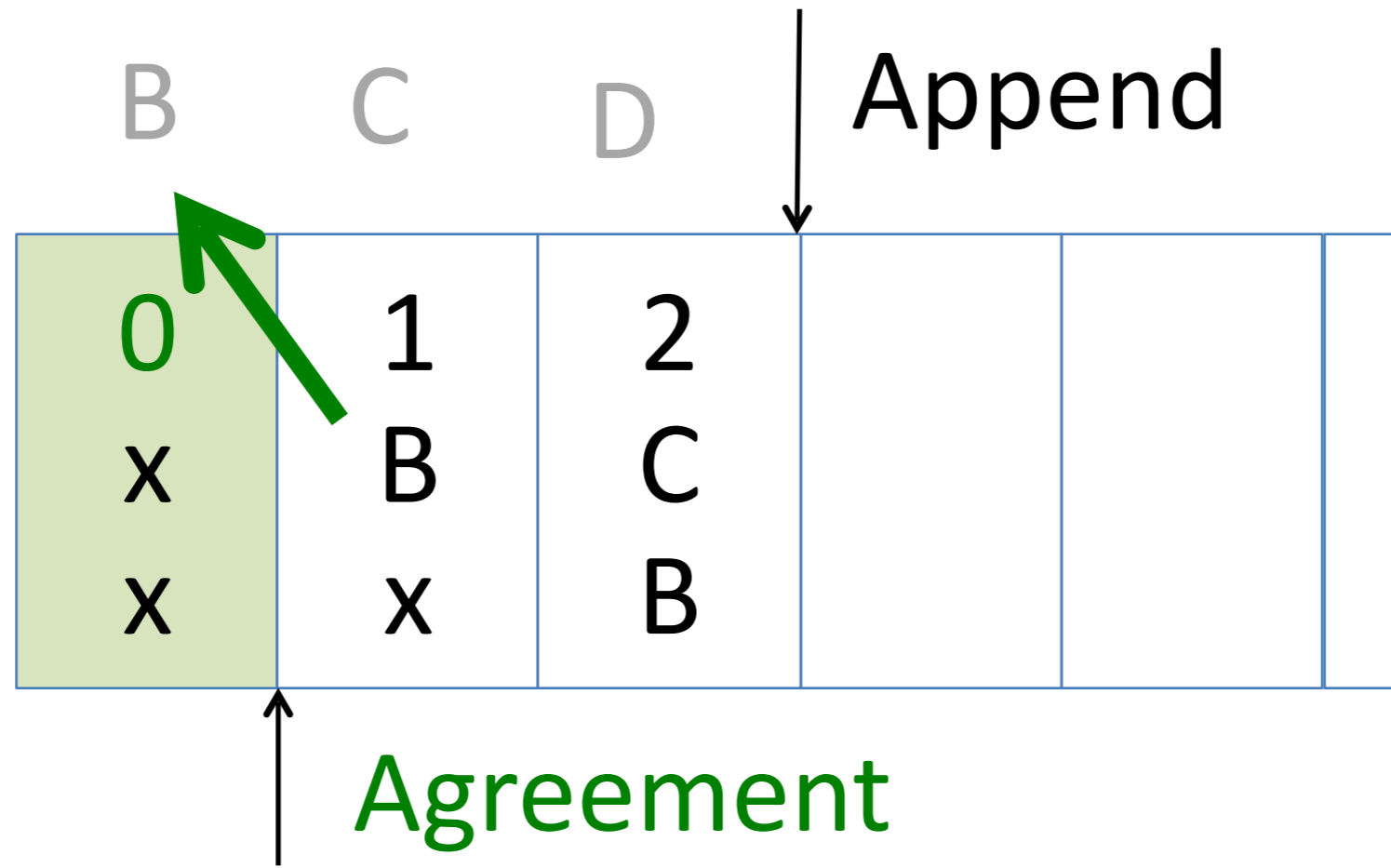




seq : 3

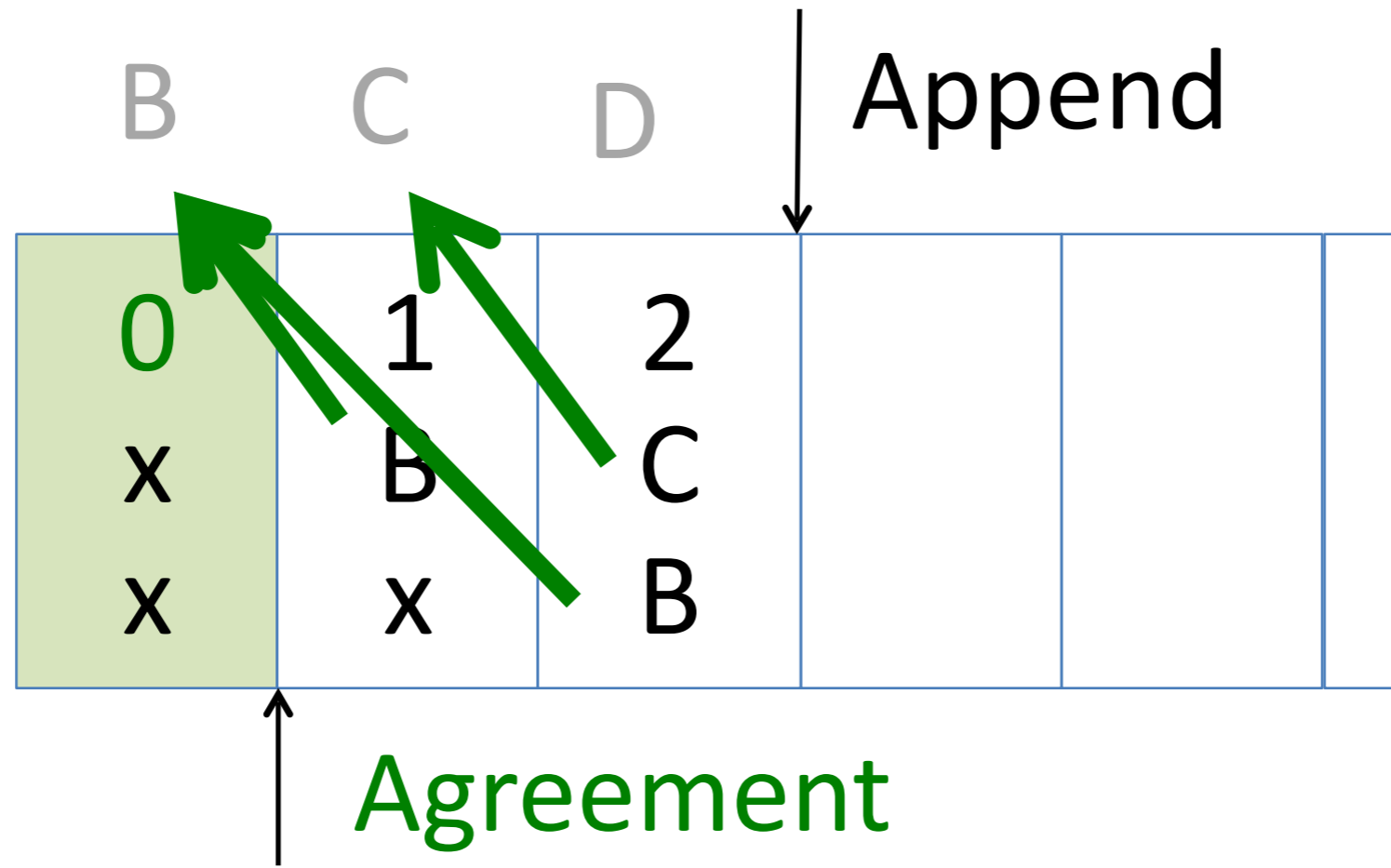


NB: Using a 64bit sequence number gives >400 years operation at 1 message per nano-second.



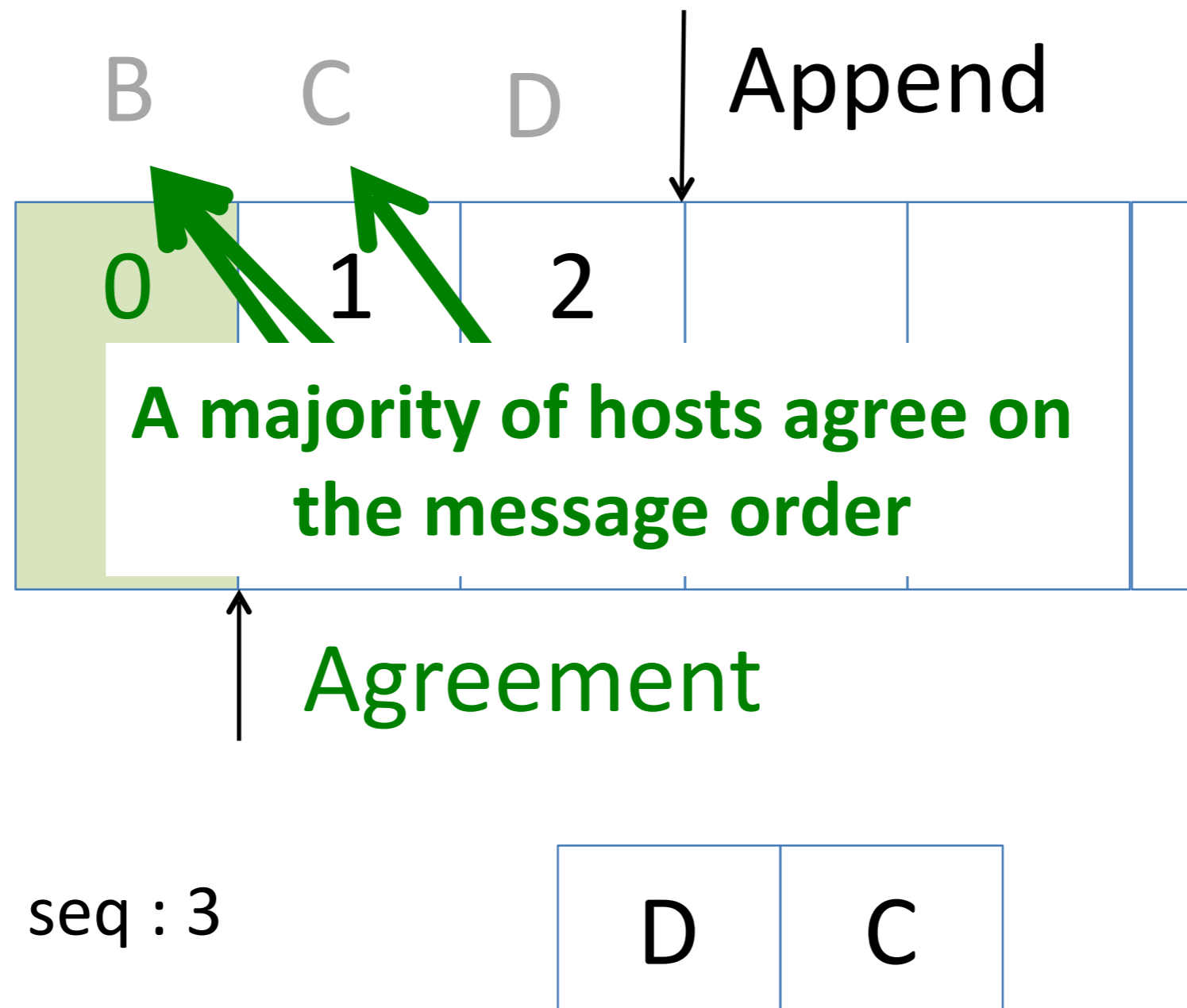
seq : 3





seq : 3



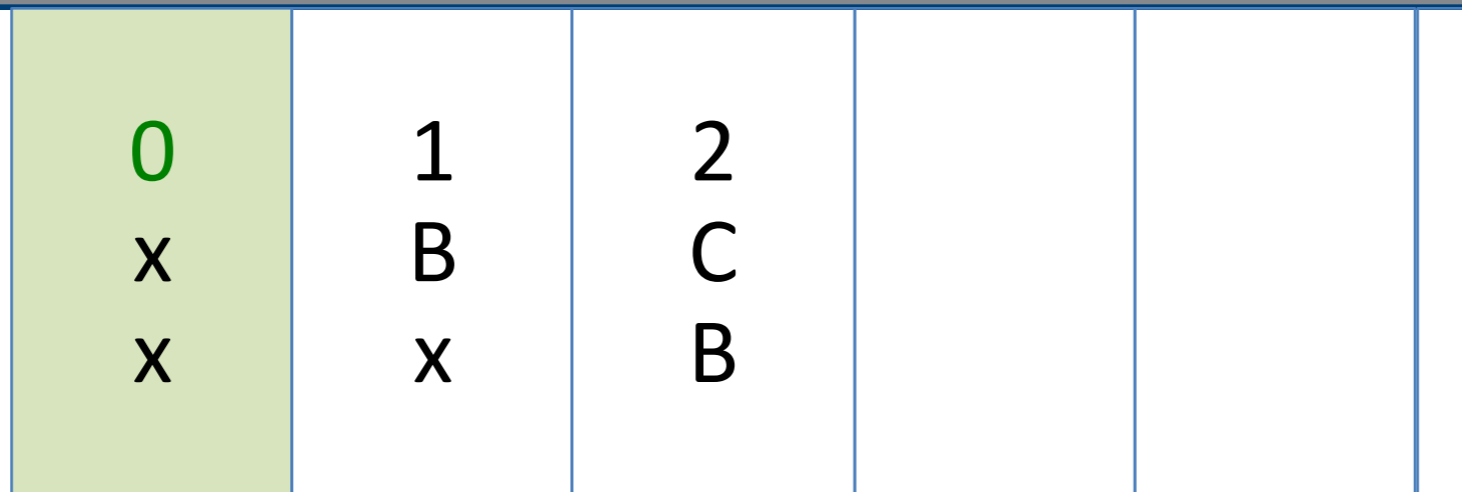


NB: Assuming at most 256 nodes, message history is 128B and can



# How to accelerate the protocol?

Host



seq : 3



NIC

seq : 3







- Protocol implemented over CamIO
  - Runs on TCP, Broadcast UDP, Raw Frames and OpenOnload (ExaSock)
- Offload engine (mostly) implemented in FPGA
  - Currently working for empty messages.



Linux TCP	36 $\mu$ s/msg	28K msg/sec)
Linux UDP	29 $\mu$ s/msg	(34K msg/sec)
Linux UDP (broadcast)	16 $\mu$ s/msg	(62K msg/sec)
Linux TCP (loopback)	9 $\mu$ s/msg	(110K msg/sec)
Linux UDP (loopback)	4 $\mu$ s/msg	(231K msg/sec)
Libexanic (raw frames)	3 $\mu$ s/msg	(359K msg/sec)
Exo HW offload	0.4 $\mu$ s/msg	(2500K msg/sec)



<b>Linux TCP</b>	<b>36<math>\mu</math>s/msg</b>	<b>28K msg/sec)</b>
Linux UDP	29 $\mu$ s/msg	(34K msg/sec)
Linux UDP (broadcast)	16 $\mu$ s/msg	(62K msg/sec)
Linux TCP (loopback)	9 $\mu$ s/msg	(110K msg/sec)
Linux UDP (loopback)	4 $\mu$ s/msg	(231K msg/sec)
Libexanic (raw frames)	3 $\mu$ s/msg	(359K msg/sec)
<b>Exo HW offload</b>	<b>0.4<math>\mu</math>s/msg</b>	<b>(2500K msg/sec)</b>

**$\approx$  100x Faster**



- Exo is a work in progress
- We use a specialised network and hardware offload to build fast scalable coordination for rack-scale architectures
- Initial results suggest at least a 100x speed improvement over existing protocols/systems



Thanks to our sponsors



Engineering and Physical Sciences  
Research Council



Research

Systems



Group

`srg@cambridge:~$`

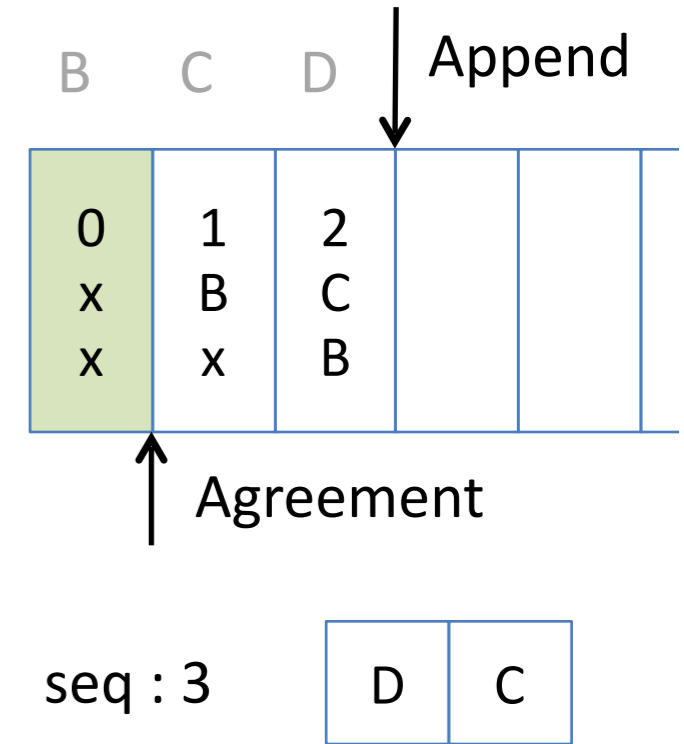
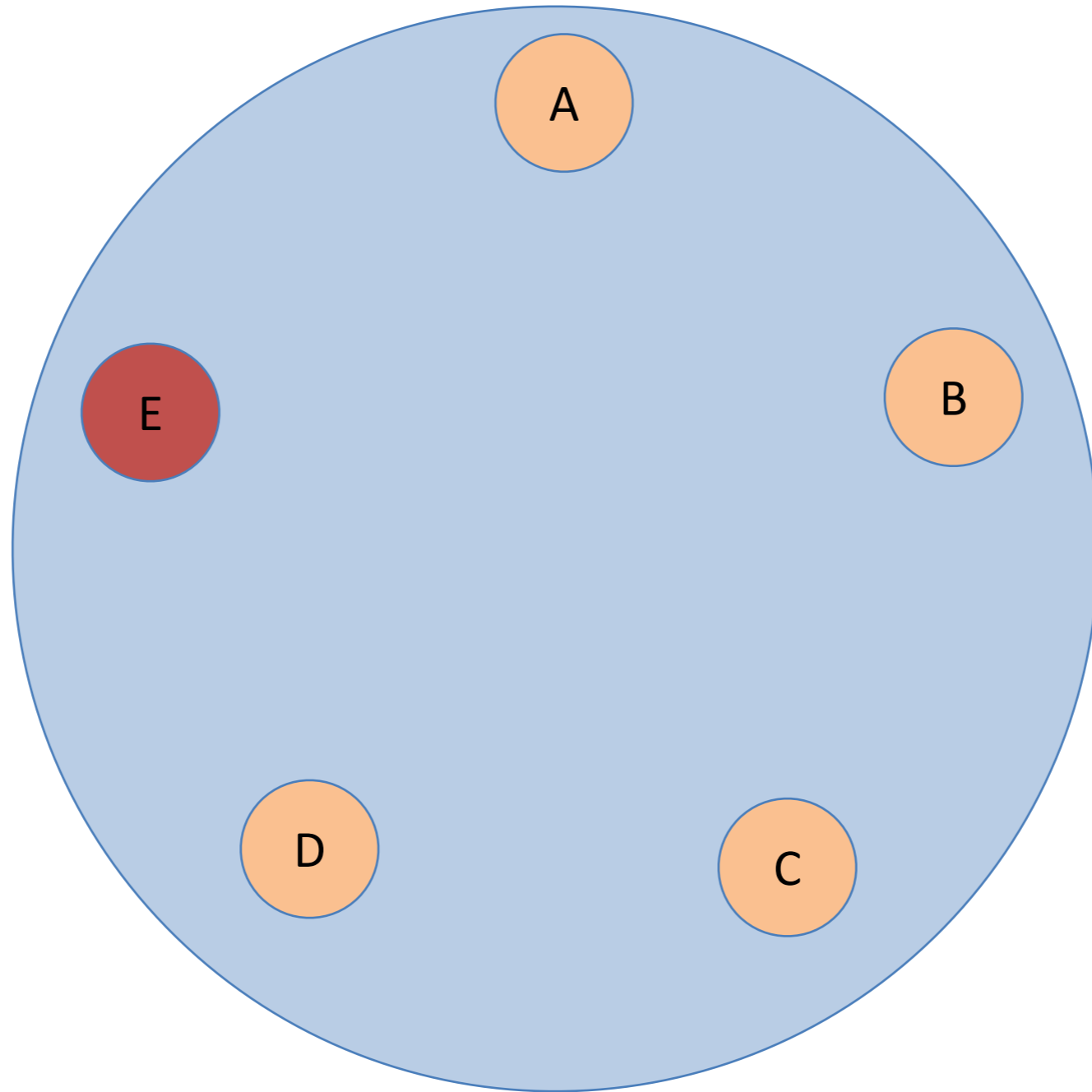
This work was jointly supported by the EPSRC INTERNET Project EP/H040536/1 and the Defense Advanced Research Projects Agency (DARPA) and the Air Force Research Laboratory (AFRL), under contract FA8750-11-C-0249. The views, opinions, and/or findings contained in this article/presentation are those of the author/presenter and should not be interpreted as representing the official views or policies, either expressed or implied, of the Defense Advanced Research Projects Agency or the Department of Defense.

Backups



# Anatomy of the Exo Protocol

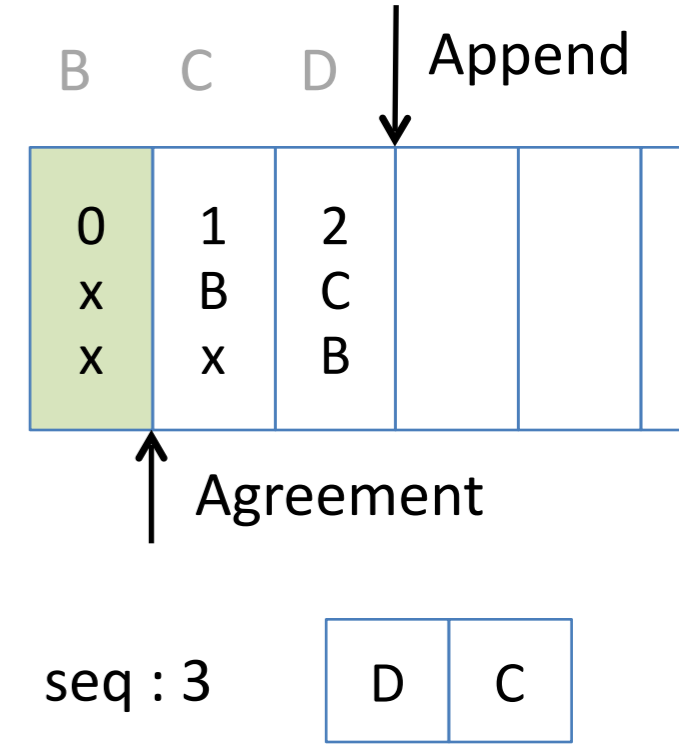
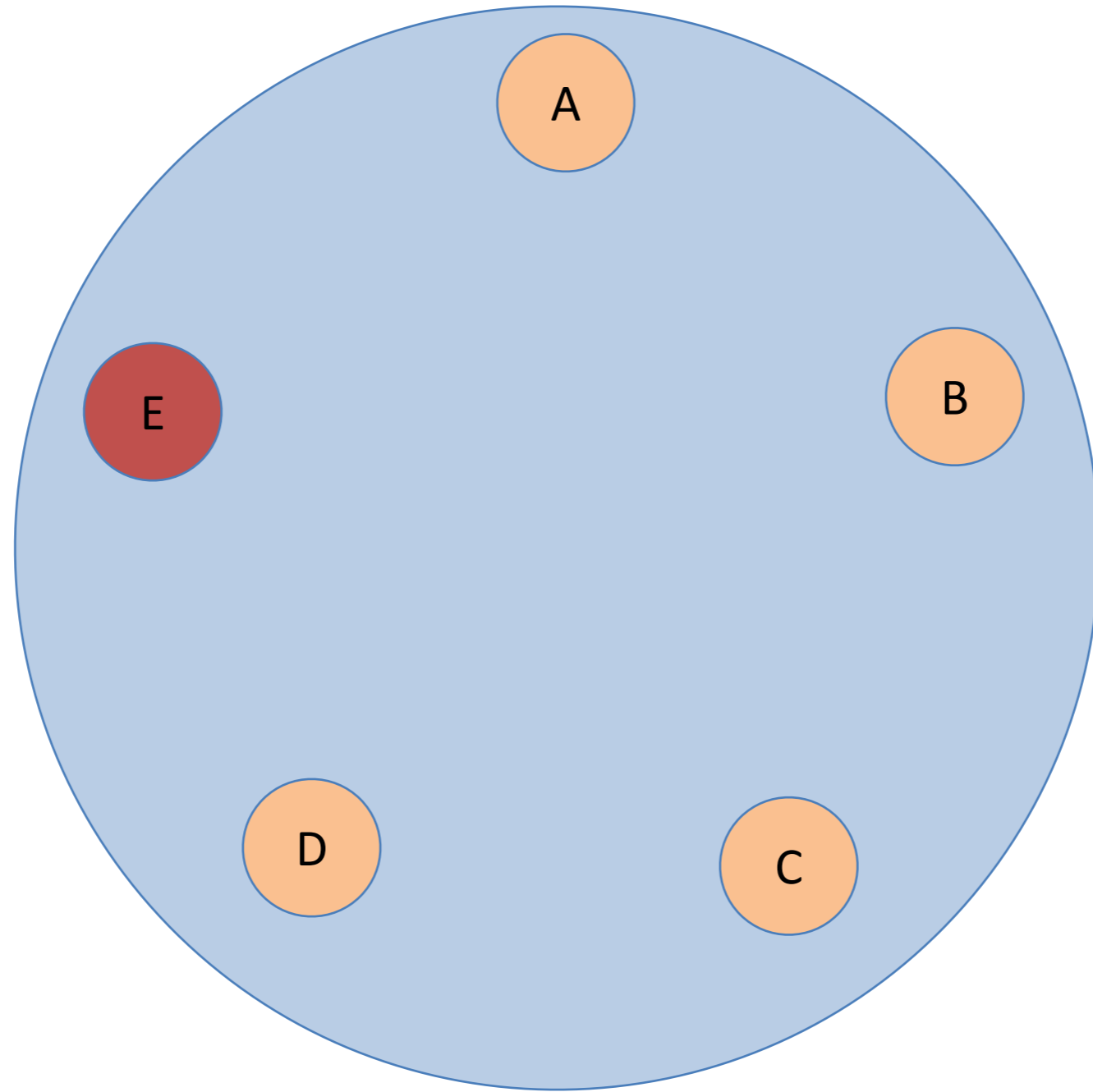
E now has the token





# Handling Errors?

E now has the token

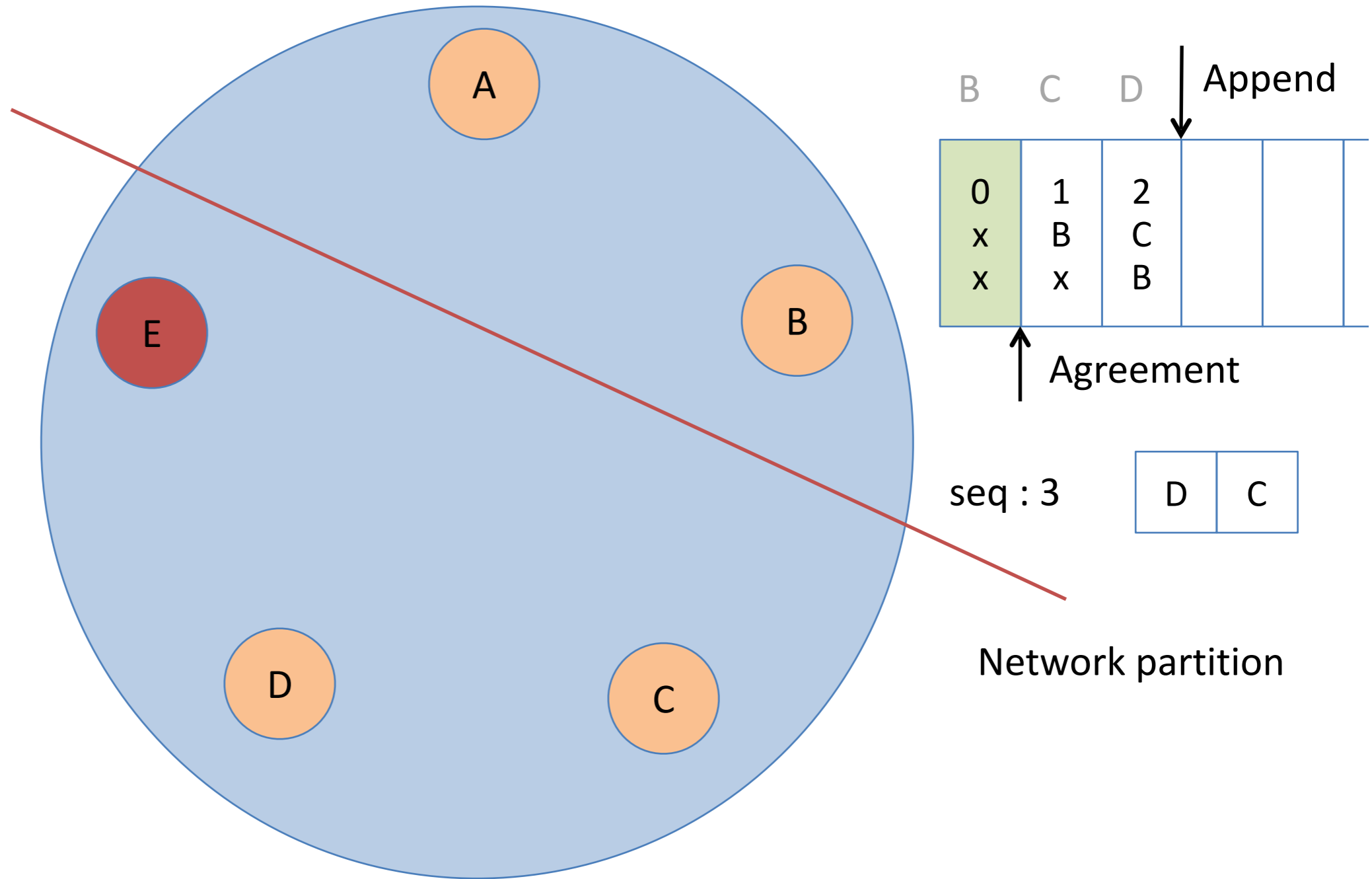






# Handling Errors?

Consider the following pathological case

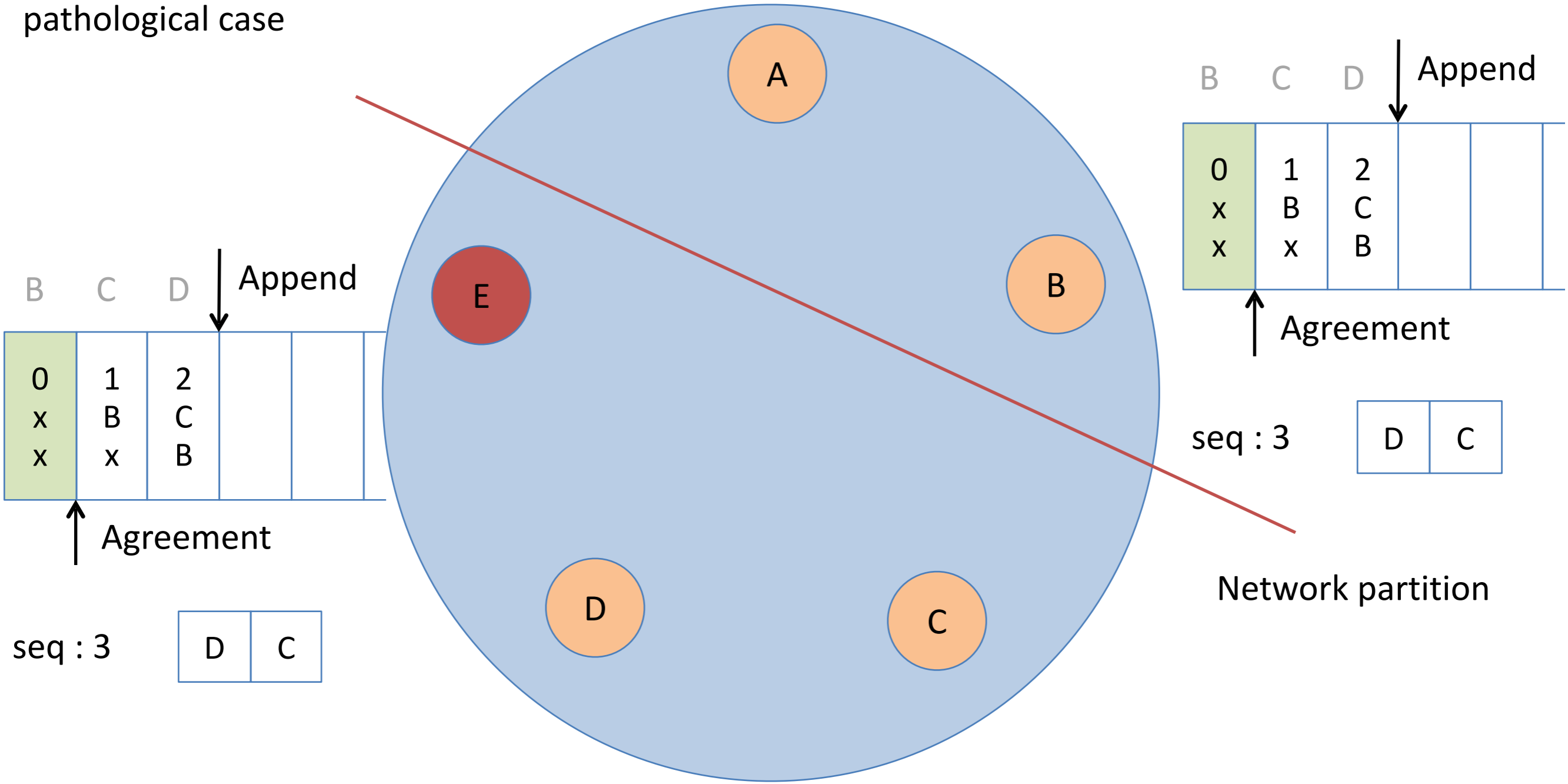


Network partition



# Handling Errors?

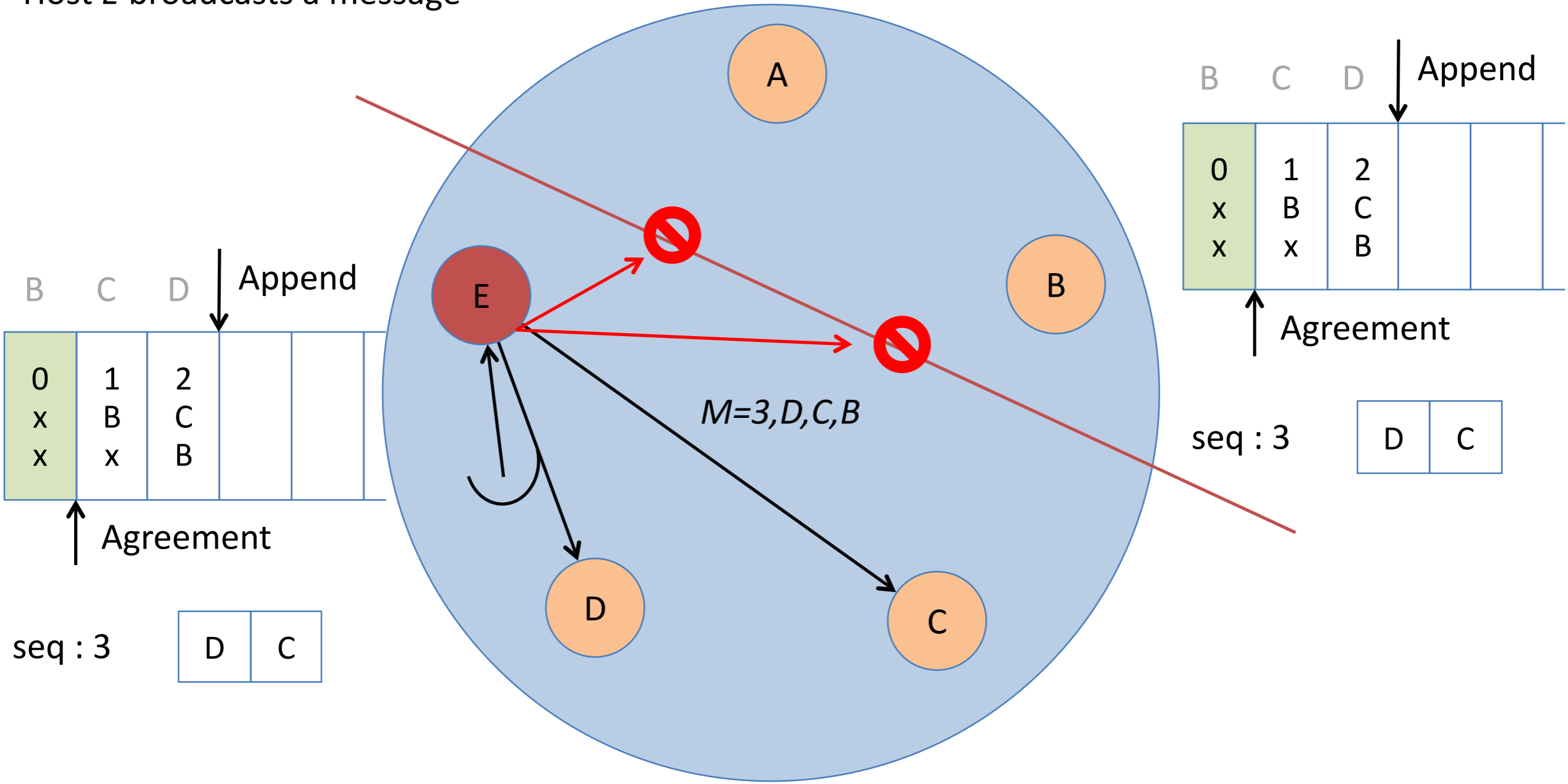
Consider the following pathological case





# Handling Errors?

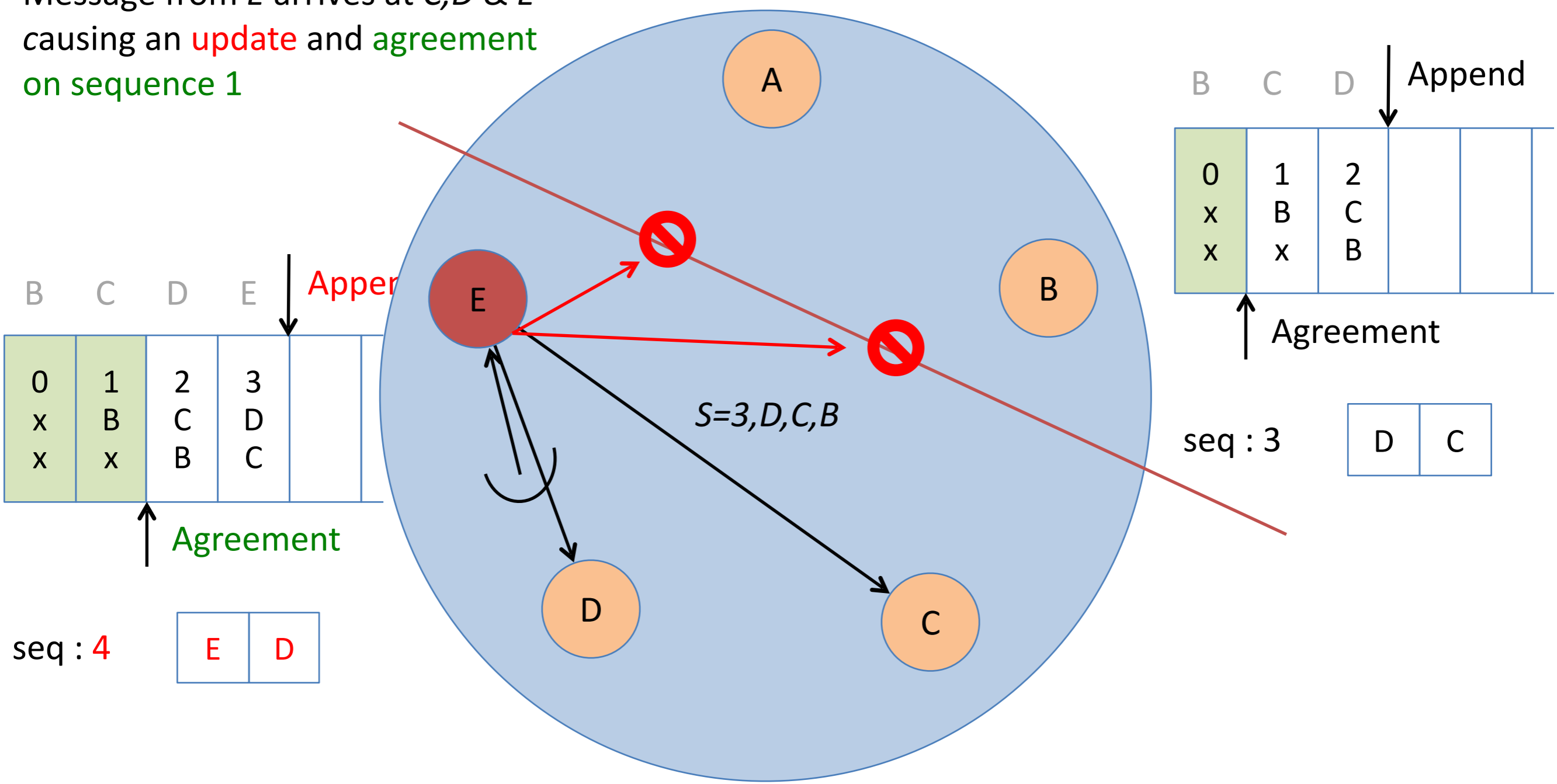
Host E broadcasts a message





# Handling Errors?

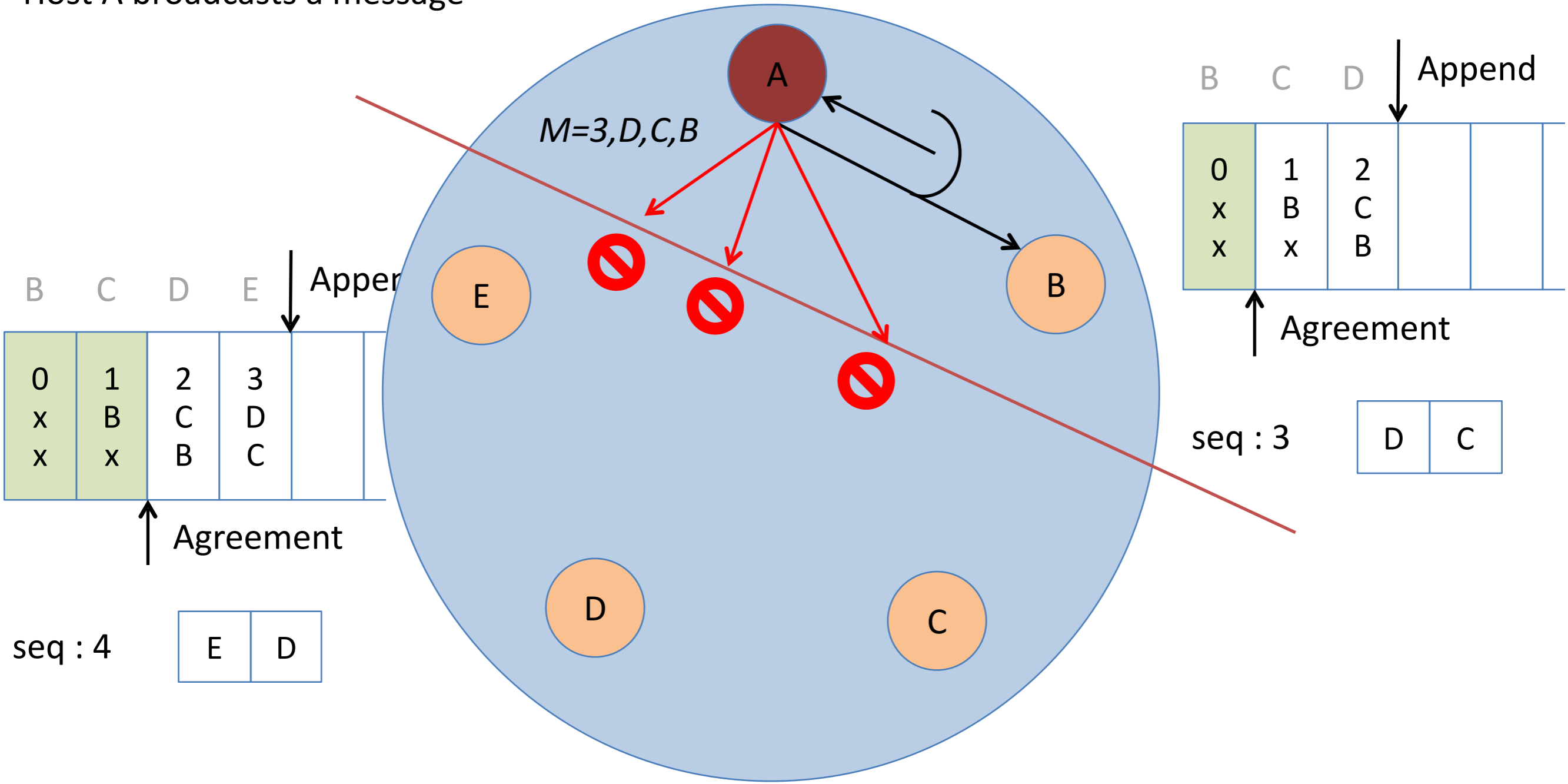
Message from *E* arrives at *C, D & E* causing an **update** and **agreement** on sequence 1





# Handling Errors?

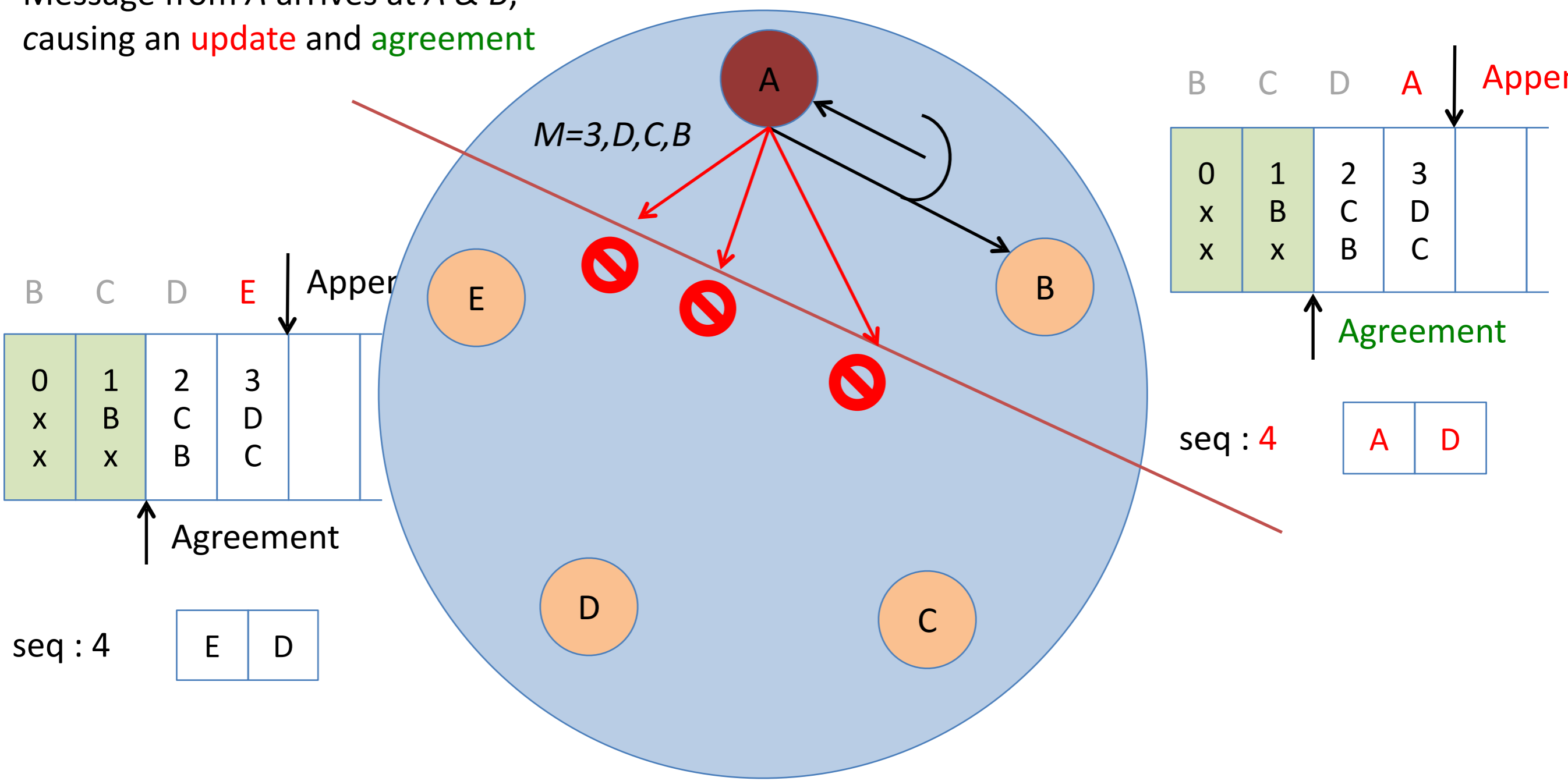
Host A broadcasts a message





# Handling Errors?

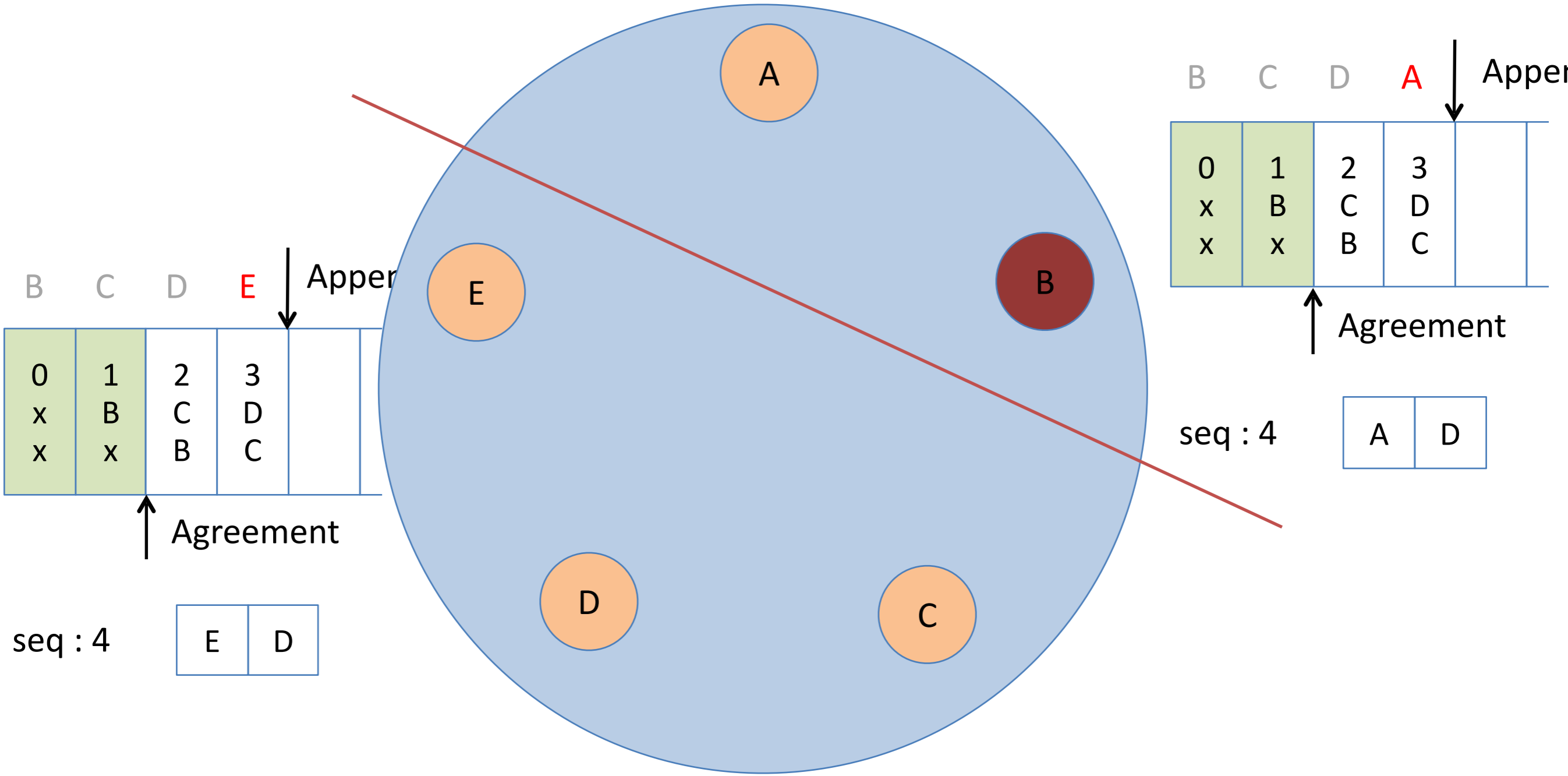
Message from A arrives at A & B, causing an **update** and **agreement**





# Handling Errors?

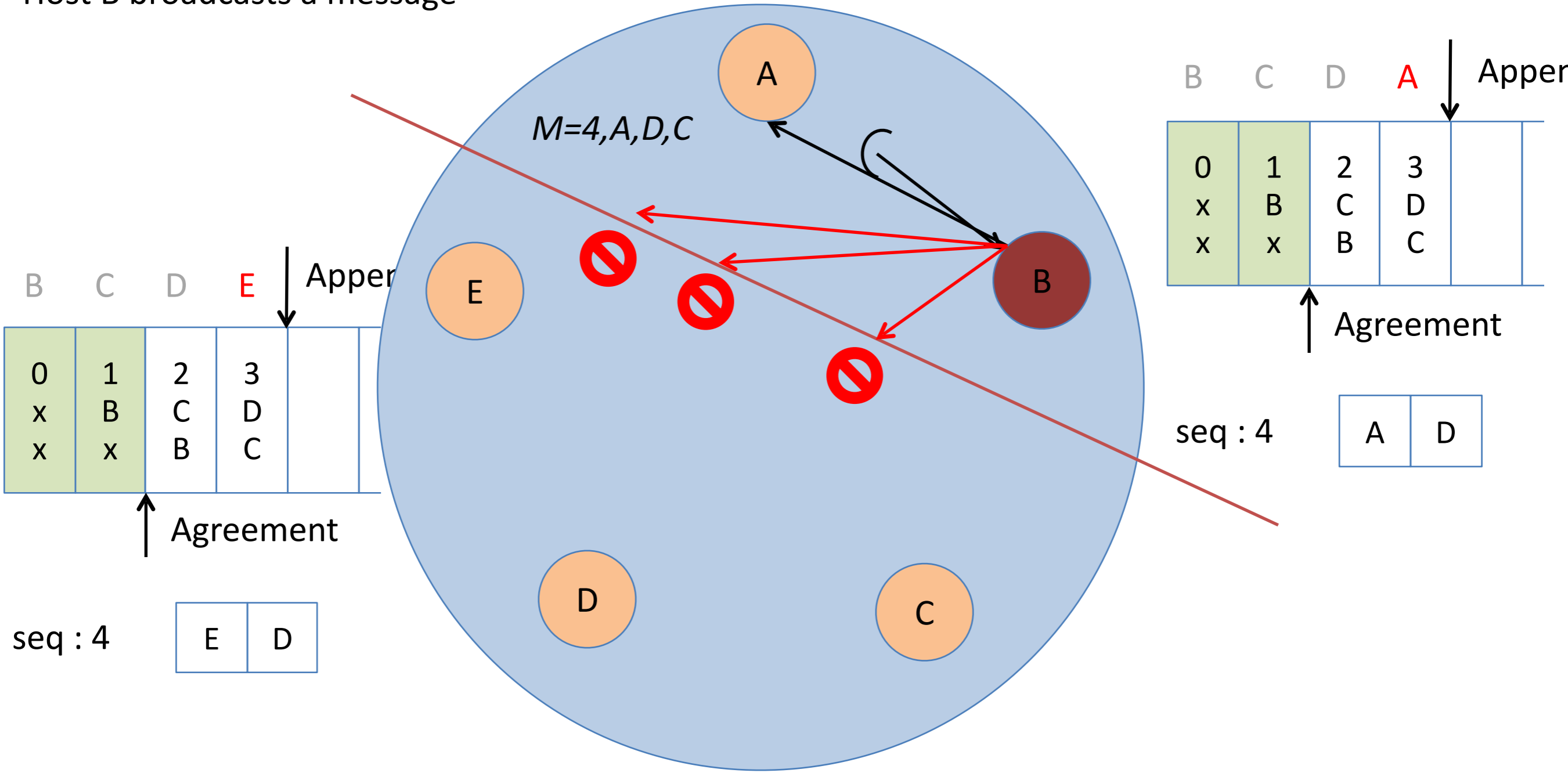
B now has the token





# Handling Errors?

Host B broadcasts a message

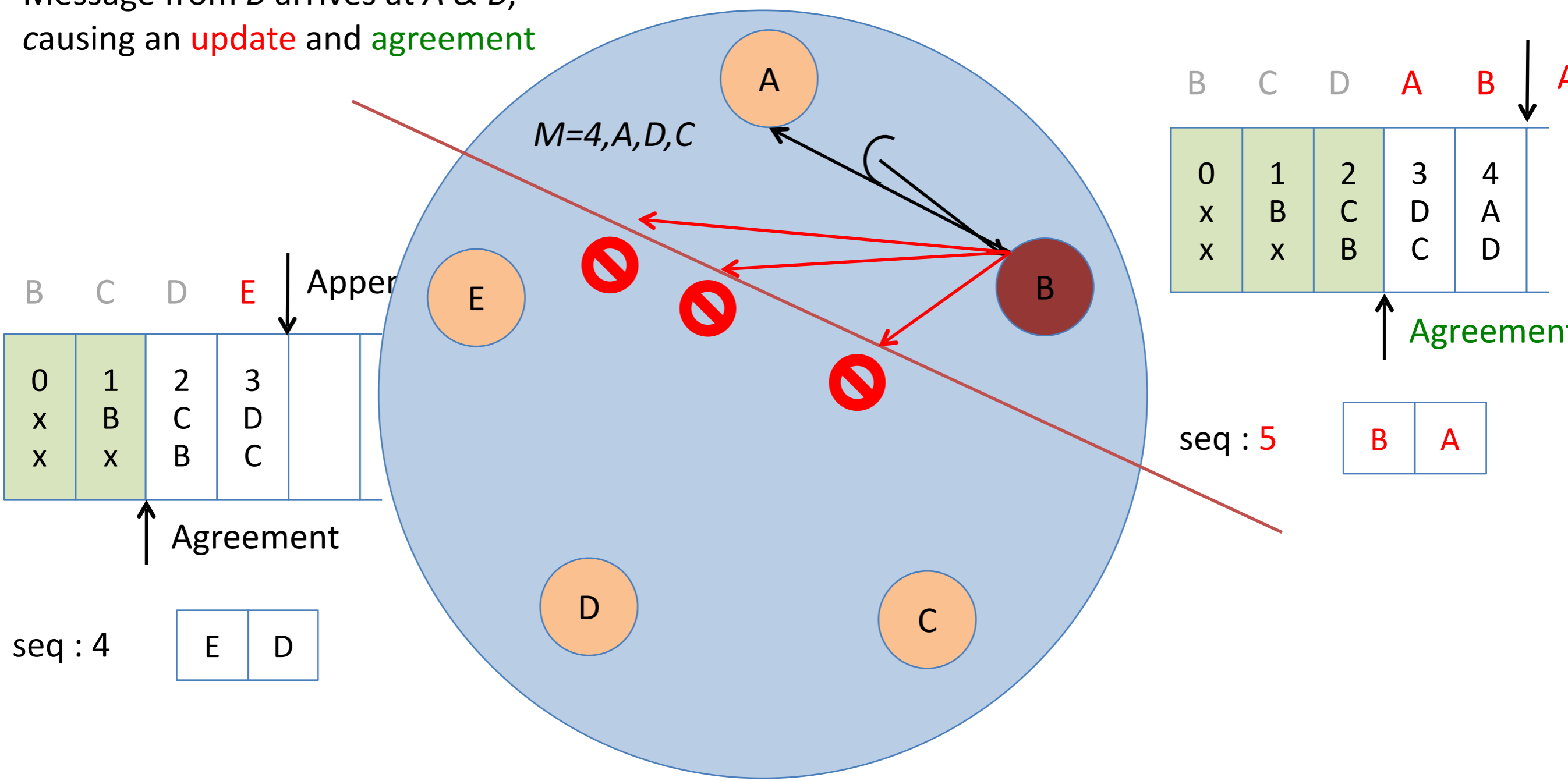






# Handling Errors?

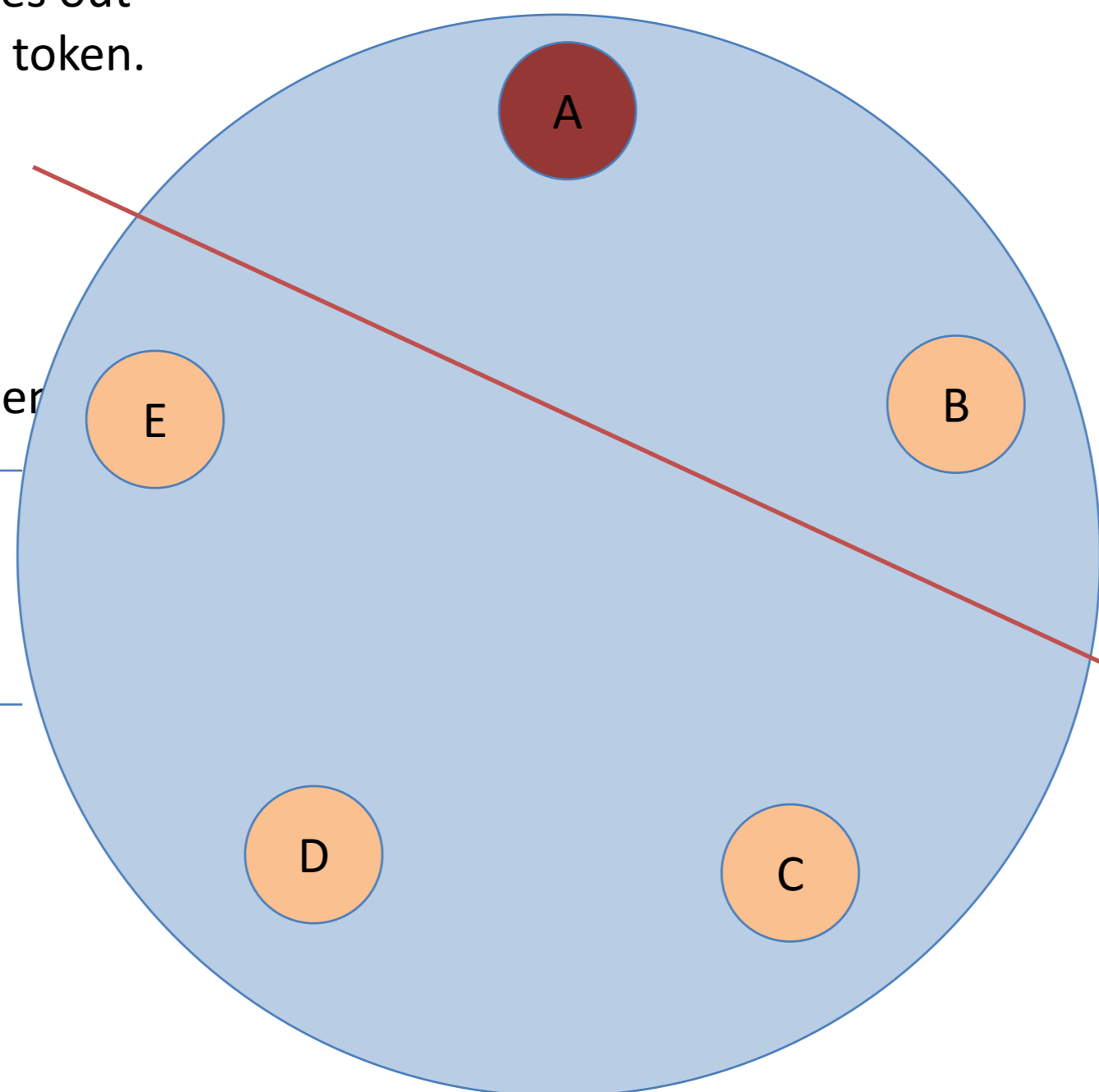
Message from *B* arrives at *A* & *B*, causing an **update** and **agreement**





# Handling Errors?

After some time, A times out and assumes it has the token.



B	C	D	E	Apper
0	1	2	3	
x	B	C	D	
x	x	B	C	
			B	

↑ Agreement

seq : 4

E	D
---	---

B	C	D	A	B	App
0	1	2	3	4	
x	B	C	D	A	
x	x	B	C	D	

↑ Agreement

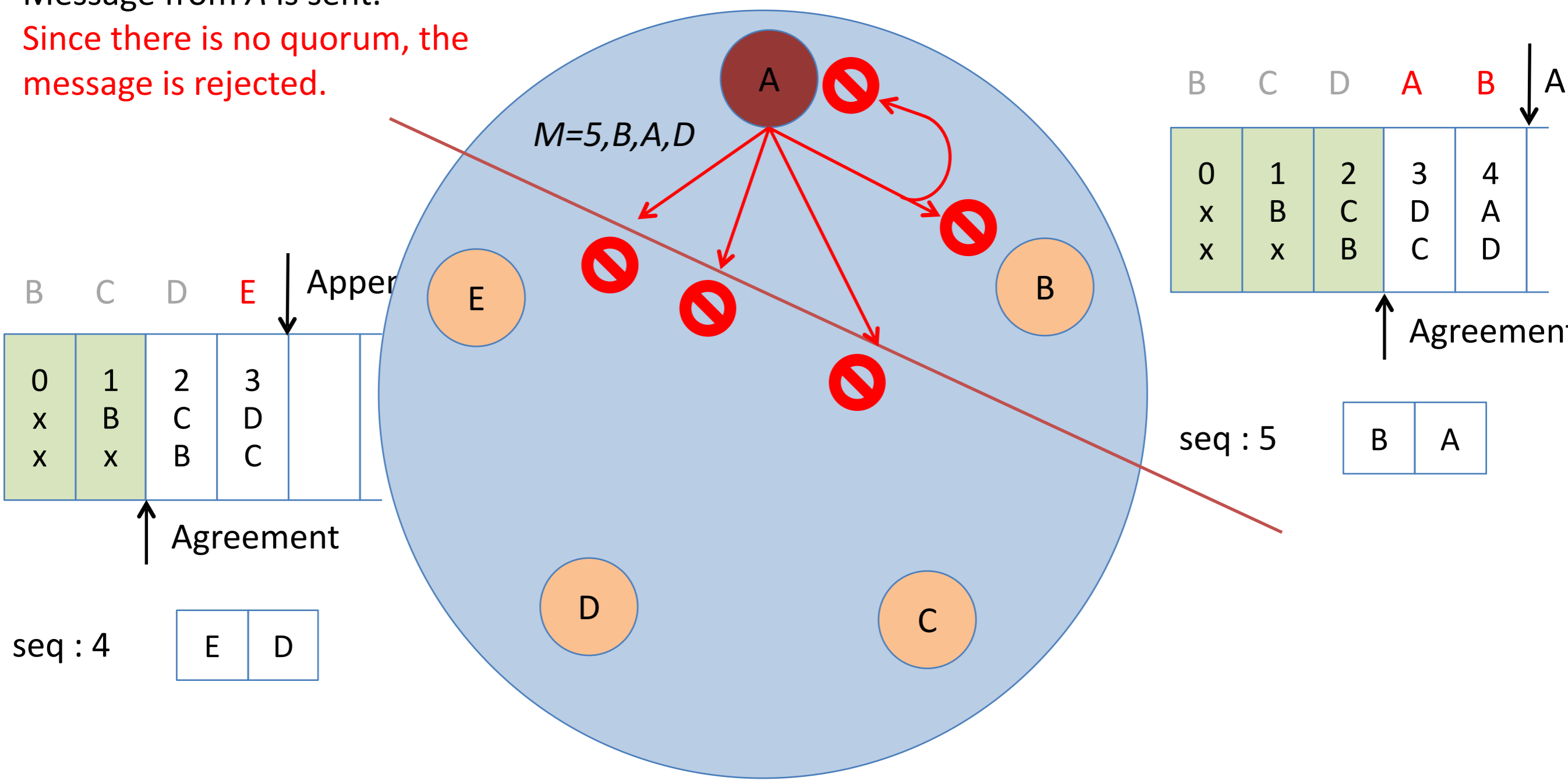
seq : 5

B	A
---	---



# Handling Errors?

Message from A is sent.  
Since there is no quorum, the message is rejected.

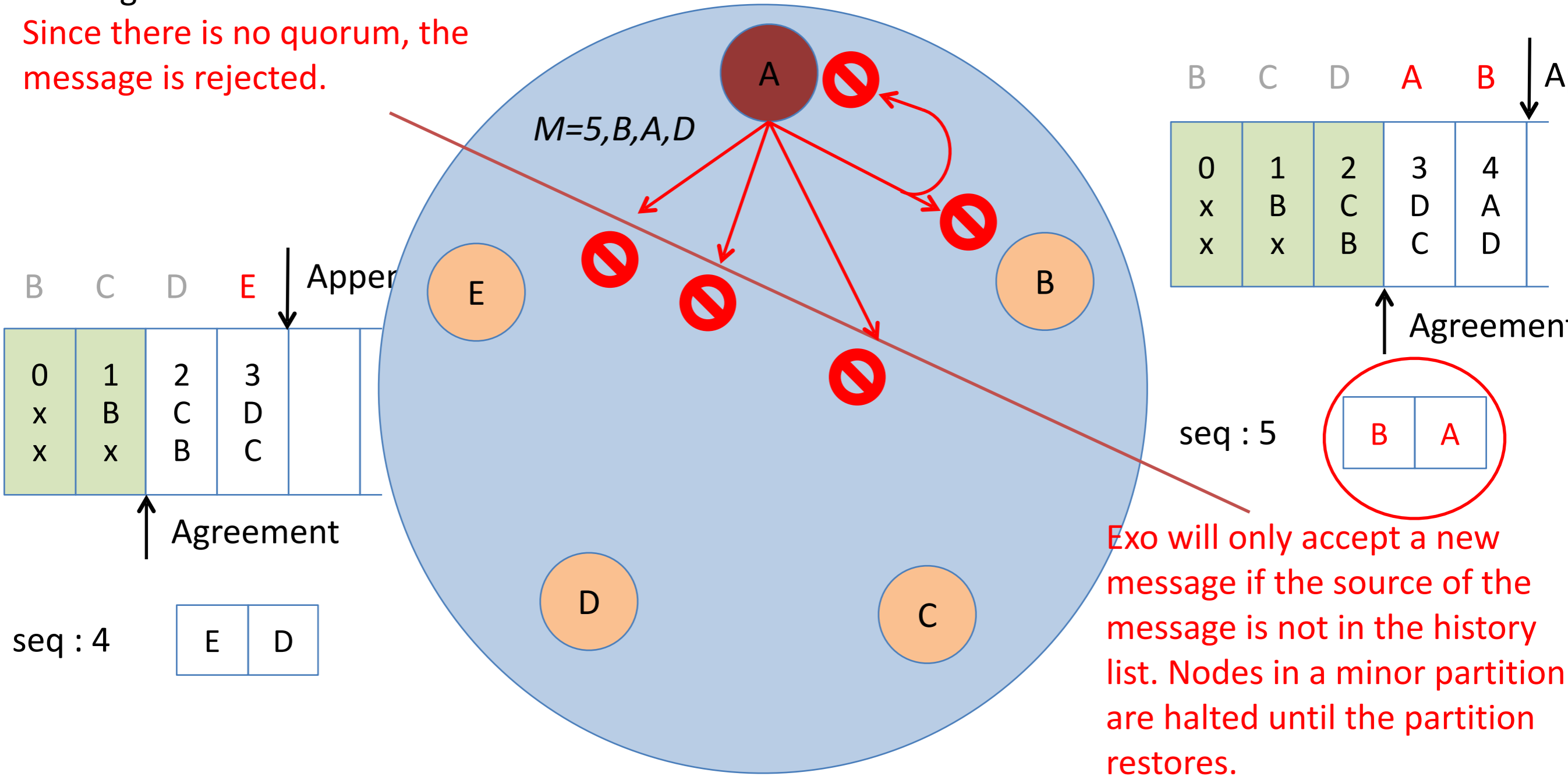




# Handling Errors?

Message from A is sent.

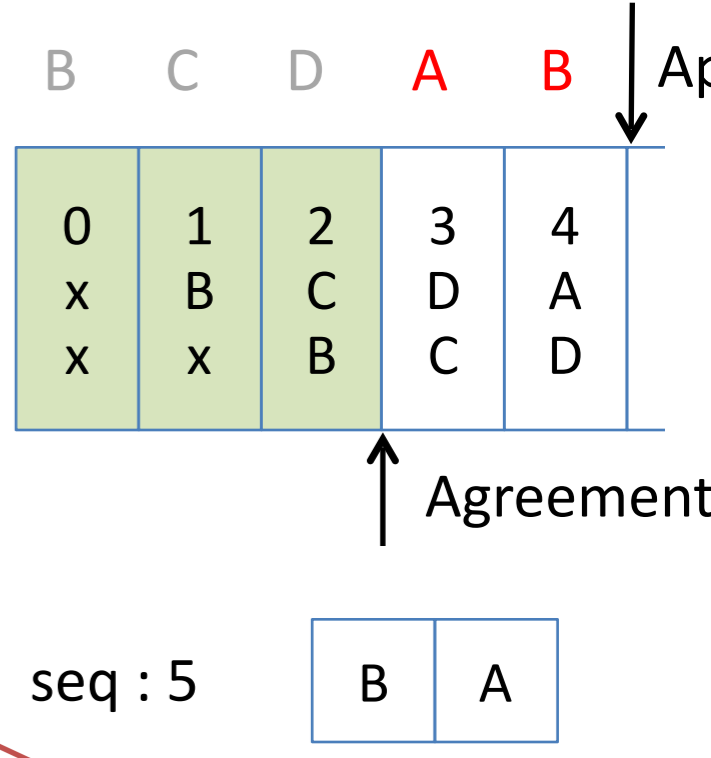
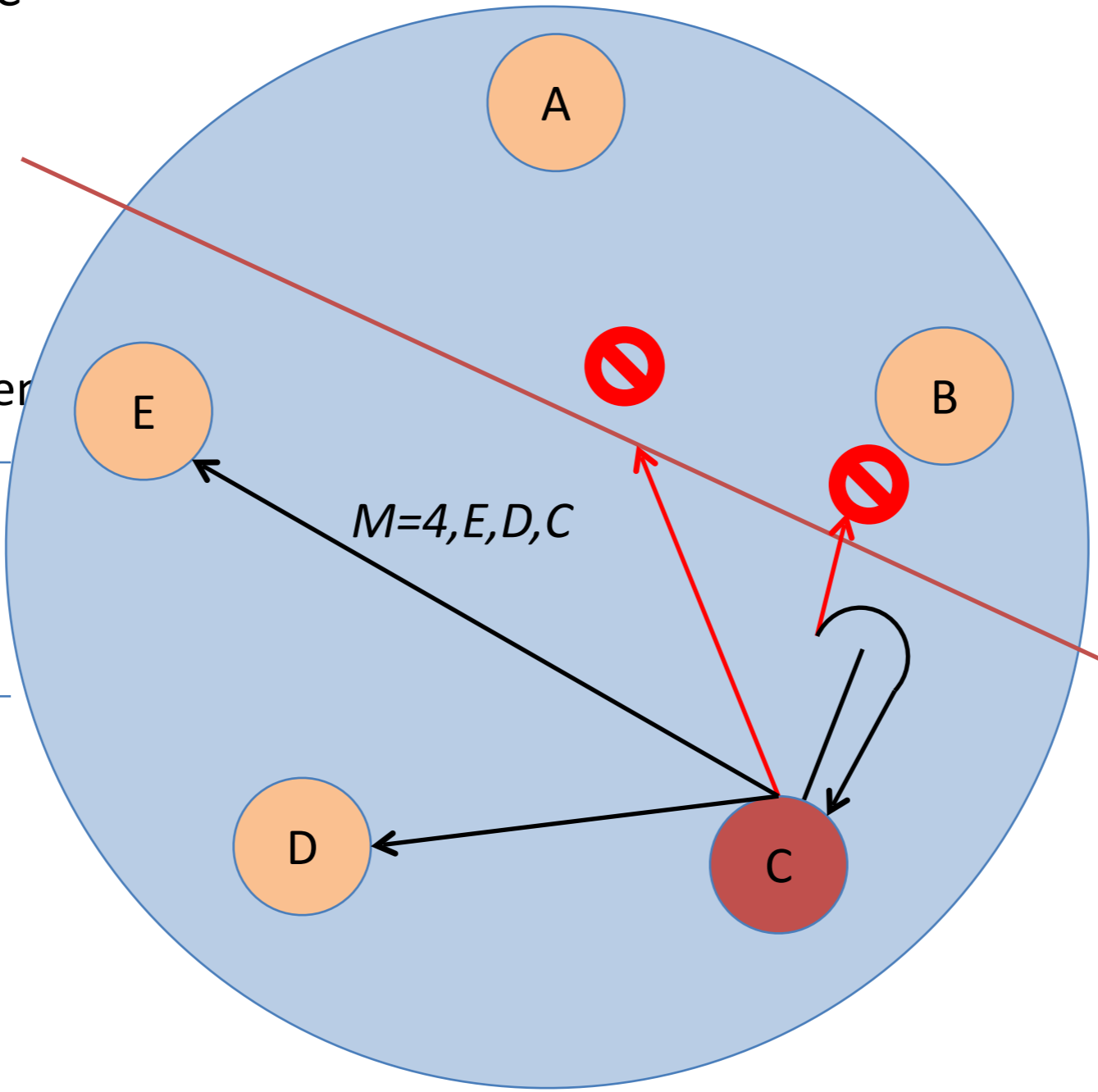
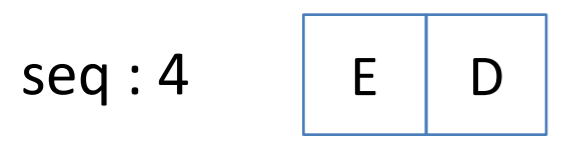
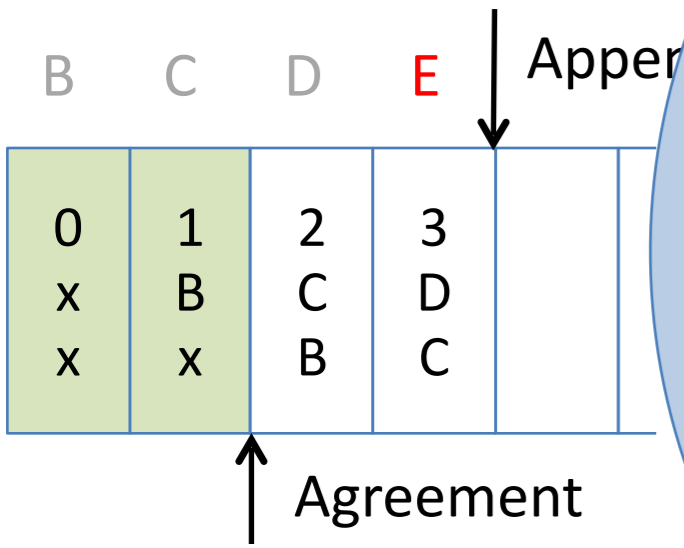
Since there is no quorum, the message is rejected.





# Handling Errors?

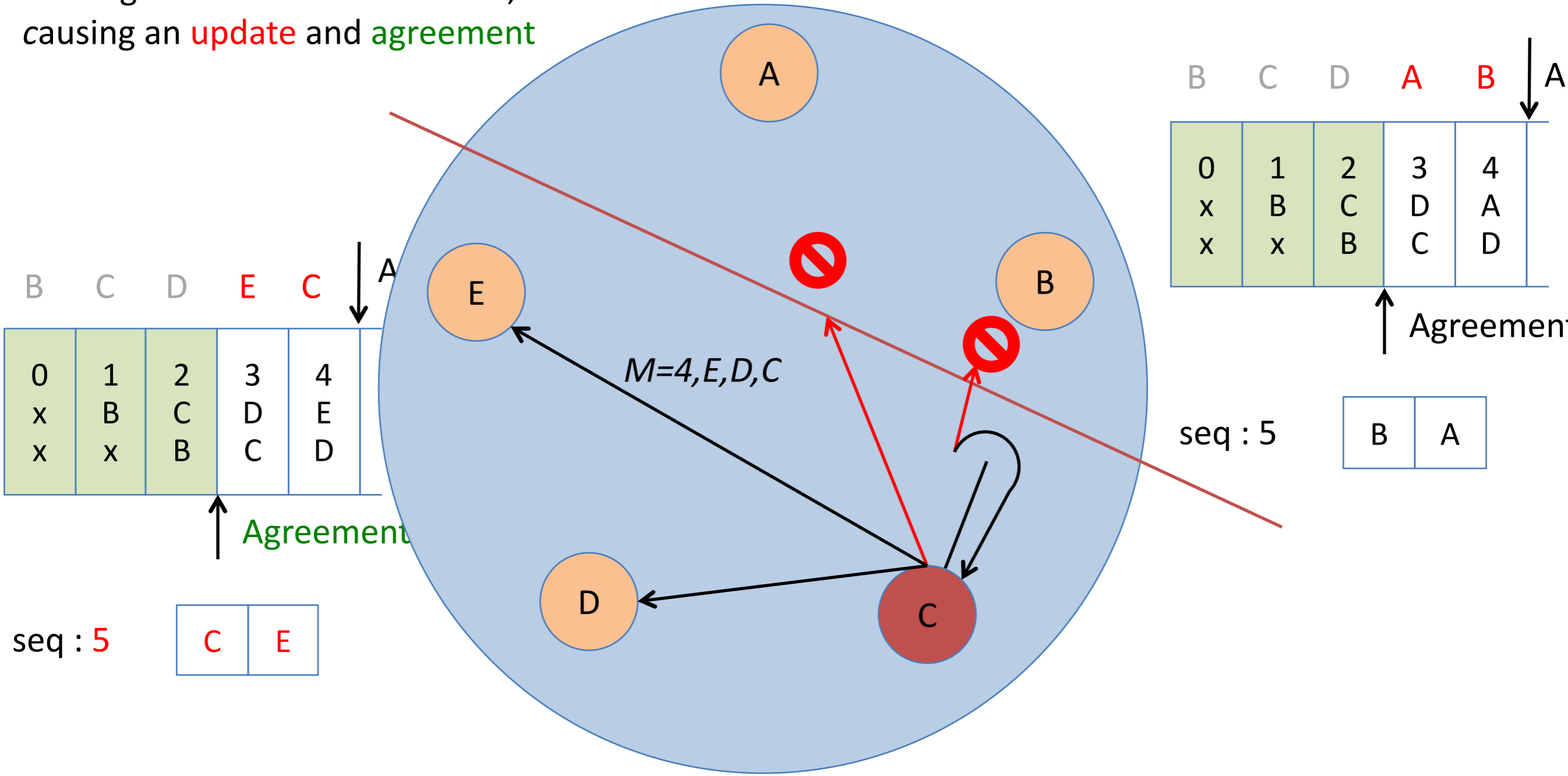
In the meantime, Host C broadcasts a message





# Handling Errors?

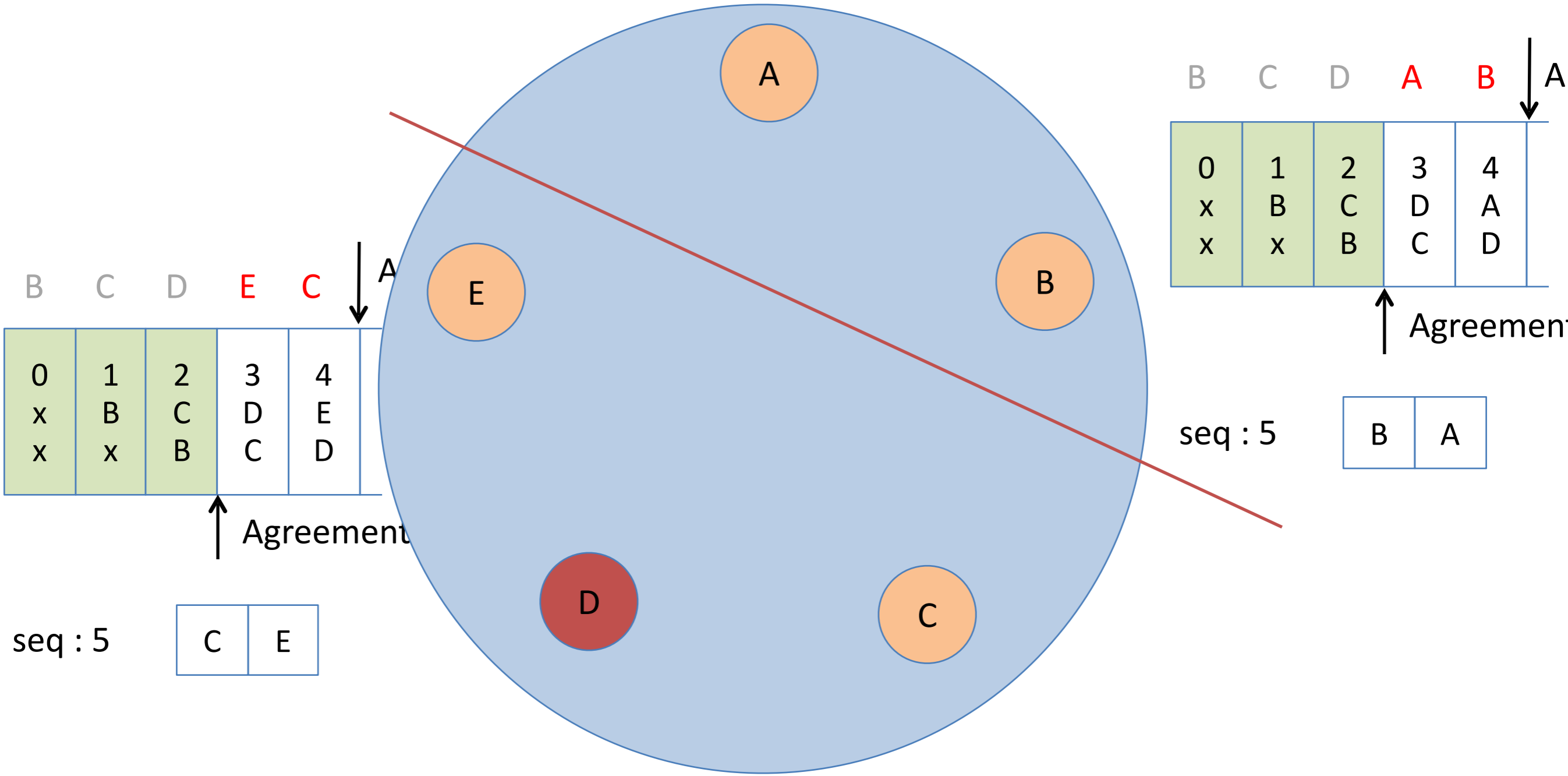
Message from B arrives at A & B, causing an **update** and **agreement**





# Handling Errors?

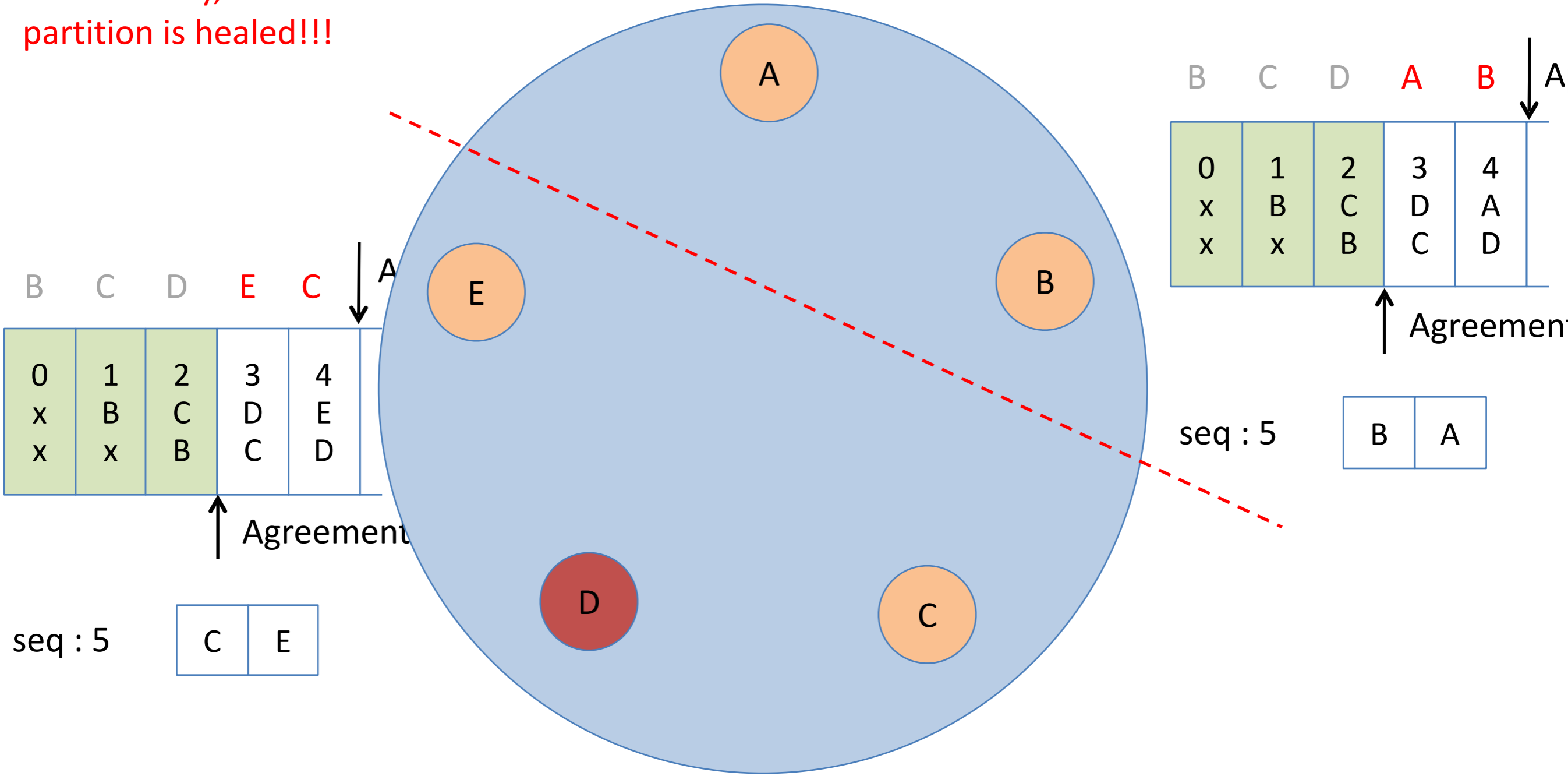
D now has the token





# Handling Errors?

Miraculously, the partition is healed!!!

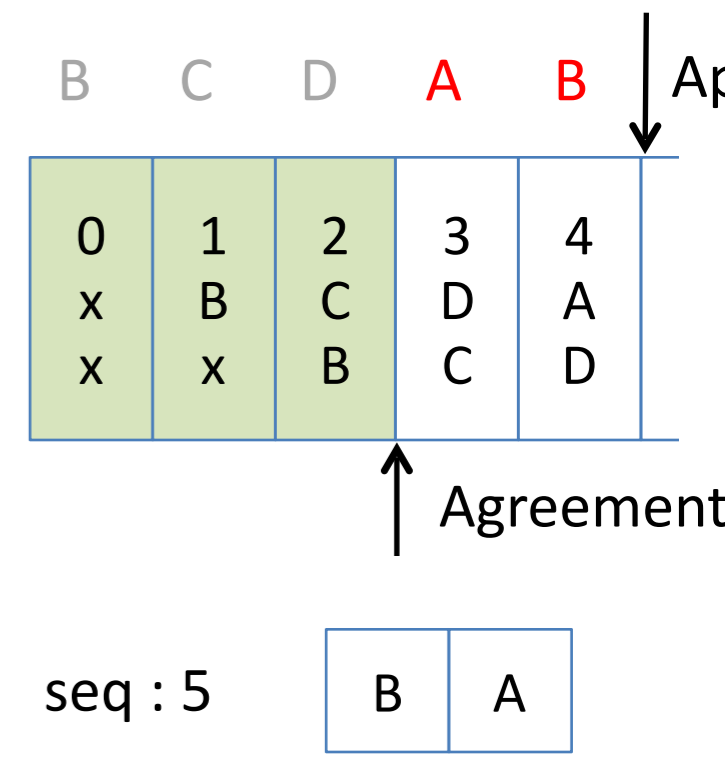
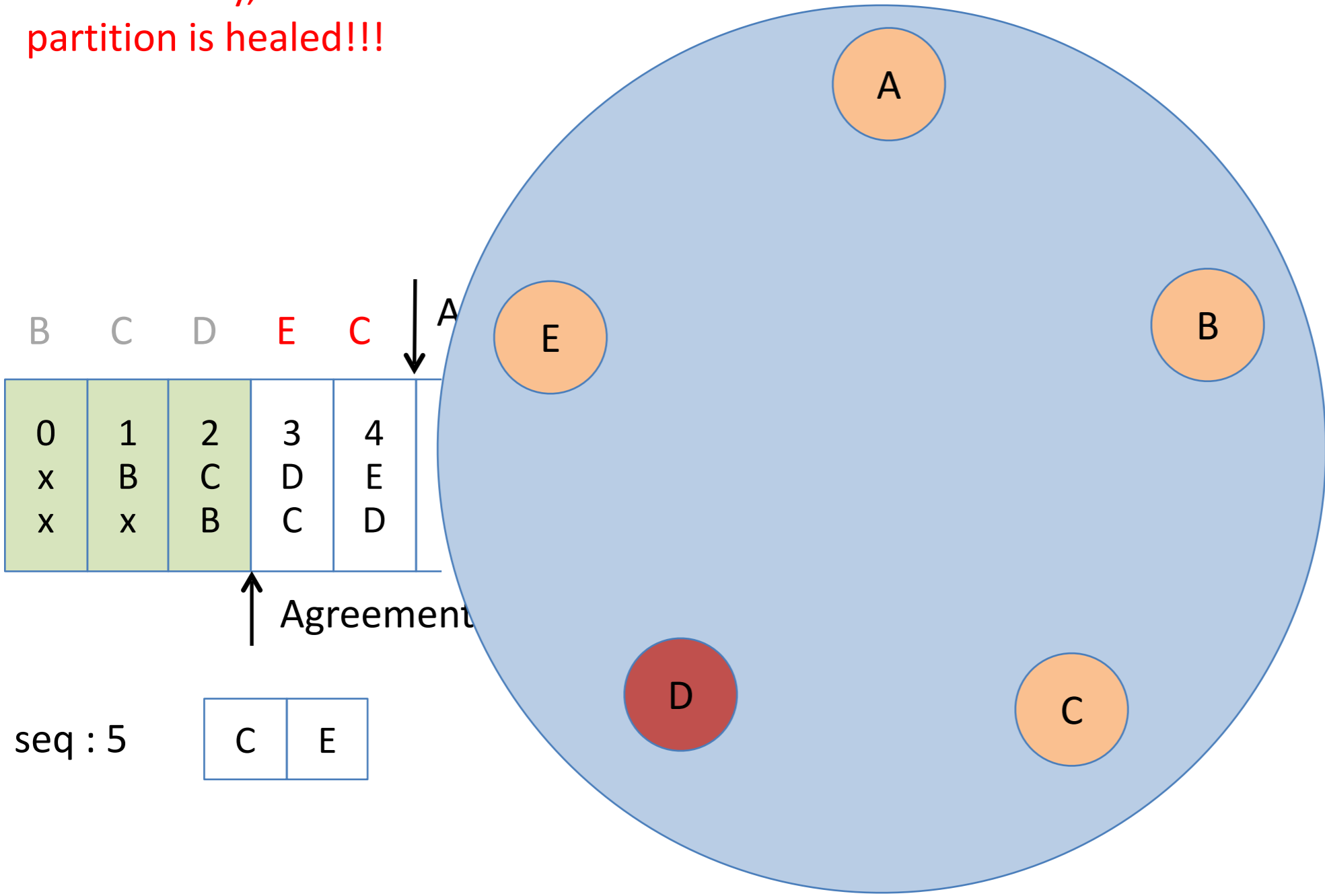






# Handling Errors?

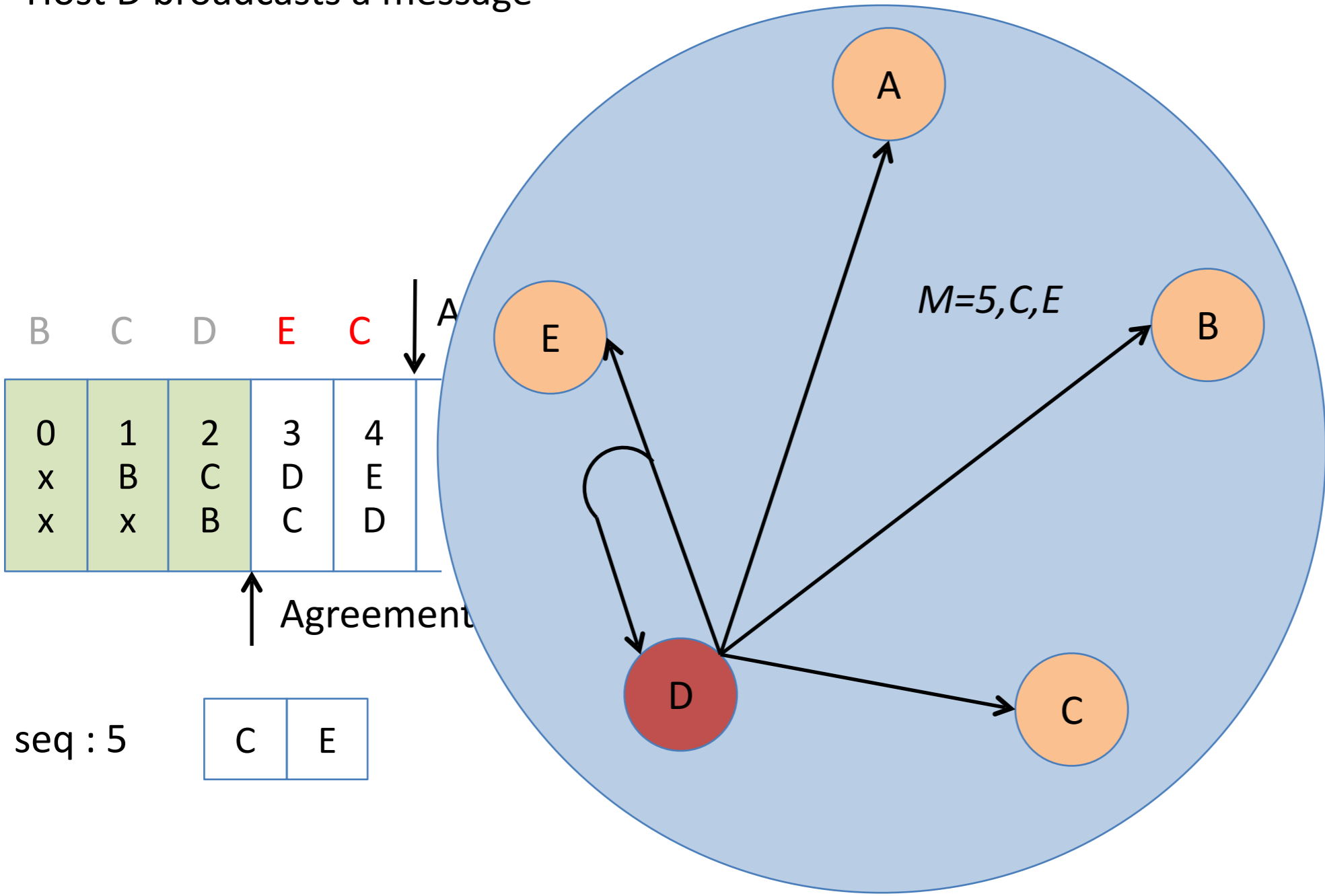
Miraculously, the partition is healed!!!





# Handling Errors?

Host D broadcasts a message



seq : 5

B	C	D	E	C	A
0	1	2	3	4	
x	B	C	D	E	
x	x	B	C	D	

↑ Agreement

C	E
---	---

seq : 5

B	C	D	A	B	A
0	1	2	3	4	
x	B	C	D	A	
x	x	B	C	D	

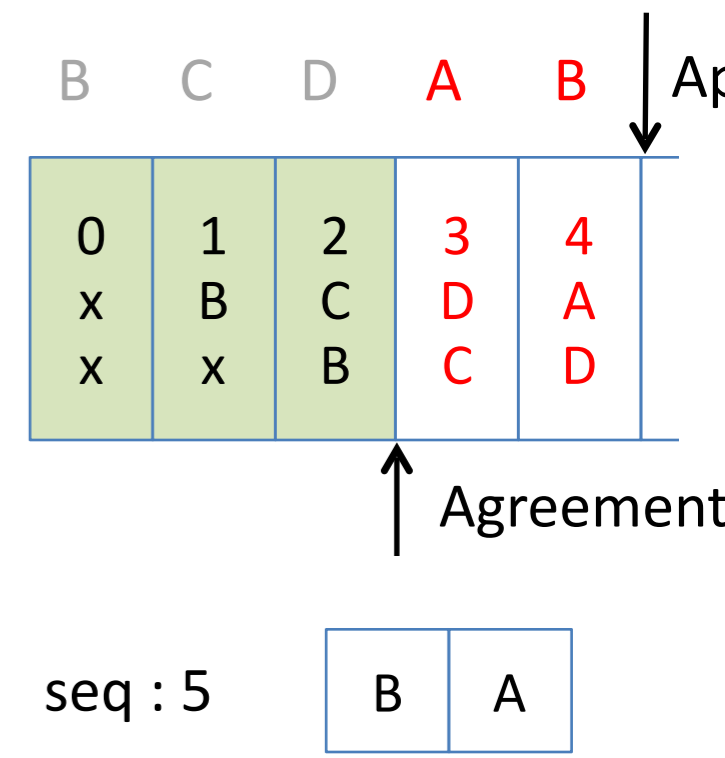
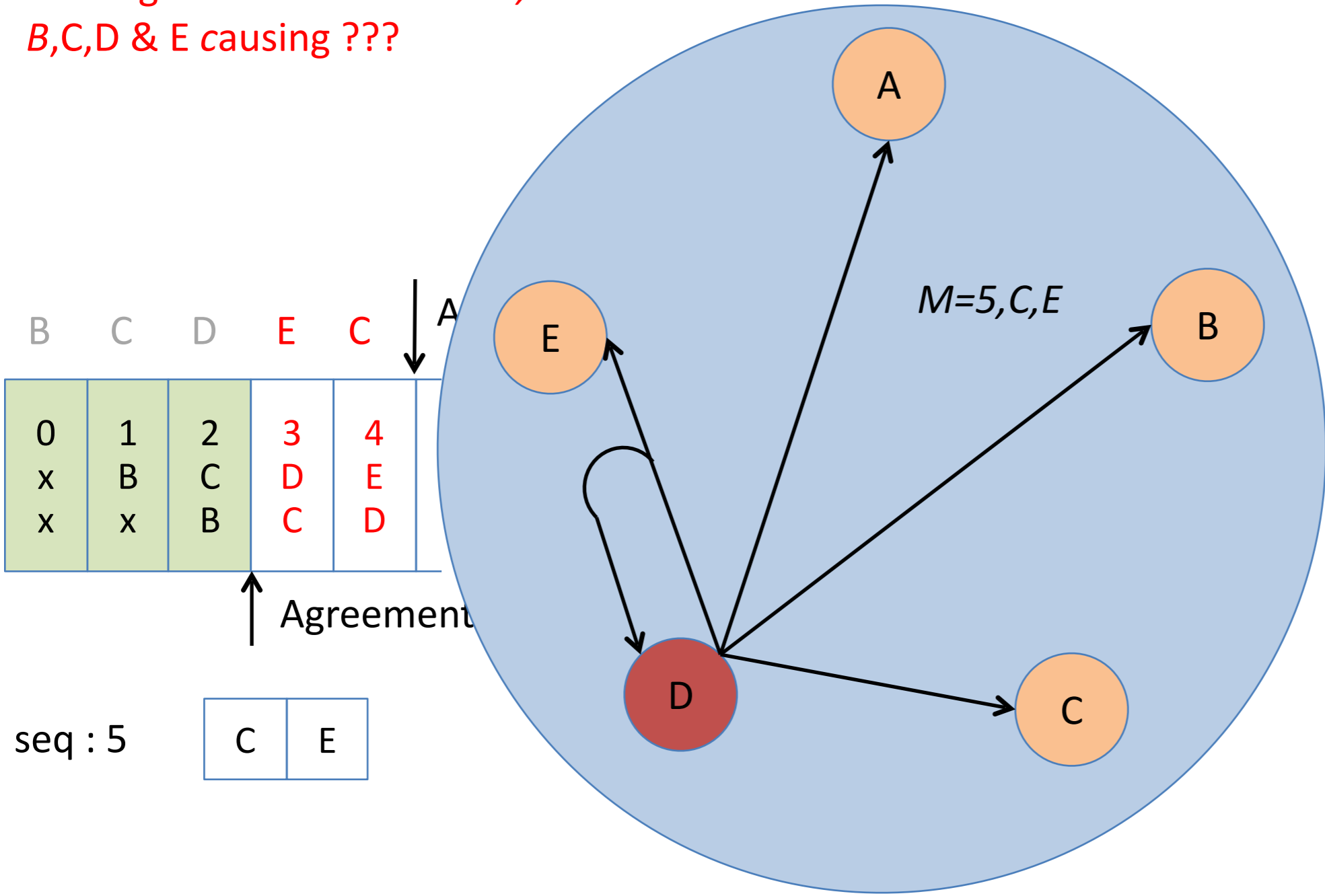
↑ Agreement

B	A
---	---



# Handling Errors?

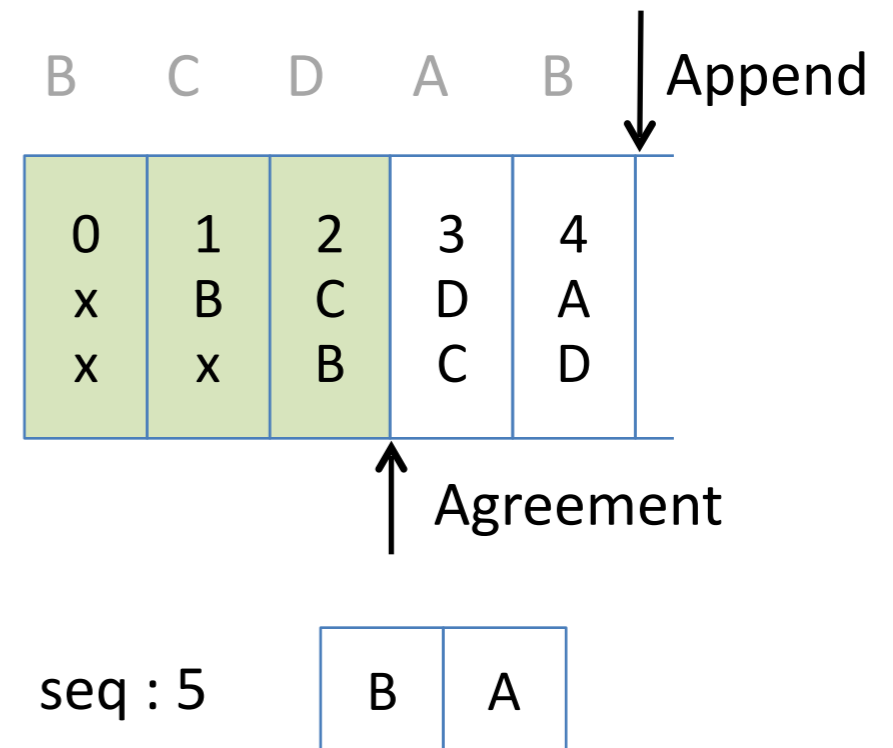
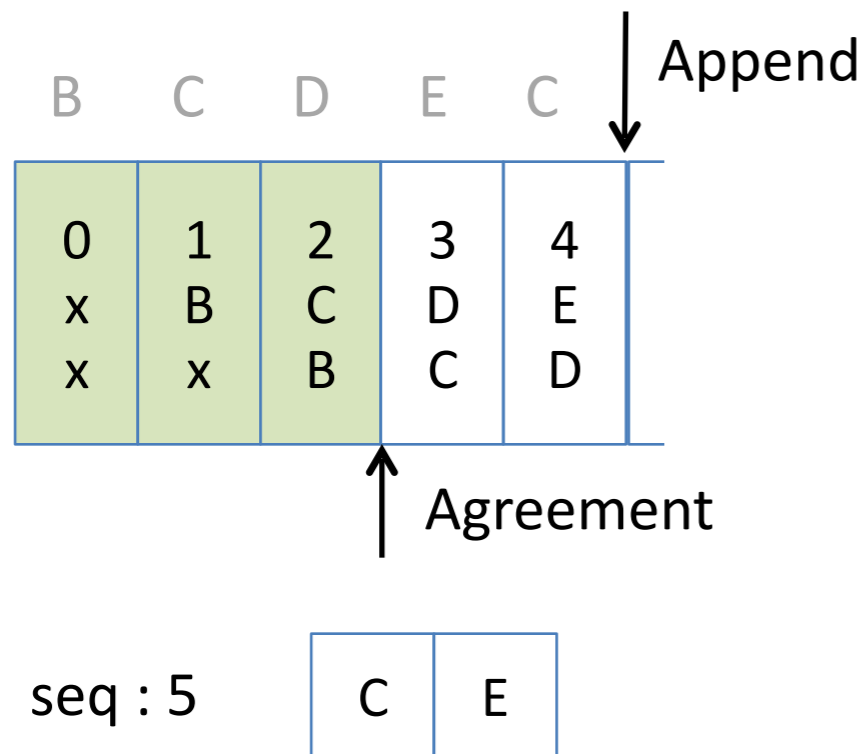
Message from *D* arrives at *A*,  
*B*, *C*, *D* & *E* causing ???





# Handling Errors?

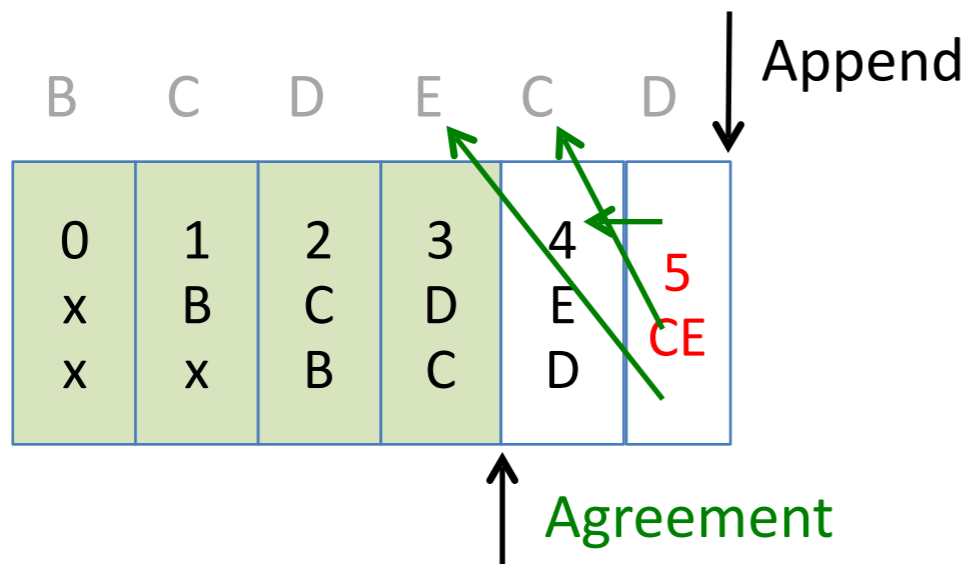
Message from *D* arrives at *A*,  
*B*,*C*,*D* & *E* causing ???



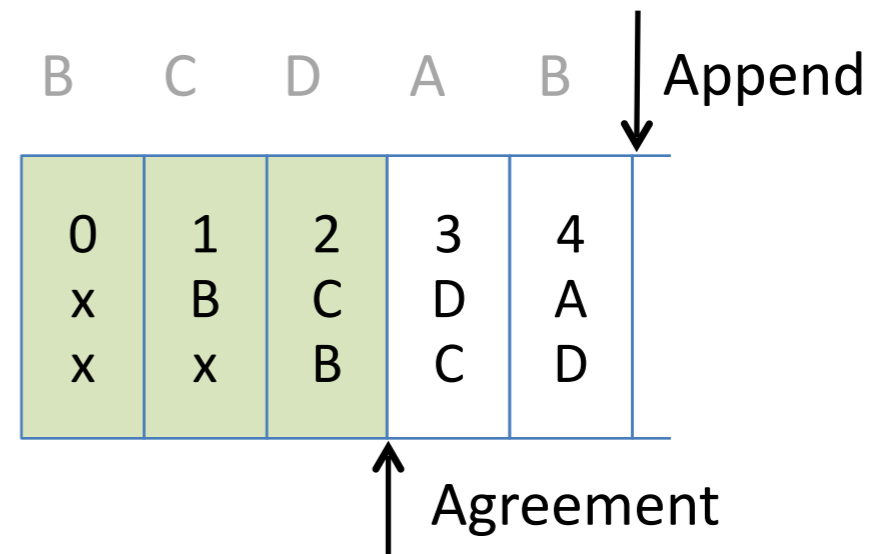


# Handling Errors?

Message from *D* arrives at *A*,  
*B*,*C*,*D* & *E* causing ???



C	E
---	---

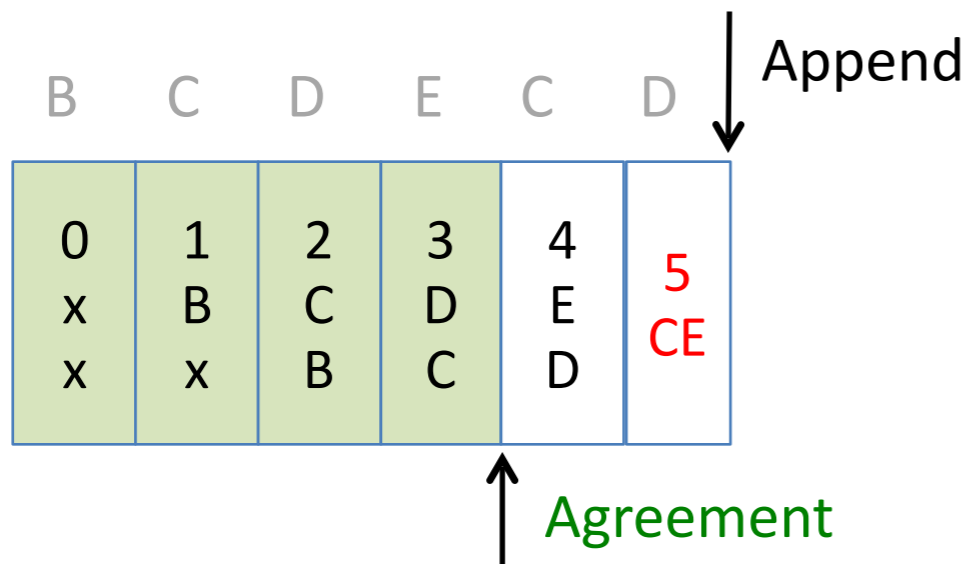


B	A
---	---

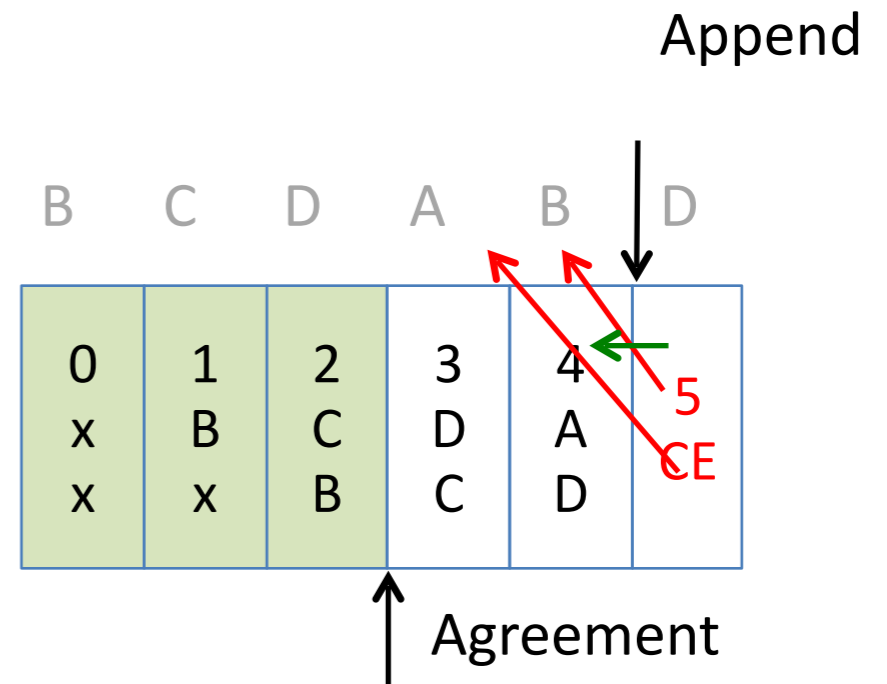
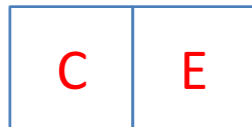


# Handling Errors?

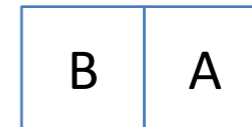
Message from *D* arrives at *A*,  
*B*,*C*,*D* & *E* causing ???



seq : 6



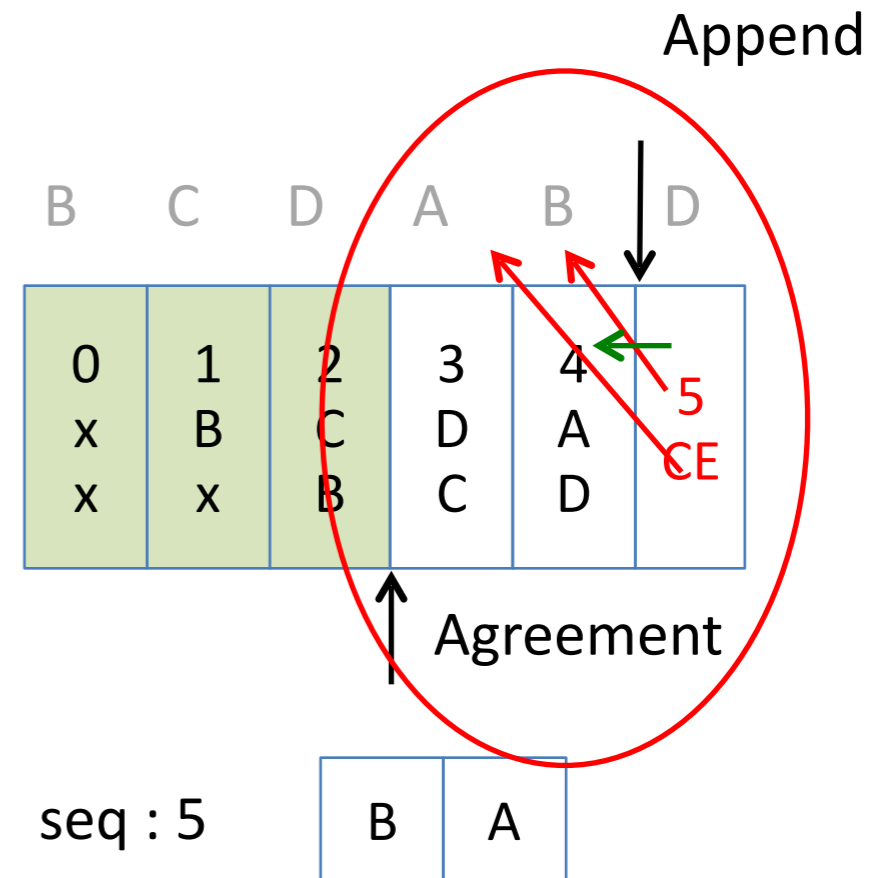
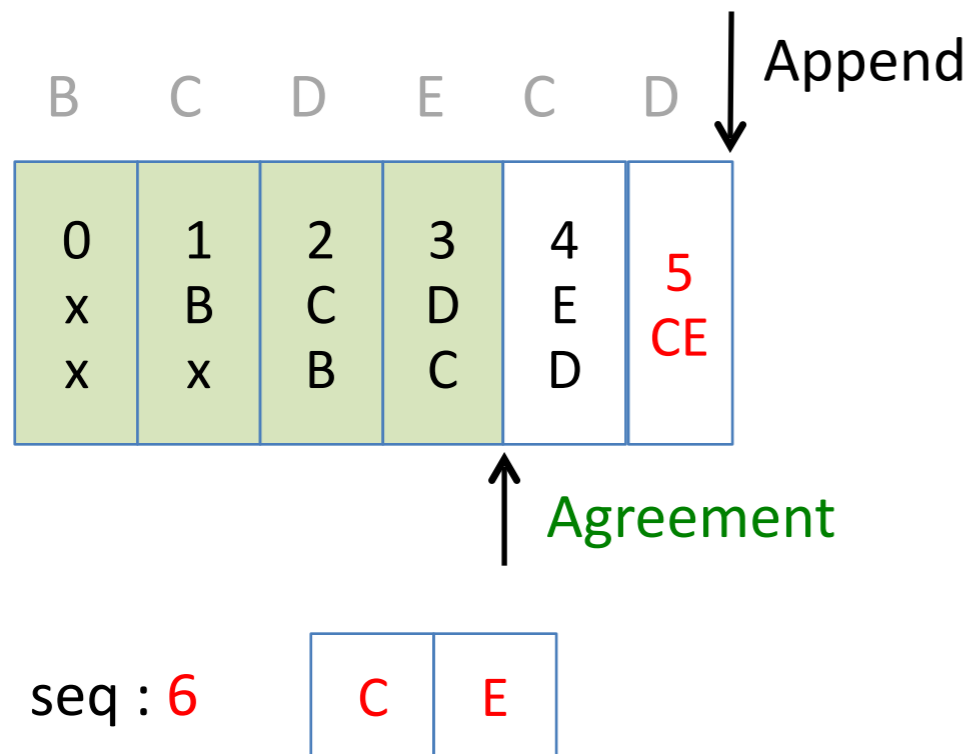
seq : 5





# Handling Errors?

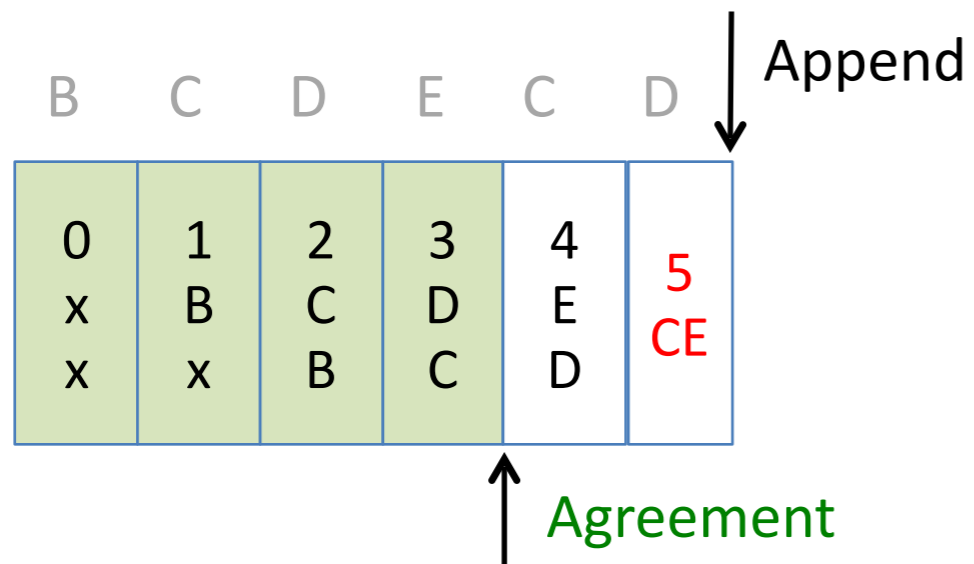
Message from *D* arrives at *A, B, C, D* & *E* causing  
 -update and agreement to *C, D* & *E*  
 -error detection for *A* & *B*





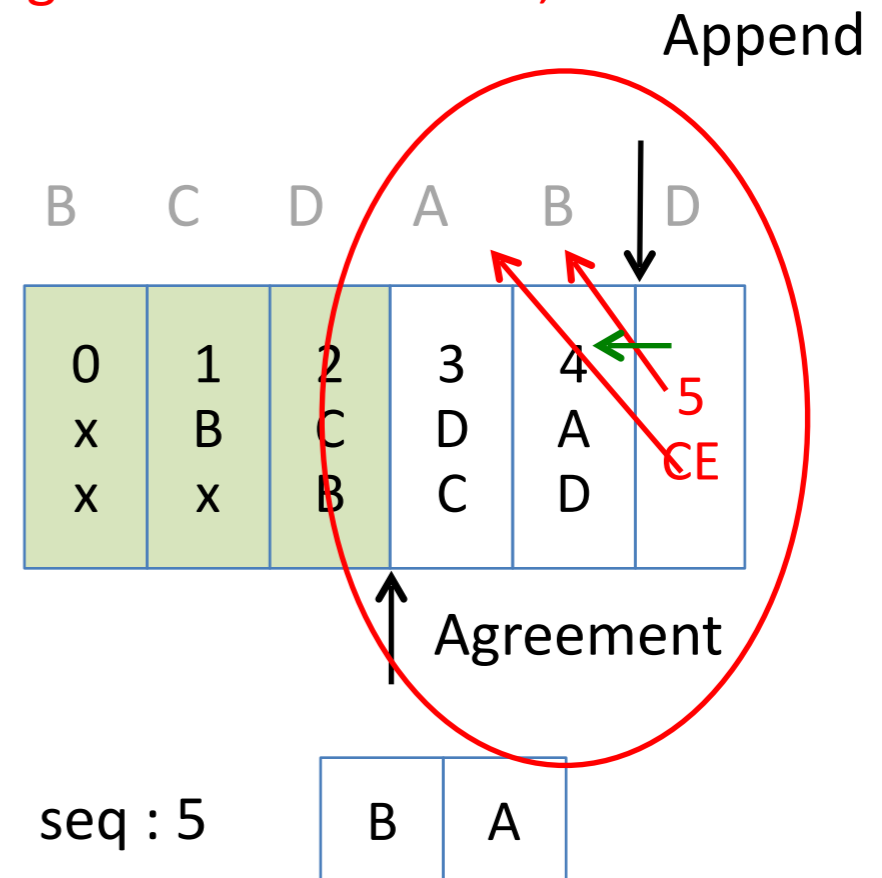
# Handling Errors?

Hosts A & B must now use the out of band network to query C, D & E about their log state for sequence 3-5. Once a majority of response are gathered is reached, host A & B can continue to move the agreement pointer



seq : 6 

C	E
---	---



seq : 5 

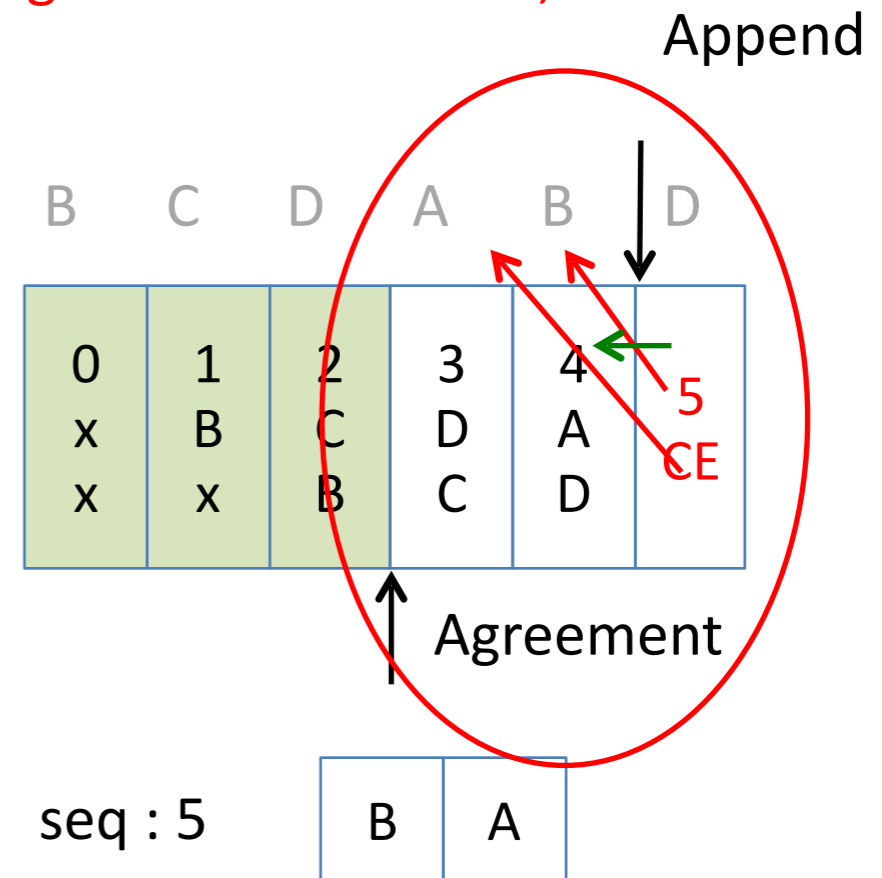
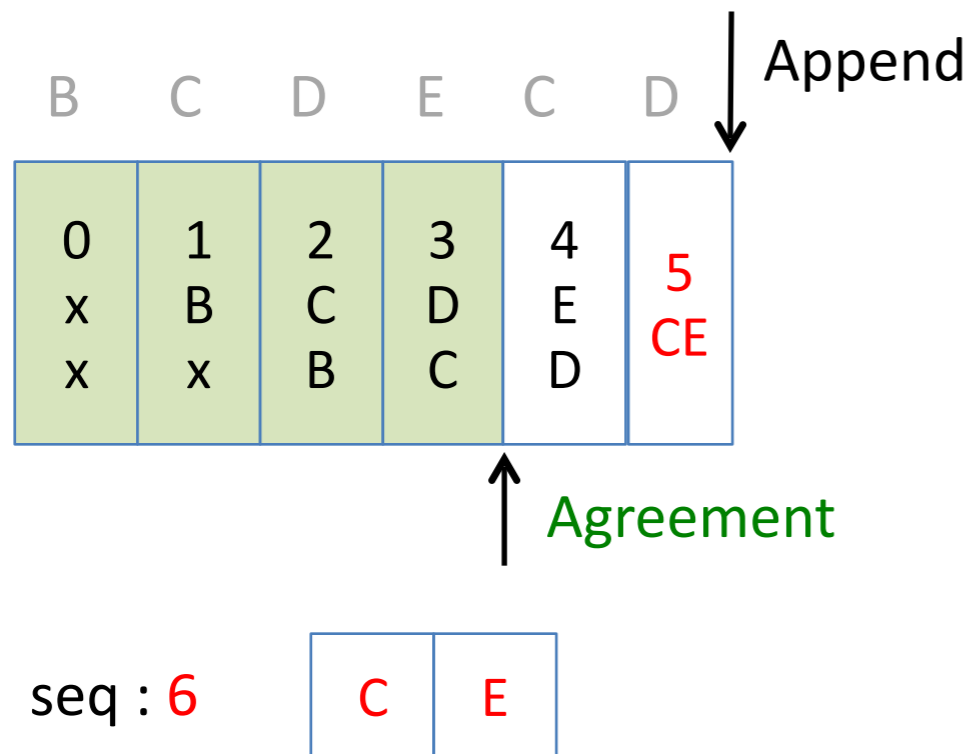
B	A
---	---





# Handling Errors?

Hosts A & B must now use the out of band network to query C, D & E about their log state for sequence 3-5. Once a majority of response are gathered is reached, host A & B can continue to move the agreement pointer



NB: New messages can be added to the log at the append pointer, but the agreement pointer cannot move until consensus has been reached