

# Rollback-Recovery for Middleboxes

Justine Sherry, SIGCOMM 2015



# Basic Concepts

Middlebox: a network appliance owning functions other than packet forwarding.

Rollback recovery: restart from recent saved state on a back-up device



---

# Correctness: Output Commit

---

Before releasing a packet: has all information reflecting that packet been committed to stable storage?

Implemented with check  
every time packet is  
released.

Gb/s flows trigger  
frequent output  
commit

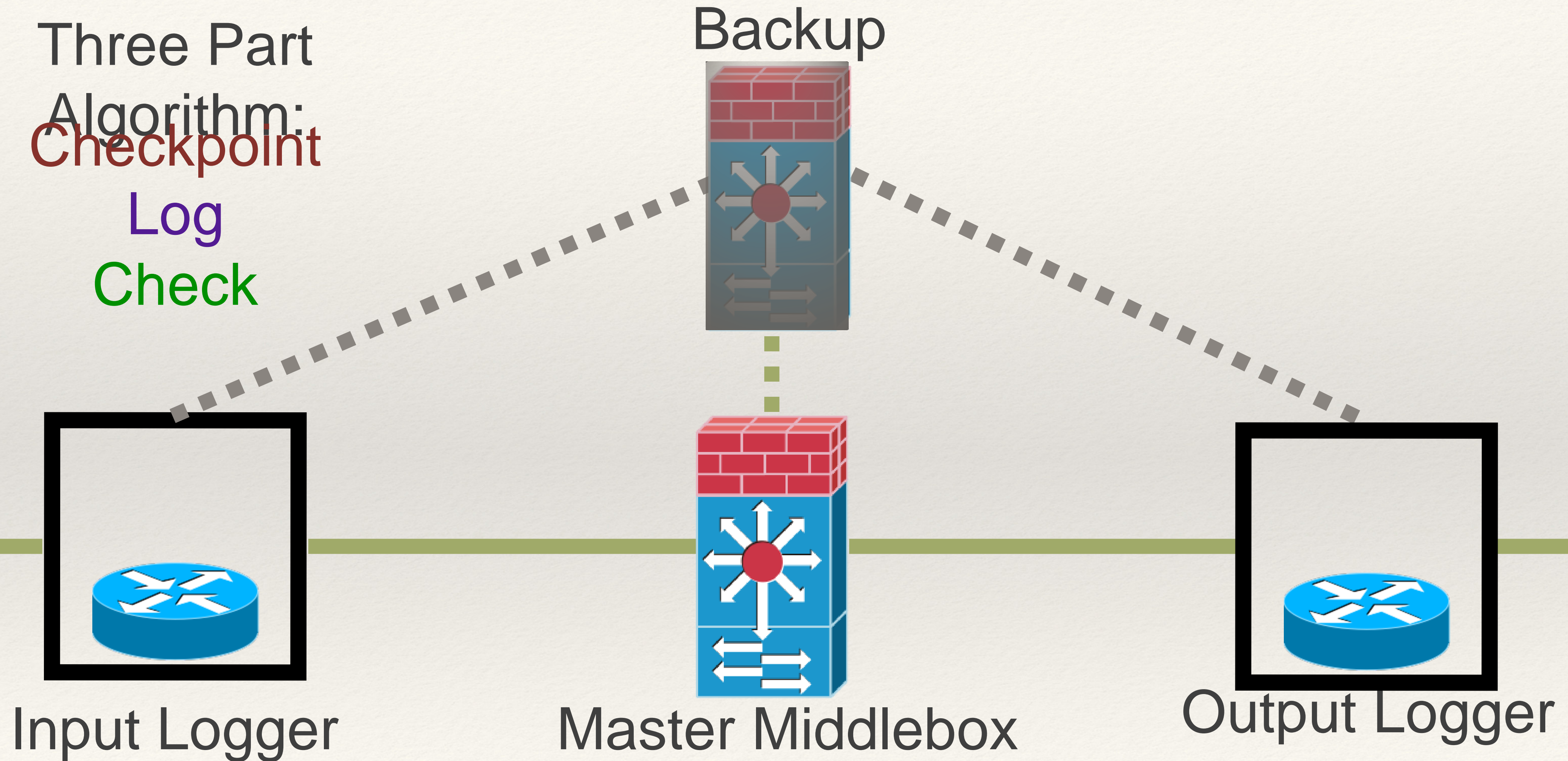


# Performance vs. Correctness



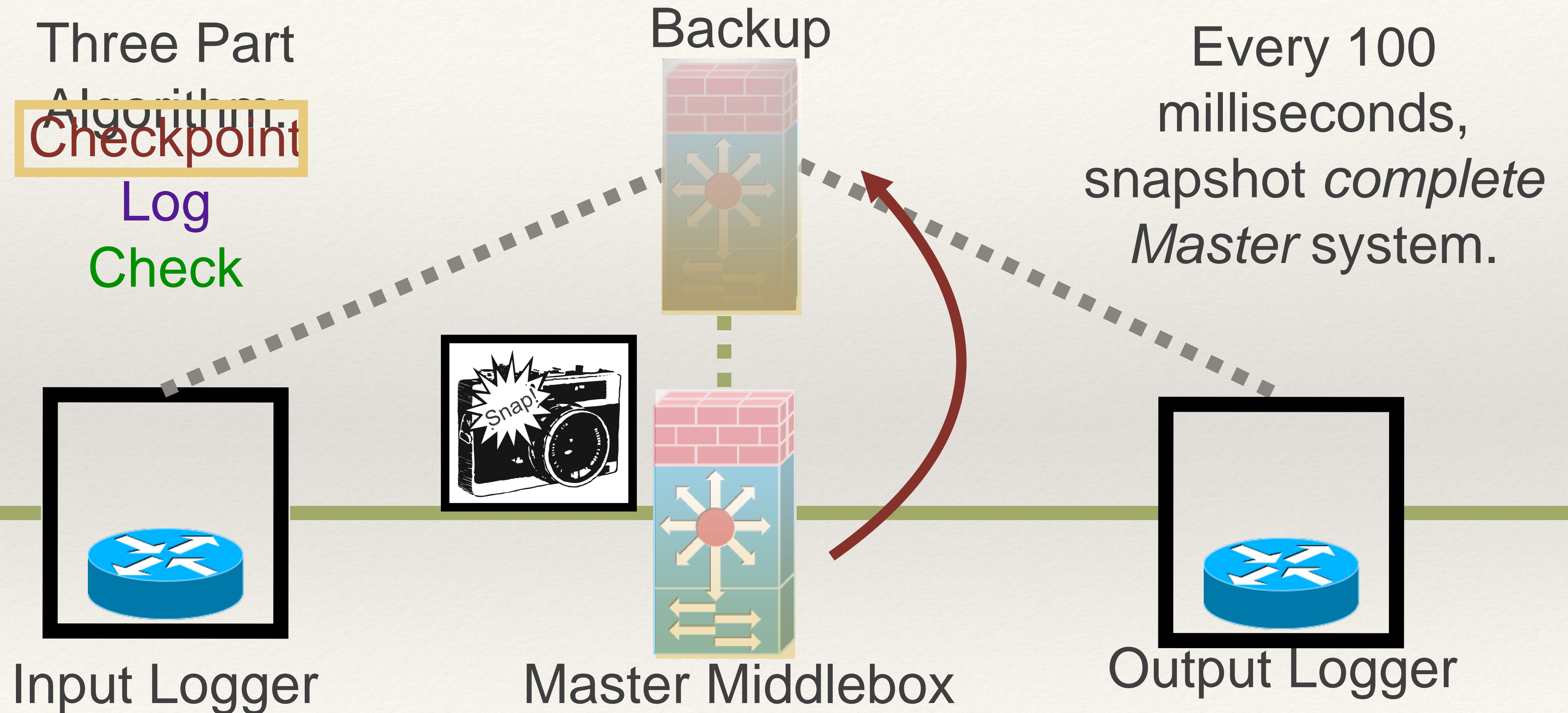
# FTMBB Implements Rollback Recovery.

Three Part  
Algorithm:  
Checkpoint  
Log  
Check



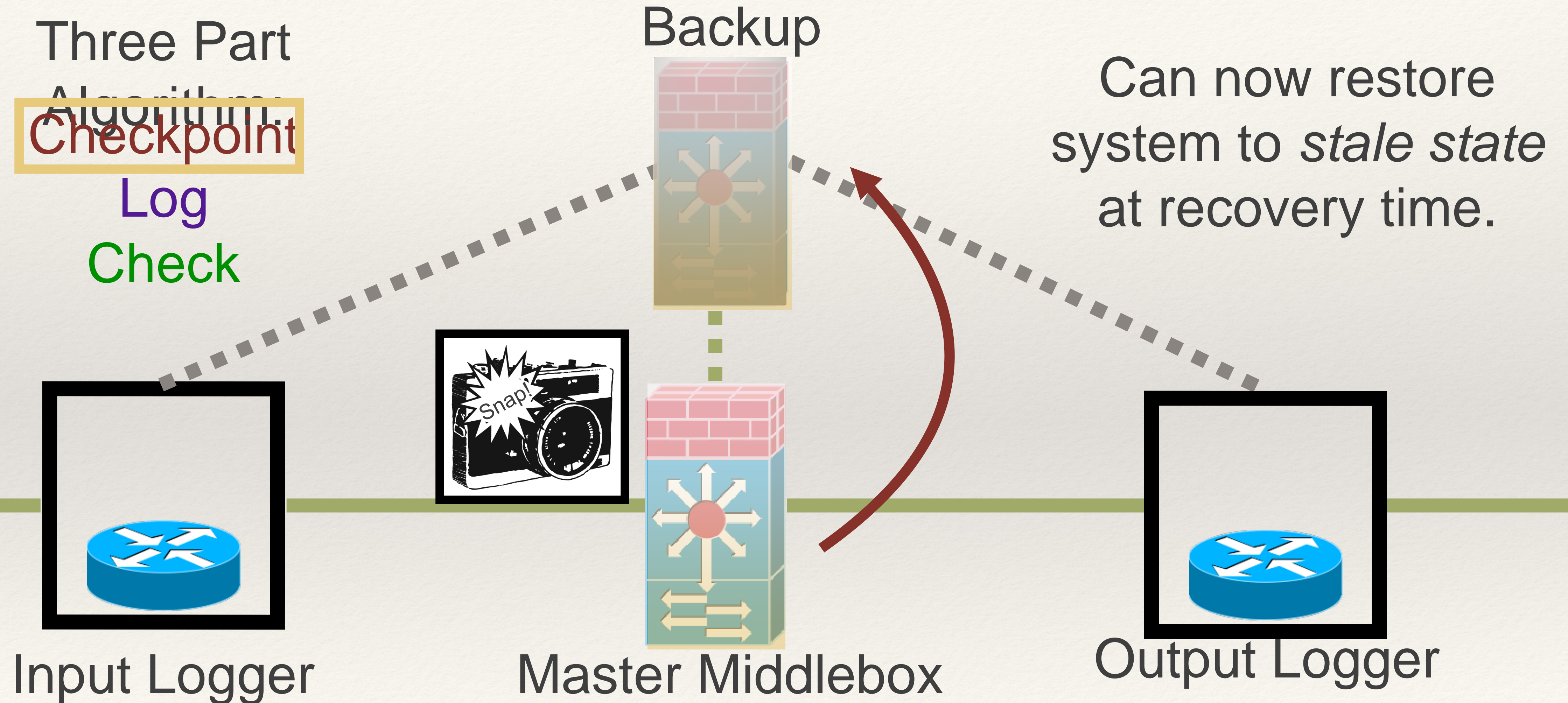


# Rollback Recovery



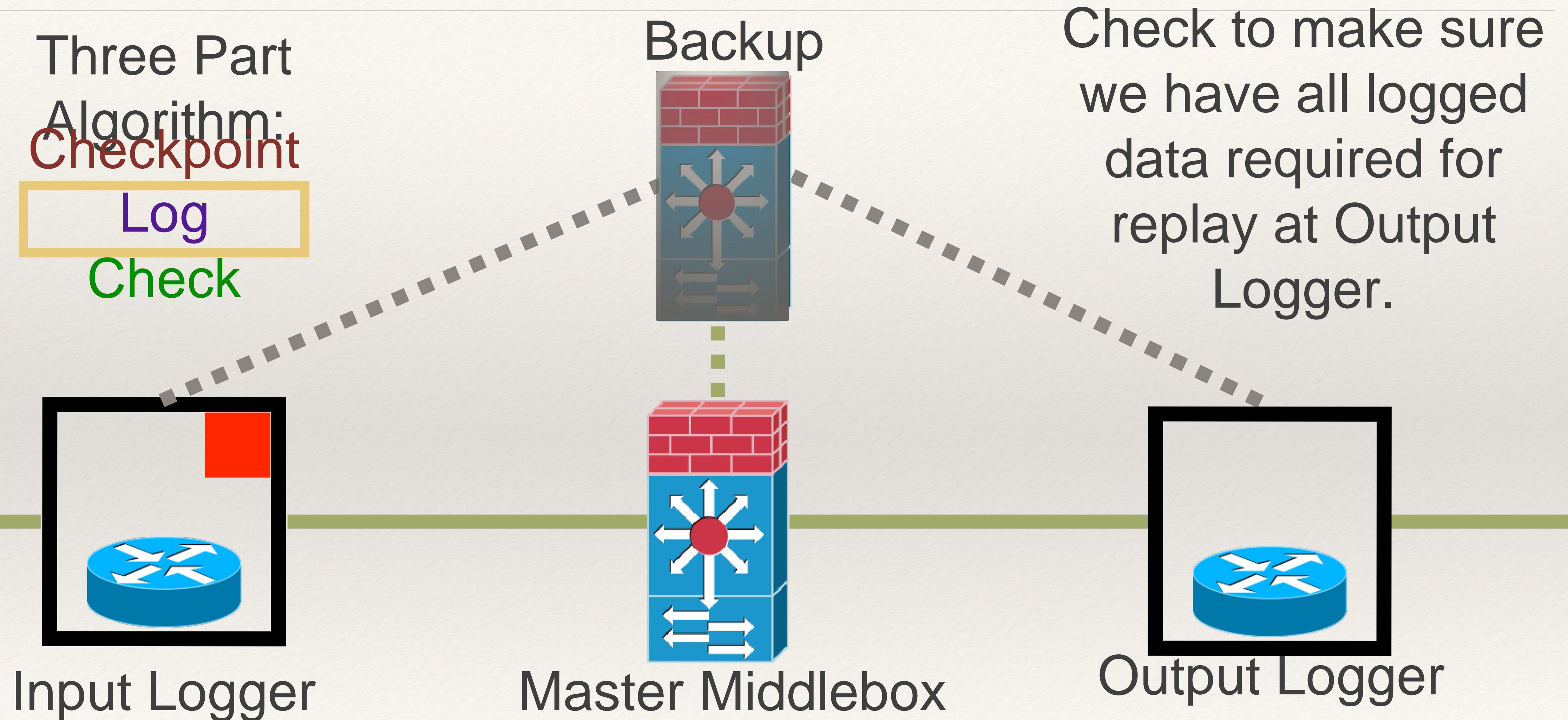


# Rollback Recovery





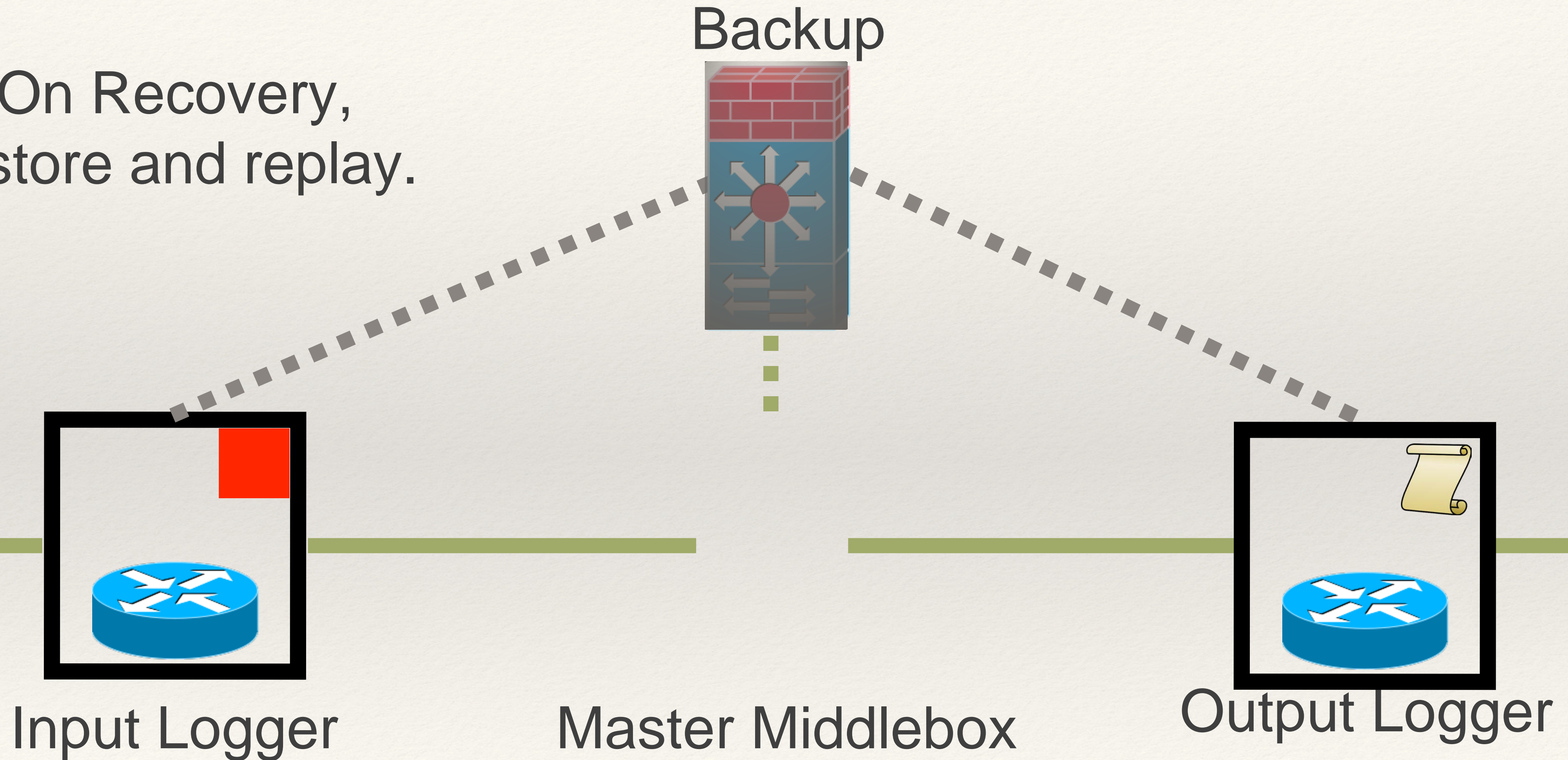
# Rollback Recovery





# Rollback Recovery

On Recovery,  
restore and replay.





# Rollback Recovery

Three Part  
Algorithm:  
Checkpoint



Open Questions:

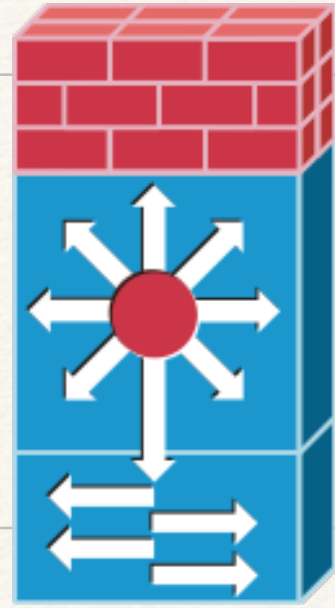
(1) What do we need to log for correct replay?

- A classically hard problem due to nondeterminism.

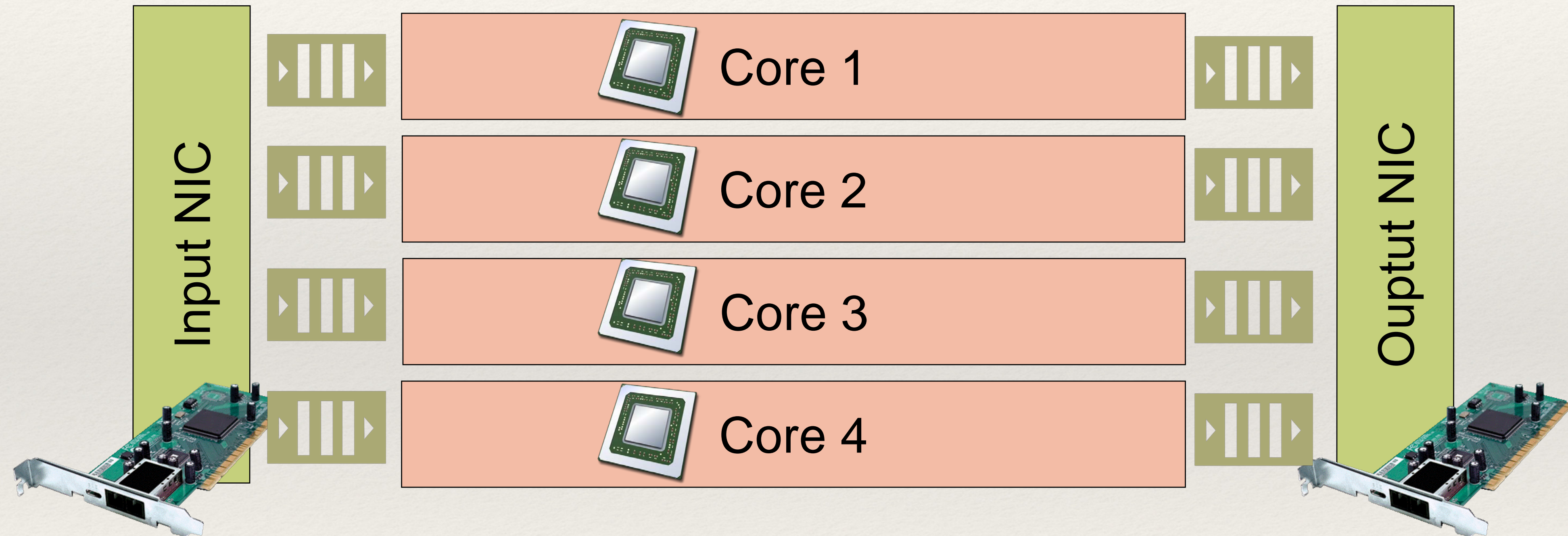
(2) How do we check that we have everything we need to replay a given packet?

- Need to monitor system state that is updated frequently and on multiple cores.

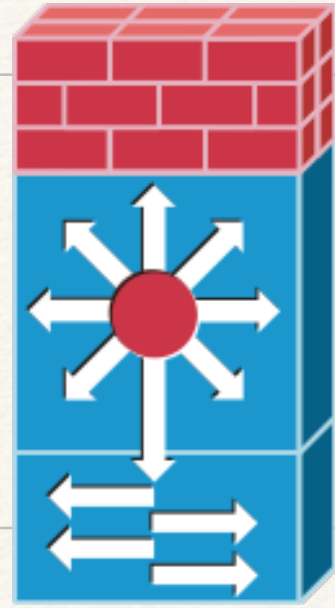




# Middlebox Architecture

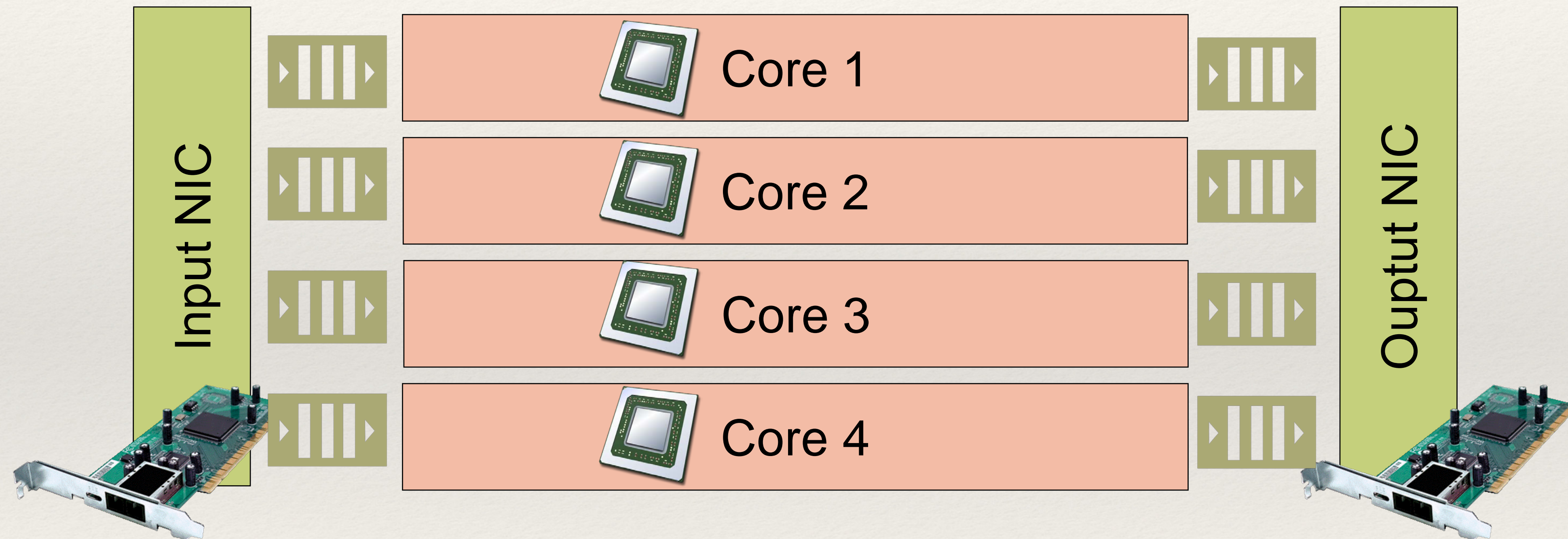






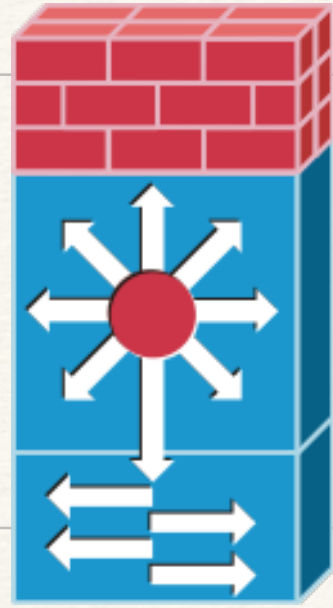
# Middlebox Architecture

Input NIC “hashes” incoming packets to cores.



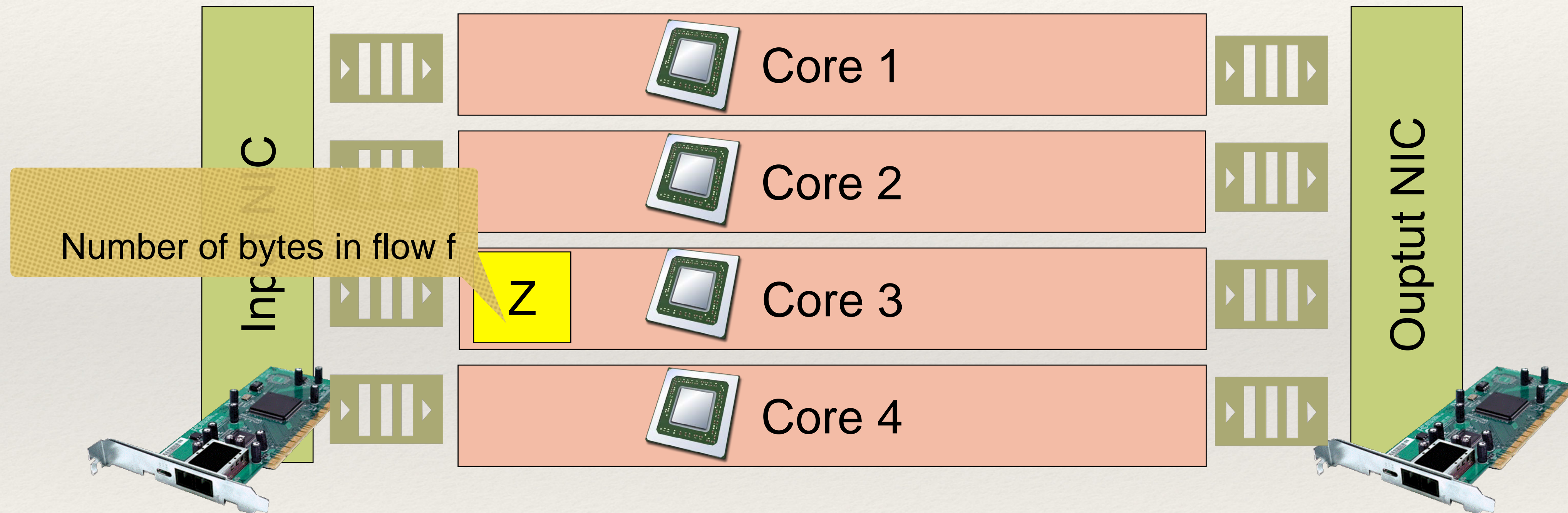
All packets from same flow are processed by same core.





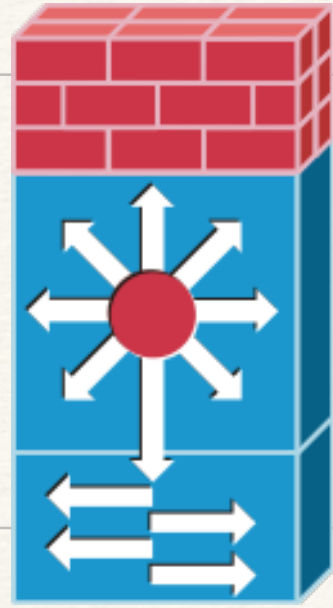
# Middlebox Architecture: State

Local state: only relevant to one connection.

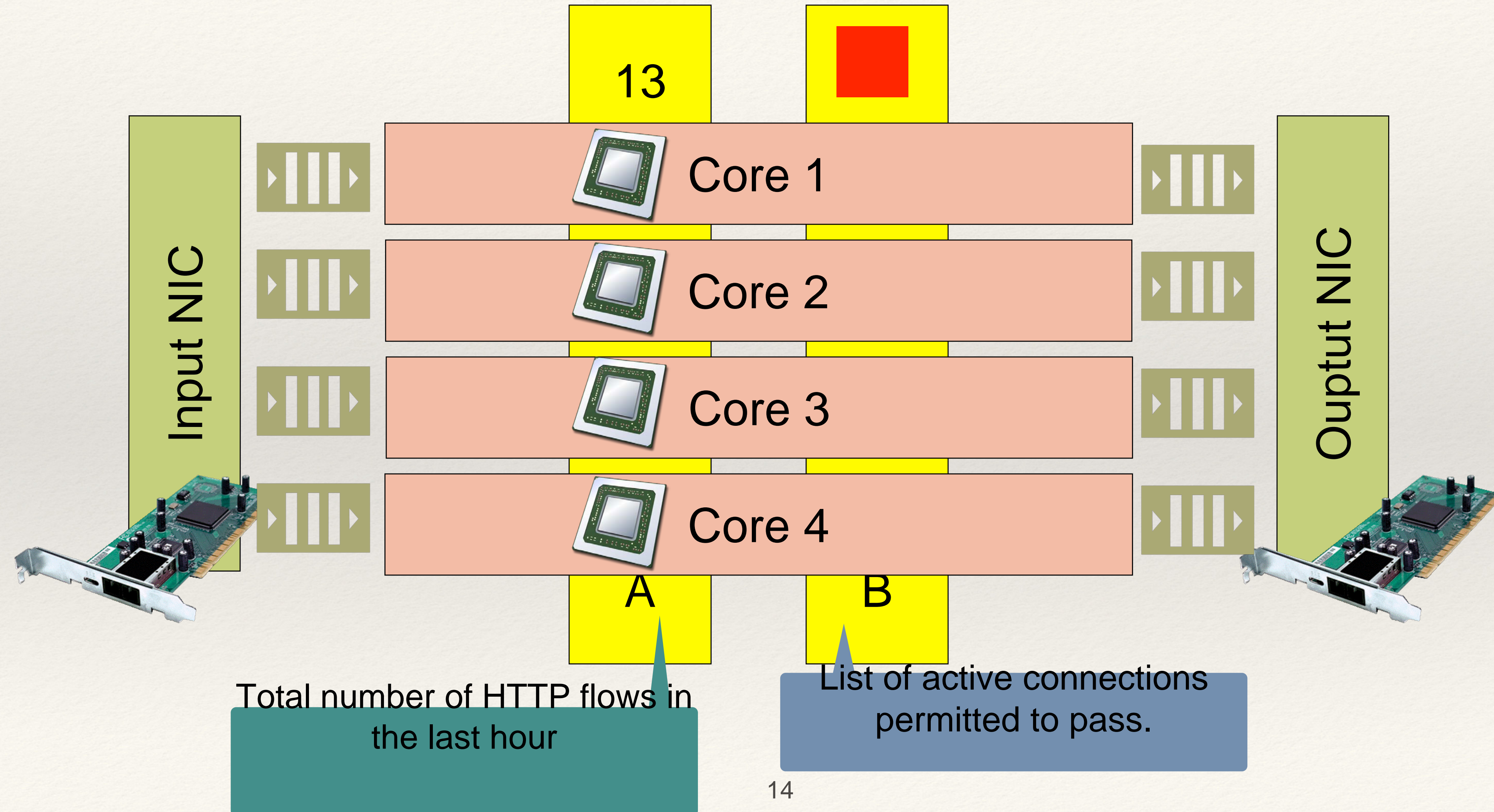


Accessing local state is *fast* because only one core “owns” the data.

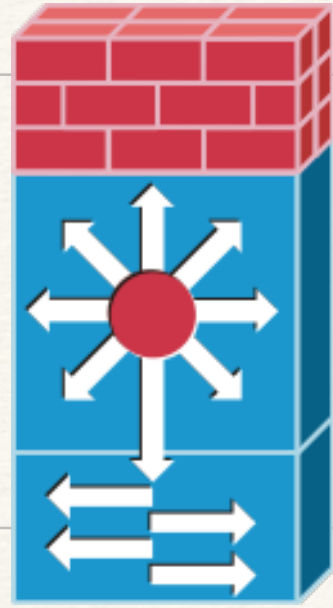




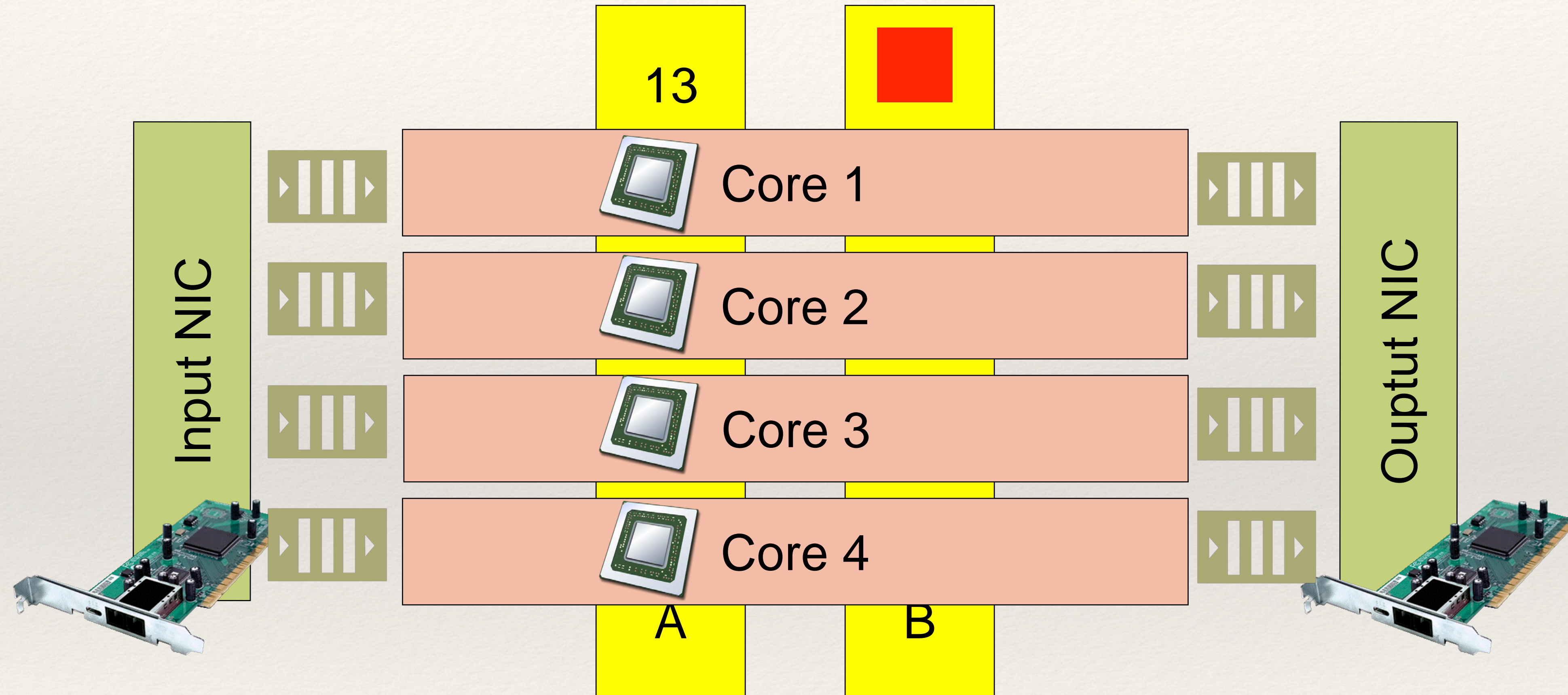
# Middlebox Architecture: State







# Middlebox Architecture: State



Reading shared state is slower.  
Writing is most expensive because it can cause *contention*!



# Rollback Recovery

Three Part  
Algorithm:

Checkpoint

Log

Rollback

Open Questions:

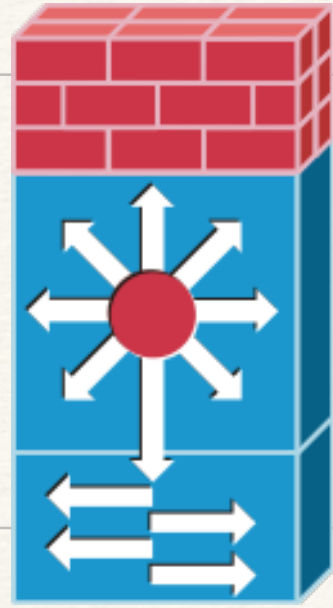
(1) What do we need to log for correct replay?

- A classically hard problem due to nondeterminism.

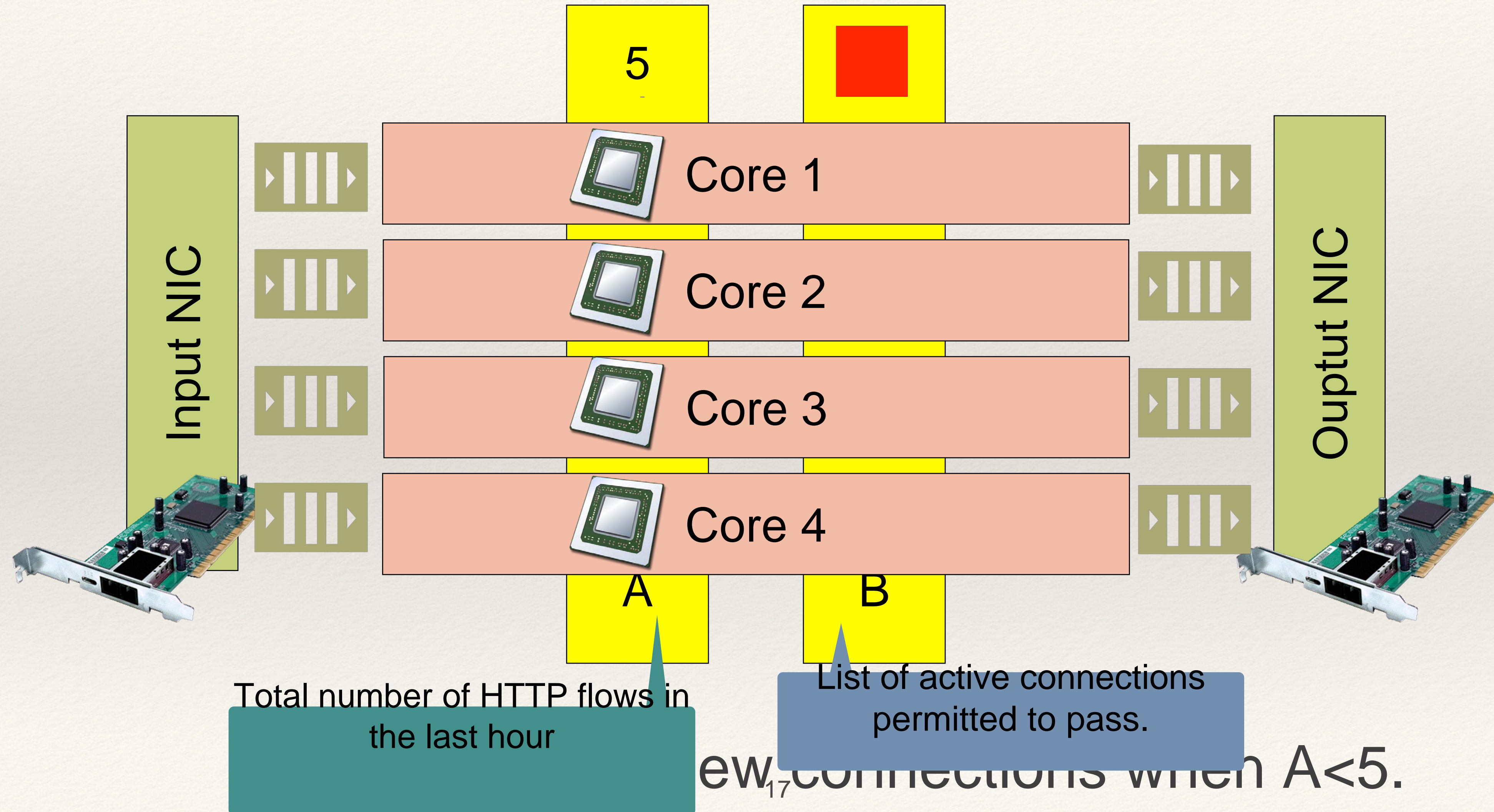
(2) How do we check that we have everything we need to replay a given packet?

- Need to monitor system state that is updated frequently and on multiple cores.

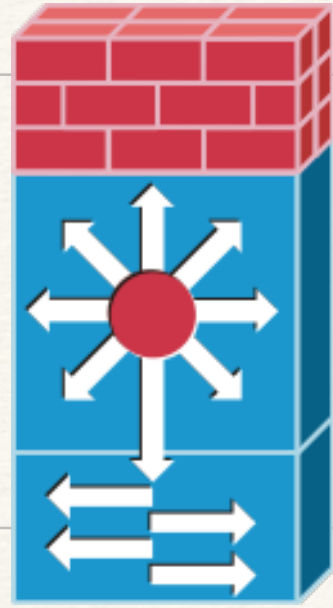




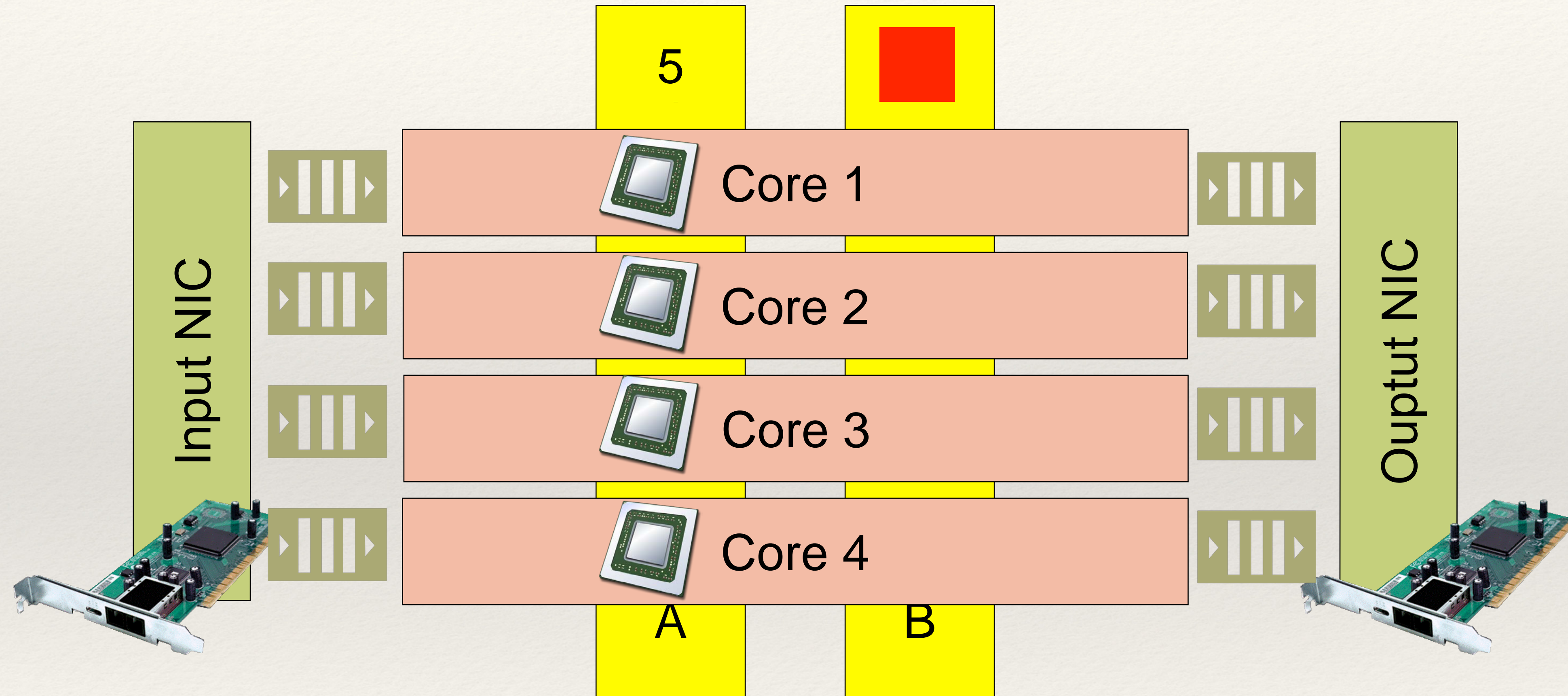
# Parallelism + Shared State





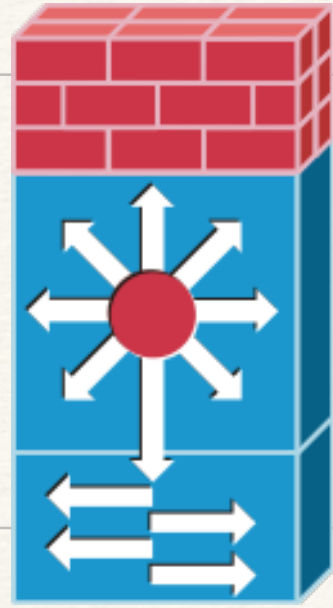


# Parallelism + Shared State

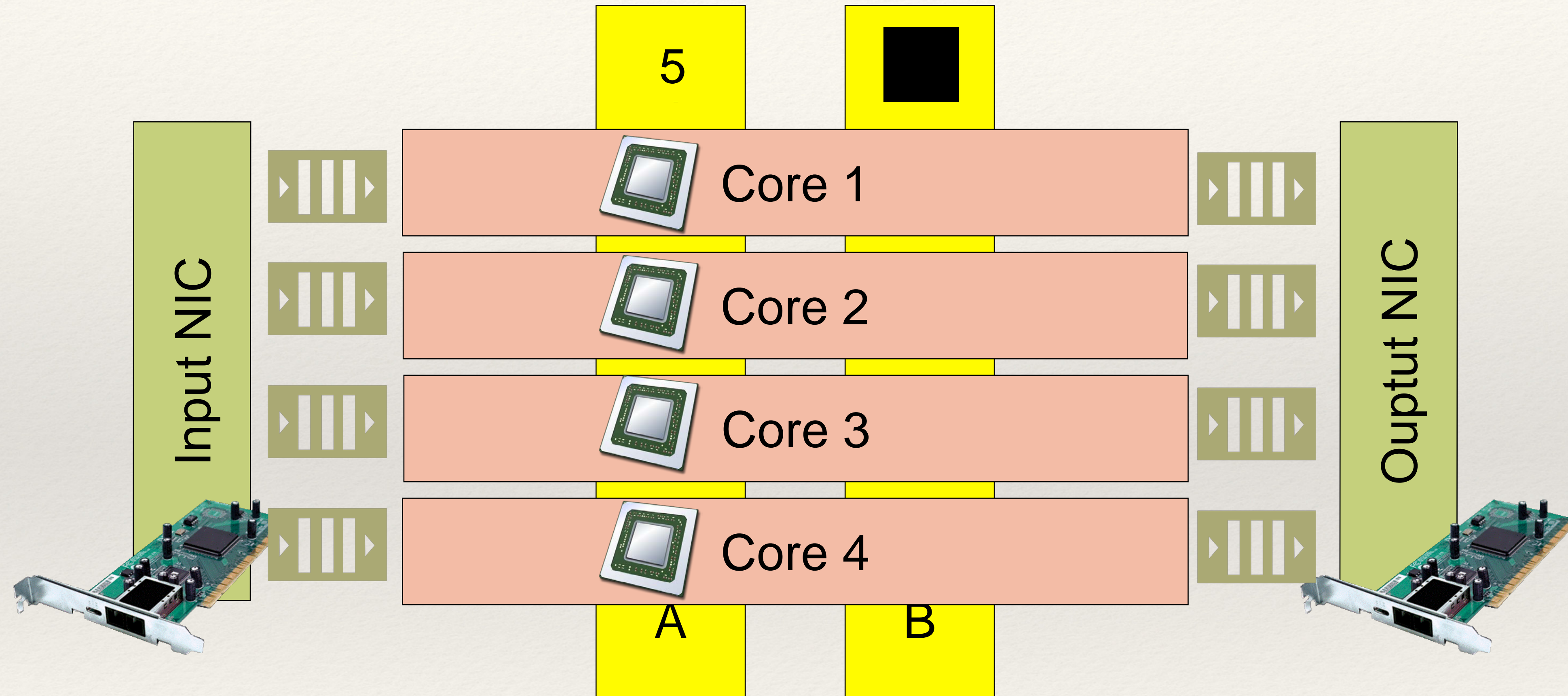


MB Rule: allow new connections, unless  $A \geq 5$ .





# Parallelism + Shared State

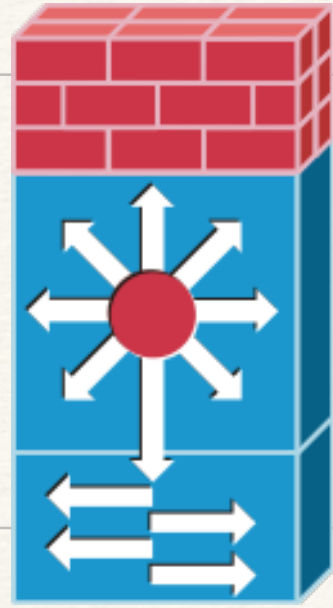


MB Rule: allow new connections, unless  $A \geq 5$ .

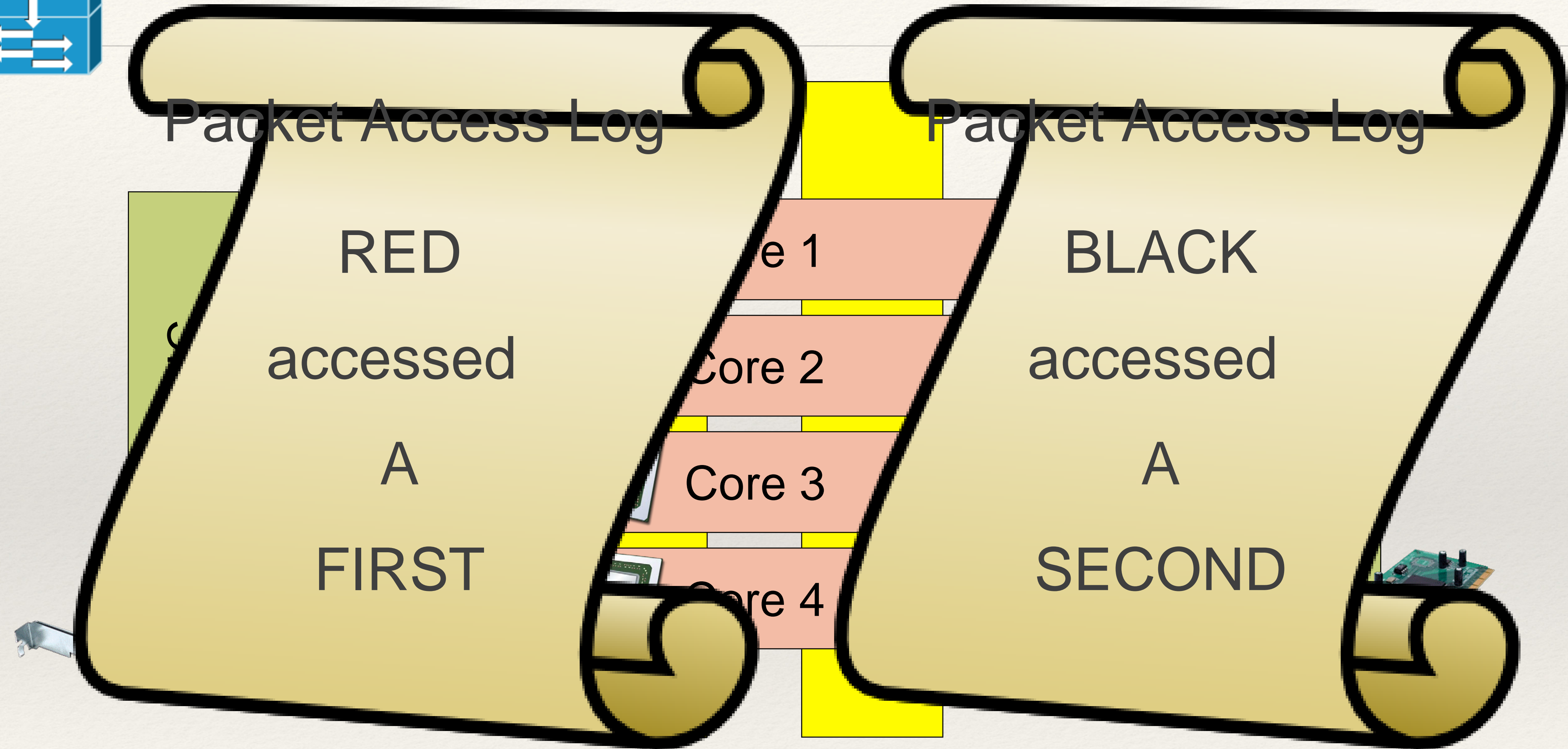


FTMB logs all accesses to shared state  
using Packet Access Logs.





# Parallelism + Shared State





# Rollback Recovery

Three Part  
Algorithm:

Checkpoint

Log  
Check

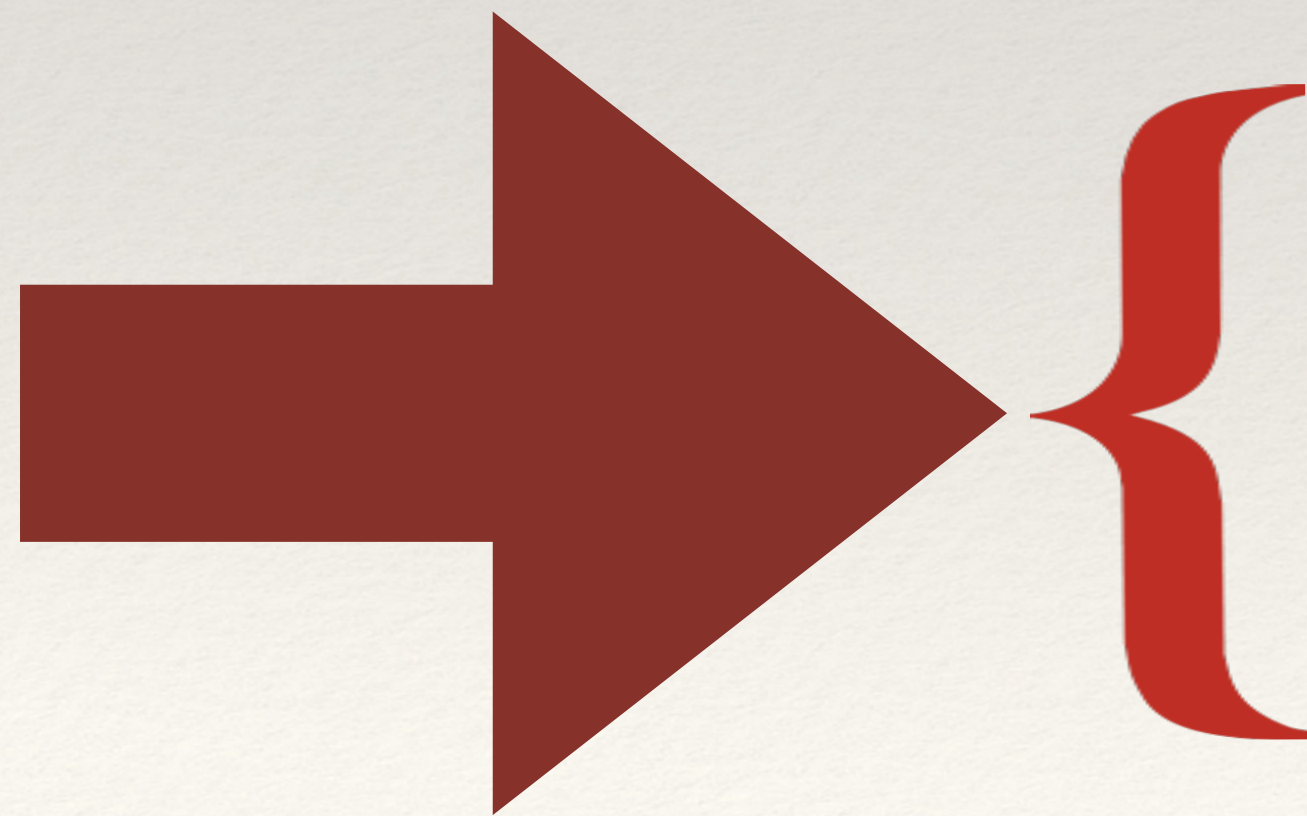
Open Questions:

(1) What do we need to log for correct replay?

- Packet Access Logs record accesses to shared state.

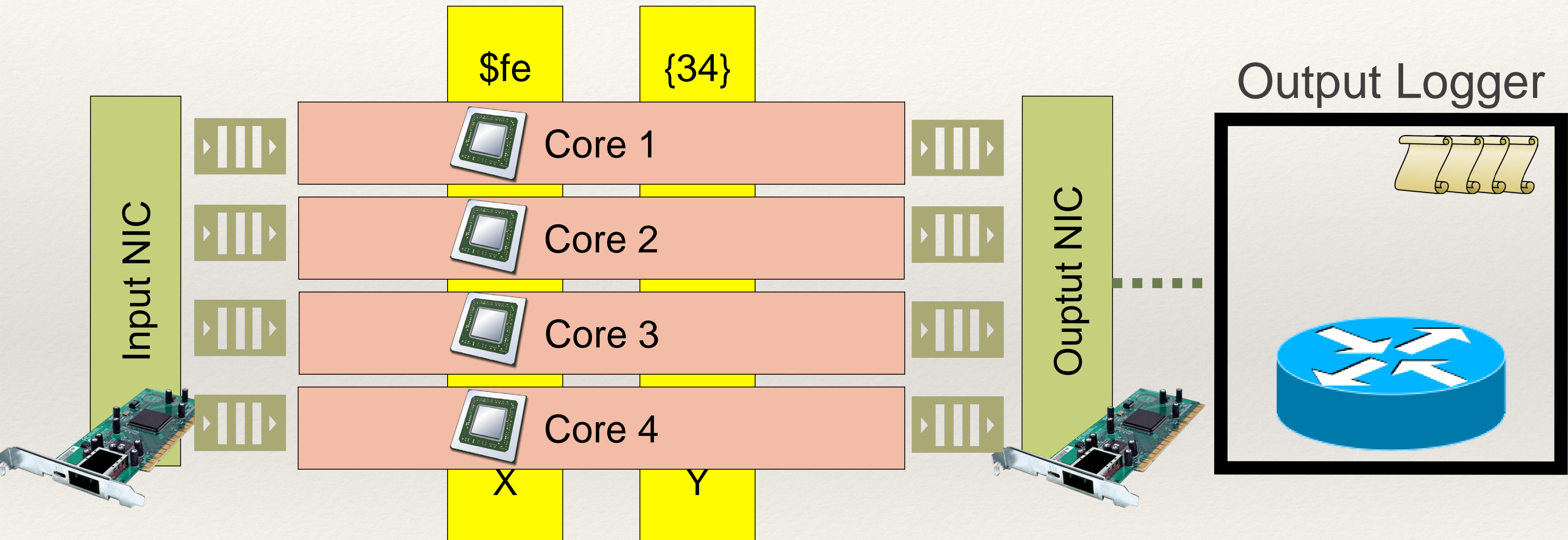
(2) How do we check that we have everything we need to replay a given packet?

- Need to monitor system state that is updated frequently and on multiple cores.



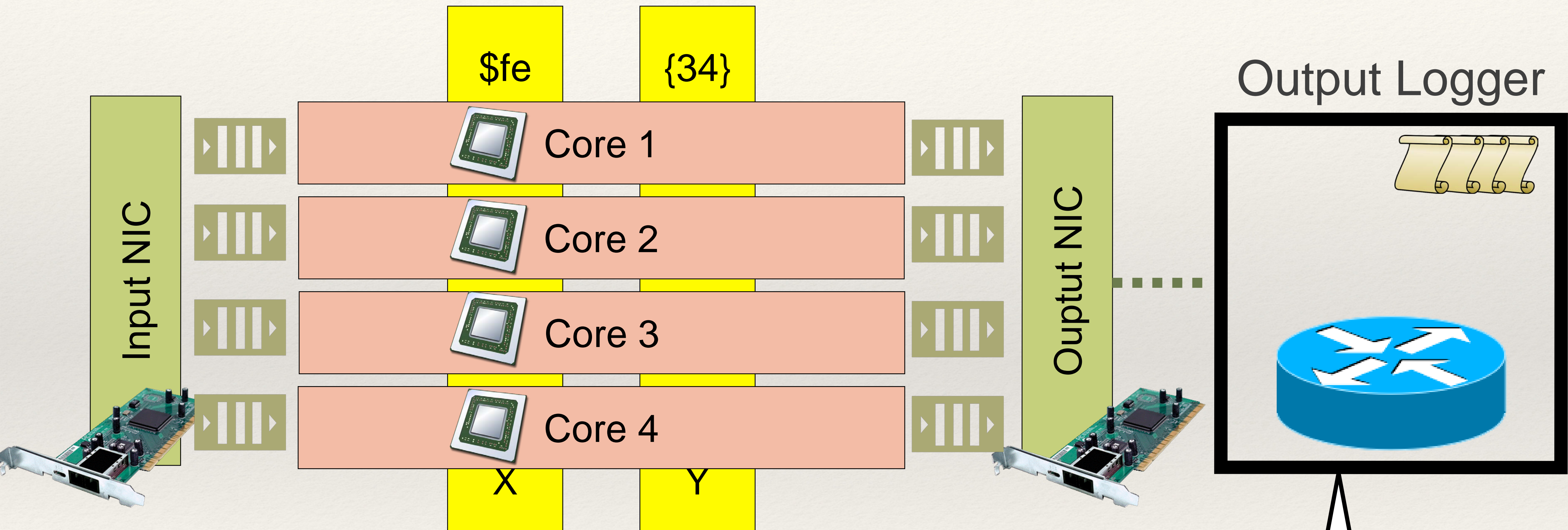


# Checking for Safe Release





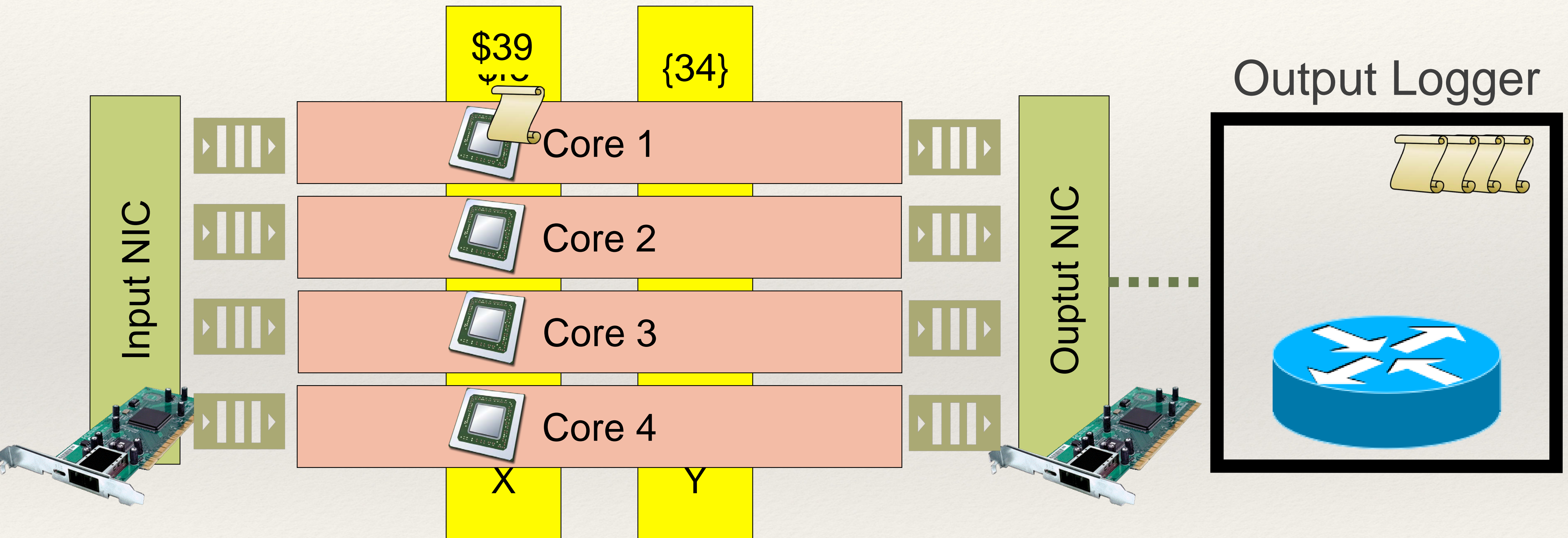
# Checking for Safe Release



Do I have all PALs so that I can replay the system?



# Checking for Safe Release

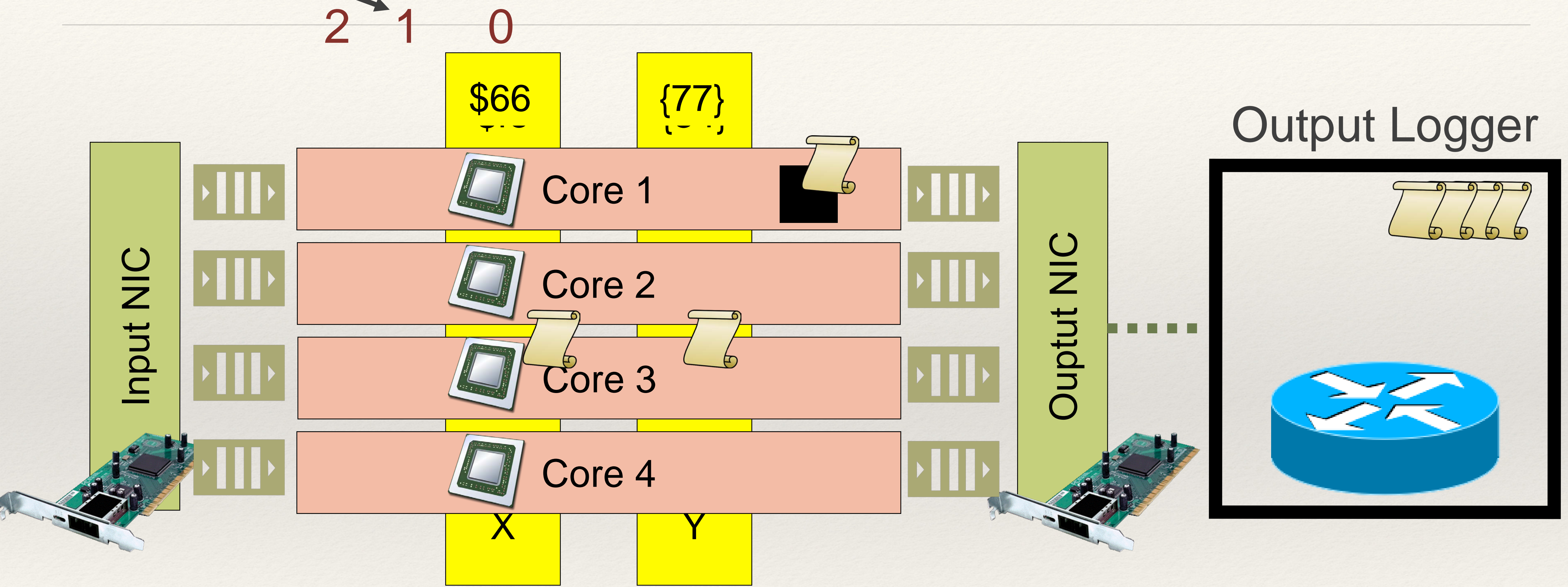


If black packet were released now, would only need PAL {X, Black, First



Need to read!

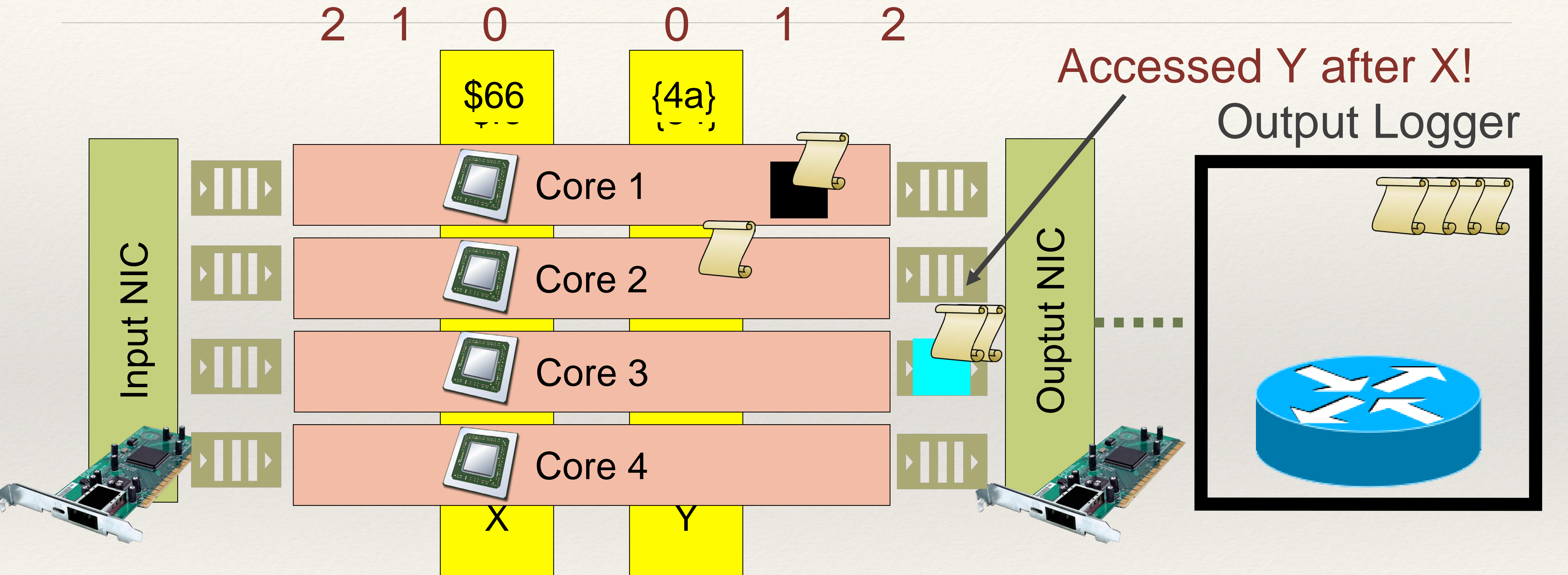
# Checking for Safe Release



If blue packet were released now, would need its own PALs, and {X, Black



# Checking for Safe Release

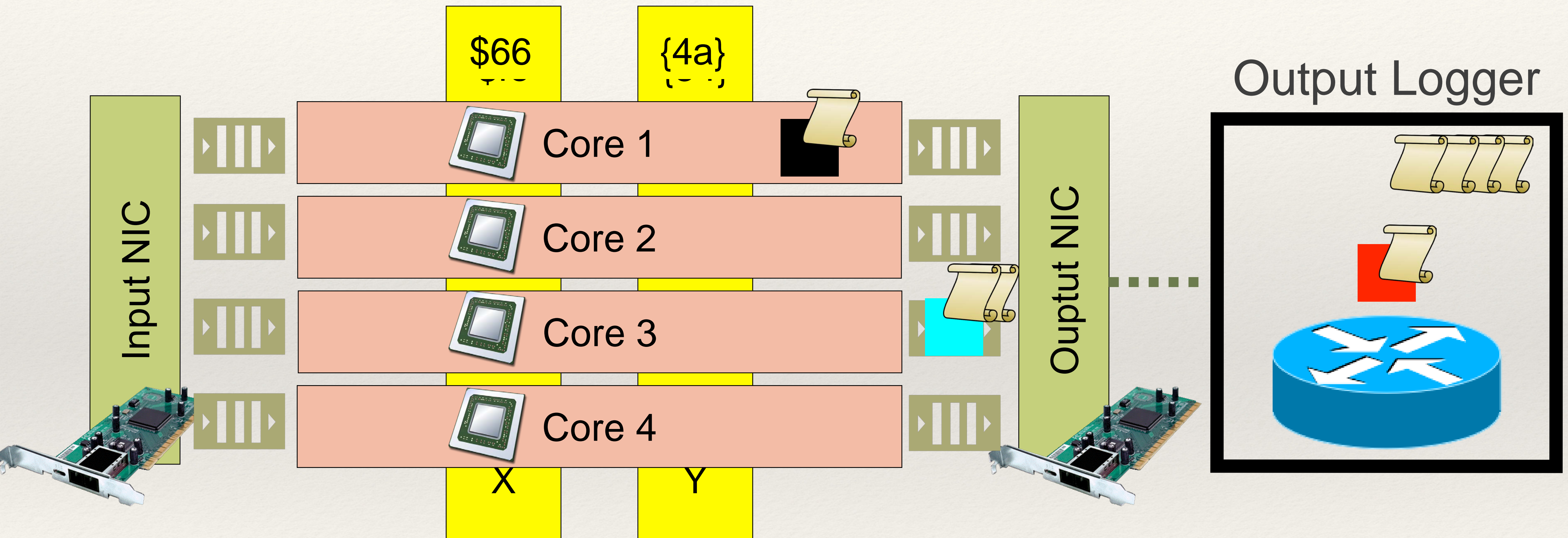


Red packet needs its own PAL, and {Blue, Y, First}

27...and {Blue, X, 2nd} and {Black, X, Fi



# Checking for Safe Release



Can depend on PALs from different cores & variables, causing a lot of ex



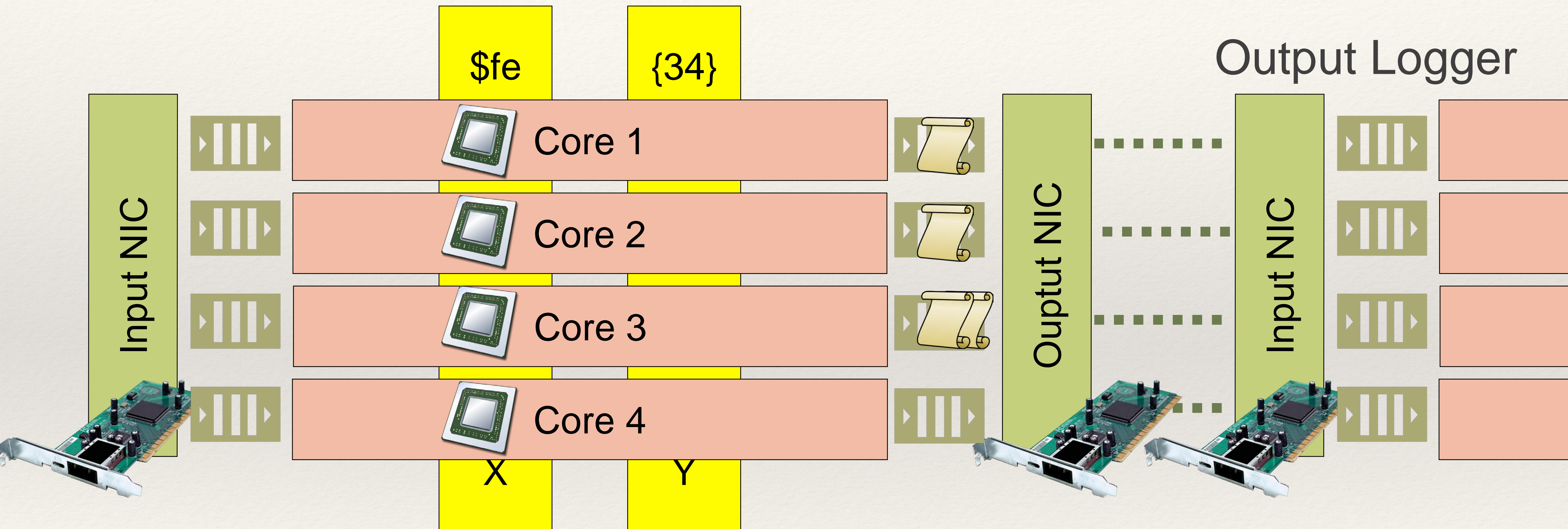
---

---

# Ordered Logging and Parallel Release



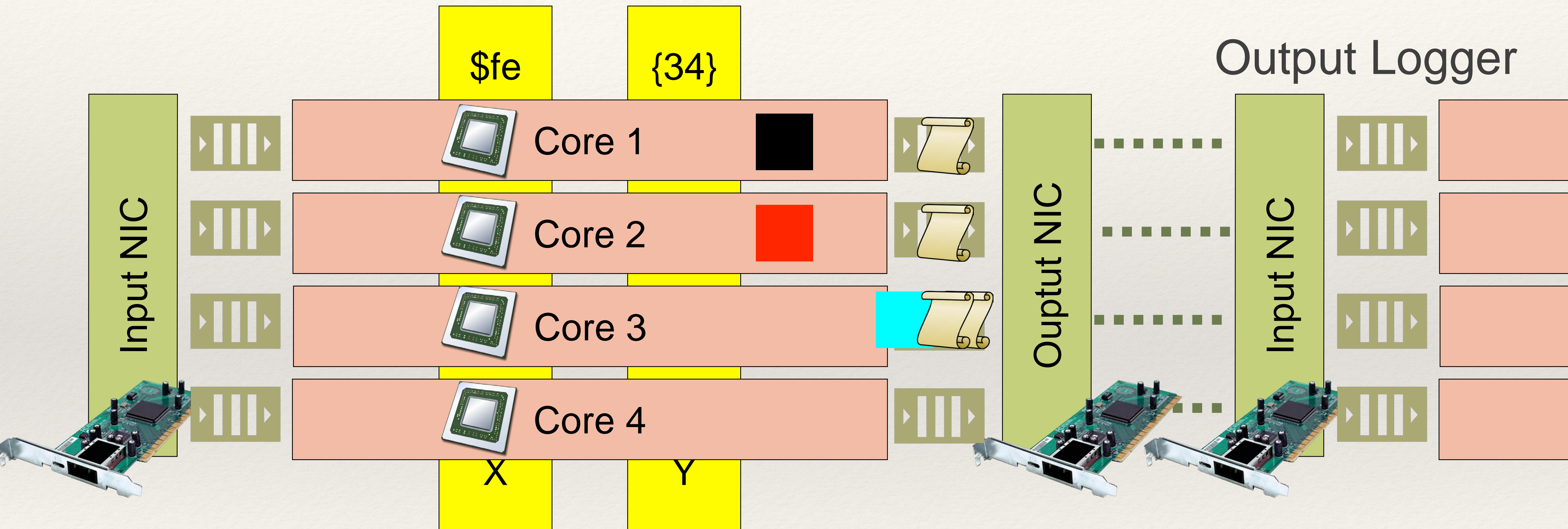
# Ordered Logging



PALs are written to output queues immediately when created.



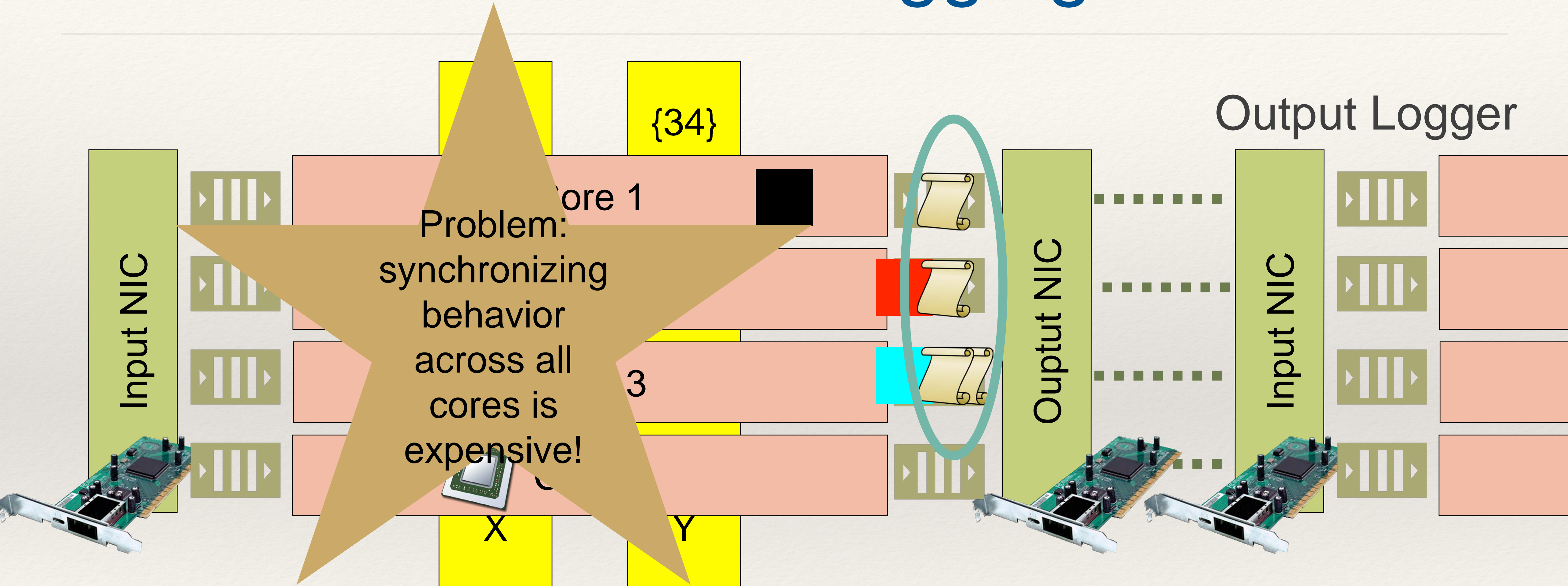
# Ordered Logging



at output queue, all PALs it depends on are already enqueued; or are already



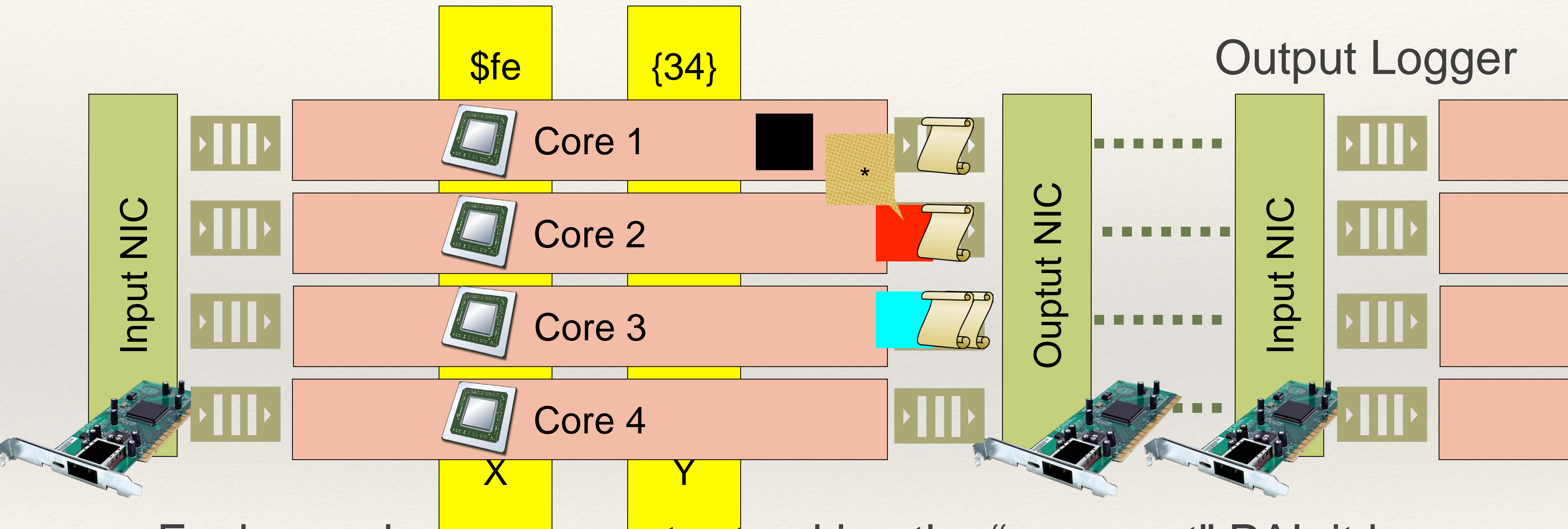
# Ordered Logging



What we want: “flush” all PALs to Output Logger. Then we’re done!



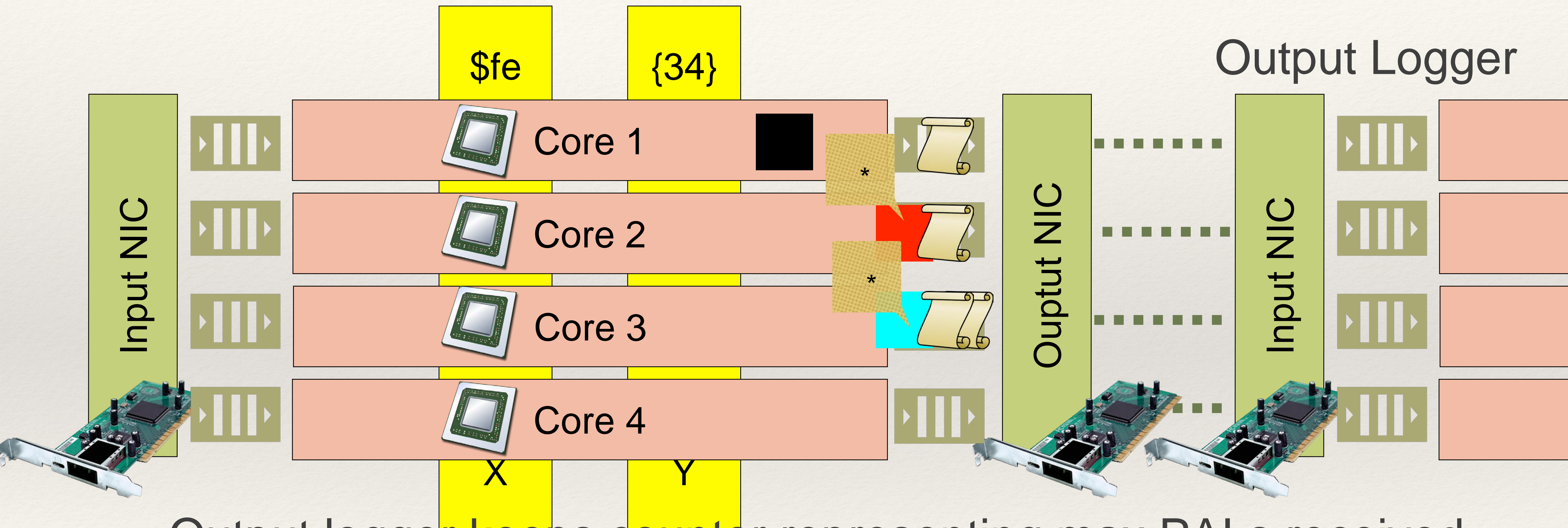
# Parallel Release



Each core keeps a counter tracking the "youngest" PAL it has created. On release, packet *reads* counters across all cores.  
( $O(\#cores^{33})$  reads)



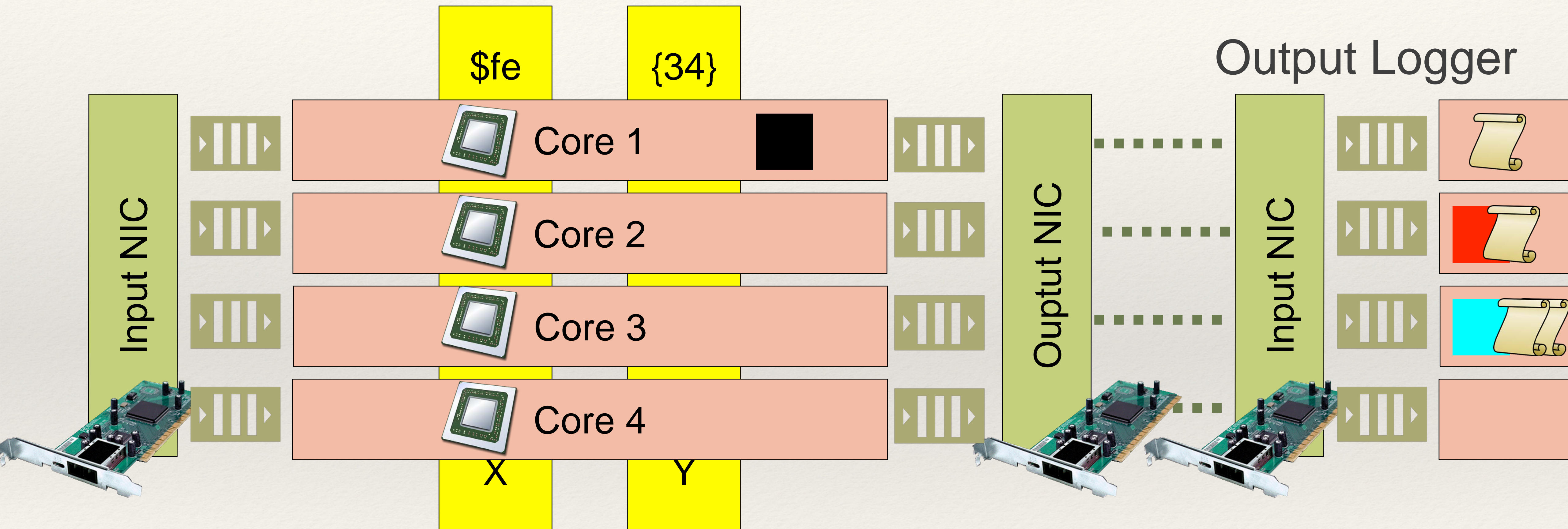
# Parallel Release



Output logger keeps counter representing max PALs received.  
Receive packet: reads all counters to compare against the  
attached<sup>34</sup> counters.



# Parallel Release



If attached counters  $\leq$  all counters, release packet!



# Recap

Three Part  
Algorithm:

Checkpoint

Log

Check

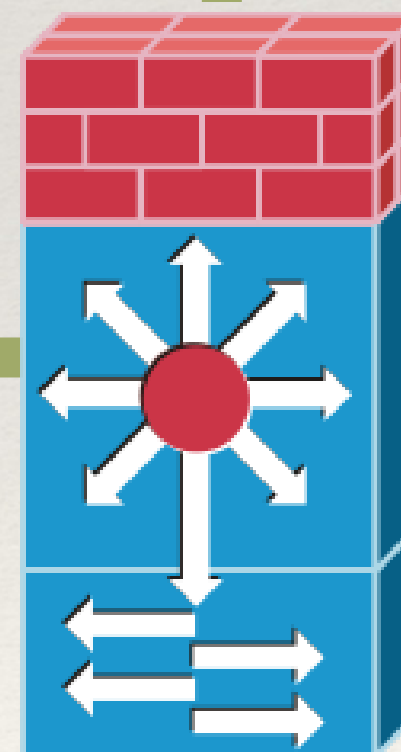
VM Snapshots

Ordered Logging and Parallel  
Release

Backup



Input Logger



Master Middlebox



Output Logger



# Thanks



# Backupland



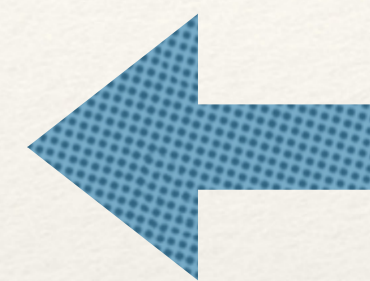
## Latency

Remus [NSDI 2008]:  
50,000us overhead

Pico [SOCC 2013]:  
8000us overhead

FTMB: 30us  
overhead

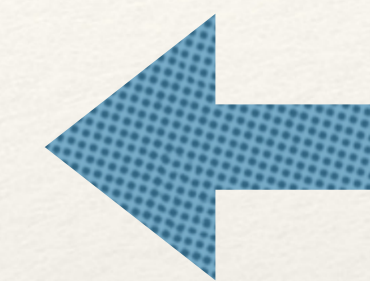
## Throughput



None higher than  
200kpps

FTMB: 1.4-4Mpps

## Recovery Time



100s of ms

FTMB: increases  
recovery time by 50-  
300ms.

Still fast enough not  
to trigger TCP  
timeouts or errors!



---

# What I'm not talking about today

---

- ❖ a prototype based on Xen, Click, and DPDK.
- ❖ a tool to instrument shared state with PALs automatically.
- ❖ How replay works in detail.
- ❖ Alternative output commit approaches.



# Ordered Logging and Parallel Release

---

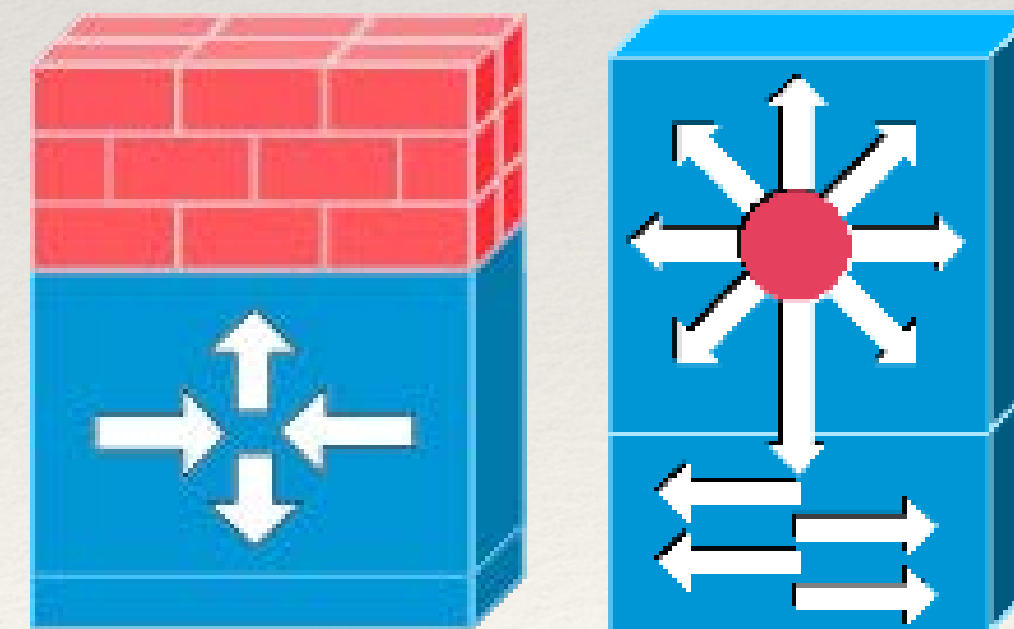
- ❖ Parallel! Threads are never blocked on each other to make progress.
- ❖ Cross-core accesses are *read only*.
  - ❖ Further amortized by batching.
- ❖ Linear: order # threads reads to perform.
- ❖ Fine-grained. Can make this decision with every packet release.



# FTMB Outside of Academia

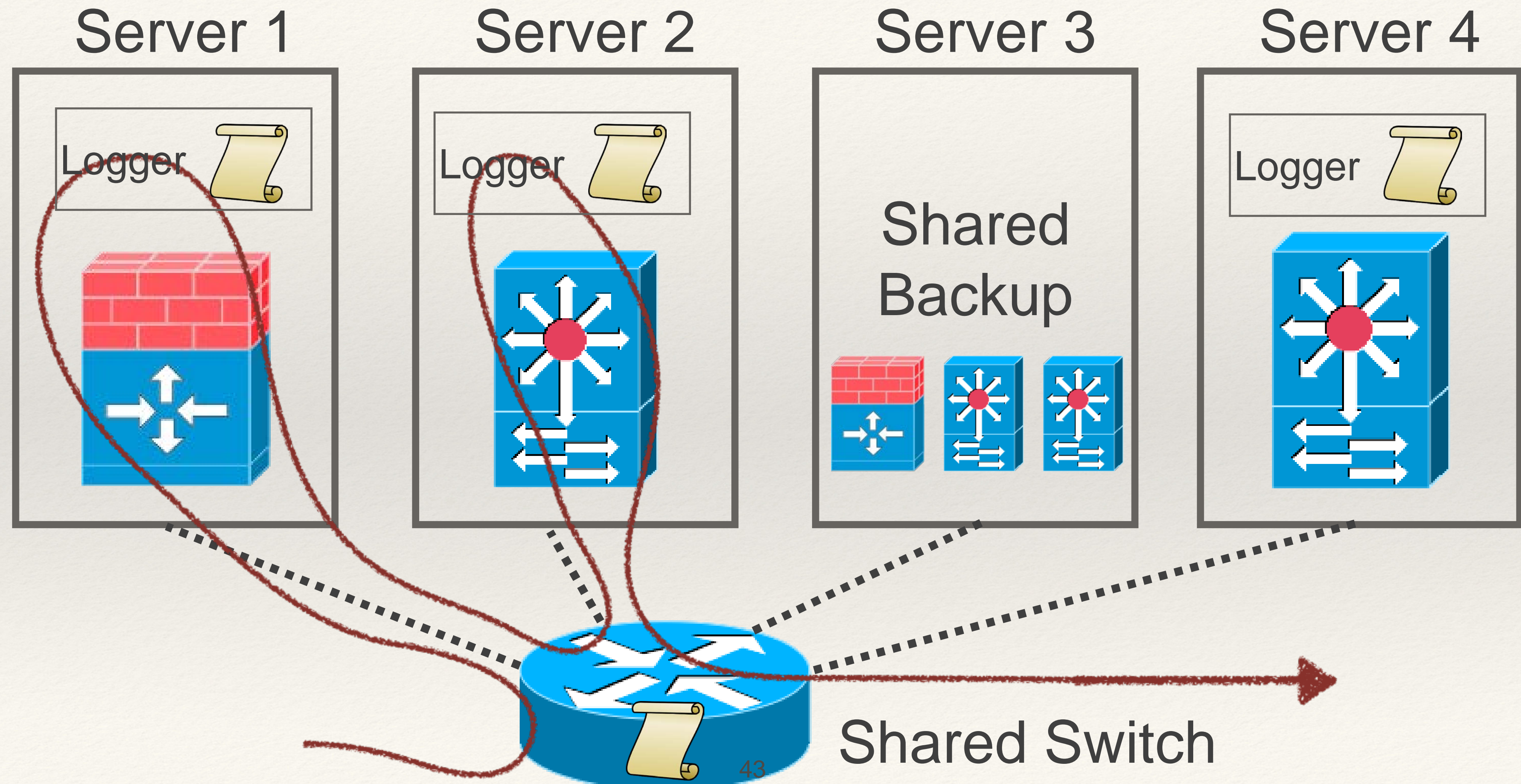
AT&T has submitted  
FTMB to ETSI:  
founders of NFV

FTMB is in trials at  
two major NFV & IDS  
vendors.





# NFV Frameworks

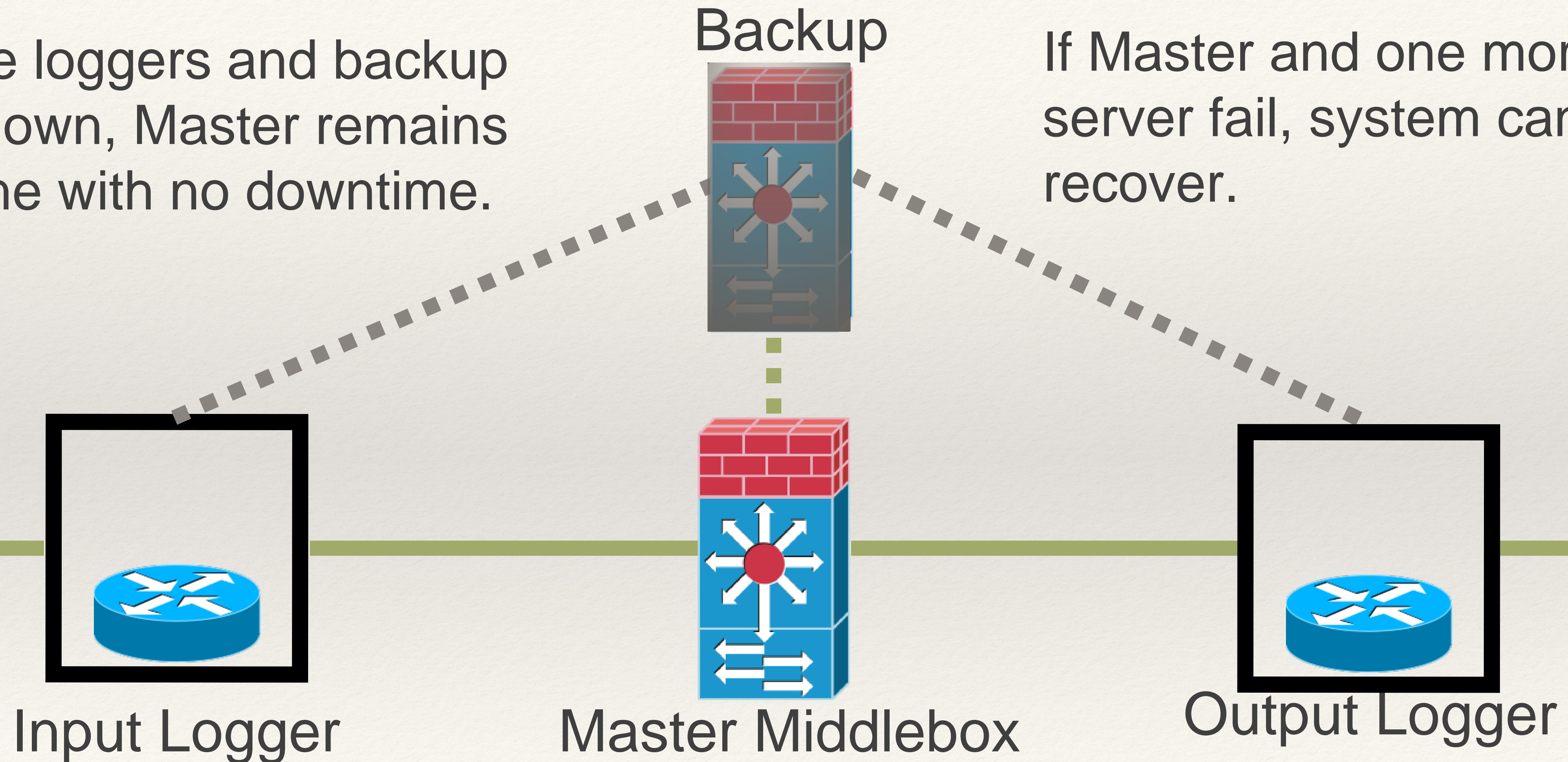




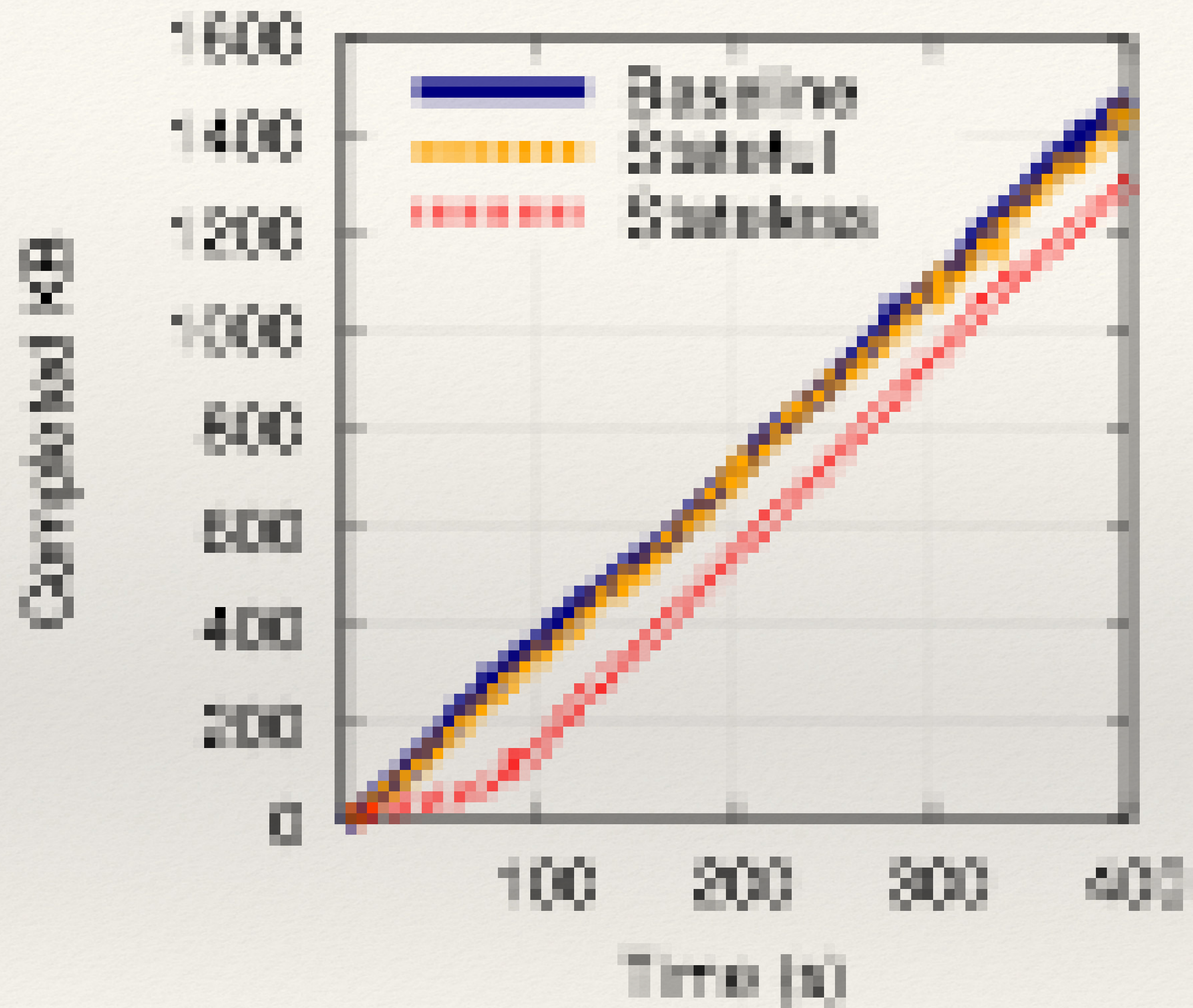
# Recovery Guarantees

If the loggers and backup go down, Master remains online with no downtime.

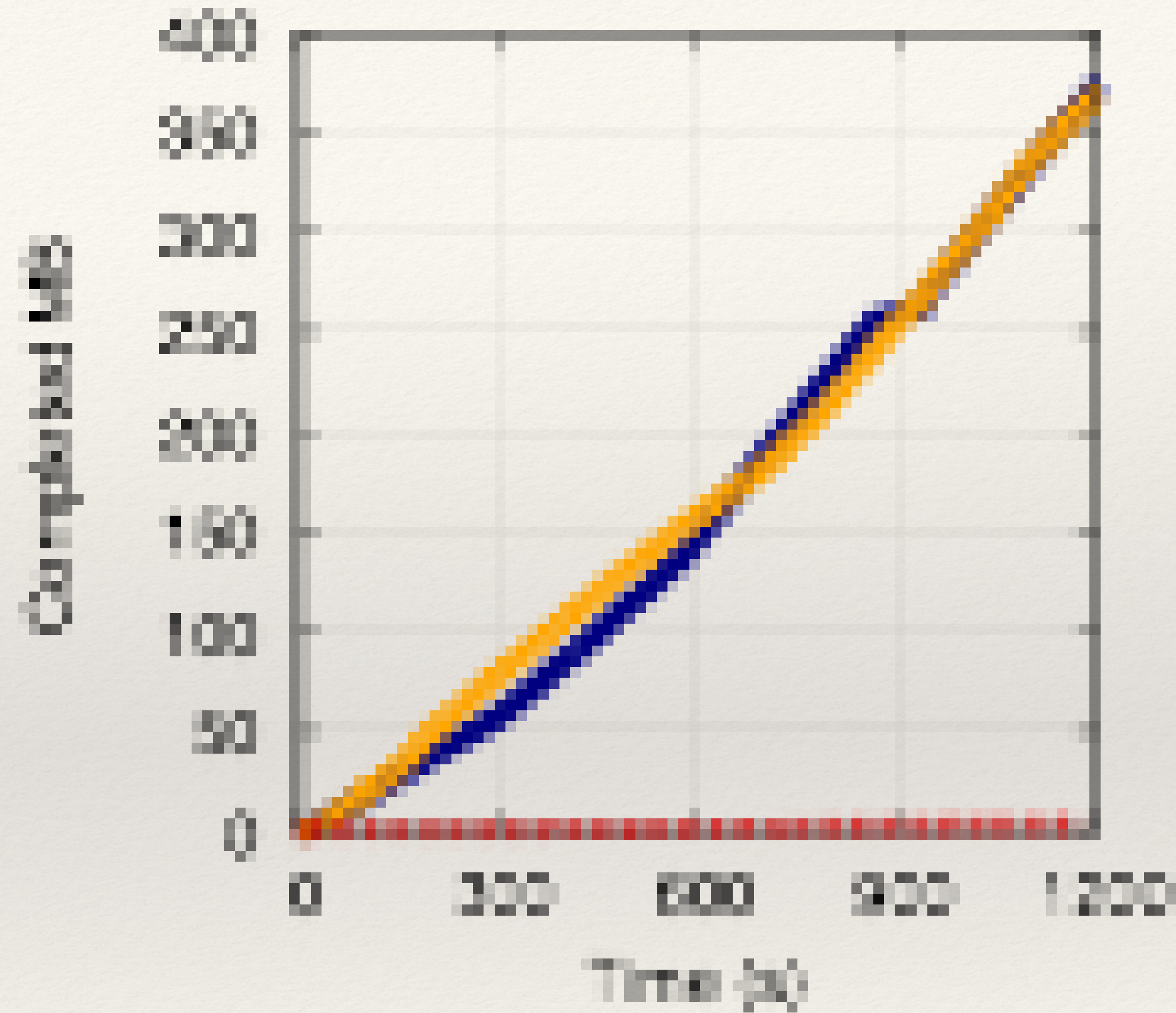
If Master and one more server fail, system cannot recover.







FTP Client



BitTorrent Client