

# Replication State Machines via Primary-Backup

---

COS 418: *Distributed Systems*

Amit Levy

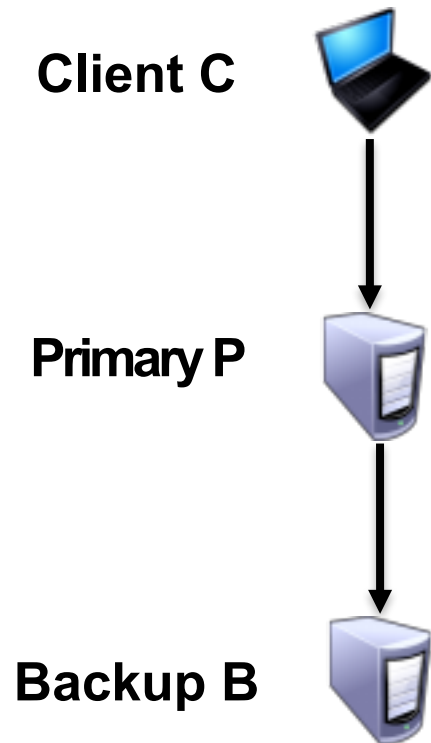
# From eventual to strong consistency

---

- Eventual consistency
  - Multi-master: Any node can accept operations
  - Asynchronous: nodes synchronize state after operations
- Eventual consistency inappropriate for many applications
  - Imagine eventually consistent file system
  - NFS clients can read/write to different masters, see different versions of files
- Stronger consistency makes applications easier to write

# Primary-Backup Replication

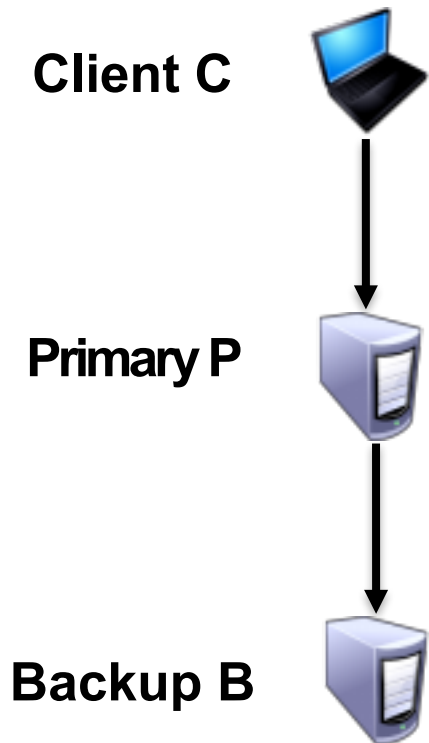
---



- **Mechanism:** Replicate and separate servers
- **Goal #1:** Reliable despite individual server failures
- **Goal #2:** Semantics of a single server

# Primary-Backup Replication

---



- Nominate one replica **primary**
- Other replicas are **backup**
  - Only **one primary at a time**
- Clients send all operations to current **primary**
- Primary **orders** clients' operations

Need to keep clients, primary, and backup in sync:  
**who is primary and who is backup**

# State machine replication

---

- **Insight:** A database is essentially a *state machine*
  - E.g. a set of (key, value) pairs is a **state**
  - Operations **transition** between states
- Each operation executed on *all* replicas, or *none at all*
  - *i.e.*, we need **distributed all-or-nothing atomicity**
- **Key assumption:** Operations are **deterministic**
- If op is deterministic, replicas will end in same state

## Implementing Fault-Tolerant Services Using the State Machine Approach: A Tutorial

FRED B. SCHNEIDER

*Department of Computer Science, Cornell University, Ithaca, New York 14853*

The state machine approach is a general method for implementing fault-tolerant services in distributed systems. This paper reviews the approach and describes protocols for two different failure models—Byzantine and fail stop. System reconfiguration techniques for removing faulty components and integrating repaired components are also discussed.

Categories and Subject Descriptors: C.2.4 [Computer-Communication Networks]: Distributed Systems—*network operating systems*; D.2.10 [Software Engineering]: Design—*methodologies*; D.4.5 [Operating Systems]: Reliability—*fault tolerance*; D.4.7 [Operating Systems]: Organization and Design—*interactive systems, real-time systems*

General Terms: Algorithms, Design, Reliability

Additional Key Words and Phrases: Client-server, distributed services, state machine approach

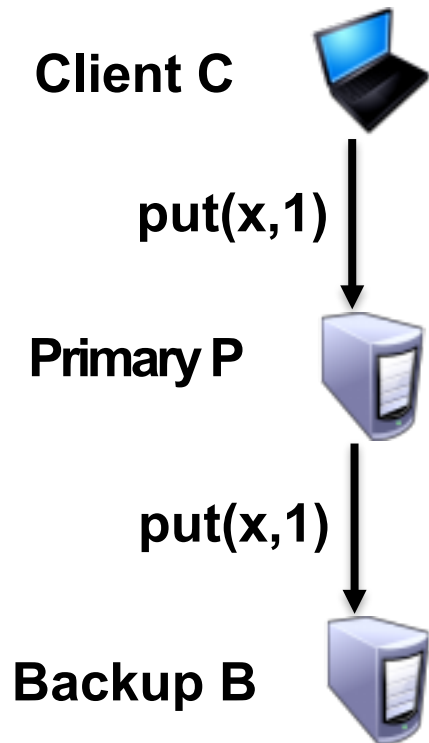
### INTRODUCTION

Distributed software is often structured in terms of *clients* and *services*. Each service comprises one or more *servers* and exports *operations* that clients invoke by making *requests*. Although using a single, centralized, server is the simplest way to imple-

service by replicating servers and coordinating client interactions with server replicas.<sup>1</sup> The approach also provides a framework for understanding and designing replication management protocols. Many protocols that involve replication of data or software—be it for masking failures or simply to facilitate cooperation without

# Primary-Backup Replication

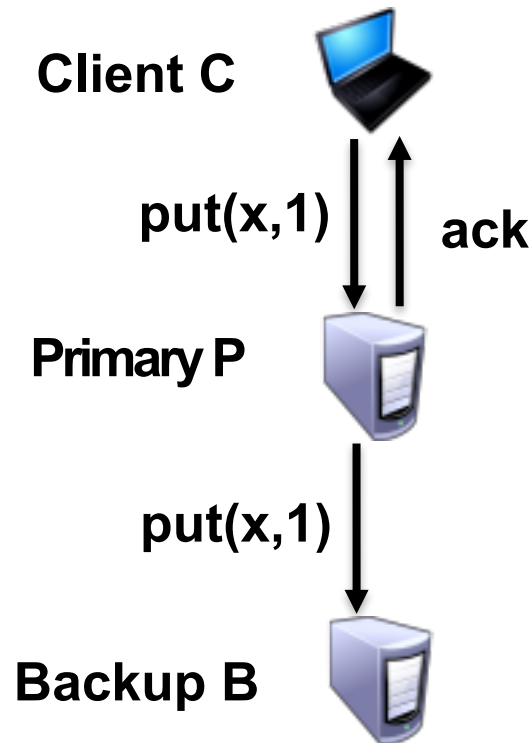
---



1. Primary gets operations
2. Primary orders ops into log
3. Replicates log of ops to backup
4. Backup exec's ops or writes to log

# Primary-Backup Replication

---



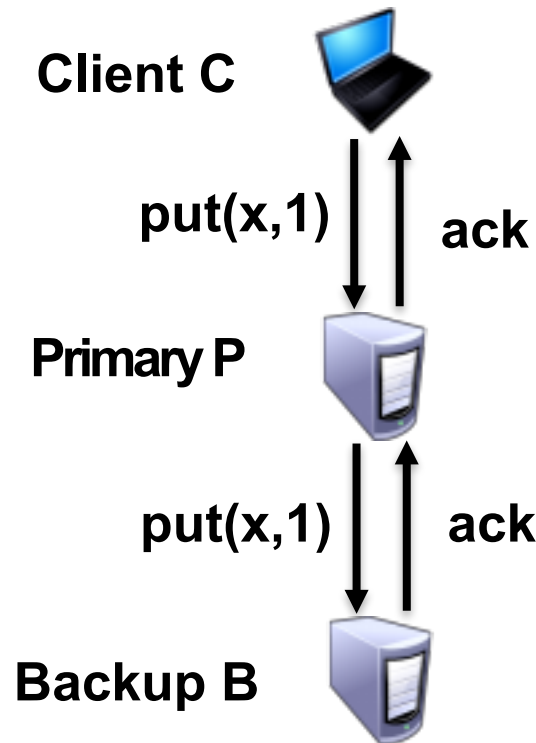
## Asynchronous Replication

1. Primary gets operations
- 2. Primary exec's ops**
3. Primary orders ops into log
4. Replicates log of ops to backup
5. Backup exec's ops or writes to log



# Primary-Backup Replication

---

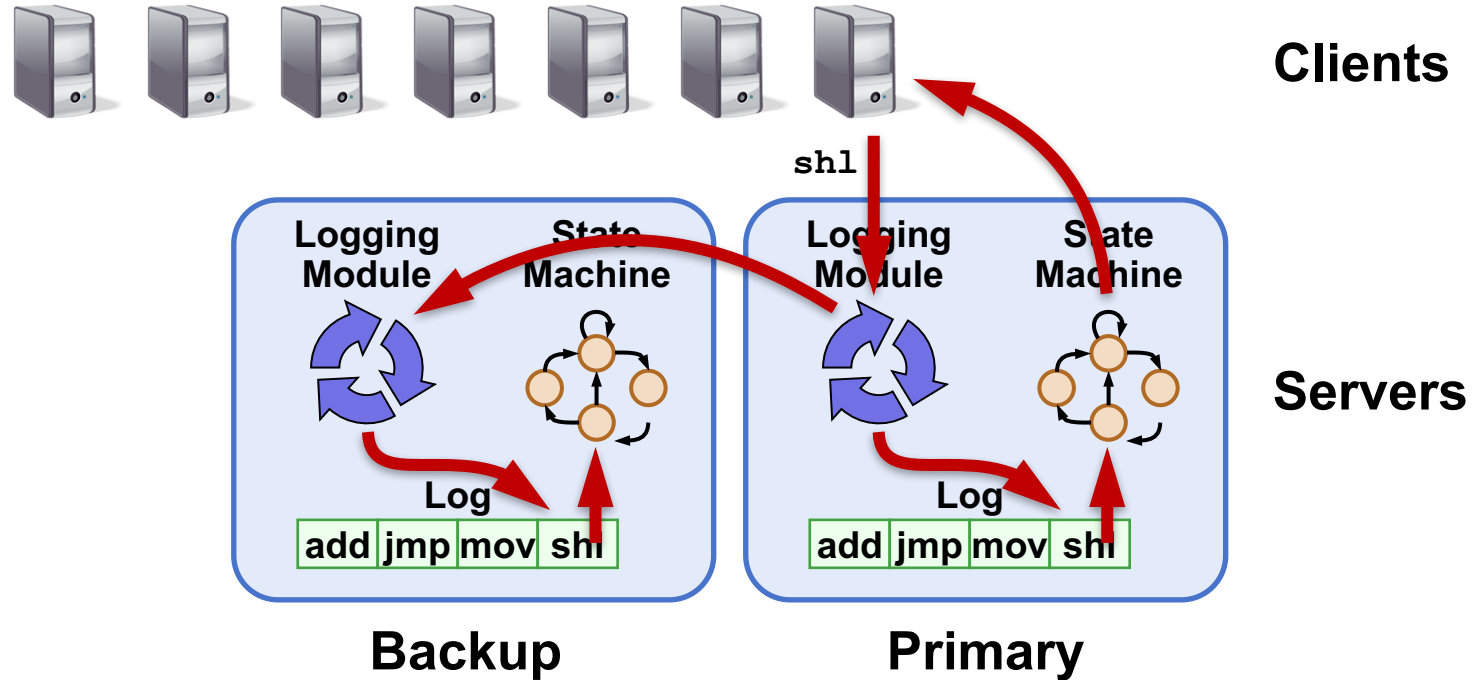


## Synchronous Replication

1. Primary gets operations
2. Primary orders ops into log
3. Replicates log of ops to backup
4. Backup exec's op or writes to log
5. **Primary gets ack, execs ops**

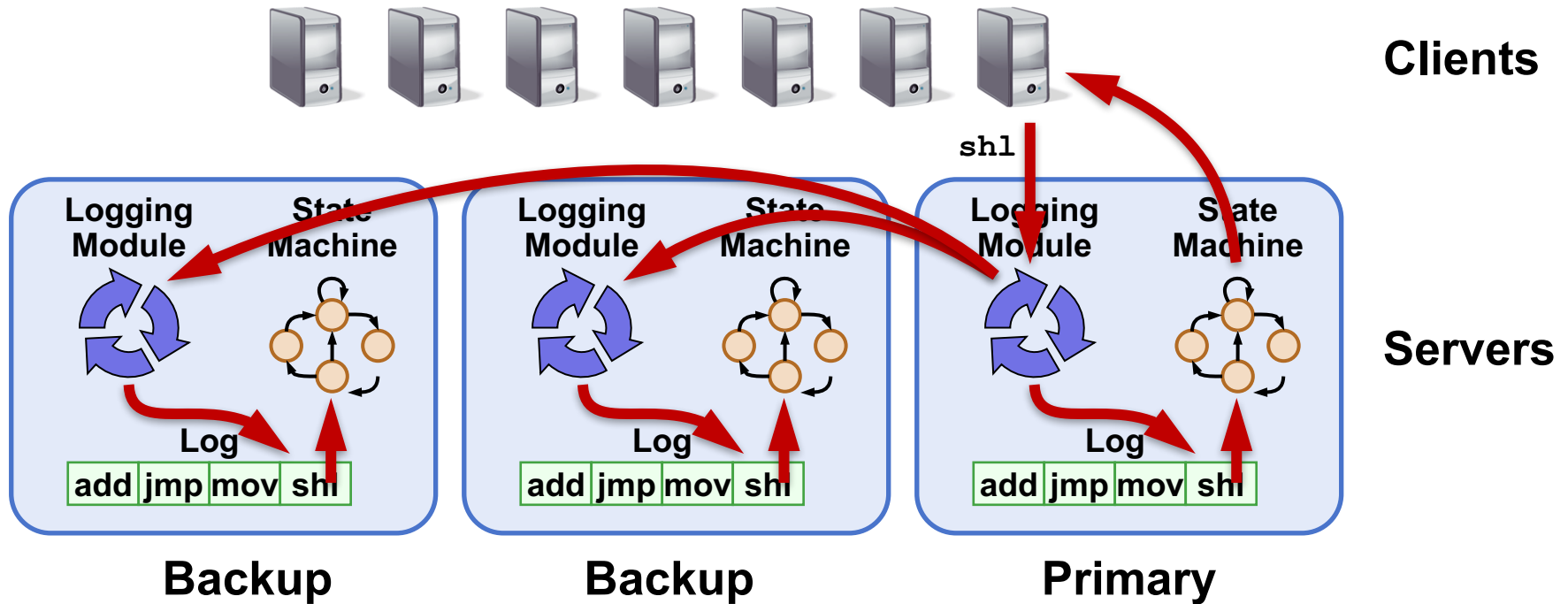
# Why does this work?

## Synchronous Replication



- Replicated log => replicated state machine
  - All servers execute same commands in same order

# Synchronous Replication



- Replicated log => replicated state machine
  - All servers execute same commands in same order

# Need determinism? Make it so!

---

- Operations are deterministic
  - No events with ordering based on local clock
    - Convert timer, network, user into logged events
  - Nothing using random inputs
- Execution order of ops is identical
  - e.g. replicated state machines (RSMs) is single threaded

# Example: Make random() deterministic

---

Almost all module functions depend on the basic function `random()`, which generates a random float uniformly in the semi-open range [0.0, 1.0). Python uses the Mersenne Twister as the core generator. It produces 53-bit precision floats and has a period of  $2^{19937}-1$ . The underlying implementation in C is both fast and threadsafe. The Mersenne Twister is one of the most extensively tested random number generators in existence. However, being completely deterministic, it is not suitable for all purposes, and is completely unsuitable for cryptographic purposes.

`random.seed(a=None)`

Initialize internal state of the random number generator.

`None` or no argument seeds from current time or from an operating system specific randomness source if available (see the `os.urandom()` function for details on availability).

`random.getstate()`

Return an object capturing the current internal state of the generator. This object can be passed to `setstate()` to restore the state.

# Example: Make random() deterministic

---

- Primary:
  - Initiates PRNG with OS-supplied randomness, gets initial seed
  - Sends initial seed to backup
- Backup
  - Initiates PRNG with seed from primary

`random.seed(a=None)`

Initialize internal state of the random number generator.

`None` or no argument seeds from current time or from an operating system specific randomness source if available (see the `os.urandom()` function for details on availability).

`random.getstate()`

Return an object capturing the current internal state of the generator. This object can be passed to `setstate()` to restore the state.

# Case study

The design of a practical system for  
fault-tolerant virtual machines

D. Scales, M. Nelson, G. Venkitachalam, VMWare

SIGOPS Operating Systems Review 44(4), Dec. 2010 ([pdf](#))

# VMware vSphere Fault Tolerance (VM-FT)

---

## Goals:

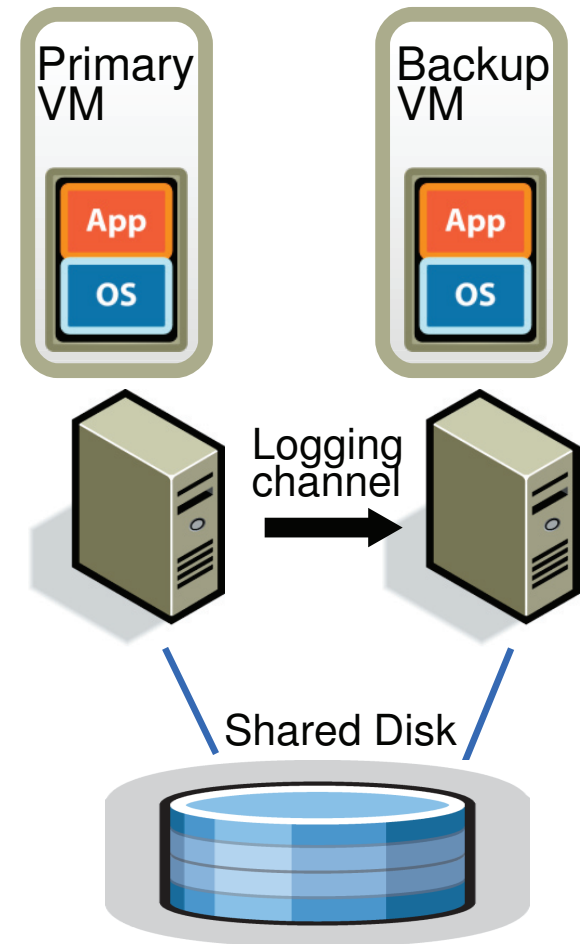
- 1. Replication of the **whole virtual machine**
- 1. **Completely transparent** to apps and clients
- 1. **High availability** for any existing software



# Overview

---

- Two virtual machines (*primary*, *backup*) on different bare metal
- *Logging channel* runs over network
- *Shared disk* via fiber channel



# Virtual Machine I/O

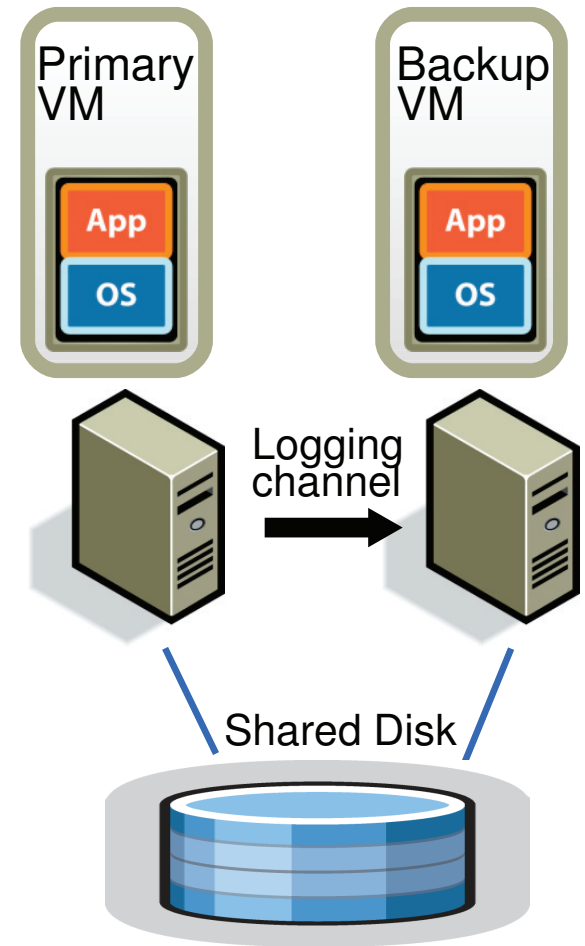
---

- **VM inputs**
  - Incoming network packets
  - Disk reads
  - Keyboard and mouse events
  - Clock timer interrupt events
- **VM outputs**
  - Outgoing network packets
  - Disk writes

# Overview

---

- **Primary** sends **inputs** to backup
- **Backup outputs** dropped
- Primary-backup **heartbeats**
  - If primary fails, backup takes over



# VM-FT: Challenges

---

1. Making the backup an exact replica of primary
1. Making the system behave like a single server
1. Avoiding two primaries (Split Brain)

# Log-based VM replication

---

- **Step 1:** Hypervisor at primary logs the causes of non-determinism
  1. Log results of **input events**
    - Including current program counter value for each
  2. Log results of **non-deterministic instructions**
    - e.g. log **result** of timestamp counter read

# Log-based VM replication

---

- **Step 2:** Primary hypervisor sends log entries to backup hypervisor
- Backup hypervisor **replays** the log entries
  - **Stops backup VM** at next input event or non-deterministic instruction
    - Delivers **same input** as primary
    - Delivers **same non-deterministic instruction result** as primary

# VM-FT Challenges

---

1. Making the backup an exact replica of primary
1. **Making the system behave like a single server**
  - **FT Protocol**
1. Avoiding two primaries (Split Brain)

# Primary to backup failover

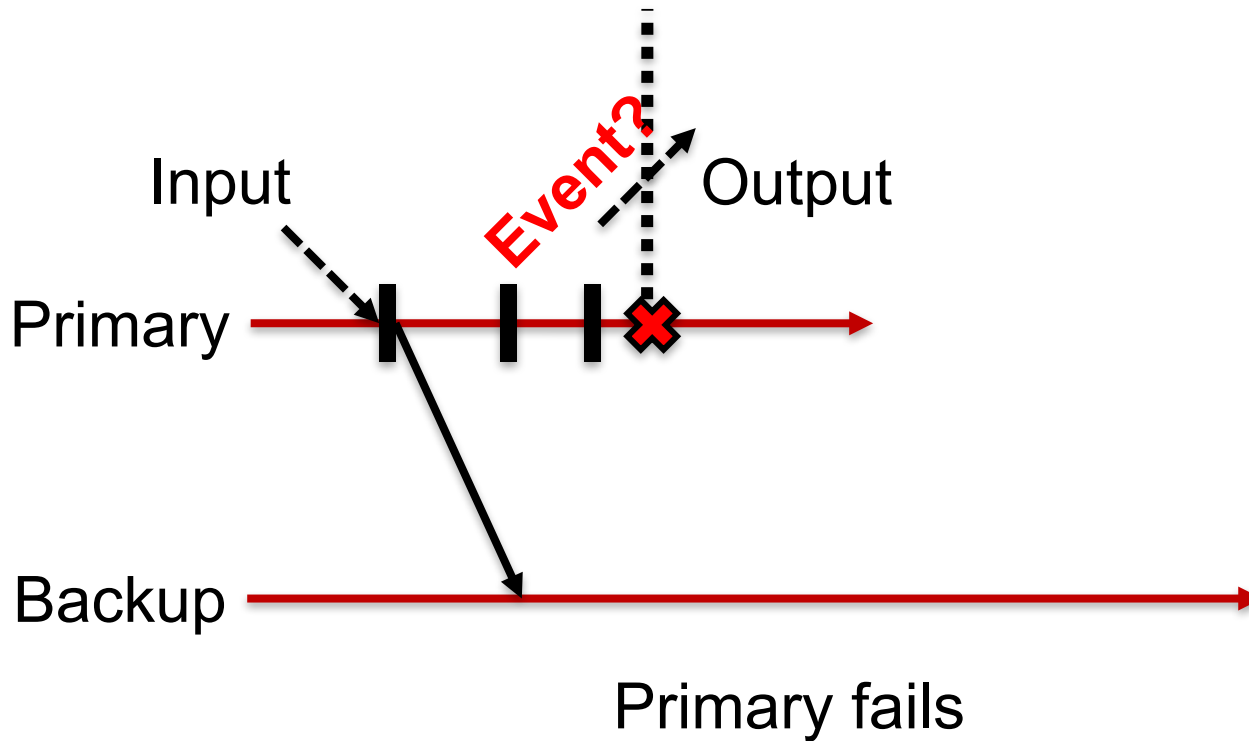
---

- When backup takes over, non-determinism makes it **execute differently** than primary would have
  - **This is okay!**
- **Output requirement**
  - When backup takes over, execution is **consistent** with **outputs** the primary has already sent



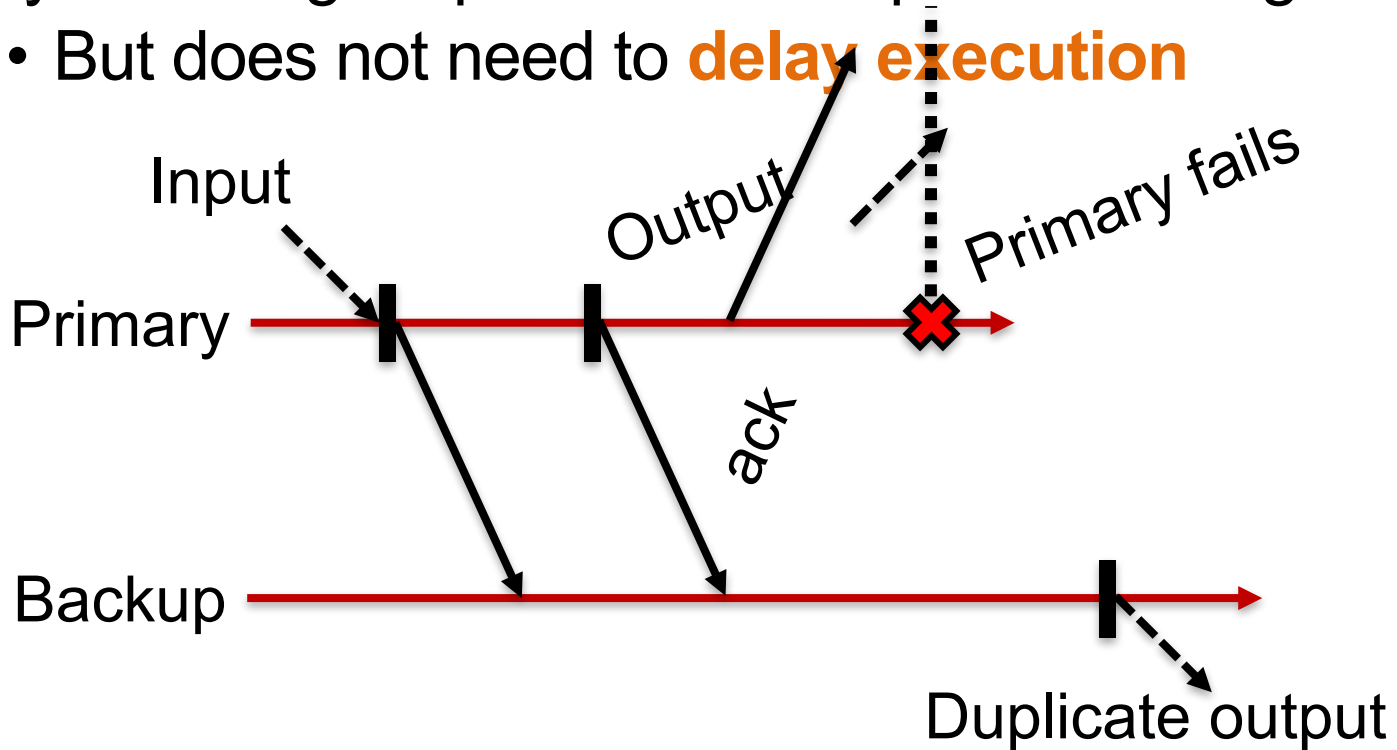
# The problem of inconsistency

---



# VM-FT protocol

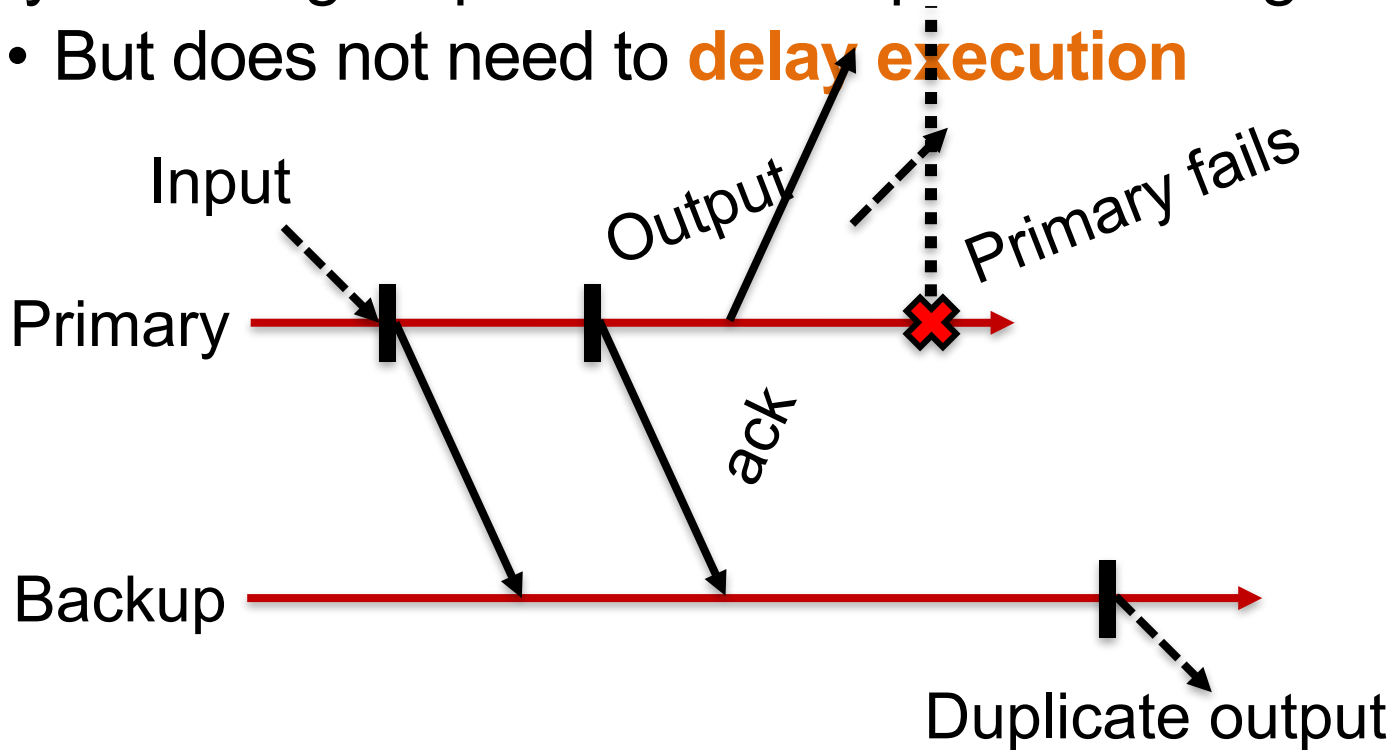
- Primary logs each output operation
- Delays sending output until Backup acknowledges it
  - But does not need to **delay execution**



# VM-FT protocol

“If a tree falls in forest” metaphor:  
If event happens and nobody sees it yet,  
did it really happen?

- Primary logs each output operation
- Delays sending output until Backup acknowledges it
  - But does not need to **delay execution**



Can restart execution at an output event

# VM-FT: Challenges

---

- 1. Making the backup an exact replica of primary
- 1. Making the system behave like a single server
- 1. **Avoiding two primaries (Split Brain)**
  - Logging channel may **break**

# Detecting and responding to failures

---

- Primary and backup each run UDP heartbeats, monitor logging traffic from their peer
- Before “going live” (backup) or finding new backup (primary), execute an **atomic test-and-set** on a variable in shared storage
- If the replica finds variable already set, it **aborts**

# VM-FT: Conclusion

---

- Challenging application of primary-backup replication
- Design for correctness and consistency of replicated VM outputs despite failures
- Performance results show generally **high performance, low logging bandwidth overhead**

# Monday

How \*do\* we detect failures?

Take over from master on failures?

“View Change Protocols”

View = Current System Configuration