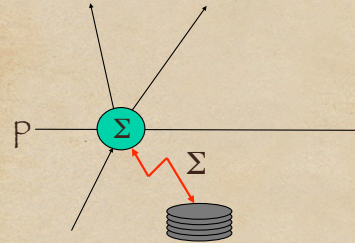


## Rollback-Recovery

## Uncoordinated Checkpointing



- Easy to understand
- No synchronization overhead
- Flexible
  - can choose **when** to checkpoint
- To recover from a crash:
  - go back to last checkpoint
  - restart

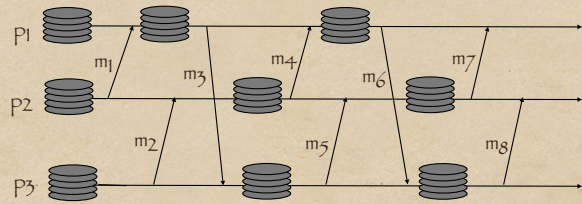
## How (not) to take a checkpoint

- Block execution, save entire process state to stable storage
  - very high overhead during failure-free execution
  - lots of unnecessary data saved on stable storage

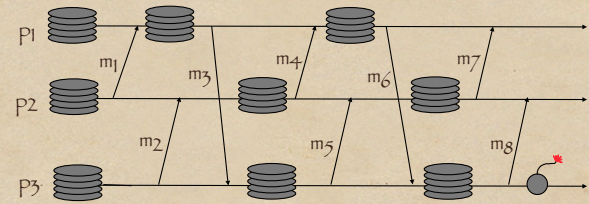
## How to take a checkpoint

- Take checkpoints incrementally
  - save only pages modified since last checkpoint
  - use “dirty” bit to determine which pages to save
- Save only “interesting” parts of address space
  - use application hints or compiler help to avoid saving useless data (e.g. dead variables)
- Do not block application execution during recovery
  - copy-on-write

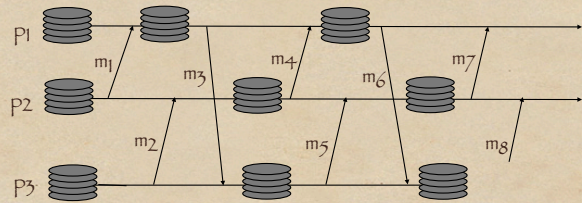
# The Domino Effect



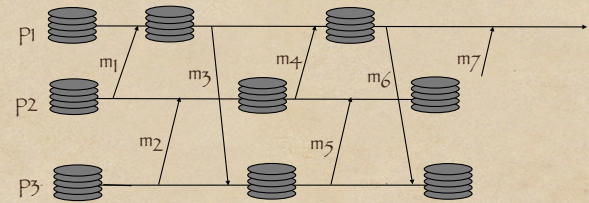
# The Domino Effect



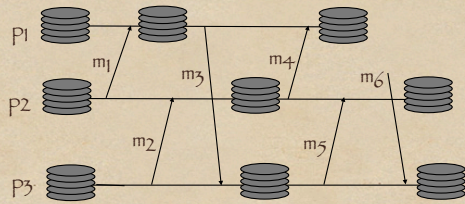
# The Domino Effect



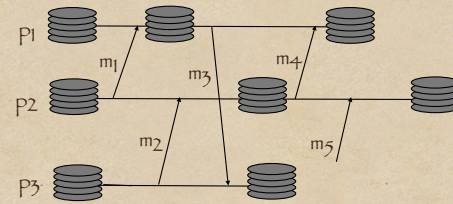
# The Domino Effect



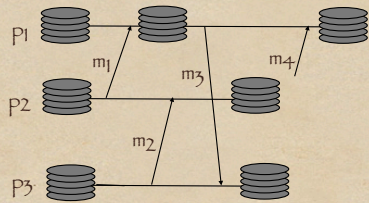
# The Domino Effect



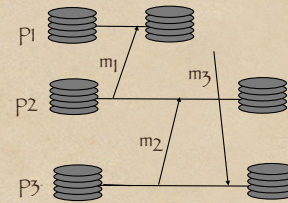
# The Domino Effect



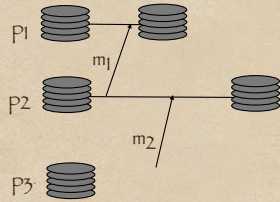
# The Domino Effect



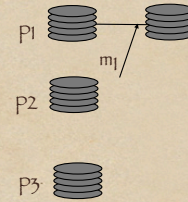
# The Domino Effect



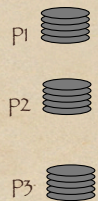
## The Domino Effect



## The Domino Effect



## The Domino Effect



## How to Avoid the Domino Effect

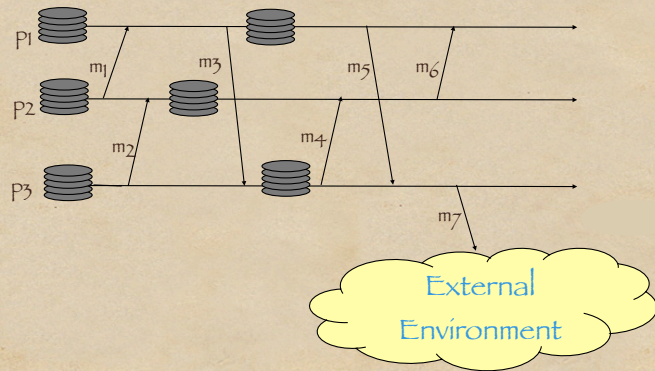
### Coordinated Checkpointing

- No independence
- Synchronization Overhead
- Easy Garbage Collection

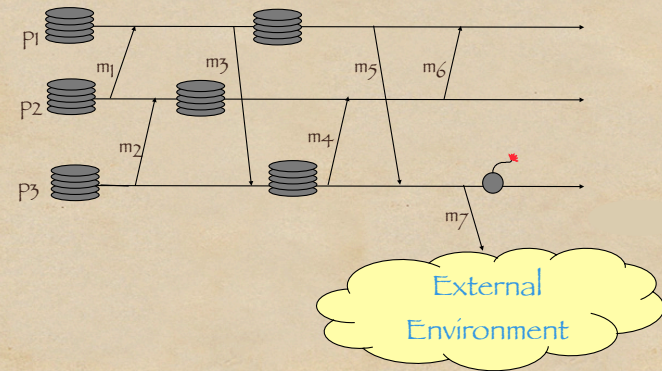
**Communication Induced Checkpointing**: detect dangerous communication patterns and checkpoint appropriately

- Less synchronization
- Less independence
- Complex

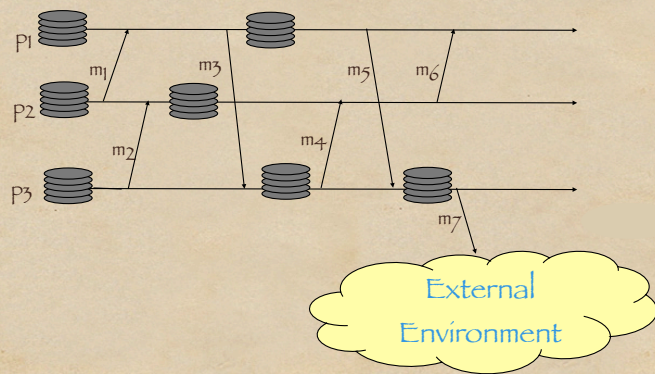
# The Output Commit Problem



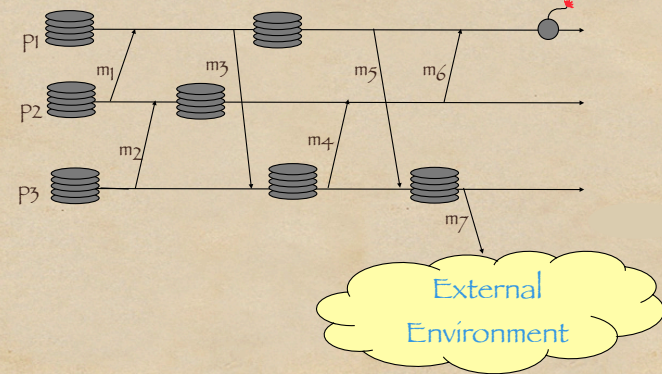
# The Output Commit Problem



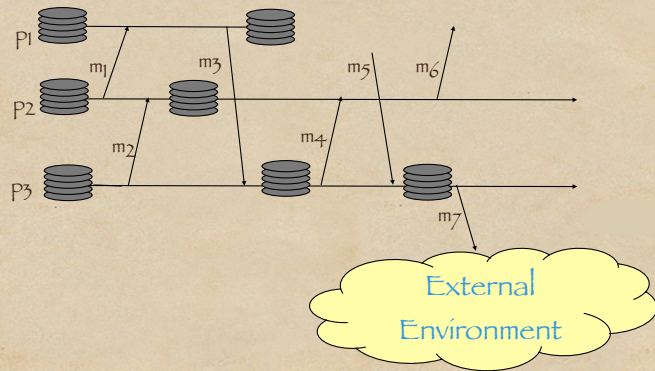
# The Output Commit Problem



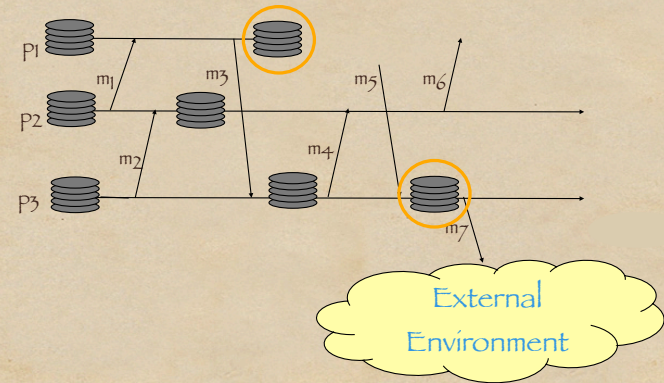
# The Output Commit Problem



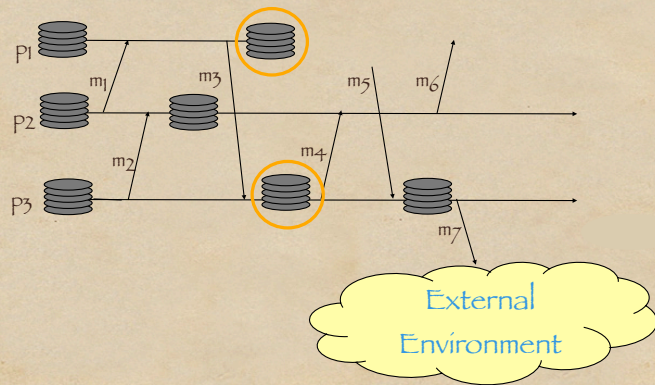
# The Output Commit Problem



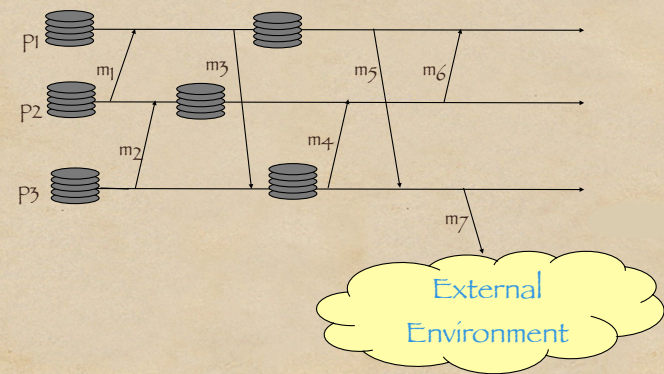
# The Output Commit Problem



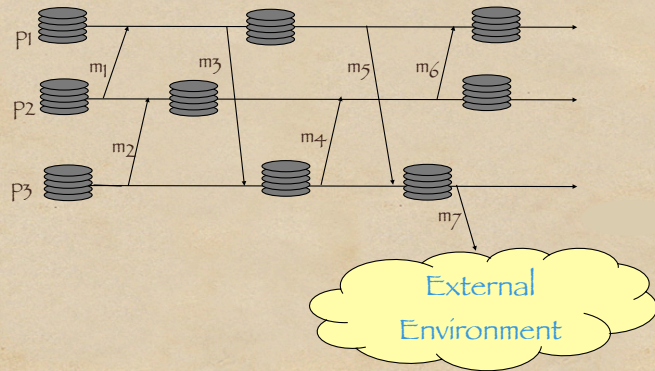
# The Output Commit Problem



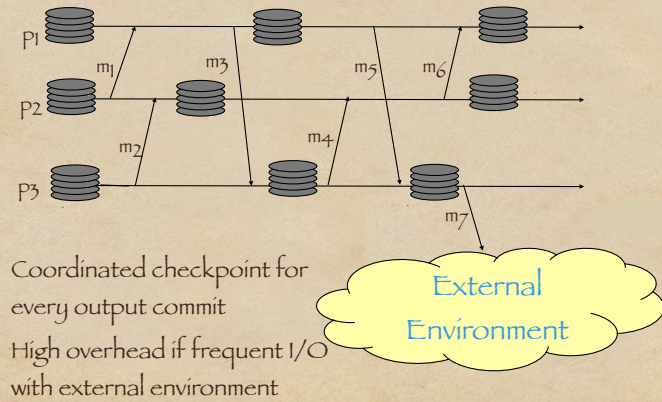
# The Output Commit Problem



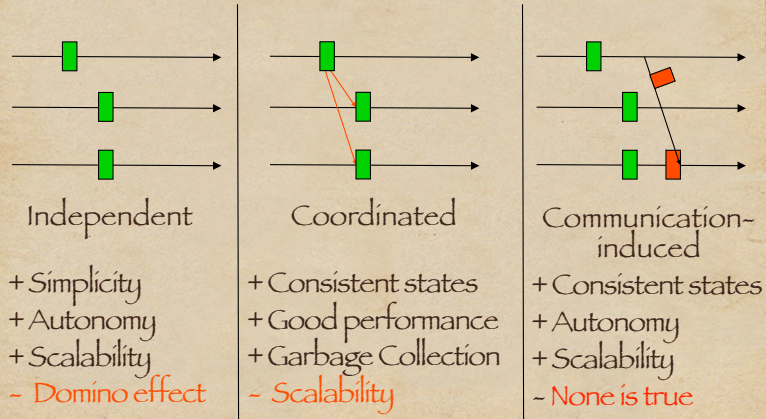
# The Output Commit Problem



# The Output Commit Problem



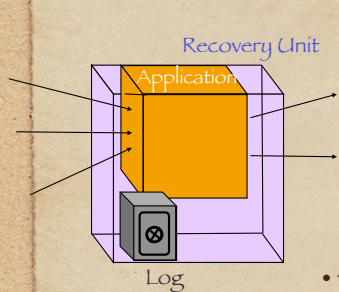
# Distributed Checkpointing at a Glance



# Message Logging

- Can avoid domino effect
- Works with coordinated checkpoint
- Works with uncoordinated checkpoint
- Can reduce cost of output commit
- More difficult to implement

# How Message Logging Works



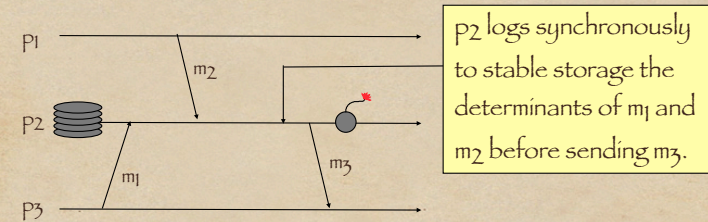
To tolerate crash failures:

- periodically checkpoint application state;
- log on stable storage **determinants** of non-deterministic events executed after checkpointed state.
- for message delivery events:  
 $\#m = (m.dest, m.rsn, m.source, m.ssn)$

Recovery:

- restore latest checkpointed state;
- replay non-deterministic events according to determinants

# Pessimistic Logging



Never creates orphans

- may incur blocking
- straightforward recovery

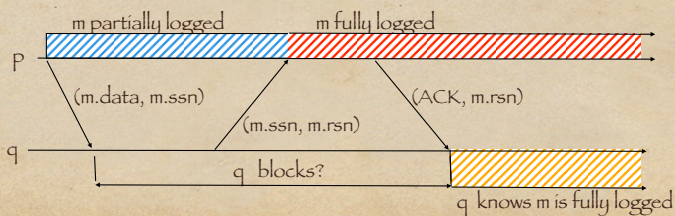
# Sender Based Logging

(Johnson and Zwaenepoel, FTCS 87)

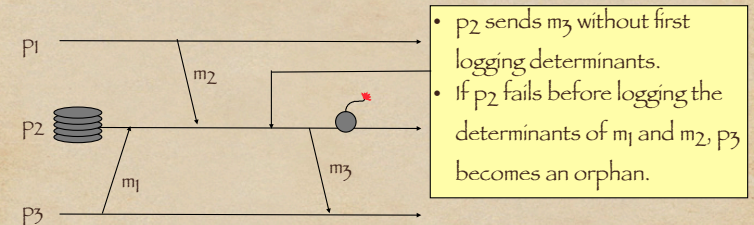
Message log is maintained in volatile storage at the sender.

A message m is logged in two steps:

- before sending m, the sender logs its content: m is **partially logged**
- the receiver tells the sender the receive sequence number of m, and the sender adds this information to its log: m is **fully logged**.



# Optimistic Logging



Eliminates orphans during recovery

- non-blocking during failure-free executions
- rollback of correct processes
- complex recovery



## Causal Logging

- No blocking in failure-free executions
- No orphans
- No additional messages
- Tolerates multiple concurrent failures
- Keeps determinant in volatile memory
- Localized output commit

## Preliminary Definitions

Given a message  $m$  sent from  $m.source$  to  $m.dest$ ,

$$\text{Depend}(m): \left\{ p \in P \mid \begin{array}{l} \forall (p = m.dest) \text{ and } p \text{ delivered } m \\ \forall (\exists e_p : (\text{deliver}_{m.dest}(m) \rightarrow e_p)) \end{array} \right\}$$

$\text{Log}(m)$ : set of processes with a copy of the determinant of  $m$  in their volatile memory

$p$  orphan of a set  $C$  of crashed processes:

$$(p \notin C) \wedge \exists m : (\text{Log}(m) \subseteq C \wedge p \in \text{Depend}(m))$$

## The “No-Orphans” Consistency Condition

No orphans after crash  $C$  if:

$$\forall m : (\text{Log}(m) \subseteq C) \Rightarrow (\text{Depend}(m) \subseteq C)$$

No orphans after any  $C$  if:

$$\forall m : (\text{Depend}(m) \subseteq \text{Log}(m))$$

The Consistency Condition

$$\forall m : (\neg \text{stable}(m) \Rightarrow (\text{Depend}(m) \subseteq \text{Log}(m)))$$

## Optimistic and Pessimistic

No orphans after crash  $C$  if:

$$\forall m : (\text{Log}(m) \subseteq C) \Rightarrow (\text{Depend}(m) \subseteq C)$$

Optimistic weakens it to:

$$\forall m : (\text{Log}(m) \subseteq C) \Rightarrow \diamond(\text{Depend}(m) \subseteq C)$$

No orphans after any crash if:

$$\forall m : (\neg \text{stable}(m) \Rightarrow (\text{Depend}(m) \subseteq \text{Log}(m)))$$

Pessimistic strengthens it to:

$$\forall m : (\neg \text{stable}(m) \Rightarrow |\text{Depend}(m)| \leq 1)$$

# Causal Message Logging

No orphans after any crash of size at most  $f$  if:

$$\forall m : (\neg \text{stable}(m) \Rightarrow (\text{Depend}(m) \subseteq \text{Log}(m)))$$

Causal strengthens it to:

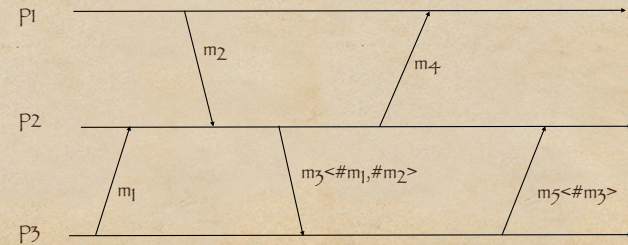
$$\forall m : \left( \neg \text{stable}(m) \Rightarrow \left( \begin{array}{l} \wedge (\text{Depend}(m) \subseteq \text{Log}(m)) \\ \wedge \diamond (\text{Depend}(m) = \text{Log}(m)) \end{array} \right) \right)$$

# An Example

Causal Logging:

$$\forall m : (\neg \text{stable}(m) \Rightarrow (\text{Depend}(m) \subseteq \text{Log}(m)))$$

If  $f=1$ ,  $\text{stable}(m) \equiv |\text{Log}(m)| \geq 2$



# Recovery for $f=1$

