



# Commit

## Distributed Systems

---

**Ali Kamandi, PH.D.**

School of Engineering Science

College of Engineering

University of Tehran

**kamandi@ut.ac.ir**

**2024**



# Distributed Database Commit

- چند سایت در اجرای یک تراکنش نقش دارند. وقتی هر یک از سایت‌ها وظیفه خود را به پایان رساند، باید تصمیم‌گیری کند که تراکنش commit شود یا abort.
- در صورتی که هر سایت توانسته باشد کارهای مربوط به خود را به صورت محلی انجام دهد، ترجیح می‌دهد تصمیم commit بگیرد و در غیر این صورت abort.
- وقتی کامیت انجام می‌شود که همه سایت‌ها روی آن توافق داشته باشند.
- ترجیحا توافق باید روی کامیت انجام شود.

# مدل مسئله

- فرض می‌کنیم خطای ارسال پیام (لینک ارتباطی) نداریم و صرفاً **خطای سایت یا نود** داریم.  
( Process Failure )
- مجموعه تصمیم‌گیری شامل  $\{0,1\}$  است: 0 به معنای abort و 1 به معنای commit خواهد بود.
- شبکه را گراف کامل در نظر می‌گیریم.

**Agreement:** No two processes decide on different values.

**Validity:**

1. If any process starts with 0, then 0 is the only possible decision value.
2. If all processes start with 1 and there are no failures, then 1 is the only possible decision value.

**Termination:**

The **weak termination** condition says that if there are no failures then all processes eventually decide.

The **strong termination** condition (also known as the non-blocking condition) says that all non-faulty processes eventually decide.

---

MODULE *TCommit*

---

CONSTANT *RM*      The set of participating resource managers

VARIABLE *rmState*      *rmState*[*rm*] is the state of resource manager *rm*.

---

$TCTypeOK \triangleq$   
 $rmState \in [RM \rightarrow \{\text{"working"}, \text{"prepared"}, \text{"committed"}, \text{"aborted"}\}]$

$TCInit \triangleq \quad rmState = [r \in RM \mapsto \text{"working"}]$

$canCommit \triangleq \forall r \in RM : rmState[r] \in \{ \text{"prepared"}, \text{"committed"} \}$

$notCommitted \triangleq \forall r \in RM : rmState[r] \neq \text{"committed"}$

---

$Prepare(r) \triangleq \begin{aligned} &\wedge rmState[r] = \text{"working"} \\ &\wedge rmState' = [rmState \text{ EXCEPT } ![r] = \text{"prepared"}] \end{aligned}$

$Decide(r) \triangleq \begin{aligned} &\vee \wedge rmState[r] = \text{"prepared"} \\ &\quad \wedge canCommit \\ &\quad \wedge rmState' = [rmState \text{ EXCEPT } ![r] = \text{"committed"}] \\ &\vee \wedge rmState[r] \in \{ \text{"working"}, \text{"prepared"} \} \\ &\quad \wedge notCommitted \\ &\quad \wedge rmState' = [rmState \text{ EXCEPT } ![r] = \text{"aborted"}] \end{aligned}$

$TCNext \triangleq \exists r \in RM : Prepare(r) \vee Decide(r)$

---

$TCConsistent \triangleq$

$\forall r1, r2 \in RM : \neg \wedge rmState[r1] = \text{"aborted"}$   
 $\wedge rmState[r2] = \text{"committed"}$

---

$TCSpec \triangleq TCInit \wedge \Box[TCNext]_{rmState}$

THEOREM  $TCSpec \Rightarrow \Box(TCTypeOK \wedge TCConsistent)$

---

# Two-Phase Commit

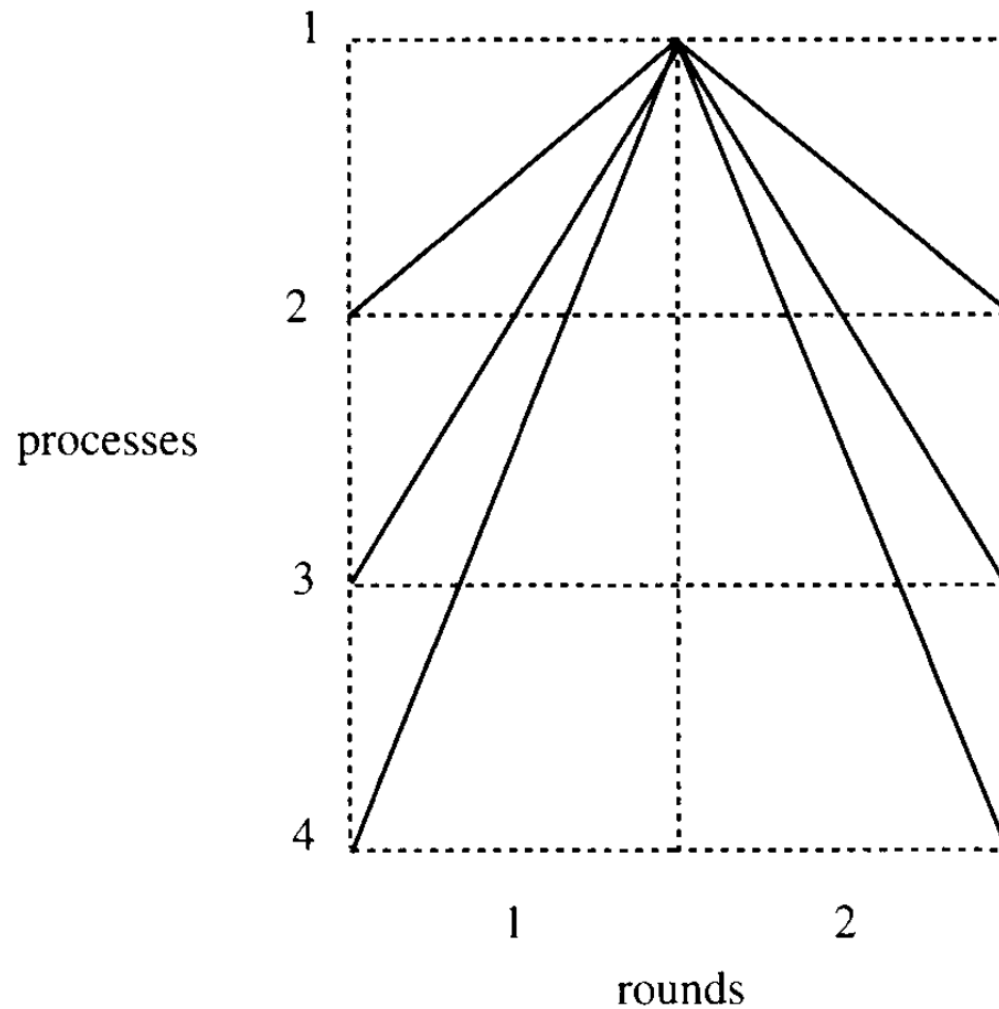
The algorithm assumes a distinguished process, say process 1.

**Round 1:** All processes except for process 1 send their initial values to process 1, and any process whose initial value is 0 decides 0. Process 1 collects all these values, plus its own initial value, into a vector. If all positions in this vector are filled in with 1s, then process 1 decides 1. Otherwise—that is, if there is some position in the vector that contains 0 or else some position that is not filled in (because no message was received from the corresponding process)—process 1 decides 0.

**Round 2:** Process 1 broadcasts its decision to all the other processes. Any process other than process 1 that receives a message at round 2 and has not already decided at round 1 decides on the value it receives in that message.



# Communication pattern in Two-PhaseCommit



- Agreement
- Validity
- Weak Termination
- Strong Termination

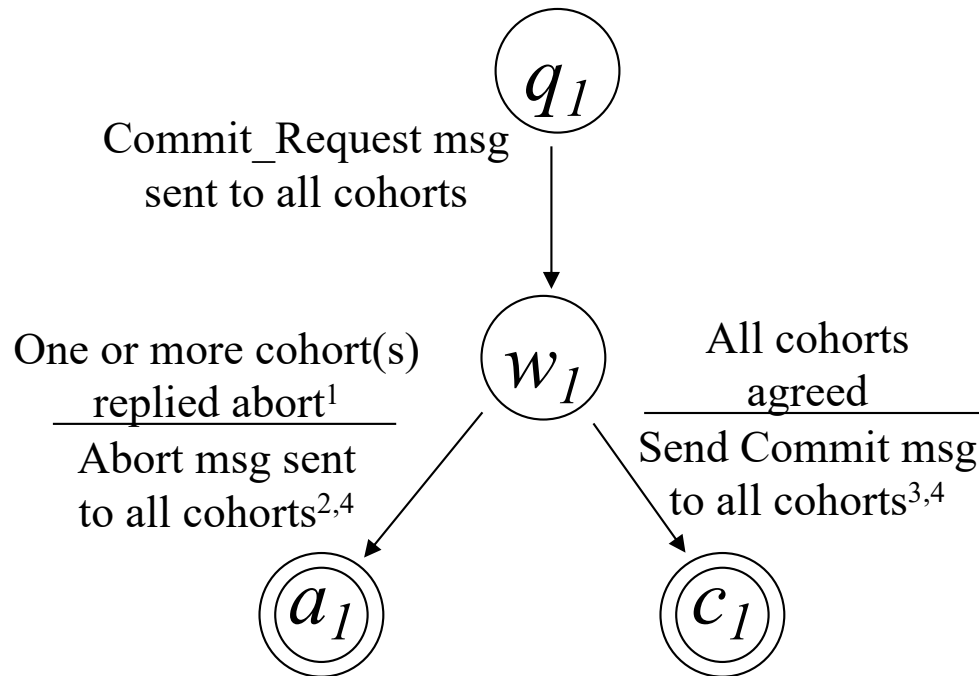
if process 1 fails before beginning its broadcast in round 2, then no nonfaulty process whose initial value is 1 ever decides.

# Strong Termination

- اگر پروسس 1 قبل از اعلام نتیجه از دسترس خارج شود؟  
راه حل: یکی از سایت‌هایی که مثلاً نظر 0 داشته به بقیه اطلاع دهد.
- اگر همه روی 1 توافق داشته باشند و 1 قبل از اینکه هر گونه پیامی ارسال کند از دسترس خارج شود؟  
بقیه سایت‌ها از نظر 1 مطلع نبوده و نمی‌توانند تصمیم‌گیری کنند.

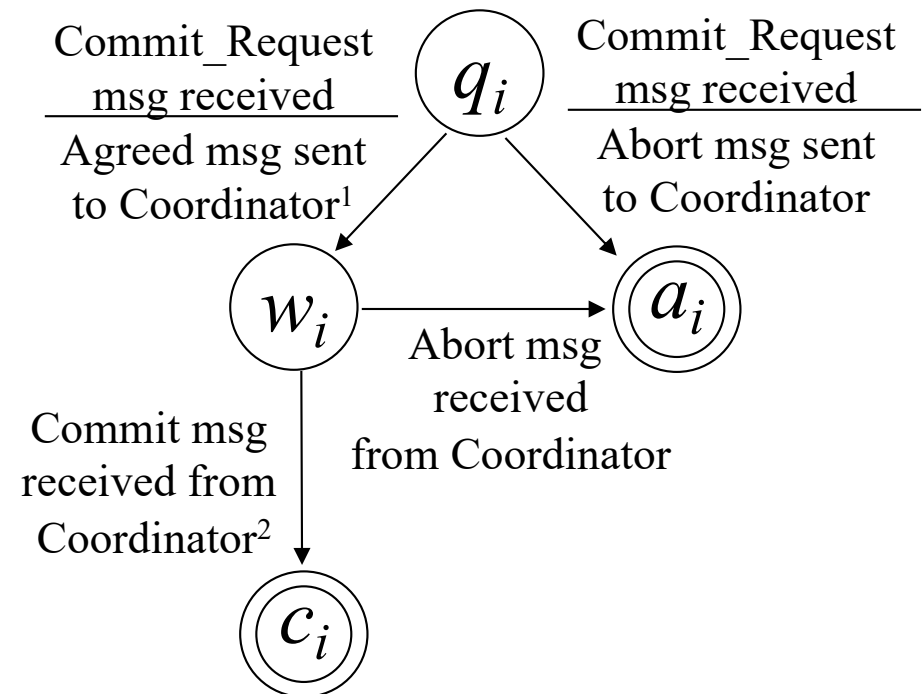
# 2-phase Commit Protocol

## Coordinator



1. Assume ABORT if there is a timeout
2. First, writes ABORT record to stable storage.
3. First, writes COMMIT record to stable storage.
4. Write COMPLETE record when all msgs confirmed.

## Cohort i (i=2,3, ..., n)



1. First, write UNDO/REDO logs on stable storage.
2. Writes COMPLETE record; releases locks

# Site Failures

Who Fails	At what point	Actions on recovery
Coordinator	before writing Commit	Send Abort messages
Coordinator	after writing Commit but before writing Complete	Send Commit messages
Coordinator	after writing Complete	None.
Cohort	before writing Undo/Redo	None. Abort will occur.
Cohort	after writing Undo/Redo	Wait for message from Coordinator.

# پیچیدگی زمانی

## ● 2 مرحله

**سوال:** چرا در حالی که الگوریتم تفاهم در شرایط خطای پروسس به  $f + 1$  مرحله نیاز داشت، پروتکل PC2 صرفاً به دو مرحله نیاز دارد؟

**پاسخ:** شرط پایان متفاوت است (weak termination)

# پیچیدگی پیامی

● در بدترین حالت  $2(n-1)$  پیام در صورتی که خطا رخ ندهد.

CONSTANT *RM*    The set of resource managers

VARIABLES

<i>rmState</i> ,	<i>rmState</i> [ <i>r</i> ] is the state of resource manager <i>r</i> .
<i>tmState</i> ,	The state of the transaction manager.
<i>tmPrepared</i> ,	The set of <i>RM</i> s from which the <i>TM</i> has received "Prepared" messages.
<i>msgs</i>	

*Messages*  $\triangleq$   
 $[type : \{ \text{"Prepared"} \}, rm : RM] \cup [type : \{ \text{"Commit"}, \text{"Abort"} \}]$

*TPTypeOK*  $\triangleq$   
 $\wedge rmState \in [RM \rightarrow \{ \text{"working"}, \text{"prepared"}, \text{"committed"}, \text{"aborted"}, \text{"Failed"} \}]$   
 $\wedge tmState \in \{ \text{"init"}, \text{"done"} \}$   
 $\wedge tmPrepared \subseteq RM$   
 $\wedge msgs \subseteq Messages$

*TPInit*  $\triangleq$   
 $\wedge rmState = [r \in RM \mapsto \text{"working"}]$   
 $\wedge tmState = \text{"init"}$   
 $\wedge tmPrepared = \{ \}$   
 $\wedge msgs = \{ \}$



---


$$\begin{aligned}
& \text{TM}RcvPrepared(r) \triangleq \\
& \quad \wedge \text{tmState} = \text{"init"} \\
& \quad \wedge [type \mapsto \text{"Prepared"}, rm \mapsto r] \in \text{msgs} \\
& \quad \wedge \text{tmPrepared}' = \text{tmPrepared} \cup \{r\} \\
& \quad \wedge \text{UNCHANGED} \langle rmState, \text{tmState}, \text{msgs} \rangle
\end{aligned}$$

$$\begin{aligned}
& \text{TMCommit} \triangleq \\
& \quad \wedge \text{tmState} = \text{"init"} \\
& \quad \wedge \text{tmPrepared} = RM \\
& \quad \wedge \text{tmState}' = \text{"done"} \\
& \quad \wedge \text{msgs}' = \text{msgs} \cup \{[type \mapsto \text{"Commit"}]\} \\
& \quad \wedge \text{UNCHANGED} \langle rmState, \text{tmPrepared} \rangle
\end{aligned}$$

$$\begin{aligned}
& \text{TMAbort} \triangleq \\
& \quad \wedge \text{tmState} = \text{"init"} \\
& \quad \wedge \text{tmState}' = \text{"done"} \\
& \quad \wedge \text{msgs}' = \text{msgs} \cup \{[type \mapsto \text{"Abort"}]\} \\
& \quad \wedge \text{UNCHANGED} \langle rmState, \text{tmPrepared} \rangle
\end{aligned}$$

$$\begin{aligned}
& \text{TMFailed} \triangleq \\
& \quad \wedge \text{tmState}' = \text{"Failed"}
\end{aligned}$$

$$\begin{aligned}
RMPrepare(r) &\triangleq \\
&\wedge rmState[r] = \text{"working"} \\
&\wedge rmState' = [rmState \text{ EXCEPT } ![r] = \text{"prepared"}] \\
&\wedge msgs' = msgs \cup \{[type \mapsto \text{"Prepared"}, rm \mapsto r]\} \\
&\wedge \text{UNCHANGED } \langle tmState, tmPrepared \rangle
\end{aligned}$$

$$\begin{aligned}
RMChooseToAbort(r) &\triangleq \\
&\wedge rmState[r] = \text{"working"} \\
&\wedge rmState' = [rmState \text{ EXCEPT } ![r] = \text{"aborted"}] \\
&\wedge \text{UNCHANGED } \langle tmState, tmPrepared, msgs \rangle
\end{aligned}$$

$$\begin{aligned}
RMRcvCommitMsg(r) &\triangleq \\
&\wedge [type \mapsto \text{"Commit"}] \in msgs \\
&\wedge rmState' = [rmState \text{ EXCEPT } ![r] = \text{"committed"}] \\
&\wedge \text{UNCHANGED } \langle tmState, tmPrepared, msgs \rangle
\end{aligned}$$

$$\begin{aligned}
RMRcvAbortMsg(r) &\triangleq \\
&\wedge [type \mapsto \text{"Abort"}] \in msgs \\
&\wedge rmState' = [rmState \text{ EXCEPT } ![r] = \text{"aborted"}] \\
&\wedge \text{UNCHANGED } \langle tmState, tmPrepared, msgs \rangle
\end{aligned}$$

$$\begin{aligned}
TPNext &\triangleq \\
&\vee TCommit \vee TMAbort \\
&\vee \exists r \in RM : \\
&\quad TMRcvPrepared(r) \vee RMPprepare(r) \vee RMChooseToAbort(r) \\
&\quad \vee RMRcvCommitMsg(r) \vee RMRcvAbortMsg(r)
\end{aligned}$$

$$TPSpec \triangleq TPInit \wedge \Box[TPNext]_{\langle rmState, tmState, tmPrepared, msgs \rangle}$$

THEOREM  $TPSpec \Rightarrow \Box TPTypeOK$

INSTANCE  $TCommit$

THEOREM  $TPSpec \Rightarrow TCSpec$

.....

### What is the behavior spec?

Initial predicate and next-state

Init: TPInit

Next: TPNext

### What is the model?

Specify the values of declared constants.

RM <- {"r1","r2","r3"}

Edit

### What to check?

☒ Deadlock

#### Invariants

Formulas true in every reachable state.

Add

Edit

Remove

#### Properties

Temporal formulas true for every possible behavior.

☐  $\Box \text{TPTypeOK}$

☐ THEOREM  $\text{TPSpec} \Rightarrow \text{TCSpec}$

☐ THEOREM  $\text{TPSpec} \Rightarrow \Box \text{TPTypeOK}$

☒  $\forall r \in \text{RM} : \langle \rangle (\text{rmState}[r] = \text{"aborted"} \vee \text{rmState}[r] = \text{"committed"})$

Add

Edit

Remove

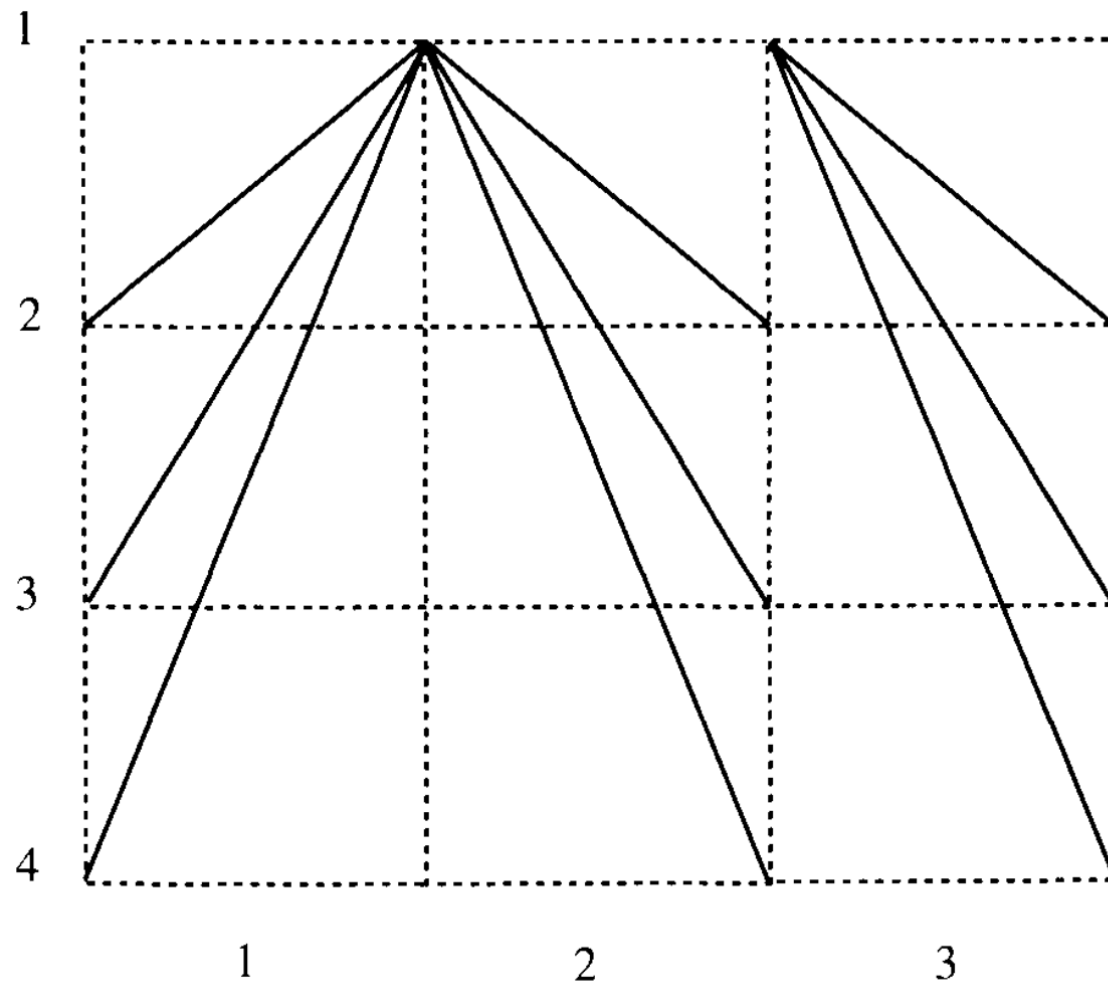
# Three-Phase Commit

**Round 1:** All processes **except for 1** send their initial values to process 1, and any process whose initial value is 0 decides 0. Process 1 collects all these values, plus its own initial value, into a vector. If all positions in this vector are filled in with 1s, then process 1 becomes ready **but does not yet decide**. Otherwise--that is, if there is some position that contains 0 or else some position that is not filled in (because no message was received from the corresponding process)--process 1 **decides 0**.

**Round 2:** If process 1 has decided 0, then it broadcasts `decide(0)`. If not, then process 1 broadcasts **ready**. Any process that receives `decide(0)` decides 0. Any process that receives **ready** becomes ready. Process 1 decides 1 if it has not already decided.

**Round 3:** If process 1 has decided 1, it broadcasts **decide(1)**. Any process that receives `decide(1)` decides 1.

# Communication pattern in Three-PhaseCommit



# Three-Phase Commit, termination protocol

**Round 4:** All (not yet failed) processes send their current status, either dec0, dec1, ready, or uncertain, to process 2. Process 2 collects all these status values, plus its own status, into a vector. Not all the positions in the vector need be filled in--process 2 just ignores those that are not. If the vector contains any dec0 values and process 2 has not already decided, then process 2 decides 0. If the vector contains any dec1 values and process 2 has not already decided, then process 2 decides 1. If all the filled-in positions in the vector contain the value uncertain, then process 2 decides 0. Otherwise--that is, if the only values in the vector are uncertain and ready and there is at least one ready--process 2 becomes ready but does not yet decide.

**Round 5:** In this and the next round, process 2 behaves similarly to process 1 in rounds 2 and 3. If process 2 has (ever) decided, then it broadcasts its decision, in a decide message. If not, then process 2 broadcasts ready. Any process that receives decide(0) or decide(1) and has not already decided, decides 0 or 1, as indicated. Any process that receives ready becomes ready. Process 2 decides 1 if it has not already decided.

**Round 6:** If process 2 has decided 1, it broadcasts decide(1). Any process that receives decide(1), and has not already decided, decides 1.

- Agreement
- Validity
- Weak Termination
- Strong Termination

در صورتی که پردازش ۱ دچار خطا نشود، برقرار است.



# تحلیل پیچیدگی

## ● پیچیدگی زمانی:

$3n$  مرحله در صورتی که  $n$  پروسس دچار خطا شود.

الگوریتم stopping agreement با حدود  $n$  مرحله به جواب می‌رسد.

سوال: چرا کامیت سه مرحله‌ای؟

پاسخ: در شرایط نرمال، الگوریتم در سه مرحله کار می‌کند که بسیار سریع‌تر خواهد بود.

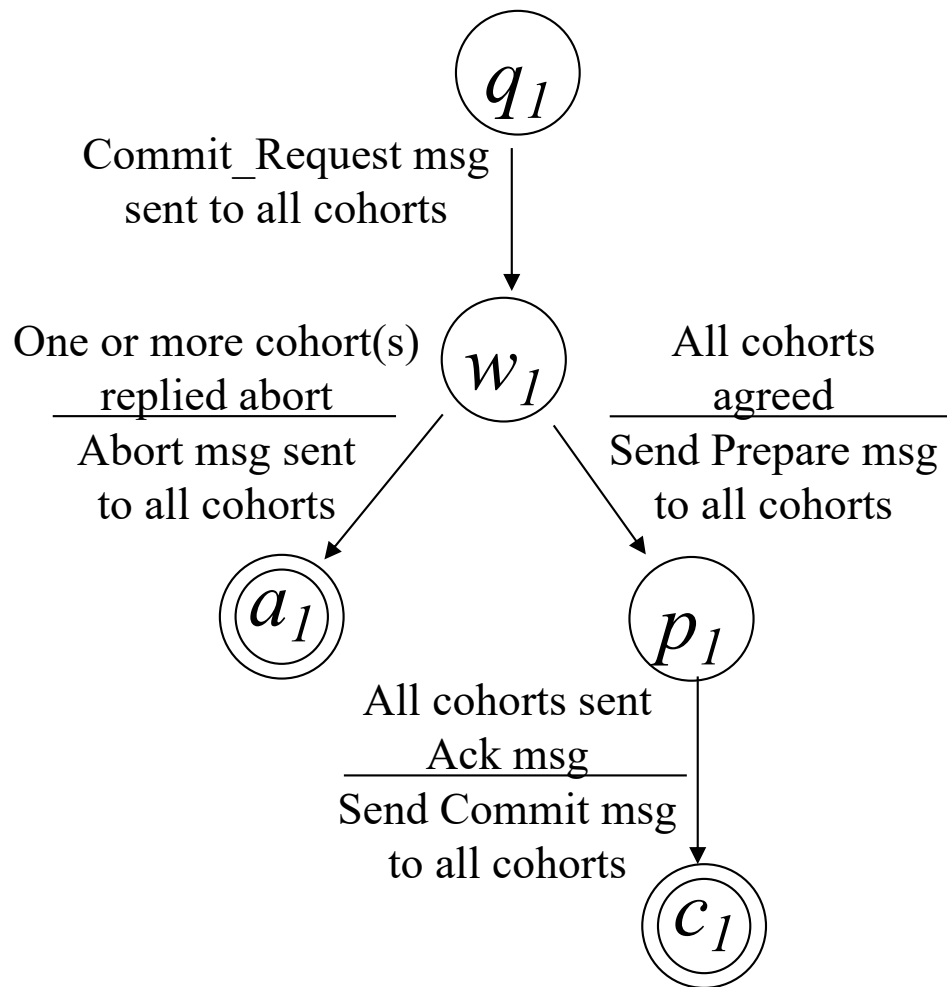
## ● پیچیدگی پیامی:

$O(n)$

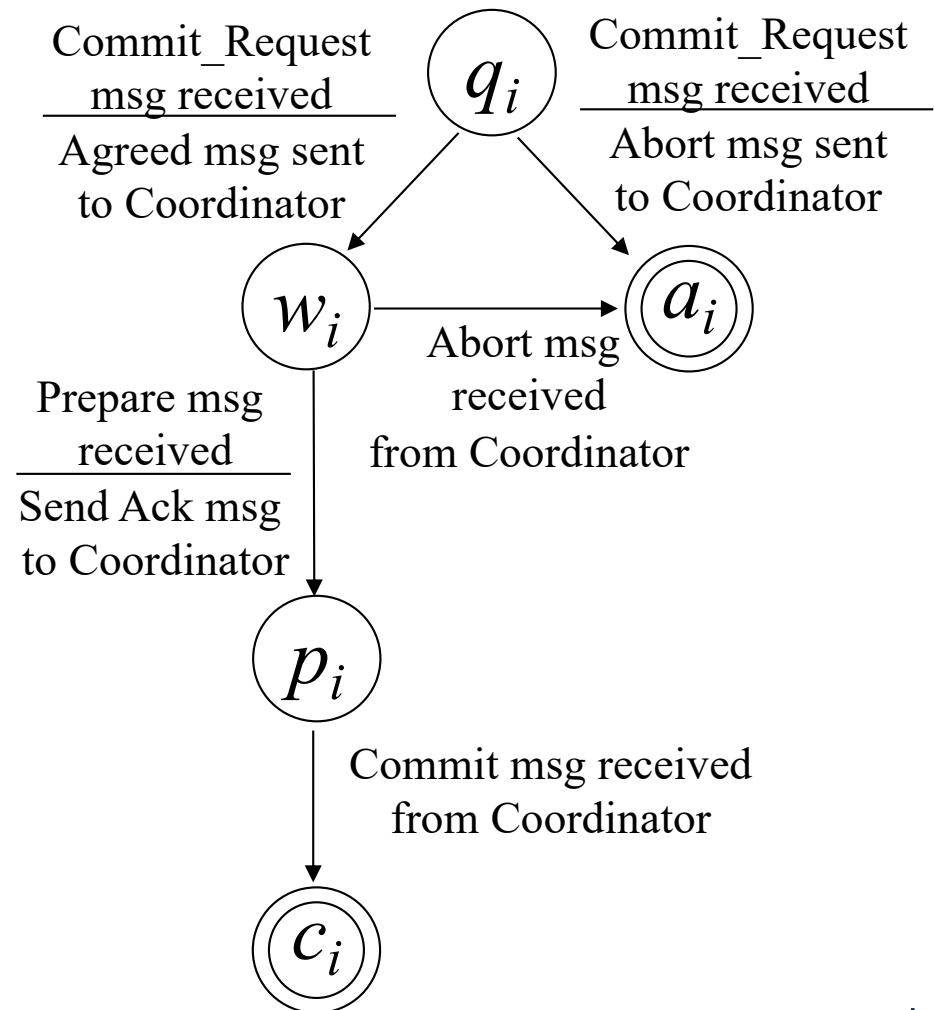
**قضیه:** هر الگوریتم که مساله کامیت را حل کند، حداقل به  $2n - 2$  پیام در صورت عدم بروز خطا و شرایطی که همه با 1 آغاز کنند، نیاز دارد.

# 3-phase Commit Protocol

## Coordinator

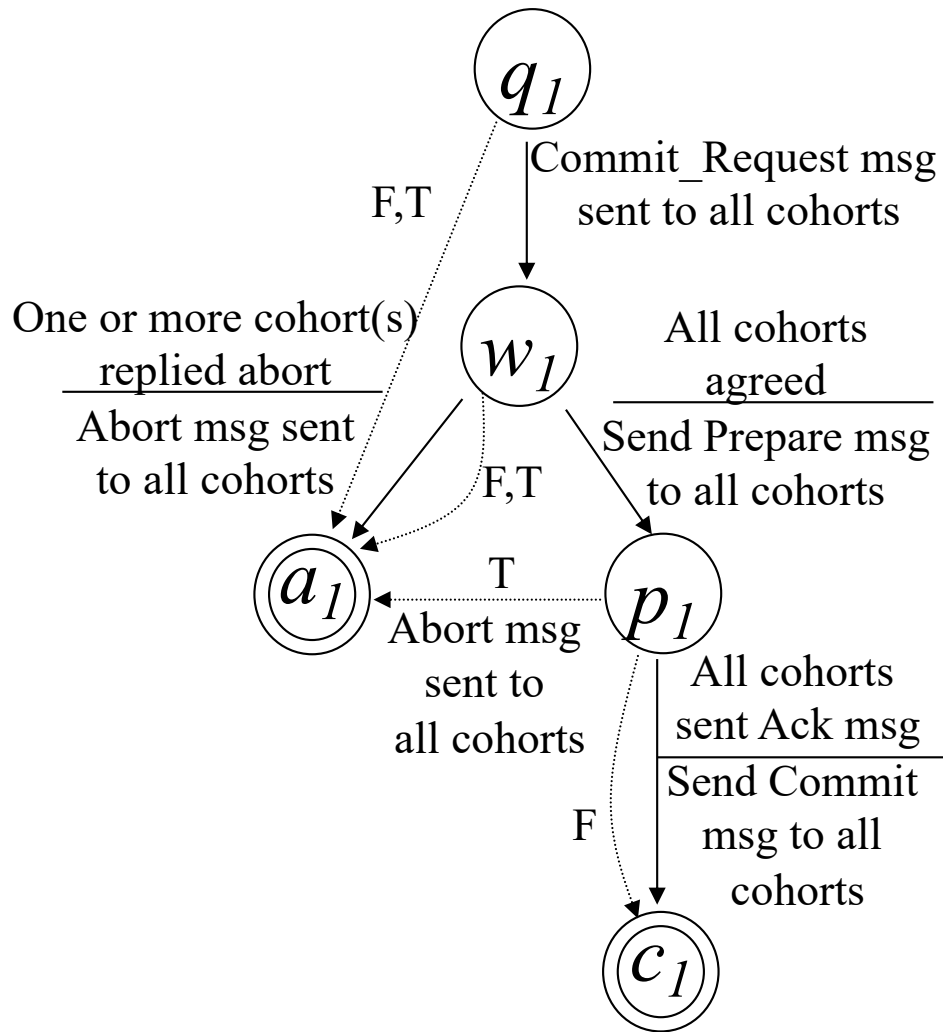


## Cohort $i$ ( $i=2,3, \dots, n$ )

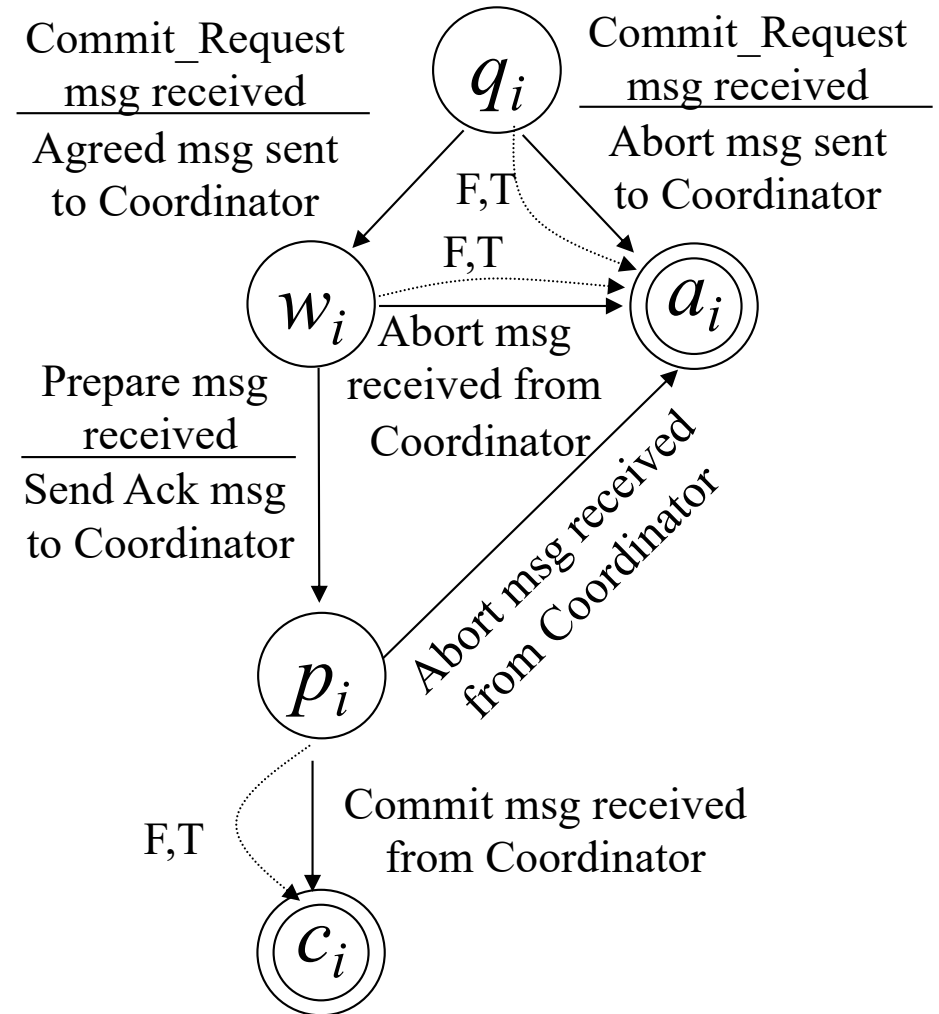


# Timeout and Failure Transitions

## Coordinator



## Cohort i ( $i=2,3, \dots, n$ )



$\text{---T---} \rightarrow$  Timeout Transition    
  $\text{---F---} \rightarrow$  Failure Transition    
  $\text{---F,T---} \rightarrow$  Failure/Timeout Transition

# References

- **Lynch, Chapter 7**