of their first waiting intervals. At the end of its second waiting interval, p broadcasts a READY message indicating that it is ready to begin the next round. However, if p receives f + 1 READY messages during its second waiting interval, it terminates its second interval early, and goes ahead and broadcasts READY. As soon as p receives n - f READY messages, it updates the clock according to the adjustment calculated earlier, and begins its next round by broadcasting its new clock value. (This algorithm uses some ideas from [3].)

A process need only keep clock differences for one round at a time. The waiting intervals are designed so that during round i a nonfaulty process p will not receive a READY message from another nonfaulty process until p has finished collecting round i clock values. Round i + 1 clock values are not broadcast until after READY is broadcast, so p will certainly not receive round i + 1 clock values until after it has finished collecting round i clock values. However, round i + 1 clock values might arrive during the second waiting interval and while the process is collecting READY messages. As a result, the adjustment is calculated at the end of the first waiting interval and the difference for any round i + 1 clock value received during round i is decremented by the amount of the adjustment.

## 5.2.2 Code for an Arbitrary Process

Global constants:  $\delta$ ,  $\epsilon$ ,  $\rho$ ,  $\eta$ , f: as usual.

Local variables (all initially arbitrary):

- T: clock time at which current round began.
- U: clock time at which the first waiting period is to end.
- V: clock time at which the second waiting period is to end.
- DIFF: array of clock differences between other processes and this one for current round.
- SENT-READY: set of processes from whom READY messages have been received in current round.
- CORR: correction variable.
- · A: adjustment to clock.

The code is in Figure 5-1.

```
beginstep(w)
 do forever /* each iteration is a round */
      T := NOW
      broadcast(T)
      U := T + (1 + \rho)(2\delta + 4\varepsilon)
      set-timer(U)
/* first waiting interval: collect clock values */
      while ~(w = TIMER & NOW = U) do
           if w = (m,q) then DIFF[q] := m + \delta - NOW endif
           endstep
           beginstep(w)
           endwhile
/* end of first waiting interval */
      A := mid(reduce(DIFF))
     V := U + (1 + \rho)(4\varepsilon + 4\rho(\delta + 2\varepsilon) + 2\rho^{2}(\delta + 2\varepsilon))
      set-timer(V)
     SENT-READY := Ø
/* second waiting interval: collect READY messages and clock values
   for next round */
     while ~(w = TIMER & NOW = V) do
           if w = (READY,q) then
                SENT-READY := SENT-READY ∪ {q}
                 if |SENT-READY| = f + 1 then exit endif
           elseif w = (m,q) then DIFF[q] := m + \delta - NOW endif
           endstep
           beginstep(w)
           endwhile
/* end of second waiting interval due to timer or f + 1 READY messages */
     broadcast(READY)
     endstep
     beginstep(w)
/* collect n - f READY messages and next round clock values */
     while true do
           if w = (READY, q) then
                SENT-READY := SENT-READY U {q}
                if |SENT-READY| = n - f then exit endif
          elseif w = (m,q) then DIFF[q] := m + \delta - NOW endif
          endstep
          beginstep(w)
          endwhile
/* update clock and begin next round */
     DIFF := DIFF - A
     CORR := CORR + A
     endstep
     beginstep(w)
     enddo
```

Figure 5-1: Algorithm 5-1, Establishing Synchronization

## 5.3 Analysis

We will use the following notation in addition to that introduced already.

- $VAL_{p}^{i}(q)$  is the value of q's round i message to p.
- DIFF<sup>i</sup><sub>p</sub>(q) = VAL<sup>i</sup><sub>p</sub>(q) +  $\delta$  ARR<sup>i</sup><sub>p</sub>(q), p's estimate of the difference between p's and q's clocks.
- DIFF $_{p}^{i}$  is the multiset of DIFF $_{p}^{i}$ (q) values.
- ti is the real time when p begins round i.
- $u_p^i$  is the real time when p begins the second waiting interval during round i.
- v<sup>i</sup><sub>p</sub> is the real time when p sends READY during round i (and thus ends the second waiting interval).
- arr<sup>i</sup><sub>p</sub>(q) is the real time when p first receives a round i clock value from q.
- $\bullet$  rdy $_{p}^{i}(q)$  is the real time when p first receives READY from q during round i.
- tmax<sup>i</sup> = max{t<sup>i</sup><sub>p</sub>} for p nonfaulty, the latest real time when a nonfaulty process begins round i.
- B<sup>i</sup> = max{|C<sub>p</sub><sup>i</sup>(tmax<sup>i</sup>) C<sub>q</sub><sup>i</sup>(tmax<sup>i</sup>)|} for p and q nonfaulty, the maximum difference between nonfaulty clock values at tmax<sup>i</sup>.

Note that tmax<sup>i</sup> has a slightly different meaning from that in Chapter 4.

From now on, terms of order  $\rho^2$  and higher will be ignored. Since  $10^{-6}$  seconds is an often quoted reasonable value for  $\rho$  [5, 7, 9], terms of order  $\rho^2$  are negligible. The second-order terms in the assignment to V in line 13 of the code are needed for strict correctness, but will not appear in the analysis.

Lemma 5-2 proves together inductively that the time between two nonfaulty processes beginning a round is bounded, and that when a nonfaulty process q receives READY from another nonfaulty process, q has already finished the first waiting period. First we show a preliminary lemma needed by Lemma 5-2.

**Lemma 5-1:** Let  $i \ge 0$ , p and q be any nonfaulty processes, and r be the first nonfaulty process to send READY at round i. Then

$$rdy^i_{\ \alpha}(p) \geq u^i_{\ \alpha} - (t^i_{\ \alpha} - t^i_{\ r}) \, + \, \delta \, + \, 3\epsilon.$$

**Proof:** Since r is the first nonfaulty process to send READY, it doesn't send until its full second waiting interval has elapsed. Then

$$rdy_{q}^{i}(p) \ge v_{p}^{i} + \delta - \epsilon$$

$$\geq v^i + \delta - \varepsilon$$

 $\geq$   $t_r^i + (2\delta + 4\epsilon) + (4\epsilon + 4\rho(\delta + 2\epsilon)) + \delta - \epsilon$ , by definition of  $v_r^i$  and the upper

$$=t^{i}_{r}+3\delta+7\varepsilon+4\rho\delta+8\rho\varepsilon,$$

$$u_{q}^{i} \leq t_{r}^{i} + (t_{q}^{i} - t_{r}^{i}) + (u_{q}^{i} - t_{q}^{i})$$

 $\leq t_r^i + (t_q^i - t_r^i) + (1 + \rho)^2 (2\delta + 4\epsilon)$ , by definition of  $u_q^i$  and the lower bound on the drift rate

$$=t_r^i+(t_q^i-t_r^i)+2\delta+4\varepsilon+4\rho\delta+8\rho\varepsilon.$$

Thus, 
$$t_r^i \ge u_q^i - (t_q^i - t_r^i) - 2\delta - 4\varepsilon - 4\rho\delta - 8\rho\varepsilon$$
, implying

$$rdy_{q}^{i}(p) \geq u_{q}^{i} - (t_{q}^{i} - t_{r}^{i}) - 2\delta - 4\varepsilon - 4\rho\delta - 8\rho\varepsilon + 3\delta + 7\varepsilon + 4\rho\delta + 8\rho\varepsilon$$

= 
$$u_q^i - (t_q^i - t_r^i) + \delta + 3\varepsilon$$
.

Lemma 5-2: For any nonfaulty processes p and q and any  $i \ge 0$ ,

(a) 
$$|t_p^i - t_q^i| \le \delta + 3\epsilon$$
, and

(b) 
$$rdy^{i}_{q}(p) \ge u^{i}_{q}$$
.

(b)  $rdy_{q}^{i}(p) \ge u_{q}^{i}$ . Proof: We proceed by induction on i.

Basis: i = 0.

- (a)  $|t_0^0 t_0^0| \le \delta + \epsilon$ , because as soon as p wakes up, it sends its round 0 message to all other processes. The receipt of this message, which occurs at most  $\delta + \varepsilon$  later, causes q to begin round 0, if it hasn't already done so.
- (b) Let r be the first nonfaulty process to send READY at round 0. By Lemma 5-1,

$$rdy_{q}^{0}(p) \ge u_{q}^{0} - (t_{q}^{0} - t_{r}^{0}) + \delta + 3\epsilon$$

$$\geq u^0 - (\delta + \epsilon) + \delta + 3\epsilon$$
, by part (a)

Induction: Assume for i - 1 and show for i.

(a) Let s be the first nonfaulty process to begin round i. Then s receives n - f READY messages during its round i - 1.(after  $u^{i-1}_s$ ). At least n - 2f of them are from nonfaulty processes by part (b) of the induction hypothesis. These n - 2f nonfaulty processes

also send READY messages to all the other processes. By  $t_s^i+2\epsilon$ , every nonfaulty process receives at least  $n-2f\geq f+1$  READY messages and broadcasts READY. Thus q receives n-f READY messages by  $t_s^i+2\epsilon+\delta+\epsilon$ . Thus,

$$t_{q}^{i} \leq t_{s}^{i} + \delta + 3\varepsilon$$

 $\leq t_{p}^{i} + \delta + 3\varepsilon$ , by choice of s,

which implies  $t_q^i - t_p^i \le \delta + 3\epsilon$ .

By reversing the roles of p and q in the above argument, we obtain  $t^i_{\ p} - t^i_{\ q} \leq \delta + 3\epsilon$ .

(b) Let r be the first nonfaulty process to send READY at round i. By Lemma 5-1,

$$rdy_{q}^{i}(p) \ge u_{q}^{i} - (t_{q}^{i} - t_{r}^{i}) + \delta + 3\varepsilon$$

$$\ge u_{q}^{i} - (\delta + 3\varepsilon) + \delta + 3\varepsilon, \text{ by part (a)}$$

$$= u_{q}^{i}. \blacksquare$$

Next we show that a process waits a sufficient length of time to receive clock values from all nonfaulty processes before beginning the second waiting interval in a round.

Lemma 5-3: Let p and q be nonfaulty, and  $i \ge 0$ . Then  $arr_{p}^{i}(q) \le u_{p}^{i}$ .

**Proof:** By the lower bound on the drift rate,  $u^i_{\ p} \geq t^i_{\ p} + 2\delta + 4\epsilon$ . Lemma 5-2 implies that q sends its round i clock value by  $t^i_{\ p} + \delta + 3\epsilon$ . Thus  $arr^i_{\ p}(q) \leq t^i_{\ p} + 2\delta + 4\epsilon \leq u^i_{\ p}$ .

The next two lemmas bound how long a round can last for one process. First we bound how long a process must wait after sending READY to receive n – f READY messages.

Lemma 5-4: For p nonfaulty and  $i \ge 0$ ,  $t^{i+1}_{p} - v^{i}_{p} \le 2\delta + 4\epsilon + 4\rho(\delta + 4\epsilon)$ .

**Proof:** The worst case occurs if p is as far ahead of the other nonfaulty processes as possible, its clock is fast, the other clocks are slow, and the slow processes' READY messages take as long as possible to arrive. However, as soon as they arrive, p begins the next round. Let q be one of the slow nonfaulty processes.

$$\begin{split} t^{i+1}_{p} - v^{i}_{p} &= (t^{i+1}_{p} - v^{i}_{q}) + (v^{i}_{q} - u^{i}_{q}) + (u^{i}_{q} - t^{i}_{q}) + (t^{i}_{q} - t^{i}_{p}) - (v^{i}_{p} - u^{i}_{p}) - (u^{i}_{p} - t^{i}_{p}) \\ &\leq (\delta + \varepsilon) + (1 + \rho)^{2} (4\varepsilon + 4\rho(\delta + 2\varepsilon)) + (1 + \rho)^{2} (2\delta + 4\varepsilon) + (\delta + 3\varepsilon) \\ &- (4\varepsilon + 4\rho(\delta + 2\varepsilon)) - (2\delta + 4\varepsilon) \end{split}$$

=  $2\delta + 4\varepsilon + 4\rho(\delta + 4\varepsilon)$ , ignoring  $\rho^2$  terms.

Lemma 5-5: For any nonfaulty process p and any  $i \ge 0$ ,

$$\begin{split} t^{i+1}_{\quad p} - t^{i}_{p} &\leq 4\delta + 12\epsilon + 4\rho(3\delta + 10\epsilon). \\ \text{Proof: } t^{i+1}_{\quad p} - t^{i}_{p} &= (t^{i+1}_{\quad p} - v^{i}_{p}) + (v^{i}_{p} - u^{i}_{p}) + (u^{i}_{p} - t^{i}_{p}) \end{split}$$

$$\leq 2\delta + 4\varepsilon + 4\rho(\delta + 4\varepsilon) + (v_p^i - u_p^i) + (u_p^i - t_p^i), \text{ by Lemma 5-4}$$

$$\leq 2\delta + 4\varepsilon + 4\rho(\delta + 4\varepsilon) + (1 + \rho)^2(4\varepsilon + 4\rho(\delta + 2\varepsilon)) + (1 + \rho)^2(2\delta + 4\varepsilon)$$

$$= 4\delta + 12\varepsilon + 4\rho(3\delta + 10\varepsilon). \blacksquare$$

Now we give an upper bound on how far apart tmax<sup>i</sup> and tmax<sup>i+1</sup> can be.

Lemma 5-6: For any  $i \ge 0$ ,

$$tmax^{i+1} - tmax^{i} \le 4\delta + 12\varepsilon + 4\rho(3\delta + 10\varepsilon).$$

**Proof:** Let p be the nonfaulty process such that  $t^{i+1}_{p} = tmax^{i+1}$ . Then

$$tmax^{i+1} - tmax^{i} = t^{i+1}_{p} - tmax^{i} \le t^{i+1}_{p} - t^{i}_{p}$$

$$\leq 4\delta + 12\varepsilon + 4\rho(3\delta + 10\varepsilon)$$
, by Lemma 5-5.

Lemma 5-7 bounds the amount of real time between the time a nonfaulty process receives a round i message from another nonfaulty process and the time the last nonfaulty process begins round i + 1.

Lemma 5-7: For any  $i \ge 0$  and nonfaulty processes p and q,

$$\begin{split} & \operatorname{tmax}^{i+1} - \operatorname{arr}^{i}_{p}(q) \leq 5\delta + 19\epsilon + 4\rho(3\delta + 10\epsilon). \\ & \operatorname{\textbf{Proof:}} \operatorname{tmax}^{i+1} - \operatorname{arr}^{i}_{p}(q) = (\operatorname{tmax}^{i+1} - \operatorname{t}^{i+1}_{p}) + (\operatorname{t}^{i+1}_{p} - \operatorname{t}^{i}_{p}) + (\operatorname{t}^{i}_{p} - \operatorname{t}^{i}_{q}) - (\operatorname{arr}^{i}_{p}(q) - \operatorname{t}^{i}_{q}) \\ & \leq (\delta + 3\epsilon) + (4\delta + 12\epsilon + 4\rho(3\delta + 10\epsilon)) + (\delta + 3\epsilon) - (\delta - \epsilon), \, \text{by Lemmas 5-2 and} \\ & 5\text{-5 and the lower bound on the message delay} \end{split}$$

= 
$$5\delta + 19\varepsilon + 4\rho(3\delta + 10\varepsilon)$$
.

The next lemma bounds the error in a nonfaulty process' estimate of another nonfaulty process' local time at a particular real time.

Lemma 5-8: Let p and r be nonfaulty. Then

$$\begin{split} &|\mathsf{DIFF}^{i}_{p}(r) + C^{i}_{p}(\mathsf{tmax}^{i+1}) - C^{i}_{r}(\mathsf{tmax}^{i+1})| \leq \varepsilon + \rho(11\delta + 39\varepsilon). \\ &\mathsf{Proof:} \, |\mathsf{DIFF}^{i}_{p}(r) + C^{i}_{p}(\mathsf{tmax}^{i+1}) - C^{i}_{r}(\mathsf{tmax}^{i+1})| \\ &= |\mathsf{VAL}^{i}_{p}(r) + \delta - \mathsf{ARR}^{i}_{p}(r) + C^{i}_{p}(\mathsf{tmax}^{i+1}) - C^{i}_{r}(\mathsf{tmax}^{i+1})|. \end{split}$$

If the quantity in the absolute value signs is negative, then this expression is equal to

$$\begin{split} &C_{r}^{i}(tmax^{i+1})-C_{p}^{i}(tmax^{i+1})+C_{p}^{i}(arr_{p}^{i}(r))-\delta-VAL_{p}^{i}(r)\\ &\leq C_{r}^{i}(tmax^{i+1})-C_{p}^{i}(tmax^{i+1})+C_{p}^{i}(arr_{p}^{i}(r))-\delta-C_{r}^{i}(arr_{p}^{i}(r)-\delta-\epsilon), \text{ since the delay is at }\\ &most\ \delta+\epsilon \end{split}$$

$$\leq C_r^i(\operatorname{tmax}^{i+1}) - C_p^i(\operatorname{tmax}^{i+1}) + C_p^i(\operatorname{arr}_p^i(r)) - \delta - C_r^i(\operatorname{arr}_p^i(r)) + (1 + \rho)(\delta + \varepsilon), \text{ since the clock drift is at most } 1 + \rho$$

$$= (C_r^i(\operatorname{tmax}^{i+1}) - C_p^i(\operatorname{tmax}^{i+1})) - (C_r^i(\operatorname{arr}_p^i(r)) - C_p^i(\operatorname{arr}_p^i(r))) - \delta + \delta + \varepsilon + \rho\delta + \rho\varepsilon$$

$$\leq 2\rho(\operatorname{tmax}^{i+1} - \operatorname{arr}_p^i(r)) + \varepsilon + \rho\delta + \rho\varepsilon, \text{ by Lemma } 4-2$$

$$\leq 2\rho(5\delta + 19\varepsilon) + \varepsilon + \rho\delta + \rho\varepsilon, \text{ by Lemma } 5-7$$

$$= \varepsilon + \rho(11\delta + 39\varepsilon).$$

If the quantity in the absolute value signs is positive, a similar argument shows that  $|\mathrm{DIFF}_{p}^{i}(r) + C_{p}^{i}(\mathrm{tmax}^{i+1}) - C_{r}^{i}(\mathrm{tmax}^{i+1})| \leq \varepsilon + \rho(11\delta + 37\varepsilon)$ .

The next lemma bounds how far apart two processes' i-th clocks are at the time when the last process begins round i + 1. The bound is in terms of how far apart the clocks are when the last process begins round i.

Lemma 5-9: For any nonfaulty p and q, and any i,

$$\begin{aligned} &|C_{p}^{i}(tmax^{i+1}) - C_{q}^{i}(tmax^{i+1})| \leq B^{i} + 8\rho(\delta + 3\epsilon). \\ &Proof: |C_{p}^{i}(tmax^{i+1}) - C_{q}^{i}(tmax^{i+1})| \\ &\leq |C_{p}^{i}(tmax^{i}) - C_{q}^{i}(tmax^{i})| + |(C_{q}^{i}(tmax^{i+1}) - C_{q}^{i}(tmax^{i+1})) - (C_{p}^{i}(tmax^{i}) - C_{q}^{i}(tmax^{i}))| \\ &\leq B^{i} + 2\rho(tmax^{i+1} - tmax^{i}), \text{ by definition of } B^{i} \text{ and Lemma 4-2} \\ &\leq B^{i} + 2\rho(4\delta + 12\epsilon), \text{ by Lemma 5-6 and ignoring } \rho^{2} \text{ terms} \end{aligned}$$

$$= B^{i} + 8\rho(\delta + 3\epsilon). \blacksquare$$

Now we can state the main result, bounding  $B^{i+1}$  in terms of  $B^{i}$ .

Theorem 5-10: 
$$B^{i+1} \le \frac{1}{2}B^i + 2\varepsilon + 2\rho(11\delta + 39\varepsilon)$$
.  
Proof:  $B^{i+1} = \max\{|C^{i+1}_p(\tan x^{i+1}) - C^{i+1}_q(\tan x^{i+1})|\}$  for nonfaulty p and q.  
Let  $x = \varepsilon + \rho(11\delta + 39\varepsilon)$ .

We now define three multisets U, V, and W that satisfy the hypotheses of Lemma A-4. Let

$$U = DIFF_{p}^{i} + C_{p}^{i}(tmax^{i+1}),$$

$$V = DIFF_{q}^{i} + C_{q}^{i}(tmax^{i+1}), \text{ and}$$

$$W = \{C_{r}^{i}(tmax^{i+1}): r \text{ is nonfaulty}\}.$$

U and V have size n; W has size n - f.

Define an injection from W to U as follows. Map each element  $C_r^i$  in W to DIFF $_p^i$ (r) +  $C_n^i$ (tmax $^{i+1}$ ) in U. Since Lemma 5-8 implies that

$$|\mathsf{DIFF}^{i}_{p}(r) + C^{i}_{p}(tmax^{i+1}) - C^{i}_{r}(tmax^{i+1})| \le x$$

for all the n – f nonfaulty processes,  $d_x(W,U) = 0$ . Similarly,  $d_x(W,V) = 0$ .

By Lemma 5-9, diam(W)  $\leq$  B<sup>i</sup> + 8 $\rho$ ( $\delta$  + 3 $\epsilon$ ). Thus, Lemma A-4 implies

 $|mid(reduce(U)) - mid(reduce(V))| \le \frac{1}{2} diam(W) + 2x$ 

$$= \frac{1}{2}B^{i} + 2\varepsilon + 2\rho(11\delta + 39\varepsilon).$$

Since mid(reduce(U)) = mid(reduce(DIFF<sub>p</sub> +  $C_p^i(tmax^{i+1})))$ 

$$= ADJ_p^i + C_p^i(tmax^{i+1})$$

$$= C^{i+1}_{p}(tmax^{i+1})$$

and similarly mid(reduce(V)) =  $C^{i+1}_{q}(tmax^{i+1})$ , the result follows.

We obtain an approximate bound on how closely this algorithm will synchronize the clocks by considering the limit of B<sup>i</sup> as the round number increases without bound.

**Theorem 5-11:** This algorithm can synchronize clocks to within  $4\varepsilon + 4\rho(11\delta + 39\varepsilon)$ .

Proof:  $\lim_{i\to\infty} B^i$ 

$$= \lim_{i \to \infty} [\mathsf{B}^0/2^i + (1 + 1/2 + \dots + 1/2^{i-1})(2\varepsilon + 2\rho(11\delta + 39\varepsilon))]$$

=  $4\varepsilon + 4\rho(11\delta + 39\varepsilon)$ , since the limit of the geometric series is 2.

As was the case for Algorithm 4-1, if the number of processes, n, increases while f, the number of faulty processes remained fixed, a greater closeness of synchronization can be achieved by modifying Algorithm 5-1 so that it computes the mean instead of the midpoint of the range of values. which approaches  $2\varepsilon + 2\rho P$  as n approaches infinity.

After modifying Algorithm 5-1, we get

$$B^j \leq B^{j-1} f/(n-2f) \, + \, 2\varepsilon \, + \, 2\rho(11\delta \, + \, 39\varepsilon).$$

This is the same as

$$B^{j} \leq B^{0}f/(n-2f) + (1 - (f/(n-2f))^{j})/(1 - f/(n-2f))(2\epsilon + 2\rho(11\delta + 39\epsilon),$$

which approaches  $2\varepsilon + 2\rho(11\delta + 39\varepsilon)$  as n approaches infinity.

## 5.4 Determining the Number of Rounds

The nonfaulty processes must determine how many rounds of this algorithm must be run to establish the desired degree of synchronization before switching to the maintenance algorithm. The basic idea is for each nonfaulty process p to estimate  $B^0$ , and then calculate a sufficient number of rounds, NROUNDS $_p$ , using the known rate of convergence.  $B^0$  is estimated by having p calculate an overestimate and an underestimate for  $C^0_q(tmax^0)$  for each q, and letting the estimated  $B^0$  be the difference between the maximum overestimate and the minimum underestimate.

Let p's overestimate for  $C_q^0(tmax^0)$  be  $OVER_p(q)$  and p's underestimate for  $C_q^0(tmax^0)$  be  $UNDER_p(q)$ .

For the overestimate, we assume that q's clock is fast, and that the maximum amount of time elapses between  $t_q^0$  (when q sent the message) and  $t_q^0$ . That maximum is  $\delta + \epsilon$  since every nonfaulty process begins round 0 as soon as it receives a message. Thus,

$$OVER_{\rho}(q) = VAL_{\rho}^{0}(q) + (1 + \rho)(\delta + \varepsilon).$$

Similarly, we can derive the underestimate. We assume that q is the last nonfaulty process to begin round 0. Thus,

$$UNDER_{p}(q) = VAL_{p}^{0}(q).$$

Process p computes its estimate of B<sup>0</sup>,

$$B_p^0 = \max_{q} \{OVER_p(q)\} - \min_{q} \{UNDER_p(q)\}.$$

Now p estimates how many rounds are needed until the spread is close enough. There is a predetermined  $\gamma \ge 4\epsilon + 4\rho(11\delta + 39\epsilon)$ , which is the desired closeness of synchronization for the start-up algorithm. After j rounds,

$$B^{j} \le B^{0}_{p}/2^{j} + (1 + 1/2 + ... + 1/2^{j-1})(2\varepsilon + 2\rho(11\delta + 39\varepsilon)).$$

Process p sets the right hand side equal to  $\gamma$  and solves for j to obtain its estimate of the required number of rounds, NROUNDS<sub>p</sub>.

Now each process executes a Byzantine Agreement protocol on the vector of NROUNDS values, one value for each process. The processes are guaranteed to have the same vector at the end of the Byzantine Agreement protocol. Each process chooses the (f + 1)-st smallest element of the resulting vector as the required number of rounds. The smallest number of rounds computed by a nonfaulty process will suffice to achieve the desired closeness of synchronization. Variations in the number of rounds computed by different nonfaulty processes are due to spurious values introduced by faulty processes and to different message delays. However, the range computed by any nonfaulty process is guaranteed to include the actual values of all nonfaulty processes at  $\max^0$ , so the range determined by the process that computes the smallest number of rounds also includes all the actual values. In order to guarantee that each process chooses a number of rounds that is at least as large as the smallest one computed by a nonfaulty process, it chooses the (f + 1)-st smallest element of the vector of values.

Any Byzantine Agreement protocol requires at least f+1 rounds. The processes can execute this algorithm in parallel with the clock synchronization algorithm, beginning at round 0. The clock synchronization algorithm imposes a round structure on the processes' communications. The Byzantine Agreement algorithm can be executed using this round structure. Each BA message can also include information needed for the clock synchronization algorithm (namely, the current clock value). However, the processes will always need to do at least f+2 rounds, one to obtain the estimated number of rounds and f+1 for the Byzantine Agreement algorithm.

## 5.5 Switching to the Maintenance Algorithm

After the processes have done the required number of rounds (denoted by r throughout this section) of the start-up algorithm, they cease executing it. The processes should begin the maintenance algorithm as soon as possible after ending the start-up algorithm in order to minimize the inaccuracy introduced by the clock drift.

In the maintenance algorithm each process broadcasts its clock value when its clock reaches  $T^i$ , for i=0,1,..., where  $T^{i+1}=T^i+P$ . Let  $T^0$  be a multiple of P. It is shown below in Lemma 5-13 that the first multiple of P reached by nonfaulty p's clock after finishing the required r rounds differs by at most one from the first multiple reached by nonfaulty q's clock after the r rounds. When a process reaches the first multiple of P after it has ended the start-up algorithm, it broadcasts its clock value as in the maintenance algorithm, but doesn't update its clock. At the next multiple of P, the process begins the full maintenance algorithm by broadcasting its clock