



DOI:10.1145/3513259

**A rigorous journey from the bakery algorithm to a distributed state machine.**

BY LESLIE LAMPORT

# Deconstructing the Bakery to Build a Distributed State Machine

IN THIS ARTICLE, the reader and I will journey between two concurrent algorithms of the 1970s that are still studied today. The journey begins at the bakery algorithm<sup>9</sup> and ends at an algorithm for implementing a distributed state machine.<sup>12</sup> I hope we enjoy the voyage and perhaps even learn something.

The bakery algorithm ensures processes execute a critical section of code one at a time. A process trying to execute that code chooses a number it believes to be higher than the numbers chosen by other such processes. The process with the lowest number goes first, with ties broken by process name. In the distributed state-machine algorithm, each process maintains a logical clock, with the clocks being synchronized by having a process

include its clock value in the messages it sends. Commands to the state machine are ordered according to the value of a process's clock when it issues a command, with ties broken by process name.

The similarity between the bakery algorithm's numbers and the state-machine algorithm's clocks has been noticed, but I know of no previous rigorous connection between them. Our trip makes this connection, going from the bakery algorithm to the state-machine algorithm through a sequence of algorithms, each (except the first) derived from the preceding one.

The first algorithm on the journey is a straightforward generalization of the bakery algorithm, mainly by allowing a process to read other processes' numbers in an arbitrary order. We then deconstruct this algorithm by having each process maintain multiple copies of its number, one for each other process. Next is a distributed version of the deconstructed algorithm obtained by having each copy of a process  $i$ 's number kept by the process that reads it, where  $i$  writes the value stored at another process by sending a message to that process. We then modify this distributed algorithm to ensure that numbers increase with each execution of the critical section. Finally, we arrive at the distributed state-machine algorithm by forgetting about critical sections and just using the numbers as logical clocks.

Not only do our algorithms date from the 1970s, but the path between them is one that could have been followed at that time. The large amount of related work done since then has neither influenced nor obviated any part of the route. At the end of our journey, a concluding section discusses that related work and why the algorithms that begin and end our path are still studied today. The correctness proofs in our journey are informal, much as they would have been in the 1970s. More modern, rigorous proofs are discussed in the concluding section.

A62



A61



A60



A59

A57

A56

## The Original Bakery Algorithm


The bakery algorithm solves the mutual-exclusion problem introduced and solved by Edsger Dijkstra.<sup>3</sup> The problem assumes a set of processes that alternate between executing a noncritical and a critical section of code. A process must eventually exit the critical section, but it may stay forever in the noncritical section. The basic requirement is that, at most, one process can be executing the critical section at any time. A solution to the mutual-exclusion problem lies at the heart of almost all multiprocess programming.

The bakery algorithm assumes processes are named by numbers from 1 through  $N$ . Figure 1 contains the code for process number  $i$ , almost exactly as it appeared in the original paper. The values of the variables *number* and *choosing* are arrays indexed by process number, with  $number[i]$  and  $choosing[i]$  initially equal to 0 for every process  $i$ . The relation  $\ll$  is lexicographical ordering on pairs of numbers, so  $(1, 3) \ll (2, 2) \ll (2, 4)$ ; it is an irreflexive total ordering on the set of all pairs of integers.


Mutual exclusion can be achieved very simply by not allowing any process to ever enter the critical section. A mutual-exclusion algorithm needs to also satisfy some progress condition. The condition Dijkstra's algorithm satisfies is *deadlock freedom*, meaning that if one or more processes try to enter the critical section, one of them must succeed. Most later algorithms satisfy the stronger requirement of *starvation freedom*, meaning that every process that tries to enter the critical section eventually does so. Before discussing mutual exclusion, we show that the bakery algorithm is starvation free. But first, some terminology.

We say that a process is *in the doorway* when it is executing statement  $M$ . After it finishes executing  $M$  until it exits its critical section, we say that it is *inside the bakery*. When it is at any other place in its code, we say that it is *outside the bakery*.

We first show that the algorithm is deadlock free. If it weren't, it would eventually reach a state in which every process is either forever in its noncritical section or forever inside the bakery. Eventually,  $choosing[i]$  would



## The bakery algorithm ensures processes execute a critical section of code one at a time.



equal 0 for all  $i$ , so every process inside the bakery would be waiting forever at statement  $L3$ . But this is impossible because the waiting process  $i$  with the smallest value of  $(number[i], i)$  would eventually enter the critical section. Hence, the algorithm is deadlock free.

To show that the algorithm is starvation free, it suffices to obtain a contradiction by assuming that a process  $i$  remains forever inside the bakery and outside the critical section. By deadlock freedom, other processes must continually enter and leave the critical section, since they cannot halt there.

However, once a process  $j$  is outside the bakery, to enter the bakery again it must execute statement  $M$  and set  $number[j]$  to be greater than  $number[i]$ . At that point, process  $j$  must remain forever inside the bakery because it will loop forever if it reaches  $L3$  with  $k = i$ . Eventually,  $number[i]$  will be less than  $number[j]$  for every process  $j$  in the bakery, so  $i$  will enter its critical section. This is the contradiction that proves starvation freedom.

Essentially, the same proof shows that the other mutual-exclusion algorithms we derive from the bakery algorithm also satisfy starvation freedom. So, we will say little more about starvation freedom. We now explain why the bakery algorithm satisfies mutual exclusion. For brevity, we abbreviate  $(number[i], i) \ll (number[j], j)$  as  $i \ll j$ .

Here is a naive proof that  $i$  and  $j$  cannot both be in their critical sections at the same time. For  $i$  to enter the critical section, it must find  $number[j] = 0$  or  $i \ll j$  when executing  $L3$  for  $k = j$ . Similarly, for  $j$  to enter the critical section, it must find  $number[i] = 0$  or  $j \ll i$  when executing  $L3$  for  $k = i$ . Since a process's number is non-zero when it executes  $L3$ , this means that for  $i$  and  $j$  both to be in their critical sections,  $i \ll j$  and  $j \ll i$  must be true, which is impossible.

This argument is flawed because it assumes that both  $i$  and  $j$  were inside the bakery when the other process executed  $L3$  for the appropriate value of  $k$ . Suppose process  $i$  read  $number[j]$  while  $j$  was in the doorway (executing  $M$ ) but had not yet set  $number[j]$ . It is possible for  $j$  to have read  $number[i] = 0$  in  $L3$  and entered the critical section, and for  $i$  then to have chosen  $number[i]$  to make  $i \ll j$  and entered the critical section.

The flaw in the argument is correct-

ed by statement  $L2$ . Since  $choosing[j]$  equals 1 when  $j$  is in the doorway, process  $i$  executed  $L3$  after  $L2$  found that  $j$  was not in the doorway; similarly,  $j$  executed  $L3$  after finding  $i$  not in the doorway. If, in both cases, the two processes were inside the bakery when  $L2$  was executed, then the naive argument is correct. If one of them, say  $j$ , was not inside the bakery, it must have been outside the bakery. Since  $i$  was then inside the bakery, with its current value of  $number[i]$ , process  $j$  must have chosen  $number[j]$  to be greater than the current value of  $number[i]$ , making  $i \ll j$  true. Hence,  $j$  could not have exited the  $L3$  loop for  $k = i$  and entered the critical section while  $i$  was still in the bakery. Therefore,  $i$  and  $j$  cannot both be in the critical section.

Observe that the  $choosing$  variable serves only to ensure that, when process  $i$  executes  $L3$  for  $k = j$ , there had been an instant when  $i$  was already inside the bakery and  $j$  was not in the doorway. This will be important later.

The most surprising property of the bakery algorithm is that it does not require reading or writing a memory register to be an atomic action. Carefully examining the proof of mutual exclusion shows that it just requires that  $number[i]$  and  $choosing[i]$  are what were later called safe registers,<sup>13</sup> ensuring only that a read not overlapping a write obtains the current register value. A read that does overlap a write can obtain any value the register might contain.

It is most convenient to describe a safe register in terms of atomic actions. We represent writing a value  $v$  to the register as two actions: the first sets its value to a special constant  $\zeta$  and the second sets it to  $v$ . We represent a read as a single atomic action that obtains the value of the register if that value does not equal  $\zeta$ . A read of  $number[i]$  when it equals  $\zeta$  can return any natural number, and a read of  $choosing[i]$  when it equals  $\zeta$  can return 0 or 1.

### Generalization of the Original Algorithm

Two generalizations of the bakery algorithm were obvious when it was published. The first is that, in statement  $M$ , it is not necessary to set  $number[i]$  to  $1 + maximum(\dots)$ . It could be set to any number greater than that maximum. It can also be set to the maximum if that

makes  $(number[j], j) \ll (number[i], i)$  for all  $j$ , but we will not bother with that generalization. We rewrite statement  $M$  using  $\succ$  to mean “is assigned a value greater than.”

The second obvious generalization is that statements  $L2$  and  $L3$  for different values of  $k$  do not have to be executed in the order specified by the **for** statement. Since the proof of mutual exclusion considers each pair of processes by themselves, the only requirement is that, for any value of  $k$ , statement  $L2$  must be executed before  $L3$ . For different values of  $k$ , those statements can be executed concurrently by different subprocesses. Also, there is no reason to execute them for  $k = i$  because their **if** tests always equal false.

These two generalizations have appeared elsewhere.<sup>5,10</sup> There is another, less obvious generalization that seems to be new: The assignment of 0 to  $number[i]$  after the process leaves

the critical section does not need to be completed before the process enters the noncritical section. In fact, that assignment need not even be completed if the process leaves the noncritical section to enter its critical section again. As long as that assignment is completed or aborted (leaving the register equal to  $\zeta$ ) before  $number[i]$  is assigned a new value in statement  $M$ , it just appears to other processes as if process  $i$  is still in the critical section or is executing the assignment statement immediately after the critical section. Therefore, mutual exclusion is still satisfied. To maintain starvation freedom, the write of 0 must eventually be completed if  $i$  remains forever in the noncritical section. There seems to be no simple way to describe in pseudo-code these requirements for setting  $number[i]$  to 0 upon completing the critical section. We simply add the mysterious keyword **asynchro-**

Figure 1. Process  $i$  of the original bakery algorithm.

```

begin integer k ;
L1: noncritical section ;
    choosing[i] := 1 ;
M: number[i] := 1 + maximum(number[1], ..., number[N]) ;
    choosing[i] := 0 ;
    for k = 1 step 1 until N do
        begin
            L2: if choosing[k] ≠ 0 then goto L2 ;
            L3: if number[k] ≠ 0 and (number[k], k) ≪ (number[i], i)
                then goto L3 ;
        end ;
    critical section ;
    number[i] := 0 ;
    goto L1
end
    
```

Figure 2. A generalization of the original bakery algorithm.

```

process i in {1, ..., N}
    variables number[i] = 0, choosing[i] = 0 ;
    while true do
        noncritical section ;
        choosing[i] := 1 ;
M: number[i] ≻ maximum(number[1], ..., number[N]) ;
        choosing[i] := 0 ;
        process j ≠ i in {1, ..., N}
            L2: await choosing[j] = 0 ;
            L3: await (number[j] = 0) ∨ ((number[i], i) ≪ (number[j], j))
        end process ;
        critical section ;
        asynchronously number[i] := 0 ; see explanation in text
    end while
end process
    
```

nously and refer to this discussion for its explanation.

The generalized algorithm is in Figure 2. Processes are explicitly declared; the outer **process** statement indicates that there are processes numbered from 1 through  $N$  and shows the code for process number  $i$ . Variables are declared with their initial values. The inner **process** statement declares that process  $i$  has  $N - 1$  subprocesses  $j$  with numbers from 1 through  $N$ , with none numbered  $i$ , and gives the code for subprocess  $j$ . That statement is executed by forking the subprocesses and continuing to the next statement (the critical section) when all subprocesses have terminated. Harmful or not, **gotos** have been eliminated. The outer loop is described as a **while** statement. The loops at  $L2$  and  $L3$  have been described with **await** statements, each of which repeatedly evaluates its predicate and terminates when it is true. The  $:\>$  in statement  $M$  and the **asynchronously** statement are explained above.

### The Deconstructed Bakery Algorithm

We have assumed that  $number[i]$  and  $choosing[i]$  are safe registers, written only by  $i$  and read by multiple readers. Such a register is easily implemented with safe registers having a single reader by keeping a copy of the register's value in a separate register for each reader.

We deconstruct the generalized bakery algorithm by implementing the safe registers  $choosing[i]$  and  $number[i]$  with single-reader registers  $localCh[j][i]$  and  $localNum[j][i]$ , for each  $j \neq i$ . Note the counterintuitive subscript order, with  $localCh[j][i]$  and  $localNum[j][i]$  containing the copies of  $choosing[i]$  and  $number[i]$  read by process  $j$ .

The pseudo-code of the deconstructed algorithm is in Figure 3. The reads of  $choosing[j]$  and  $number[j]$  by process  $i$  in the generalized algorithm are replaced by reads of  $localCh[j][i]$  and  $localNum[j][i]$ . The variable  $number[i]$  is now read only by process  $i$ , and we have eliminated  $choosing[i]$  because process  $i$  never reads it. Ad hoc notation is used in statement  $M$  to indicate that  $number[i]$  is set to be greater than the values of all  $localNum[j][i]$ .

Figure 3. The deconstructed bakery algorithm.

```

process  $i$  in  $\{1, \dots, N\}$ 
  variables  $number[i] = 0, localNum[*][i] = 0, localCh[*][i] = 0$  ;
  while true do
    noncritical section ;
    process  $j \neq i$  in  $\{1, \dots, N\}$ 
       $localCh[j][i] := 1$ 
    end process ;
     $M: number[i] := \mathbf{any} \ v > 0 \ \mathbf{with} \ \forall j \neq i : v > localNum[i][j]$  ;
     $localNum[*][i] := i$  ;
    process  $j \neq i$  in  $\{1, \dots, N\}$ 
       $localNum[j][i] := number[i]$  ;
       $localCh[j][i] := 0$  ;
       $L2: \mathbf{await} \ localCh[i][j] = 0$  ;
       $L3: \mathbf{await} \ (localNum[i][j] = 0) \vee ((number[i], i) \ll (localNum[i][j], j))$ 
    end process ;
    critical section ;
     $number[i] := 0$  ;
     $localNum[*][i] := i$  ;
    asynchronously process  $j \neq i$  in  $\{1, \dots, N\}$  see explanation in text
       $localNum[j][i] := 0$ 
    end process
  end while
end process
    
```

We have explicitly indicated the two atomic actions that represent writing a value  $v$  to the safe register  $localNum[j][i]$ , first setting its value to  $i$  and then to  $v$ . We have not bothered to do that for the writes to  $localCh[j][i]$ . The  $localCh[j][i]$  and  $localNum[j][i]$  writes are performed by subprocesses of process  $i$ , except that the  $N - 1$  separate writes of  $i$  to all the registers  $localNum[j][i]$  are represented by an assignment statement

$$localNum[*][i] := i$$

of the main process  $i$ . (This will be more convenient for our next version of the bakery algorithm.) To set  $number[i]$  to 0 after  $i$  exits the critical section, all the registers  $localNum[j][i]$  are set to  $i$  by the main process, and each is set to 0 by a separate process. We require that the setting of  $localNum[j][i]$  to 0 has been either completed or aborted when  $localNum[j][i]$  is set to  $number[i]$  by subprocess  $(i, j)$ . Again, this is not made explicit in the pseudo-code.

A proof of correctness for the deconstructed algorithm can be obtained by simple modifications to the proof for the original algorithm. For the original algorithm, we defined process  $i$  to be in the doorway while executing statement  $M$ , which ended with assigning the value of  $number[i]$ .

Since  $number[i]$  has been replaced by the registers  $localNum[j][i]$ , process  $i$  now has a separate doorway for each other process  $j$ . We say that  $i$  is in the doorway with respect to  $j$  from when it begins executing statement  $M$  until its subprocess  $j$  assigns  $number[i]$  to  $localNum[j][i]$ . We say that  $i$  is inside the bakery with respect to  $j$  from when it leaves the doorway with respect to  $j$  until it exits the critical section. The definition of  $i$  outside the bakery is the same as before.

To transform the proof of correctness of the original bakery algorithm to a proof of correctness of the deconstructed algorithm, we replace every statement that  $i$  or  $j$  is in the doorway or inside the bakery with the statement that it is there with respect to the other process. The modified proof shows that the function of statement  $L2$  is to ensure some time between  $i$  coming inside the bakery with respect to  $j$  and executing  $L3$  for  $j$ , process  $j$  was not in the doorway with respect to  $i$ .

### The Distributed Bakery Algorithm

We now implement the deconstructed bakery algorithm with a distributed algorithm. Each main process  $i$  is executed at a separate node, which we call node  $i$ , in a network of processes

that communicate by message passing. The variable  $localNum[j][i]$ , which is process  $j$ 's copy of  $number[i]$ , is kept at node  $j$ . It is set by process  $i$  to the value  $v$  by sending the message  $v$  to  $j$ . The setting of  $localNum[j][i]$  to  $i$  in the deconstructed bakery algorithm is implemented by the action of sending that message, and  $localNum[j][i]$  is set to  $v$  by process  $j$  when it receives the message. Thus, we are implementing the deconstructed algorithm by having process  $j$  obtain a previous value of  $localNum[j][i]$  on a read when  $localNum[j][i]$  equals  $i$ . Since the deconstructed algorithm allows such a read to obtain any value, this is a correct implementation.

For now, we assume that process  $i$  can write the value of  $localCh[j][i]$  atomically by a magical action at a distance. We will remove this magic later.

We assume that messages sent from a process  $i$  to any other process  $j$  are received in the order that they are sent. We represent the messages in transit from  $i$  to  $j$  by a first-in, first-out (FIFO) queue  $q[i][j]$ . We let  $\emptyset$  be the empty queue, and we define the following operations on a queue  $Q$ :

- ▶  $Append(Q, val)$  appends the element  $val$  to the end of  $Q$ .
- ▶  $Head(Q)$  is the value at the beginning of  $Q$ .
- ▶  $Behead(Q)$  removes the element at the beginning of  $Q$ .
- ▶  $Head(Q)$  and  $Behead(Q)$  are undefined if  $Q$  equals  $\emptyset$ .

The complete algorithm is in Figure 4. The shading highlights uses of  $localCh$ , whose magical properties need to be dealt with. Along with the main process  $i$ , there are concurrently executed processes  $(i, j)$  at node  $i$ , for each  $j \neq i$ . Process  $(i, j)$  receives and acts upon the messages sent to  $i$  by  $j$ .

The main process  $i$  of the distributed algorithm is obtained directly from the deconstructed algorithm by replacing the assignments of  $i$  to each  $localNum[j][i]$  with the sending of a message to  $j$ , except for two changes. The first is that statement  $M$  and the following sending of messages to other processes (represented by appending  $number[i]$  to all the message queues  $q[i][j]$ ) have been made a single atomic action. We can do this because we can view the end of each message queue  $q[i][j]$ , onto which messages are

appended, to be part of process  $i$ 's local state. A folk theorem<sup>4</sup> says that, for reasoning about a multiprocess algorithm, we can combine any number of actions that access only a process's local state into a single atomic action. That folk theorem has been formalized in a number of results starting with one by Lipton,<sup>15</sup> and perhaps the most directly applicable being Lamport.<sup>14</sup> In our algorithm, making this action appear atomic just requires preventing other processes at node  $i$  from acting on any incoming messages while the action is being executed.

The other significant change to the deconstructed algorithm is that the **asynchronously** statement has disappeared. The setting of  $localNum[j][i]$  is performed by the receipt of messages sent by  $i$  to  $j$ . FIFO message de-

livery ensures that it is set to 0 before its subsequent setting to a non-zero value. Also, since  $localNum[j][i]$  is now set by process  $(j, i)$  upon receipt of the message, the assignment to it in subprocess  $j$  of  $i$  has been removed.

Correctness of the deconstructed algorithm also depends on the assignment to  $localNum[j][i]$  being performed before process  $i$  sets  $localCh[j][i]$  to 0. Since the assignment to  $localNum[j][i]$  is now performed at node  $j$ , the ordering of those two operations is no longer trivially implied by the code. To maintain that ordering, subprocess  $j$  of  $i$  must learn that process  $(j, i)$  has set  $localNum[j][i]$  to  $number[i]$  before it can set  $localCh[j][i]$  to 0. This is done by having  $(j, i)$  send a message to  $i$  with some value  $ack$  that is not a natural number. Process  $(j, i)$  sets the value of

**Figure 4. The Distributed Bakery Algorithm, with magic.**

```

process  $i$  in  $\{1, \dots, N\}$ 
variables  $number[i] = 0, localNum[*][i] = 0, localCh[*][i] = 0,$ 
           $ackRcvd[i][*] = 0, q[*][i] = \phi;$ 
while true do
  noncritical section ;
  process  $j \neq i$  in  $\{1, \dots, N\}$ 
     $localCh[j][i] := 1$ 
  end process ;
  atomic  $M: number[i] := \mathbf{any} \ v > 0 \ \mathbf{with} \ \forall j \neq i : v > localNum[i][j];$ 
     $Append(q[i][*], number[i])$ 
  end atomic ;
  process  $j \neq i$  in  $\{1, \dots, N\}$ 
     $L0: \mathbf{await} \ ackRcvd[i][j] = 1;$ 
     $localCh[j][i] := 0;$ 
     $L2: \mathbf{await} \ localCh[i][j] = 0;$ 
     $L3: \mathbf{await} \ (localNum[i][j] = 0) \vee (number[i], i) \ll (localNum[i][j], j)$ 
  end process ;
  critical section ;
   $ackRcvd[i][*] := 0;$ 
   $number[i] := 0;$ 
   $Append(q[i][*], 0)$ 
end while
end process

process  $i, j \neq i$  in  $\{1, \dots, N\}$ 
while true do
  atomic  $\mathbf{await} \ q[j][i] \neq \phi;$ 
  if  $Head(q[j][i]) = ack$ 
  then  $ackRcvd[i][j] := 1$ 
  else  $localNum[i][j] := Head(q[j][i]);$ 
    if  $Head(q[j][i]) \neq 0$  then  $Append(q[j][i], ack)$ 
    end if
  end if ;
   $Behead(q[j][i])$ 
end atomic ;
end while
end process
    
```

$localNum[j][i]$  and sends the *ack* message to  $i$  as a single atomic action. When process  $(i, j)$  at node  $i$  receives the *ack* message, it sets  $ackRcvd[i][j]$  to 1 to notify subprocess  $j$  of process  $i$  that the *ack* has arrived. The setting of  $localNum[j][i]$  to  $number[i]$  in the deconstructed algorithm is replaced by statement  $L0$  that waits for  $ackRcvd[i][j]$  to equal 1.

The rest of the code for the main process  $i$  is the same as that of the corresponding process of the deconstructed algorithm, except that after  $i$  leaves the critical section, the asynchronous setting of all the registers  $localNum[j][i]$  to 0 is replaced by sending the message 0 to all the processes  $j$ , and  $ackRcvd[i][j]$  is reset to 0 for all  $j$ .

The asynchronously executed process  $(i, j)$  receives messages sent by  $j$  via  $q[j][i]$ . For an *ack* message, it sets  $ackRcvd[i][j]$  to 1; for a message with a value of  $number[j]$  it sets  $localNum[i][j]$  and, if the value is non-zero, sends an *ack* to  $j$ .

The one remaining problem is the magical atomic reading and writing of the register  $localCh[i][j]$ . The value of that register is used only in statement  $L2$ . The purpose of  $L2$  is to ensure that, before the execution of  $L3$ , there existed a time  $T$  when  $i$  was in the bakery with respect to  $j$  and  $j$  was not in the doorway with respect to  $i$ . We now show that statement  $L2$  is unnecessary, because executing  $L0$  ensures the existence of such a time  $T$ .

The execution of statement  $M$  by  $j$  and the sending of  $number[j]$  in a message to  $i$  are part of a single atomic action, and  $j$  enters the bakery with respect to  $i$  when that message is received at node  $i$ . Therefore,  $j$  is in the doorway with respect to  $i$  exactly when there is a message with a non-zero integer in  $q[j][i]$ . Let's call that message a *doorway* message. Process  $i$  enters the bakery with respect to  $j$  when its message containing  $number[i]$  is received at node  $j$ , an action that appends to  $q[j][i]$  the *ack* that  $L0$  is waiting to arrive. If there is no doorway message in  $q[j][i]$  at that time, then immediately after execution of that action is the time  $T$  whose existence we need to show, since it occurred before the receipt of the *ack* that  $L0$  was waiting for. If there is a doorway message in  $q[j][i]$ , then the required time  $T$  is right after that message was received at node  $i$ . Because of FIFO

message delivery, that time was also before the receipt of the *ack* that  $L0$  is waiting for. In both cases, executing  $L0$  ensures there was some time  $T$  after  $i$  entered inside the bakery with respect to  $j$  when  $j$  was not in the doorway with respect to  $i$ . Hence, statement  $L2$  is redundant.

Because  $L2$  is the only place where the value of  $localCh[i][j]$  is read, we can eliminate  $localCh$  and all statements that set it. Removing all the grayed statements in Figure 4 gives us the distributed bakery algorithm, with no magic.

The first paper devoted to distributed mutual exclusion was apparently that of Ricart and Agrawala.<sup>19</sup> Their algorithm can be viewed as an optimization and simplification of our algorithm. It delays the sending of *ack* messages in such a way that a process can enter its critical section when it receives an *ack* from every other process, so it does not have to keep track of other processes' numbers. The number 0 messages sent upon exiting the critical section can therefore be eliminated, yielding an algorithm with fewer messages. Although nicer than our algorithm, the Ricart-Agrawala algorithm is not directly on the path we are traveling.

### A Distributed State Machine

In a distributed state machine,<sup>12</sup> there is a set of processes at separate nodes in a network, each wanting to execute state-machine commands. The processes must agree on the order in which all the commands are executed. To execute a command, a process must know the entire sequence of preceding commands.

A distributed mutual-exclusion algorithm can be used to implement a distributed state machine by having a process execute a single command in the critical section. The order in which processes enter the critical section determines the ordering of the commands. It is easy to devise a protocol that has a process in its critical section send its current command to all other processes, which order it after all preceding commands. Starting with this idea and the distributed bakery algorithm, we will obtain the distributed state-machine algorithm<sup>12</sup> by eliminating the critical section.

The bakery algorithm is based on

the idea that if two processes are trying to enter the critical section at about the same time, then the process  $i$  with the smaller value of  $(number[i], i)$  enters first. We now make that true no matter when the two processes enter the critical section. Define a version of the bakery algorithm to be *number-ordered* if it satisfies this condition: If process  $i$  enters the critical section with  $number[i] = n_i$  and process  $j$  later enters the critical section with  $number[j] = n_j$ , then  $(n_i, i) \ll (n_j, j)$ . We now make the distributed bakery number-ordered. We can do that because we have generalized the bakery algorithm to set  $number[i]$  to any number greater than the maximum value of the values of  $number[j]$  it reads, not just to the next-largest number.

We add to the distributed bakery algorithm a variable  $maxNum$ , where  $maxNum[i][j]$  is the largest value  $localNum[i][j]$  has equaled, for  $j \neq i$ . We let  $maxNum[i][i]$  be the largest value  $number[i]$  has equaled. We then make two changes to the algorithm. First, we replace statement  $M$  with the statement in Figure 5.

Second, in process  $(i, j)$ , if  $localNum[i][j]$  is assigned a non-zero value, then  $maxNum[i][j]$  is assigned that same value. The FIFO ordering of messages assures the new value of  $maxNum[i][j]$  will be greater than its previous value. Clearly,  $localNum[i][j]$  always equals  $maxNum[i][j]$  or 0. The value of  $number[i]$  chosen this way is therefore allowed by statement  $M$  of the distributed algorithm, so this is a correct implementation of that algorithm. We now show that it is number-ordered.

Suppose  $i$  enters the critical section with  $number[i] = n_i$  and  $j$  later enters the critical section with  $number[j] = n_j$ . It's evident that  $(n_i, i) \ll (n_j, j)$  if  $i = j$ , so we can assume  $i \neq j$ . The proof of mutual exclusion for the deconstructed algorithm shows that either (i)  $(n_i, i) \ll (n_j, j)$  or (ii)  $j$  chose  $n_j$  after reading a value of  $localNum[i][j]$  written after  $i$  set it to  $n_i$ . In our modified version of the distributed algorithm,  $j$  reads  $maxNum[j][i]$  not  $localNum[i][j]$  to set  $number[j]$ , and  $maxNum[j][i]$  never decreases. Therefore,  $(n_i, i) \ll (n_j, j)$  is true also in case (ii), so the algorithm is number-ordered.

Since the algorithm is number-ordered, we don't need the critical

**Figure 5. A new version of statement M.**

```

M: number[i] := maximum(maxNum[i][1], ..., maxNum[i][N]);
    maxNum[i][i] := number[i]
    
```

**Figure 6. A newer version of statement M.**

```

atomic
M: maxNum[i][i] := maximum(maxNum[i][1], ..., maxNum[i][N]);
    Append((maxNum[i][i], Cmd), q[i][*])
end atomic
    
```

section to implement a distributed state machine. We can order the commands by the value ( $number[i], i$ ) would have had when  $i$  entered the critical section to execute the command. Process  $i$  can send the command it is executing in the messages containing the value of  $number[i]$  that it sends to other processes. In fact, we don't need  $number[i]$  at all. When we send that message,  $number[i]$  has the same value as  $maxNum[i][i]$ . We can eliminate everything in the main process  $i$  except the atomic statement containing statement  $M$ , which can now be written as in Figure 6, where  $Cmd$  is process  $i$ 's current command.

There is one remaining problem. Process  $i$  saves the messages containing commands that it sends and receives, accumulating a set of triples  $(v, j, Cmd)$  indicating that process  $j$  issued a command  $Cmd$  with  $number[j]$  having the value  $v$ . It knows that those commands are ordered by  $(v, j)$ . However, to execute the command in  $(v, j, Cmd)$ , it has to know that it has received all commands  $(w, k, Dcmd)$  with  $(w, k) \ll (v, j)$ . Process  $i$  knows that, for each process  $k$ , it has received all commands  $(w, k, Dcmd)$  with  $w \leq maxNum[i][k]$ . However, suppose  $i$  has received no commands from  $k$ . How can  $i$  be sure that  $k$  hasn't sent a command in a message that  $i$  hasn't yet received? The answer is to use the distributed bakery algorithm's *ack* messages. Here's how.

For convenience, we let process  $i$  keep  $maxNum[i][i]$  always equal to the maximum of the values  $maxNum[i][j]$  (including  $j = i$ ). It does this by increasing  $maxNum[i][i]$ , if necessary, when receiving a message with the value of  $maxNum[i][j]$  from another process  $j$ . Upon receiving a message  $(v, Cmd)$  from process  $j$ , process  $i$  sets

$maxNum[i][j]$  to  $v$  (possibly increasing  $maxNum[i][i]$ ) and sends back to  $j$  the message  $(maxNum[i][i], ack)$ . When that message is received,  $j$  sets  $maxNum[j][i]$  accordingly, (increasing  $maxNum[j][j]$  if necessary). When  $i$  has received all the *ack* messages for a command it issued with  $maxNum[i][i]$  equal to  $v$ , all its values of  $maxNum[i][j]$  will be  $\geq v$ , so process  $i$  knows it has received all commands ordered before its current command. It can therefore execute all of them, in the appropriate order, and then execute its current command.

This is almost identical to the distributed state-machine algorithm,<sup>12</sup> where  $maxNum[i][i]$  is called process  $i$ 's clock. (The sketch of the algorithm given there is not detailed enough to mention the other registers  $maxNum[i][j]$ .) The one difference is that, when process  $i$  receives a message from  $j$  with a new value  $v$  of  $maxNum[i][j]$ , the algorithm requires  $maxNum[i][i]$  to be set to a value  $> v$ , whereas  $\geq v$  suffices. The algorithm remains correct if the value of  $maxNum[i][i]$  increases by any amount at any time. Thus, the registers  $maxNum[i][i]$  could be logical clocks that are also used for other purposes.

We have described all the pieces of a distributed state-machine algorithm but have not put them together into pseudo-code. "The precise algorithm is straightforward, and we will not bother to describe it."<sup>12</sup>

### Ancient and Recent History

In addition to being the author of this article, I am the author of the starting and ending algorithms of our journey. The bakery algorithm is among hundreds of algorithms that implement mutual exclusion using only read and

write operations to shared memory.<sup>22</sup> A number of them improve the bakery algorithm, the most significant improvement being a bound on the chosen numbers.<sup>6,21</sup> But all improvements seem to add impediments to our path, except for one: Moses and Patin<sup>17</sup> optimized the bakery algorithm by allowing process  $i$  to stop waiting for process  $j$  at statement  $L3$  if it reads two different values of  $number[j]$ . However, it is irrelevant to our path because it optimizes a case that cannot occur in the distributed bakery algorithm.

Mutual-exclusion algorithms based on read and write operations have been of no practical use for decades, since modern computers provide special instructions to implement mutual exclusion more efficiently. Now, they are studied mainly as concurrent programming exercises. The bakery algorithm is of interest because it was the first mutual-exclusion algorithm not to assume lower-level mutual exclusion, which is implied by atomic reads and writes of shared memory. The distributed state-machine algorithm is interesting because it preserves causality. But it too is less important than the problem it solves.

The most important contribution of my state-machine paper was the observation that any desired form of cooperation in a network of computers can be obtained by implementing a distributed state machine. The obvious next step was to make the implementation fault tolerant. The work addressing that problem is too extensive to discuss here. Fault-tolerant state-machine algorithms have become the standard building block for implementing reliable distributed systems.<sup>20</sup>

There was no direct connection between the creation of the bakery algorithm and of the state-machine algorithm. The bakery algorithm was inspired by a bakery in the neighborhood where I grew up. A machine dispensed numbers to its customers that determined the order in which they were served. The state-machine algorithm was inspired by an algorithm of Paul Johnson and Robert Thomas.<sup>7</sup> They used the  $\ll$  relation and process identifiers to break ties, but I don't know if that was inspired by the bakery algorithm.



The path between the two algorithms that we followed is not the one I originally took. That journey began when I was looking for an example of a distributed algorithm for notes I was writing. Stephan Merz suggested the mutual-exclusion algorithm I had used to illustrate the state-machine algorithm. I found it to be too complicated, so I simplified it. (I did not remember the Ricart-Agrawala algorithm and was only later reminded of it by a referee). After stripping away things that were not needed for that particular state machine, I arrived at the distributed bakery algorithm. It was obviously related to the original bakery algorithm, but it was still not clear exactly how.

I wanted to make the distributed algorithm an implementation of the bakery algorithm. I started with the generalization of having subprocesses of each process interact independently with the other processes; that was essentially how I had been describing the bakery algorithm for years. Delaying the setting of  $number[i]$  to 0 was required because the distributed algorithm's message that accomplished it could be arbitrarily delayed. It took me a while to realize that I should deconstruct the multi-reader register  $number[i]$  into multiple single-reader registers, and that both the original bakery algorithm and the distributed algorithm implemented that deconstructed algorithm.

The path back from the distributed bakery algorithm to the distributed state-machine algorithm was easy. It may have helped that I had previously used the idea of modifying the bakery algorithm to make values of  $number[i]$  keep increasing. Paradoxically, that was done to keep those values from getting too large.<sup>10</sup>

Correctness of a concurrent algorithm is expressed with two classes of properties: *safety* properties, such as mutual exclusion, that assert what the algorithm may do, and *liveness* properties, such as starvation freedom, that assert what the algorithm must do.<sup>1</sup> Safety properties depend on the actions the algorithm can perform; liveness properties depend as well on assumptions, often implicit, about what actions the algorithm must perform.

The kind of informal behavioral

reasoning I have used here is notoriously unreliable. I believe the best rigorous proofs of safety properties are usually based on invariants—predicates that are true of every state of every possible execution.<sup>2</sup> Invariance proofs that the bakery algorithm satisfies mutual exclusion have often been used to illustrate formalisms or tools.<sup>5,11</sup> An informal sketch of such a proof for the decomposed bakery algorithm is in an expanded version of this article, which is available on the Web.<sup>8</sup> Elegant rigorous proofs of progress properties can be written using temporal logic.<sup>18</sup>

Rigorous proofs are longer than informal ones and can intimidate readers not used to them. I almost never write one until I believe that what I want to prove is true. For the correctness of our algorithms, that belief was based on the reasoning embodied in the informal proofs I presented—the same kind of reasoning I used when I discovered the bakery and distributed state-machine algorithms.

I understood the two algorithms well enough to be confident in the correctness of the non-distributed versions of the bakery algorithm and of the derivation of the state-machine algorithm from the distributed bakery algorithm. Model checking convinced me of the correctness of the distributed bakery algorithm and confirmed the confidence my informal invariance proof had given me that the deconstructed algorithm satisfies mutual exclusion.

More recently, Stephan Merz wrote a formal, machine-checked version of my informal invariance proof. He also wrote a machine-checked proof that the actions of the distributed bakery algorithm implement the actions of the deconstructed bakery algorithm under a suitable data refinement. These two proofs show that the deconstructed algorithm satisfies mutual exclusion. The proofs are available on the Web.<sup>16</sup> □

**References**

1. Alpern, B. and Schneider, F. Defining liveness. *Information Processing Letters* 21, 4 (Oct. 1985), 181–185.
2. Ashcroft, E. Proving assertions about parallel programs. *Journal of Computer and System Sciences* 10, 1 (Feb. 1975), 110–135.
3. Dijkstra, E. Solution of a problem in concurrent programming control. *Commun. ACM* 8, 9 (Sept. 1965), 569.
4. Harel, D. On folk theorems. *Commun. ACM* 23, 7 (July 1980), 379–389.
5. Hesselink, W. Mechanical verification of Lamport's

- bakery algorithm. *Science of Computer Programming* 78, 9 (2013), 1622–1638.
6. Jayanti, P., Tan, K., Friedland, G., and Katz, A. Bounding Lamport's bakery algorithm. In L. Pacholski and P. Ruzicka, eds., *SOFSEM 2001: 28<sup>th</sup> Conference on Current Trends in Theory and Practice of Informatics 2234, Lecture Notes in Computer Science*, Springer (2001), 261–270.
7. Johnson, P. and Thomas, R. The maintenance of duplicate data bases. Request for Comment RFC #677, NIC #31507, ARPANET Network Working Group, (January 1975).
8. Lamport, L. Online supplemental material for Deconstructing the bakery to build a distributed state machine. <http://lamport.azurewebsites.net/pubs/bakery/deconstruction.html>.
9. Lamport, L. A new solution of Dijkstra's concurrent programming problem. *Commun. ACM* 17, 8 (Aug. 1974), 453–455.
10. Lamport, L. Concurrent reading and writing. *Commun. ACM* 20, 11 (Nov. 1977), 806–811.
11. Lamport, L. Proving the correctness of multiprocess programs. *IEEE Transactions on Software Engineering SE-3*, 2 (Mar. 1977), 125–143.
12. Lamport, L. Time, clocks, and the ordering of events in a distributed system. *Commun. ACM* 21, 7 (July 1978), 558–565.
13. Lamport, L. On interprocess communication. *Distributed Computing* 1 (1986), 77–101.
14. Lamport, L. A theorem on atomicity in distributed algorithms. *Distributed Computing* 4, 2 (1990), 59–68.
15. Lipton, R. Reduction: A method of proving properties of parallel programs. *Commun. ACM* 18, 12 (Dec. 1975), 717–721.
16. Merz, S. Online TLA+ specifications and proofs for "deconstructing the bakery to build a distributed state machine." <https://members.loria.fr/SMerz/papers/distributed-bakery.html>.
17. Moses, Y. and Patkin, K. Mutual exclusion as a matter of priority. *Theoretical Computer Science* 751 (2018), 46–60.
18. Pnueli, A. The temporal logic of programs. In *Proceedings of the 18<sup>th</sup> Annual Symposium on the Foundations of Computer Science*, IEEE (Nov. 1977), 46–57.
19. Ricart, G. and Agrawala, A. An optimal algorithm for mutual exclusion in computer networks. *Commun. ACM* 24, 1 (1981), 9–17.
20. Schneider, F. Implementing fault-tolerant services using the state machine approach: A tutorial. *ACM Computing Surveys* 22, 4 (December 1990), 299–319.
21. Taubenfeld, G. The black-white bakery algorithm and related bounded-space, adaptive, local-spinning and FIFO algorithms. In R. Guerraoui, (Ed.), *Proceedings of Distributed Computing, 18<sup>th</sup> Intern. Conf. 3274, Lecture Notes in Computer Science*, Springer (Oct. 4, 2004), 56–70.
22. Taubenfeld, G. Concurrent programming, mutual exclusion. In M-Y Kao, (Ed.) *Encyclopedia of Algorithms—2016 Edition*, 421–425. Springer, 2016.

Leslie Lamport is the recipient of the 2013 ACM A.M. Turing Award.



This work is licensed under a <http://creativecommons.org/licenses/by/4.0/>



Watch the author discuss this work in the exclusive *Communications* video. <https://cacm.acm.org/videos/deconstructing-the-bakery>