

# 6.852: Distributed Algorithms

## Fall, 2015

### Lecture 16

# Lamport's Bakery Algorithm



# Lamport's Bakery Algorithm

- Like taking tickets for service in a bakery.
- Nice features:
  - Uses only single-writer, multi-reader registers.
  - Extends to even weaker (“safe”) registers, in which operations have durations, and a read that overlaps a write receives an arbitrary response.
  - Guarantees lockout-freedom, in fact, almost-FIFO behavior.
- But:
  - Registers are of unbounded size.
    - Algorithm can be simulated using bounded registers, but not easily (uses “Bounded Concurrent Timestamps”).
- Shared variables:
  - For each process  $i$ :
    - $\text{choosing}(i)$ , a Boolean, written by  $i$ , read by all, initially 0
    - $\text{number}(i)$ , a natural number, written by  $i$ , read by all, initially 0

# Bakery Algorithm

- First part, up to  $\text{choosing}(i) := 0$  (the “Doorway”, D):
  - Process  $i$  chooses a number  $>$  all the numbers it reads for the other processes; writes this in  $\text{number}(i)$ .
  - While doing this, keeps  $\text{choosing}(i) = 1$ .
  - Two processes could choose the same number (unlike in a real bakery).
  - Break ties with process ids.
- Second part:
  - Wait to see that no others are choosing, and no one else has a smaller number.
  - That is, wait to see that your ticket is the smallest.
  - Never go back to the beginning of this part---just proceed step by step, waiting when necessary.

# Code

## Shared variables:

for every  $i \in \{1, \dots, n\}$ :

$\text{choosing}(i) \in \{0, 1\}$ , initially 0, writable by  $i$ , readable by all  $j \neq i$

$\text{number}(i)$ , a natural number, initially 0, writable by  $i$ , readable by  $j \neq i$ .

$\text{try}_i$

$\text{choosing}(i) := 1$

$\text{number}(i) := 1 + \max_{j \neq i} \text{number}(j)$

$\text{choosing}(i) := 0$

for  $j \neq i$  do

    waitfor  $\text{choosing}(j) = 0$

    waitfor  $\text{number}(j) = 0$  or  $(\text{number}(i), i) < (\text{number}(j), j)$

$\text{crit}_i$

$\text{exit}_i$

$\text{number}(i) := 0$

$\text{rem}_i$

# Correctness: Mutual exclusion

- **Key invariant:** If process  $i$  is in  $C$ , and process  $j \neq i$  is in  $(T - D) \cup C$ ,

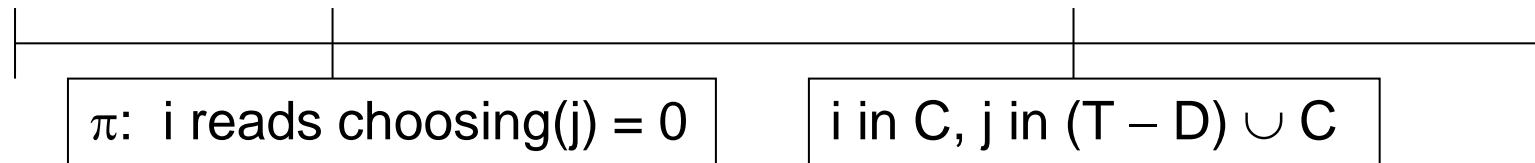
Trying region after doorway, or critical region

then  $(\text{number}(i), i) < (\text{number}(j), j)$ .

- **Proof:**
  - Could prove by induction.
  - Instead, give argument based on events in executions.

# Correctness: Mutual exclusion

- **Invariant:** If  $i$  is in  $C$ , and  $j \neq i$  is in  $(T - D) \cup C$ , then  $(\text{number}(i), i) < (\text{number}(j), j)$ .
- **Proof:**
  - Consider a point where  $i$  is in  $C$  and  $j \neq i$  is in  $(T - D) \cup C$ .
  - Then before  $i$  entered  $C$ , it must have read  $\text{choosing}(j) = 0$ , event  $\pi$ .



- **Case 1:**  $j$  sets  $\text{choosing}(j) := 1$  (starts choosing) after  $\pi$ .
  - Then  $\text{number}(i)$  is set before  $j$  starts choosing.
  - So  $j$  sees the “correct”  $\text{number}(i)$  and chooses something bigger.
  - That suffices.
- **Case 2:**  $j$  sets  $\text{choosing}(j) := 0$  (finishes choosing) before  $\pi$ .
  - Then when  $i$  reads  $\text{number}(j)$  in its second waitfor loop, it gets the “correct”  $\text{number}(j)$ .
  - Since  $i$  decides to enter  $C$ , it must see  $(\text{number}(i), i) < (\text{number}(j), j)$ .

# Correctness: Mutual exclusion

- **Invariant:** If  $i$  is in  $C$ , and  $j \neq i$  is in  $(T - D) \cup C$ , then  $(\text{number}(i), i) < (\text{number}(j), j)$ .
- **Proof of mutual exclusion:**
  - Apply invariant both ways.
  - Contradictory requirements.



# Liveness Conditions

- **Progress:**
  - By contradiction.
  - If not, eventually region changes stop, leaving everyone in T or R, and at least one process in T.
  - Everyone in T eventually finishes choosing.
  - Then nothing blocks the smallest (number, index) process from entering C.
- **Lockout-freedom:**
  - Consider any  $i$  that enters T.
  - Suppose for contradiction that  $i$  never reaches C.
  - Eventually it finishes the doorway.
  - Thereafter, any newly-entering process picks a bigger number.
  - Progress implies that some processes continue to enter C, as long as  $i$  is still in T.
  - In fact, this must happen infinitely many times!
  - But those with bigger numbers can't get past  $i$ , contradiction.

# FIFO Condition

- Not really FIFO ( $\rightarrow T$  vs.  $\rightarrow C$ ), but almost:
  - **FIFO after the doorway**: if  $j$  leaves  $D$  before  $i \rightarrow T$ , then  $j \rightarrow C$  before  $i \rightarrow C$ .
- But the “doorway” is an artifact of this algorithm, so this isn’t a meaningful way to evaluate the algorithm!
- Maybe say “there exists a doorway such that”...
- But then we could take  $D$  to be the entire trying region, making the property trivial.
- To make the property nontrivial:
  - Require  $D$  to be “wait-free”: a process is guaranteed to complete  $D$  if it keeps taking steps, regardless of what other processes do.
  - $D$  in the Bakery Algorithm is wait-free.
- The algorithm is **FIFO after a wait-free doorway**.

# Impact of Bakery Algorithm

- Originated some important ideas:
  - Wait-freedom
    - Fundamental notion for theory of fault-tolerant asynchronous distributed algorithms.
  - Weakly coherent memories
    - Beginning of formal study: definitions, and some algorithmic strategies for coping with them.



# Rest of today

- Mutual exclusion with read/write memory:
  - Burns's algorithm
  - Lower bound on the number of read/write registers
- Mutual exclusion with read-modify-write operations
- Generalized resource allocation problems
- **Reading:** Sections 10.6-10.8, 10.9; Chapter 11
- **Next:**
  - Asynchronous shared-memory systems with failures.
  - Consensus in asynchronous shared-memory systems.
  - Impossibility of consensus [Fischer, Lynch, Paterson]
  - Reading: Chapter 12

# Burns's Algorithm

# Space and memory considerations

- All the mutual exclusion algorithms we have seen (Dijkstra, Peterson, Lamport Bakery) use **more than  $n$  shared variables**.
  - Bakery algorithm could use just  $n$  variables. (Why?)
- All but the Bakery algorithm use **multi-writer variables**.
  - These can be expensive to implement in terms of lower-level primitives like 1W1R variables.
- Bakery uses **infinite-size variables**.
  - Difficult (but possible) to adapt it to use finite-size variables [Dolev, Shavit].
- Q: Can we do better?

# Burns's algorithm

- Uses just  $n$  single-writer Boolean read/write variables.
- Simple.
- Guarantees:
  - Mutual exclusion,
  - Progress,
  - But not lockout-freedom.

# Code

## Shared variables:

for every  $i \in \{1, \dots, n\}$ :

$\text{flag}(i) \in \{0, 1\}$ , initially 0, writable by  $i$ , readable by all  $j \neq i$

## Process $i$ :

$\text{try}_i$

L:  $\text{flag}(i) := 0$

for  $j \in \{1, \dots, i-1\}$  do

if  $\text{flag}(j) = 1$  then go to L

$\text{flag}(i) := 1$

for  $j \in \{1, \dots, i-1\}$  do

if  $\text{flag}(j) = 1$  then go to L

M: for  $j \in \{i+1, \dots, n\}$  do

if  $\text{flag}(j) = 1$  then go to M

$\text{crit}_i$

$\text{exit}_i$

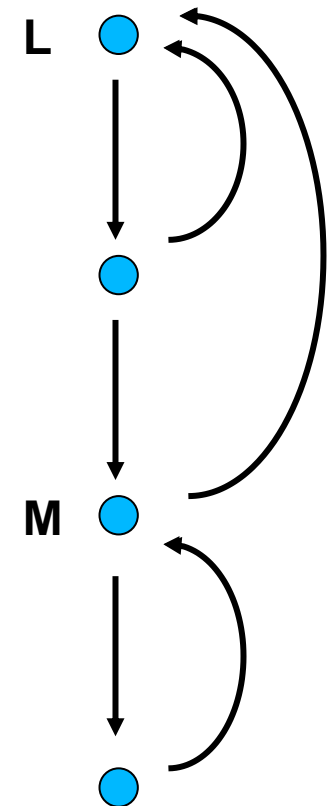
$\text{flag}(i) := 0$

$\text{rem}_i$



# That is,...

- Each process goes through 3 loops, sequentially:
  1. Check flags of processes with **smaller** indices.
  2. Check flags of processes with **smaller** indices.
  3. Check flags of processes with **larger** indices.
- If it passes all tests,  $\rightarrow$  C.
- Otherwise, drops back:

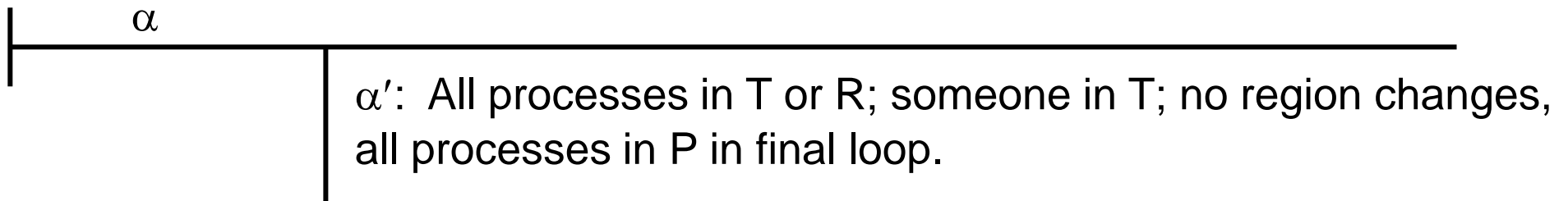


# Correctness of Burns's algorithm

- Mutual exclusion + progress
- **Mutual exclusion:**
  - Like the proof for Dijkstra's algorithm, but now with flags set to 1 rather than 2.
  - If processes  $i$  and  $j$  are ever in  $C$  simultaneously, both must have set their flags  $:= 1$ .
  - Assume WLOG that process  $i$  sets  $\text{flag}(i) := 1$  (for the last time) first.
  - Keeps  $\text{flag}(i) = 1$  until process  $i$  leaves  $C$ .
  - After  $\text{flag}(i) := 1$ , must have  $\text{flag}(j) := 1$ , then  $j$  must see  $\text{flag}(i) = 0$ , before  $j \rightarrow C$ .
  - Impossible!

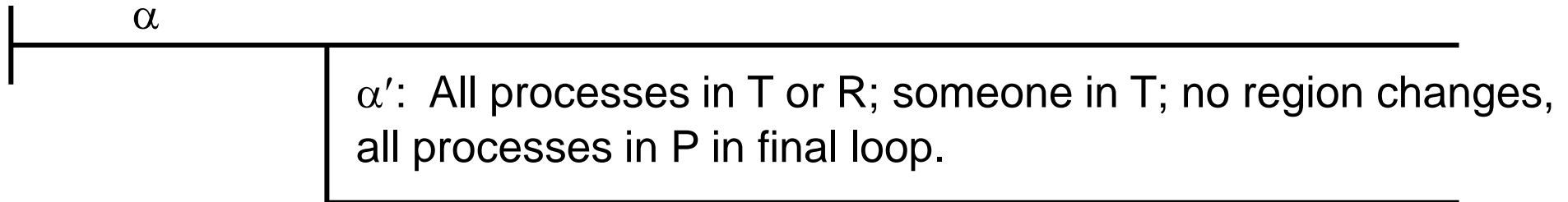
# Progress for Burns's algorithm

- Consider fair execution  $\alpha$  (each process keeps taking steps).
- Assume for contradiction that, after some point in  $\alpha$ , some process is in T, no one is in C, and no one  $\rightarrow$  C later.
- WLOG, we can assume that every process is in T or R, and no region changes occur after that point in  $\alpha$ .
- Call the processes in T the **contenders**.
- Divide the contenders into two sets:
  - P, the contenders that ever reach M, and
  - Q, the contenders that never reach M.
- After some point in  $\alpha$ , all contenders in P have reached M; they never drop back thereafter to before M.



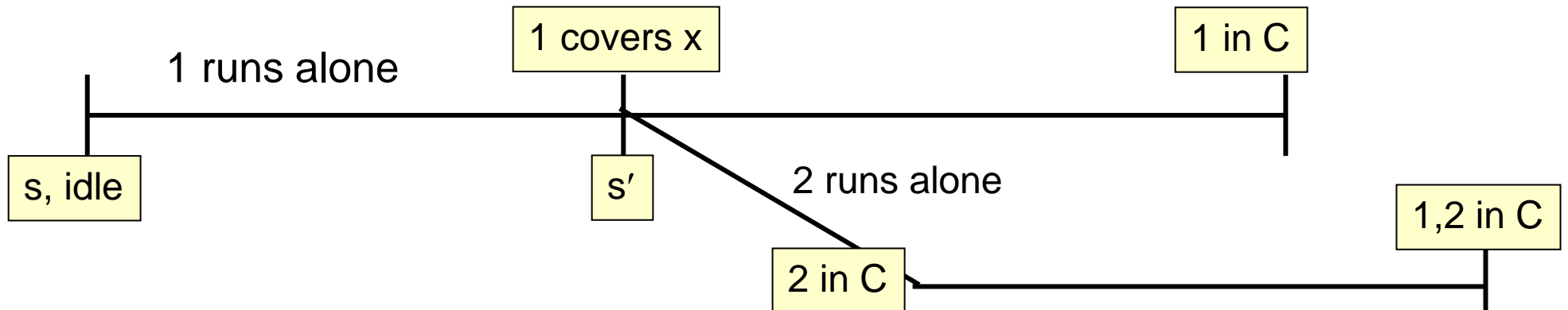
# Progress for Burns's algorithm

- P, the contenders that reach label M, and
- Q, the contenders that never reach M.



- Claim P contains at least one process:
  - Process with the lowest index among all the contenders is not blocked from reaching M.
- Let  $i$  = largest index of a process in P.
- Claim process  $i$  eventually  $\rightarrow C$ : All others with larger indices eventually see a smaller-index contender and drop back to L, setting their flags  $:= 0$  (and these stay  $= 0$ ).
- So  $i$  eventually sees all these  $= 0$  and  $\rightarrow C$ .
- Contradiction.

# A Lower Bound on the Number of Registers for Mutual Exclusion



# Lower Bound on the Number of Registers

- All the mutual exclusion algorithms we've studied:
  - Use read/write shared memory, and
  - Use at least  $n$  read/write shared variables.
- That's one variable per **potential** contender.
- **Q:** Can we use fewer than  $n$  read/write shared variables?
- Not single-writer. (Why?)
- Not even multi-writer!

# Lower bound on number of registers

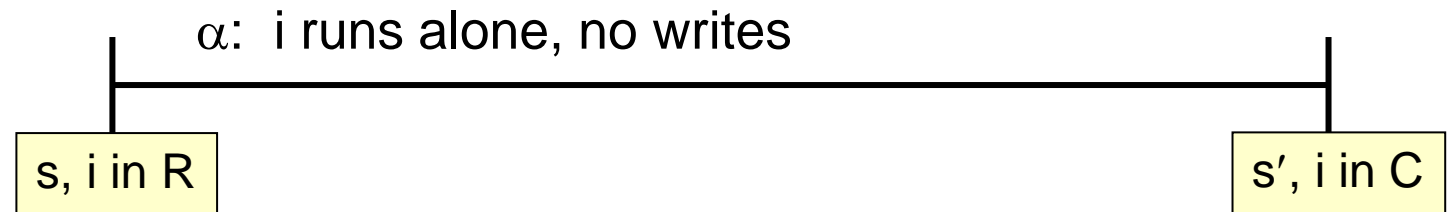
- Lower bound of  $n$  holds even if:
  - We require **only mutual exclusion + progress** (no stronger liveness properties).
  - The variables can be **any size**.
  - All variables can be **read and written by all processes**.
- Start with basic facts about any mutex algorithm  $A$  that uses only read/write shared variables.
- **Lemma 1:** If  $s$  is a reachable, idle system state (meaning all processes are in  $R$ ), and if process  $i$  runs alone from  $s$ , then eventually  $i \rightarrow C$ .
- **Proof:** By the progress requirement.
- **Corollary:** If  $i$  runs alone from a system state  $s'$  that is **indistinguishable** from  $s$  by  $i$ ,  $s' \sim^i s$ , then eventually  $i \rightarrow C$ .
- **Indistinguishable:** Same state of  $i$  and same shared variable values.

# Lower bound on registers

- **Lemma 2:** Suppose that  $s$  is a reachable system state in which process  $i \in R$ . Suppose  $i \rightarrow C$  on its own, from  $s$ . Then along the way, process  $i$  writes to some shared variable.

- **Proof:**

- By contradiction; suppose it doesn't.
- Then:



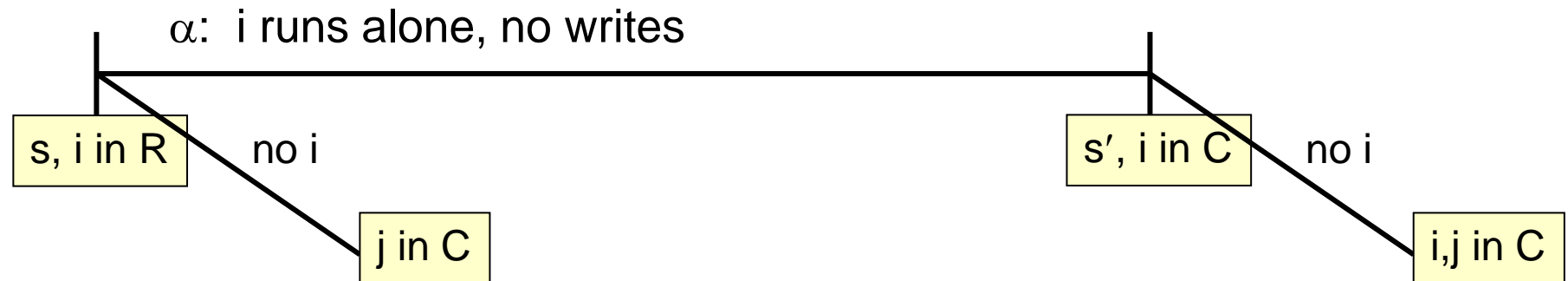
- Then  $s' \sim^j s$  for every  $j \neq i$ .
- There is some execution fragment from  $s$  in which process  $i$  takes no steps, and in which some other process  $j \rightarrow C$ .
  - By the progress requirement.





# Lower bound on registers

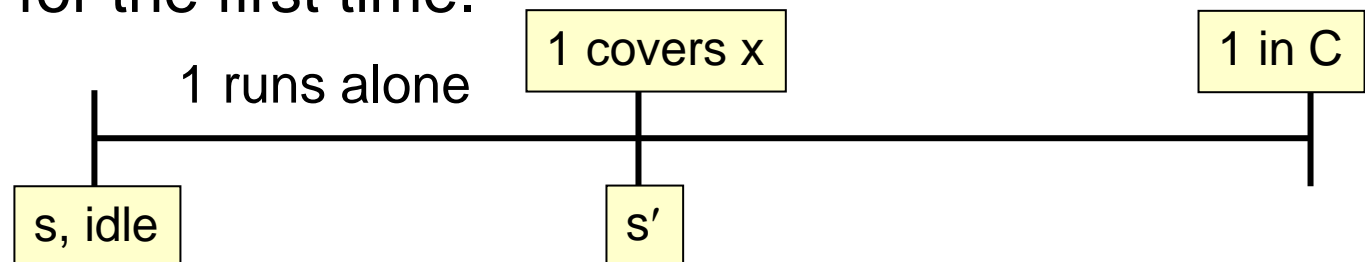
- **Lemma 2:** Suppose that  $s$  is a reachable system state in which  $i \in R$ . Suppose process  $i \rightarrow C$  on its own, from  $s$ . Then along the way, process  $i$  writes to some shared variable.
- **Proof, cont'd:**
  - There is some execution fragment from  $s$  in which process  $i$  takes no steps, and in which some other process  $j \rightarrow C$ .



- Then there is also such a fragment from  $s'$ .
- Yields a counterexample execution:
  - System gets to  $s$ , then  $i$  alone takes it to  $s'$ , then others get  $j \in C$ .
  - Contradiction because  $i, j$  are in  $C$  at the same time.

# Lower bound on registers

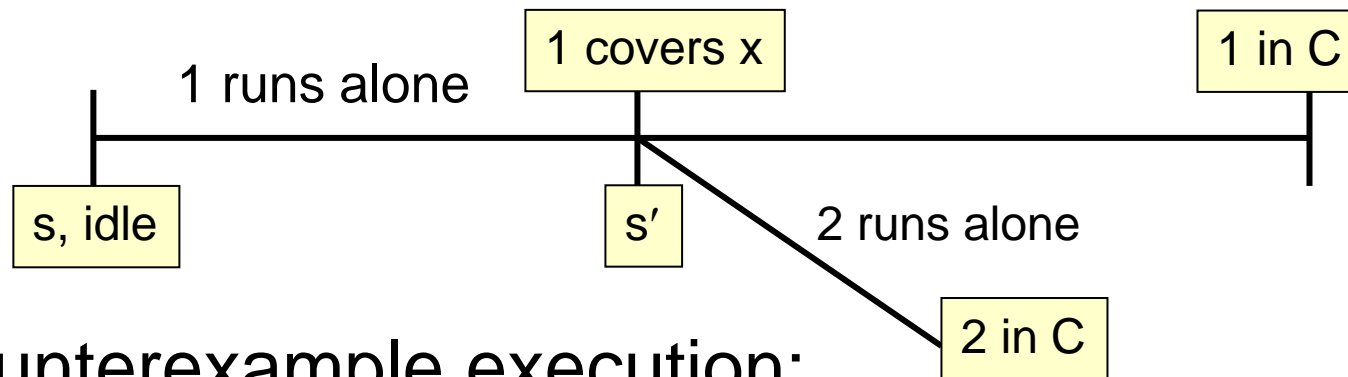
- Back to showing  $\geq n$  shared variables needed...
- Special case: **2 processes and 1 variable**:
  - Suppose A is a 2-processes mutex algorithm using 1 read/write shared variable x.
  - Start in initial (idle) state s.
  - Run process 1 alone,  $\rightarrow C$ , writes x on the way.
    - By Lemmas 1 and 2.
  - Consider the point where process 1 is just about to write x, i.e., it **covers** x, for the first time.



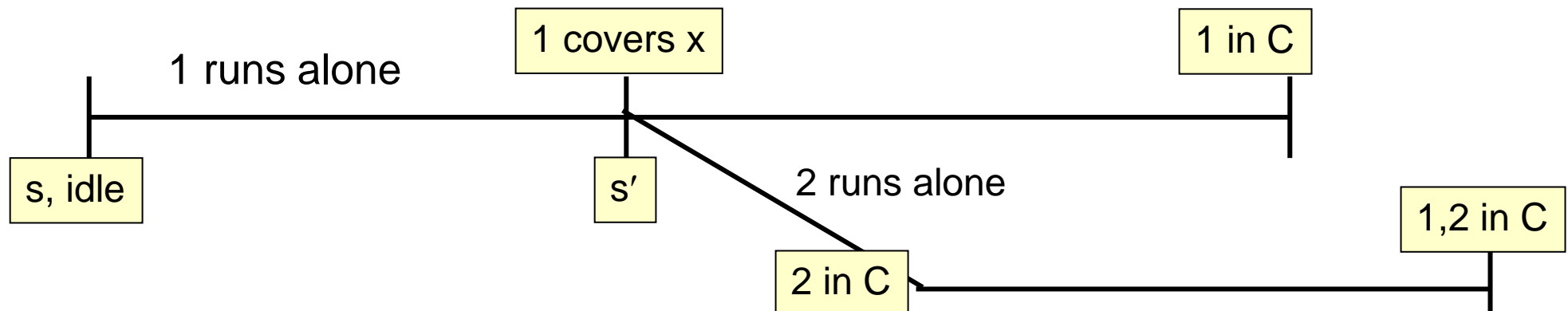
- Note that  $s' \sim^2 s$ , because 1 doesn't write between s and s'.
- So process 2 can reach C on its own from s'.
  - By Corollary to Lemma 1.

# 2 processes, 1 variable

- Process 2 can reach C on its own from  $s'$ :



- Counterexample execution:
  - Run 1 until it covers  $x$ , then let 2 reach  $C$ .
  - Then resume 1, letting it write  $x$  and then  $\rightarrow C$ .
  - When it writes  $x$ , it overwrites anything 2 might have written there on its way to  $C$ ; so 1 never sees any evidence of 2.

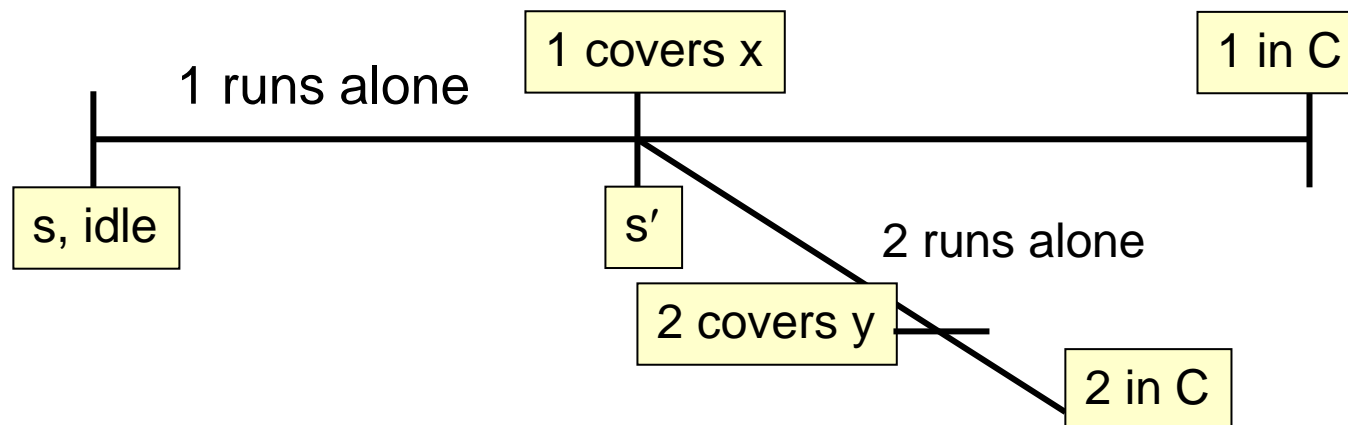


# Another special case: 3 processes, 2 variables

- Processes 1, 2, 3; variables x,y.
- Similar construction, with a couple of twists.
- Start in initial (idle) state s.
- Run processes 1 and 2 until:
  - Each covers one of x,y---both variables covered.
  - Resulting state is indistinguishable by 3 from a reachable idle state.
- **Q:** How to do this?
- For now, assume we can.
- Then run 3 alone,  $\rightarrow C$ .
- Then let 1 and 2 take one step each, overwriting both variables, and obliterating all traces of 3.
- Continue running 1 and 2; they run as if 3 were still in R.
- By the progress requirement, one eventually  $\rightarrow C$ .
- Contradicts mutual exclusion.

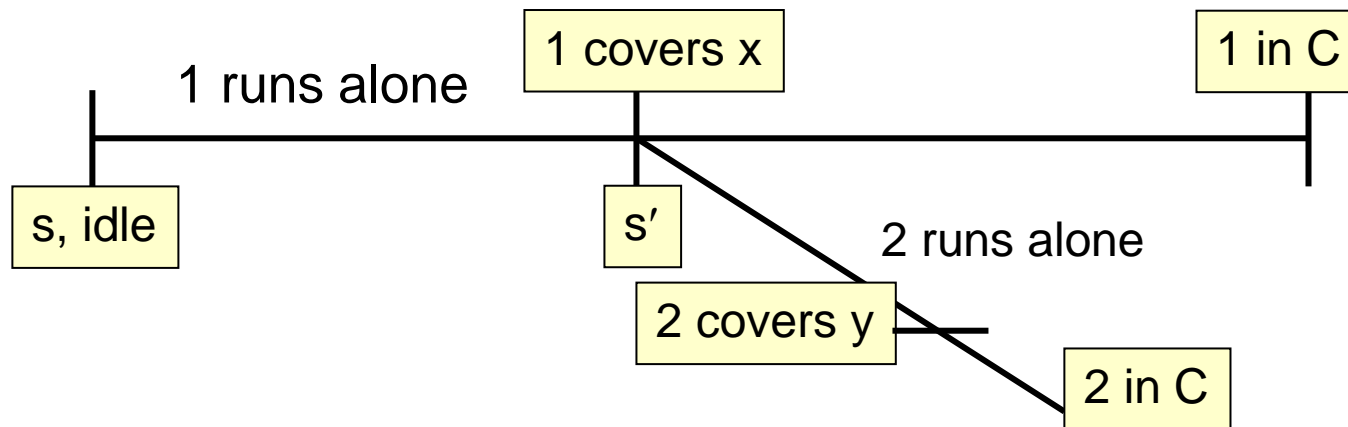
# 3 processes, 2 variables

- It remains to show how to maneuver 1 and 2 so that:
  - Each covers one of x,y.
  - Resulting state is indistinguishable by 3 from a reachable idle state.
- **First try:**
  - Run 1 alone until it first covers a shared variable, say x.
  - Then run 2 alone until it  $\rightarrow C$ .
  - **Claim:** Along the way, it must write the **other** shared variable y.
    - If not, then after  $2 \rightarrow C$ , 1 could take one step, overwriting whatever 2 wrote to x, and thus obliterating all traces of 2.
    - Then 1 continues  $\rightarrow C$ , violating mutual exclusion.
  - Stop 2 just when it first covers y; then 1 and 2 cover x and y.



# 3 processes, 2 variables

- Recall the goal: Maneuver 1 and 2 so that:
  - Each covers one of  $x, y$ .
  - Resulting state is indistinguishable by 3 from a reachable idle state.



- But this is not quite right...** resulting state might not be indistinguishable by 3 from an idle state.
- 2 could have written  $x$  before writing  $y$ .

# 3 processes, 2 variables

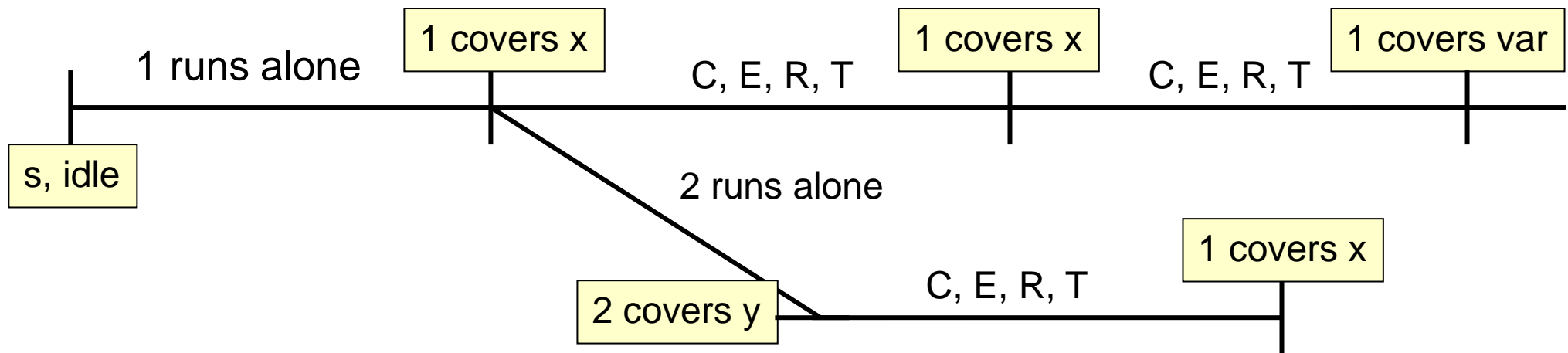
- Maneuver 1 and 2 so that:
  - Each covers one of x,y.
  - Resulting state is indistinguishable by 3 from a reachable idle state.
- **Second (successful) try:**
  - Run 1 alone until it first covers a shared variable.
  - Continue running 1, through C, E, R, back in T, until it again first covers a variable.
  - And once again.



- In two of the three covering states, 1 must cover the same variable.
- E.g., suppose in first two states, 1 covers x (other cases analogous).

# 3 processes, 2 variables

- Counterexample execution:
  - Run 1 until it covers x the first time.
  - Then run 2 until it first covers y (must do so).



- Then let 1 write x and continue until it covers x again.
- Now both variables are (again) covered.
- This time, the final state is indistinguishable by 3 from an idle state.
- As needed.



# General case:

## $n$ processes, $n-1$ variables

- Extends 3-process 2-variable case, using induction.
- Need strengthened version of Lemma 2:
- **Lemma 2'**: Suppose that  $s$  is a reachable system state in which  $i \in R$ . Suppose process  $i \rightarrow C$  on its own, from  $s$ . Then along the way, process  $i$  writes to some shared variable **that is not covered (in  $s$ ) by any other process.**
- **Proof:**
  - Similar to Lemma 2.
  - Contradictory execution fragment begins by overwriting all the covered variables, obliterating any evidence of  $i$ .

# n processes, n-1 variables

- **Definition:**  $s'$  is **k-reachable** from  $s$  if there is an execution fragment from  $s$  to  $s'$  involving only steps by processes 1 to  $k$ .

# n processes, n-1 variables

- Now suppose (for contradiction) that A solves mutual exclusion for n processes, with n-1 shared variables.
- **Main Lemma:** For any  $k \in \{1, \dots, n-1\}$  and from any idle state, there is a k-reachable state in which processes  $1, \dots, k$  cover k distinct shared variables, and that is indistinguishable by processes  $k+1, \dots, n$  from some k-reachable idle state.
- **Proof:** In a minute...
- First assume we have this, for  $k = n-1$ .
- Then run n alone,  $\rightarrow C$ .
  - Can do this, by Corollary to Lemma 1.
- Along the way, it must write some variable that isn't covered by  $1, \dots, n-1$ .
  - By Lemma 2'.
- But all n-1 variables are covered, contradiction.
- It remains to prove the Main Lemma...

# Proof of the Main Lemma

- **Main Lemma:** For any  $k \in \{1, \dots, n-1\}$  and from any idle state, there is a  $k$ -reachable state in which processes 1 to  $k$  cover  $k$  distinct shared variables, and that is indistinguishable by processes  $k+1$  to  $n$  from some  $k$ -reachable idle state.
- **Proof:** Induction on  $k$ .
  - **Base case ( $k=1$ ):**
    - Run process 1 alone until just before it first writes a shared variable.
    - 1-reachable state, process 1 covers a shared variable, indistinguishable by the other processes from initial state.
  - **Inductive step (Assume for  $k \leq n-2$ , show for  $k+1$ ):**
    - By inductive hypothesis, get a  $k$ -reachable state  $t_1$  in which processes 1, ...,  $k$  cover  $k$  variables, and that is indistinguishable by processes  $k+1, \dots, n$  from some  $k$ -reachable idle state.

# Proof of the Main Lemma

- **Main Lemma:** For any  $k \in \{1, \dots, n-1\}$  and from any idle state, there is a  $k$ -reachable state in which processes 1 to  $k$  cover  $k$  distinct shared variables, and that is indistinguishable by processes  $k+1$  to  $n$  from some  $k$ -reachable idle state.
- **Proof:** Inductive step (Assume for  $k \leq n-2$ , show for  $k+1$ ):
  - By I.H., get a  $k$ -reachable state  $t_1$  in which  $1, \dots, k$  cover  $k$  variables, and that is indistinguishable by  $k+1, \dots, n$  from some  $k$ -reachable idle state.
  - Let each of  $1, \dots, k$  take one step, overwriting covered variables.
  - Run  $1, \dots, k$  until all are back in  $R$ ; resulting state is idle.
  - By I.H. get another  $k$ -reachable state  $t_2$  in which  $1, \dots, k$  cover  $k$  variables, and that is indistinguishable by  $k+1, \dots, n$  from some  $k$ -reachable idle state.
  - Repeat, getting  $t_3, t_4, \dots$ , until we get  $t_i$  and  $t_j$  ( $i < j$ ) that cover the same set  $X$  of variables. (Why is this guaranteed to happen?)
  - Run  $k+1$  alone from  $t_i$  until it first covers a variable not in  $X$ .
  - Then run  $1, \dots, k$  as if from  $t_i$  to  $t_j$  (they can't tell the difference).
  - Now processes  $1, \dots, k+1$  cover  $k+1$  different variables.
  - And the result is indistinguishable by  $k+2, \dots, n$  from an idle state.

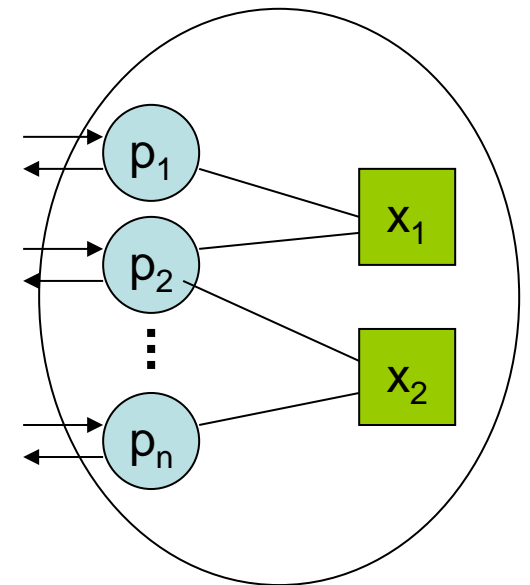
# Discussion

- Bell Labs research failure:
  - At Bell Labs (many years ago), Gadi Taubenfeld learned that the Unix group was trying to develop an asynchronous mutual exclusion algorithm for many processes that used only a few read/write shared registers.
  - He told them it was impossible.

# Discussion

- Newer research direction:
  - Develop “space-adaptive” algorithms that **potentially** use many variables, but are **guaranteed** to use only a few if only a few processes are contending.
  - Also “time-adaptive” algorithms.
  - See work by [Moir, Anderson], [Attiya, Friedman]
  - Time-adaptive and space-adaptive algorithms often yield better performance, lower overhead, in practice.

# Mutual Exclusion with Read-Modify-Write Shared Variables





# Mutual exclusion with RMW shared variables

- Stronger memory operations (synchronization primitives):
  - Test-and-set, fetch-and-increment, swap, compare-and-swap, load-linked/store-conditional,...
  - All are special types of read-modify-write operations.
- All modern computer architectures include shared memory that supports some of these operations, in addition to read and write operations.
- Also include read/write registers, which support just read and write operations.
- More powerful variables are more expensive (cost of hardware, time to access) than read/write registers.
- Not all are equally powerful; we'll come back to this later.
- Q: Do such stronger memory primitives enable better algorithms, e.g., for mutual exclusion?

# Mutual exclusion with RMW: Test-and-set algorithm

- **test-and-set** operation: Sets value to 1, returns previous value.
  - Usually for binary variables.
- **Test-and-set mutual exclusion algorithm (trivial):**
  - One shared binary variable  $x$ , 0 when no one has been granted the resource (initial state), 1 when someone has.
  - **Trying protocol:** Repeatedly test-and-set  $x$  until get 0.
  - **Exit protocol:**  $x := 0$ .

try<sub>i</sub>  
  waitfor(test-and-set( $x$ ) = 0)  
crit<sub>i</sub>

exit<sub>i</sub>  
   $x := 0$   
rem<sub>i</sub>

- Guarantees mutual exclusion + progress.
- No fairness. To get fairness, we can use a more expensive queue-based algorithm:

# Mutual exclusion with RMW: Queue-based algorithm

- **queue** shared variable
  - Supports enqueue, dequeue, head operations.
  - Can be quite large!
- **Queue mutual exclusion algorithm:**
  - One shared variable **Q**: FIFO queue.
  - **Trying protocol**: Add self to **Q**, wait until you're at the head.
  - **Exit protocol**: Remove self from **Q**.

try<sub>i</sub>

enqueue(**Q**,i)

waitfor(head(**Q**) = i)

crit<sub>i</sub>

exit<sub>i</sub>

dequeue(**Q**)

rem<sub>i</sub>

- **Fairness**: Guarantees bounded bypass (indeed, no bypass = 1-bounded bypass).

# Mutual exclusion with RMW: Ticket-based algorithm

- **Modular fetch-and-increment operation,  $f\&i_n$** 
  - Variable values are integers mod  $n$ .
  - Increments variable mod  $n$ , returns the previous value.
- **Ticket mutual exclusion algorithm:**
  - Like Bakery algorithm: Take a number, wait till it's your turn.
  - Guarantees bounded bypass (no bypass).
  - Shared variables: **next**, **granted**: integers mod  $n$ , initially 0
    - Support modular fetch-and-increment.
  - **Trying protocol**: Increment **next**, wait till your ticket is granted.
  - **Exit protocol**: Increment **granted**.

```
tryi  
  ticket :=  $f\&i_n(\text{next})$   
  waitfor(granted = ticket)  
criti
```

```
exiti  
   $f\&i_n(\text{granted})$   
remi
```

# Ticket-based algorithm

- **Space complexity:**

- Each shared variable takes on at most  $n$  values.
- Total number of variable values:  $n^2$
- Total size of variables in bits:  $2 \log n$

- **Compare with queue:**

- Total number of variable values:  
$$n! + (n \text{ choose } (n-1)) (n-1)! + (n \text{ ch } (n-2)) (n-2)! + \dots + (n \text{ ch } 1) 1!$$
$$= n! (1 + 1/1! + 1/2! + 1/3! + \dots + 1/(n-1)!)$$
$$\leq n! e = O(n^n)$$
- Size of variable in bits:  $O(n \log n)$

```
tryi  
  ticket := f&in(next)  
  waitfor(granted = ticket)  
criti
```

```
exiti  
  f&in(granted)  
remi
```

# Variable Size for Mutual Exclusion with RMW

- **Q:** How small could we make the RMW variable?
- 1 bit, for just mutual exclusion + progress (simple test and set algorithm).
- With fairness guarantees?
- $O(n)$  values ( $O(\log n)$  bits) for bounded bypass.
  - Can get  $n+k$  values, for small  $k$ .

In practice, on a real shared-memory multiprocessor,  
we want a few variables of size  $O(\log n)$ .  
So ticket algorithm is pretty good (in terms of space).

- **Theoretical lower bounds:**
  - $\Omega(n)$  values needed for bounded bypass,  $\Omega(\sqrt{n})$  for lockout-freedom.

# Variable Size for Mutual Exclusion with RMW

- Theoretical lower bound:
  - $\Omega(n)$  values needed for bounded bypass,  $\Omega(\sqrt{n})$  for lockout-freedom.
- **Significance:**
  - Achieving **mutual exclusion + lockout freedom** is not **trivial**, even though we assume that the processes get fair access to the shared variables.
  - Thus, fair access to the shared variables does not immediately translate into fair access to higher-level critical sections.
- For example, consider bounded bypass:

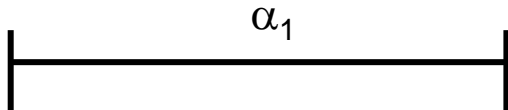
# Lower bound for mutual exclusion + bounded bypass

- **Theorem:** In any 1-variable mutual exclusion algorithm guaranteeing progress and bounded bypass, using a single RMW shared variable, the variable must be able to take on at least  $n$  distinct values.
- Essentially, we need enough space to keep a process index, or a counter of the number of active processes, in shared memory.
- General RMW shared variable: Allows read, arbitrary computation, and write, all in one step.
- **Proof:** By contradiction.
  - Suppose Algorithm A achieves mutual exclusion + progress +  $k$ -bounded bypass, using one RMW variable with  $< n$  values.
  - Construct a bad execution, which violates  $k$ -bounded bypass:

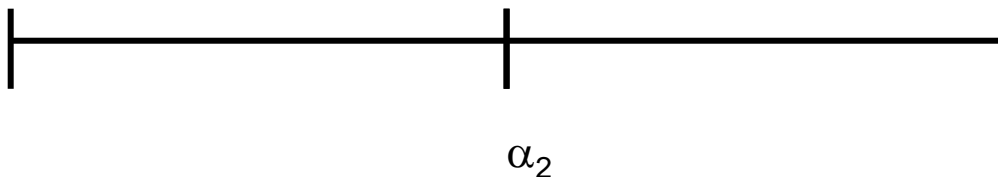


# Lower bound for mutual exclusion + bounded bypass

- **Theorem:** In any 1-variable mutual exclusion algorithm guaranteeing progress and bounded bypass, using a single RMW shared variable, the variable must be able to take on at least  $n$  distinct values.
- **Proof:** By contradiction.
  - Suppose Algorithm A achieves mutual exclusion + progress +  $k$ -bounded bypass, using one RMW variable with  $< n$  values.
  - Run process 1 from initial state, until  $\rightarrow C$ , execution  $\alpha_1$ :



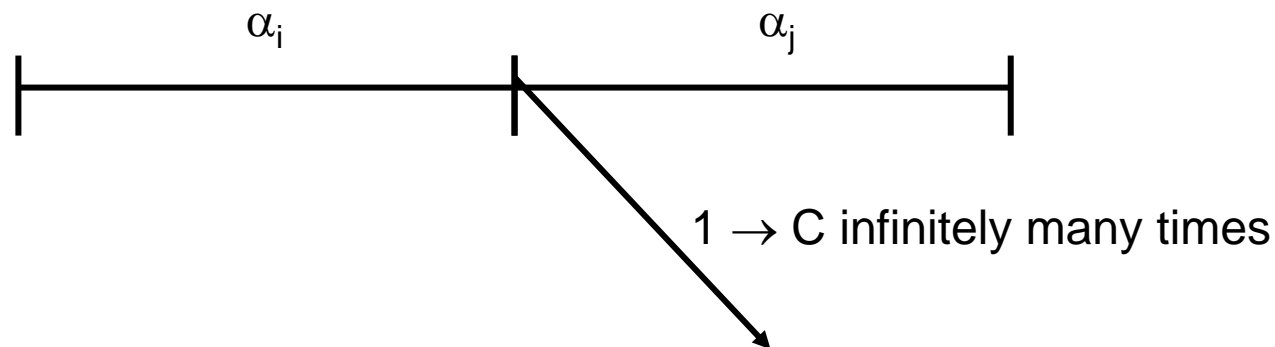
- Run process 2 until it accesses the variable,  $\alpha_2$ :



- Continue by running each of 3, 4, ...,  $n$ , obtaining  $\alpha_3, \alpha_4, \dots, \alpha_n$ .

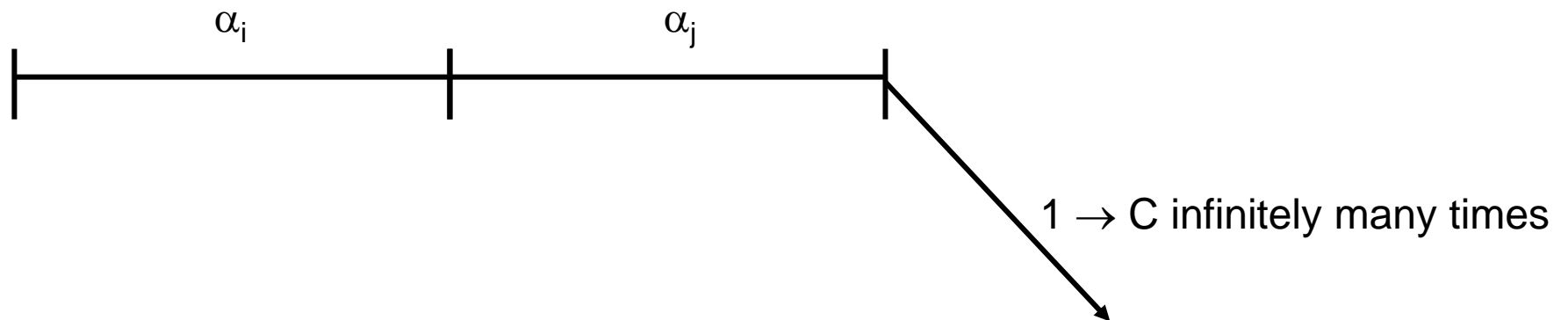
# Lower bound for mutual exclusion + bounded bypass

- **Theorem:** In any 1-variable mutual exclusion algorithm guaranteeing bounded bypass, using a single RMW shared variable, the variable must be able to take on at least  $n$  distinct values.
- **Proof, cont'd:**
  - Since the variable takes on  $< n$  values, there must be two processes,  $i$  and  $j$ ,  $i < j$ , for which  $\alpha_i$  and  $\alpha_j$  leave the variable with the same value  $v$ .
  - Now extend  $\alpha_i$  so that  $1, \dots, i$  exit, then  $1$  reenters repeatedly,  $\rightarrow C$  infinitely many times.
    - Possible since progress is required in a fair execution.



# Lower bound for mutual exclusion + bounded bypass

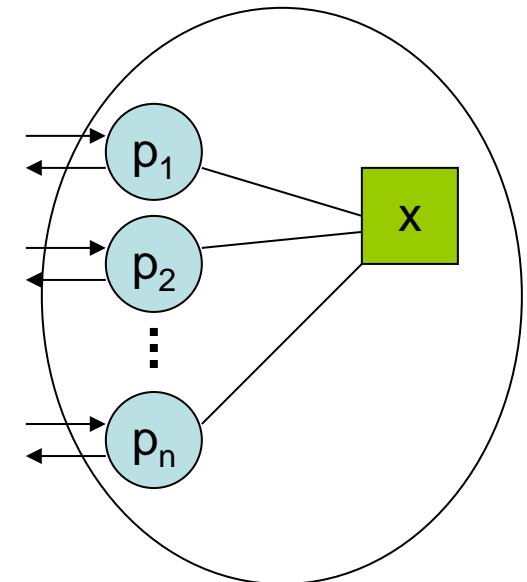
- **Theorem:** In any 1-variable mutual exclusion algorithm guaranteeing bounded bypass, using a single RMW shared variable, the variable must be able to take on at least  $n$  distinct values.
- **Proof, cont'd:**
  - Now apply the same steps after  $\alpha_j$ .
  - Result is an execution in which process 1  $\rightarrow C$  infinitely many times, while process  $j$  remains in  $T$ .
  - Violates bounded bypass.



- Note: The extension of  $\alpha_j$  isn't a fair execution; this is OK since fairness isn't required for a violation of bounded bypass.

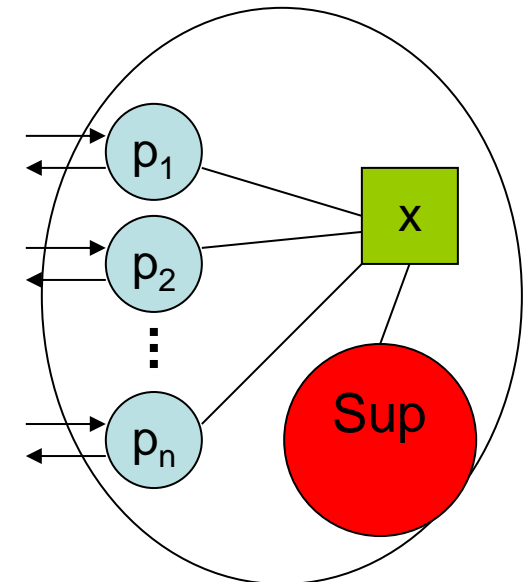
# $O(n)$ algorithm for mutual exclusion + bounded bypass

- It's possible to achieve fair mutual exclusion using a single shared variable that can hold (essentially) only one process UID, plus a constant amount of other information.
- **Algorithm idea:**
  - Divide Trying region into two parts:
    - An unordered Buffer
    - An unordered Main region
  - Processes in Main region go to Critical one at a time, in any order, until Main is empty.
  - When Main is empty, all processes in Buffer go to Main, one at a time, in any order, until Buffer is empty.
  - Prevents unbounded bypass.
- But, how can we manage all this in a distributed way?



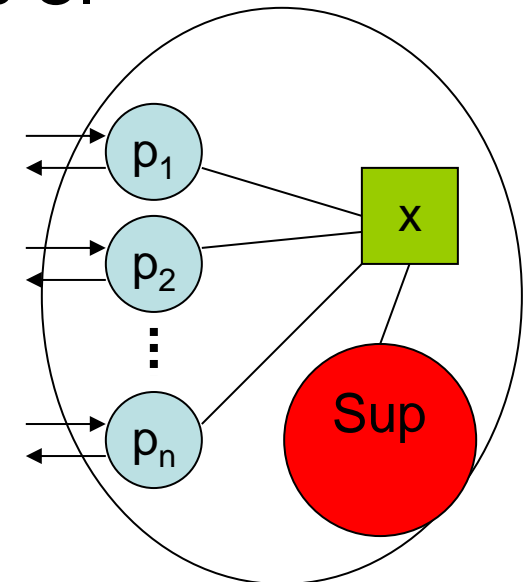
# $O(n)$ algorithm for mutual exclusion + bounded bypass

- First, pretend we have a dedicated Supervisor process (we will have to simulate this with ordinary processes).
- Divide Trying region into Buffer, Main regions
- Supervisor keeps a local count of the processes in Buffer, Main.
- Shared variable has **count** component, which a newly-entering process increments.
- Sup adds this **count** to its local Buffer-count, while resetting the shared **count** to 0.
- Sup has enough info to decide when someone should go from Buffer to Main, or from Main to Critical.
- Tells processes by putting special tokens in the shared variable, which an arbitrary qualified process can pick up.



# $O(n)$ algorithm for mutual exclusion + bounded bypass

- Divide Trying region into Buffer, Main regions
- Pretend we have a dedicated Supervisor process, which manages all the process movement, using tokens in the shared variable for communication.
- Eliminating the Supervisor: Let the process in the Critical region emulate it.
- Passes the Supervisor role to the next entrant to C.
- Does so using an exit protocol, using small messages in the shared variable for communication.
- If no new entrants are around, just leave the (small) Supervisor state in the shared variable.



# Mutual Exclusion + Lockout-freedom

- Can solve with  $O(n)$  values.
  - Actually, can achieve  $n/2 + k$ , small constant  $k$ .
- Lower bound of  $\Omega(\sqrt{n})$  values.
  - Actually, about  $\sqrt{n}$ .
  - Uses a more complicated version of the construction for the bounded bypass lower bound.

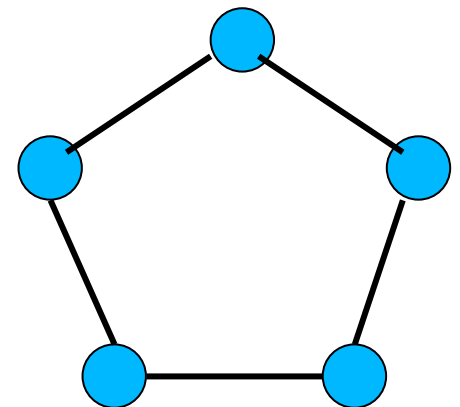
# Generalized Resource Allocation

- A very quick tour
- Textbook, Chapter 11



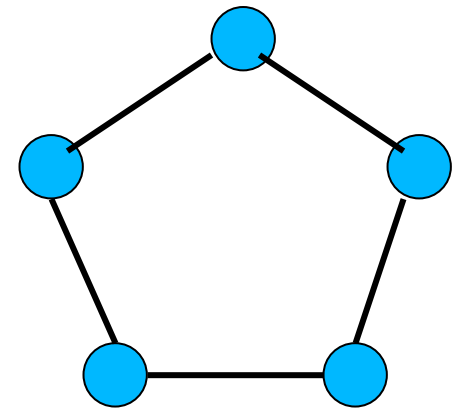
# Generalized resource allocation

- Mutual exclusion: Problem of allocating a single non-sharable resource.
- Can generalize to more resources, some sharing.
- Exclusion specification  $\mathcal{E}$  (for a given set of users):
  - Any collection of sets of users, closed under superset.
  - Expresses which users are incompatible, can't coexist in the critical section.
- **Example: k-exclusion** (any  $k$  users are OK, but not  $k+1$ )
  - $\mathcal{E} = \{ E : |E| > k \}$
- **Example: Reader-writer locks**
  - Relies on classification of users as readers vs. writers.
  - $\mathcal{E} = \{ E : |E| > 1 \text{ and } E \text{ contains a writer} \}$
- **Example: Dining Philosophers (Dijkstra)**
  - $\mathcal{E} = \{ E : E \text{ includes a pair of neighbors} \}$



# Resource specifications

- Some exclusion specs can be described conveniently in terms of requirements for concrete resources.
- Resource specification: Each user needs a certain subset of the resources.
- Can't share: Users with intersecting sets exclude each other.



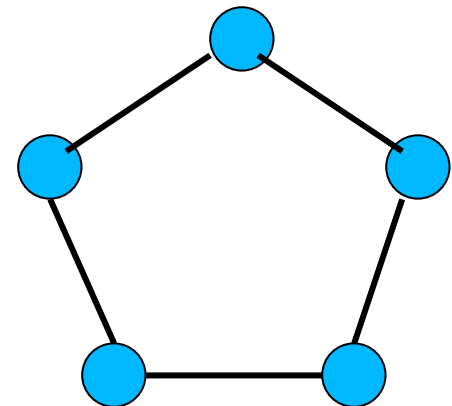
- **Example: Dining Philosophers (Dijkstra)**
  - $\mathcal{E} = \{ E : E \text{ includes a pair of neighbors } \}$
  - Forks (resources) between adjacent philosophers; each needs both adjacent forks in order to eat.
  - Only one can hold a particular fork at a time, so adjacent philosophers must exclude each other.
- Not every exclusion problem can be expressed in this way.
  - k-exclusion cannot.

# Resource allocation problem, for a given exclusion specification $\mathcal{E}$

- Same shared-memory architecture as for mutual exclusion (processes and shared variables).
- **Well-formedness:** As before.
- **Exclusion:** There is no reachable state in which the set of users in  $C$  is a set in  $\mathcal{E}$ .
- **Progress:** As before.
- **Lockout-freedom:** As before.
- But these don't capture concurrency requirements.
  - Any lockout-free mutual exclusion algorithm also satisfies all these conditions (provided that  $\mathcal{E}$  doesn't contain any singleton sets).
- Can add concurrency conditions, e.g.:
  - Independent progress: If  $i \in T$  and every  $j$  that could conflict with  $i$  remains in  $R$ , then eventually  $i \rightarrow C$ .
  - Time bound: Obtain better bounds from  $i \rightarrow T$  to  $i \rightarrow C$ , even in the presence of conflicts, than we can for mutual exclusion.

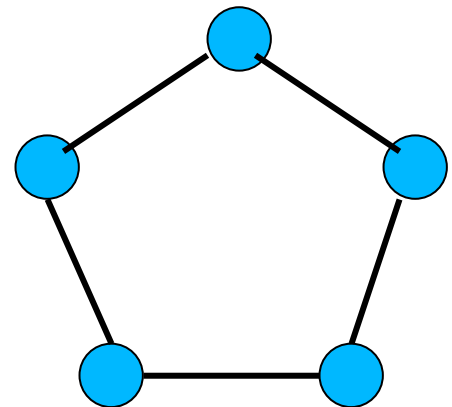
# Dining Philosophers Problem

- Like Mutual Exclusion, with a different exclusion condition:
  - No two neighbors are in C at once (exclusion specification), or
  - Forks on edges, each philosopher needs both adjacent forks to eat (explicit resource specification).
- Can use progress and fairness conditions as for Mutex.
- Can add new conditions to capture concurrent access to C.
- Dijkstra posed the problem, gave a solution for a strong shared-memory model.
  - Globally-shared variables, atomic access to all of shared memory.
- Distributed version: Assume the only shared variables are RMW variables corresponding to the forks, accessible only by processes at the endpoints.



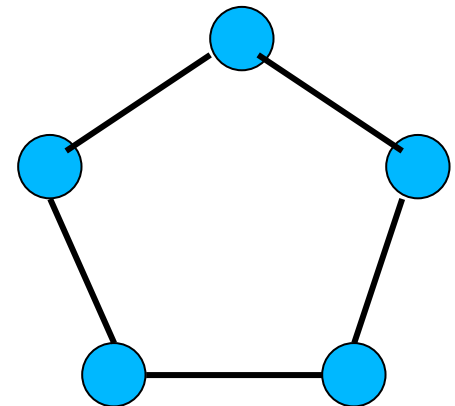
# Impossibility Result

- **Theorem 1:** If all processes are identical and refer to forks by local names L and R, and all shared vars have the same initial value, then we can't guarantee DP exclusion + progress.
- **Proof:** Can't break symmetry:
  - Consider only executions that work in synchronous rounds, prove by induction on rounds that symmetry is preserved.
  - By progress, someone  $\rightarrow C$ .
  - By symmetry, all do, violating DP exclusion.
- **Example symmetric algorithm:** Wait for R fork first, then L fork.
  - Guarantees DP exclusion.
  - Progress fails---all might get R, then wait forever for L (deadlock).

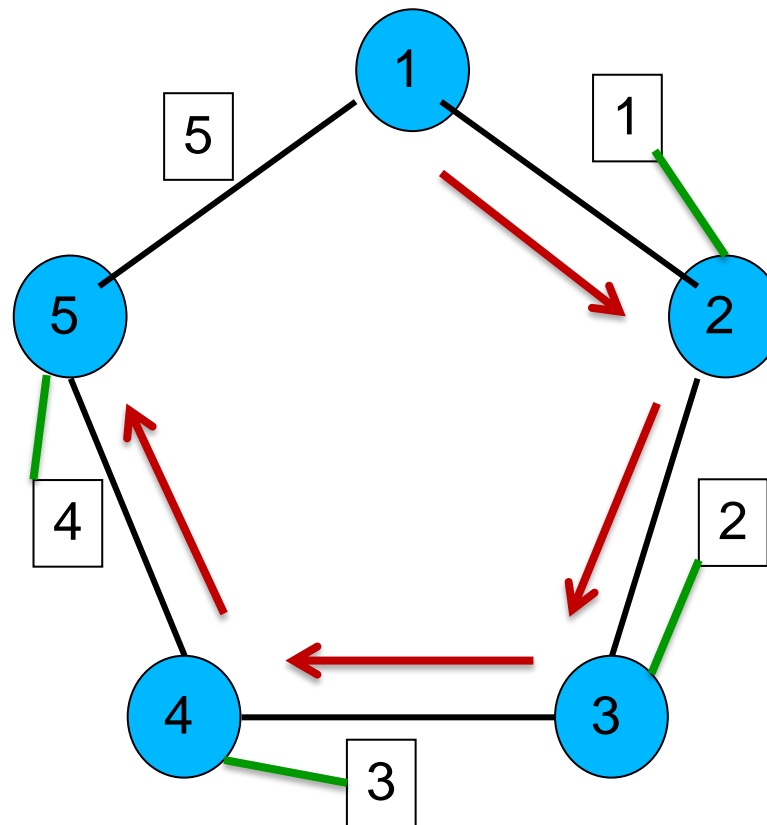


# DP Algorithms

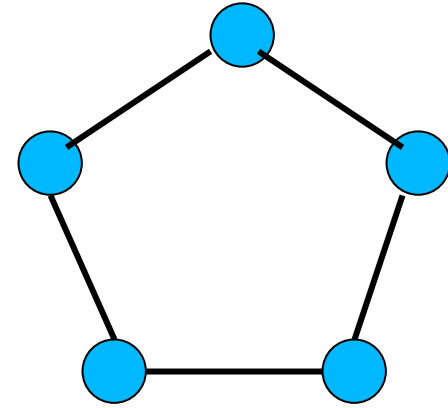
- So we need a mechanism to break symmetry.
- **Solution 1:** Number forks in increasing order around the table; every process picks up its smaller numbered fork first.
  - Yields DP exclusion, progress, lockout freedom, independent progress.
  - But the time isn't good---we can have a long “waiting chain” of processes waiting for neighbors to release forks.



# Creating a long waiting chain



# DP Algorithms



- **Solution 2:** Right/Left algorithm (**Burns**):
  - Classify processes as R or L (at least one of each).
  - R processes pick up right fork first, L processes pick up left fork first.
  - Yields DP exclusion, progress, lockout freedom, independent progress.
  - In even-sized rings in which R and L alternate, the lengths of waiting chains are limited to 2.
  - Yields a good (constant) time bound, LTTR.
- Generalize to solve any resource problem:
  - Represent the problem as an undirected graph.
  - Nodes = resources.
  - Edge between two resources if some user wants both.
  - Color the nodes of the graph; order colors.
  - All processes acquire resources in order of colors.



# Next time

- Impossibility of consensus in the presence of failures.
- Reading:
  - Lynch, Chapter 12
  - [Fischer, Lynch, Paterson] 2<sup>nd</sup> Dijkstra prize