# Deadlock

## Chapter 32 in "Three Easy Steps"
## Chapter 19 in Harmony Book

### CS 4410
### Operating Systems

# Dining Philosophers [Dijkstra 68]

Pi: **do forever**
    acquire( left(i) );
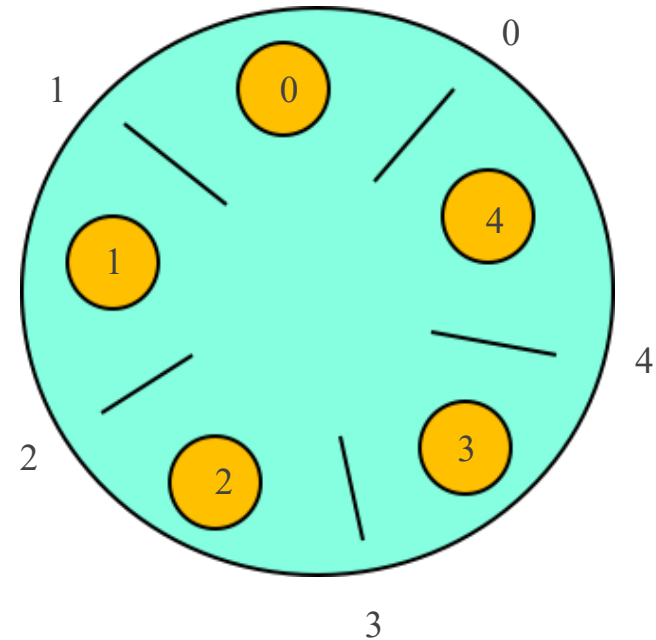    acquire( right(i) );
    eat
    release( left(i) );
    release( right(i) );
    think
  **end**

left(i):    i
right(i):   i+1 mod 5

# Dining Philosophers in Harmony

```
1   from synch import Lock, acquire, release
2
3   const N = 5
4
5   forks = [Lock(),] * N
6
7   def diner(which):
8       let left, right = (which, (which + 1) % N):
9           while choose({ False, True }):
10              acquire(?forks[left])
11              acquire(?forks[right])
12              # dine
13              release(?forks[left])
14              release(?forks[right])
15              # think
16
17  for i in {0..N-1}:
18      spawn diner(i)
```

# Dining Philosophers in Harmony

```
 1  from synch import Lock, acquire, release
 2
 3  const N = 5
 4
 5  forks = [Lock(),] * N
 6
 7  def diner(which):
 8      let left, right = (which, (which + 1) % N):
 9          while choose({ False, True }):
10              acquire(?forks[left])
11              acquire(?forks[right])
12              # dine
13              release(?forks[left])
14              release(?forks[right])
15              # think
16
17  for i in {0..N-1}:
18      spawn diner(i)
```

harmony –cN=3

**Summary: some execution cannot terminate**

Here is a summary of an execution that exhibits the issue:

- Schedule thread T0: **init**()
  - Line 5: Initialize forks to [ False, False, False ]
  - **Thread terminated**
- Schedule thread T2: diner(1)
  - Line 9: Choose True
  - Line synch/36: Set forks[1] to True (was False)
  - Preempted in diner(1) --> acquire(?forks[2]) about to execute atomic section in line synch/35
- Schedule thread T3: diner(2)
  - Line 9: Choose True
  - Line synch/36: Set forks[2] to True (was False)
  - Preempted in diner(2) --> acquire(?forks[0]) about to execute atomic section in line synch/35
- Schedule thread T1: diner(0)
  - Line 9: Choose True
  - Line synch/36: Set forks[0] to True (was False)
  - Preempted in diner(0) --> acquire(?forks[1]) about to execute atomic section in line synch/35

---

**Final state** (all threads have terminated or are blocked):

- Threads:
  - T1: (blocked) diner(0) --> acquire(?forks[1])
    - about to execute atomic section in line synch/35
  - T2: (blocked) diner(1) --> acquire(?forks[2])
    - about to execute atomic section in line synch/35
  - T3: (blocked) diner(2) --> acquire(?forks[0])
    - about to execute atomic section in line synch/35
- Variables:
  - forks: [ True, True, True ]

4

# Problematic Emergent Properties

**Starvation**:   Process waits forever

**Deadlock**:  A set of processes exists, where each is <span style="color:red">blocked</span> and can become unblocked only by actions of another process in the set.

- Deadlock implies Starvation (but not *vice versa*)

- Starvation often tied to **fairness**:  A process is not forever blocked awaiting a condition that (i) becomes continuously true or (ii) infinitely-often becomes true.

*Testing for starvation or deadlock is difficult in practice*

# More Examples of Deadlock

*Example (initially in1 = in2 = False):*
    in1 = True;  **await not** in2;  in1 = False
    //
    in2 := True;  **await not** in1;  in2 = False

*Example (initially lk1 = lk2 = released):*
    acquire(lk1); acquire(lk2);  release(lk2); release(lk1);
    //
    acquire(lk2); acquire(lk1); release(lk1); release(lk2);

# System Model

- Set of resources requiring "exclusive" access
    - Might be "k-exclusive access" if resource has capacity for k
    - Examples:  buffers, packets, I/O devices, processors, …

- Protocol to access a resource causes blocking:
    - If resource is free, then access is granted;  process proceeds
    - If resource is in use, then process blocks
    – Use resource
    – Release resource

# When is deadlock possible?

# Necessary Conditions for Deadlock

1. **Mutual Exclusion**.  Acquire can block invoker

2. **Hold & wait**.  A process can be blocked while holding resources

3. **No preemption**.  Allocated resources cannot be reclaimed.   Explicit release operation needed

4. **Circular waits** are possible

   *Let p →q denote "p waits for q to release a resource".  Then*

   P1 → P2 → … → Pn → P1

# Deadlock is Undesirable

- Deadlock <u>prevention</u>:  Ensure that a necessary condition cannot hold

- Deadlock <u>avoidance</u>:  System does not allocate resources that will lead to a deadlock

- Deadlock <u>detection</u>:  Allow system to deadlock; detect it; recover

# Deadlock Prevention: Negate 1

**#1: Eliminate mutual exclusion / bounded resources:**

- Make resources sharable without locks
  - E.g., time-shared CPU
  - Harmony book has examples of non-blocking concurrent data structures
- Have sufficient resources available, so acquire never delays
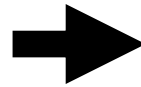  - E.g., make sure bounded queue is "large enough"

# Deadlock Prevention: Negate 2

## #2: **Eliminate hold and wait**

Don't hold some resources when requesting others

- Re-write code:

```
acquire(?foo_lock);
foo1();
acquire(?bar_lock);
bar();
release(?bar_lock);
foo2();
release(?foo_lock);
```

➡

```
acquire(?foo_lock);
foo1();
release(?foo_lock);
acquire(?bar_lock);
bar();
release(?bar_lock);
acquire(?foo_lock);
foo2();
release(?foo_lock);
```

- *Assuming* bar() *does not access shared variables protected by foo_lock, are these the same?*
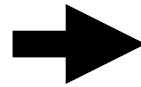
# Deadlock Prevention: Negate 2

## #2: **Eliminate hold and wait**

Don't hold some resources when requesting others

- Re-write code:

```
acquire(?foo_lock);
foo1();
acquire(?bar_lock);
bar();
release(?bar_lock);
foo2();
release(?foo_lock);
```

➡

```
acquire(?foo_lock);
foo1();
release(?foo_lock);
acquire(?bar_lock);
bar();
release(?bar_lock);
acquire(?foo_lock);
foo2();
release(?foo_lock);
```

- *Answer: no.*  *The state that foo_lock protects may change between foo1() and foo2() in code on the right*

# Deadlock Prevention: Negate 2

## #2: **Eliminate hold and wait**

Don't hold some resources when requesting others

- Re-write code

- Another approach: request all resources at once
  - Problems:
    - Processes don't know what they need ahead of time
    - No mechanism to request all resources at the same time
    - Starvation (if waiting on many popular resources)
    - Low utilization (need resource only for a bit)

# Simultaneous Acquisition in Harmony

```
5    mutex = synch.Lock()
6    forks = [False,] * N
7    conds = [synch.Condition(),] * N
8
9    def diner(which):
10       let left, right = (which, (which + 1) % N):
11           while choose({ False, True }):
12               synch.acquire(?mutex)
13               while forks[left] or forks[right]:
14                   if forks[left]:
15                       synch.wait(?conds[left], ?mutex)
16                   if forks[right]:
17                       synch.wait(?conds[right], ?mutex)
18               assert not (forks[left] or forks[right])
19               forks[left] = forks[right] = True
20               synch.release(?mutex)
21               # dine
22               synch.acquire(?mutex)
23               forks[left] = forks[right] = False
24               synch.notify(?conds[left])
25               synch.notify(?conds[right])
26               synch.release(?mutex)
27               # think
```

wait for both forks and then grab them both

release both forks

14

# Simultaneous Acquisition in Harmony

```
5   mutex = synch.Lock()
6   forks = [False,] * N
7   conds = [synch.Condition(),] * N
8
9   def diner(which):
10      let left, right = (which, (which + 1) % N):
11          while choose({ False, True }):
12              synch.acquire(?mutex)
13              while forks[left] or forks[right]:
14                  if forks[left]:
15                      synch.wait(?conds[left], ?mutex)
16                  if forks[right]:
17                      synch.wait(?conds[right], ?mutex)
18              assert not (forks[left] or forks[right])
19              forks[left] = forks[right] = True
20              synch.release(?mutex)
21              # dine
22              synch.acquire(?mutex)
23              forks[left] = forks[right] = False
24              synch.notify(?conds[left])
25              synch.notify(?conds[right])
26              synch.release(?mutex)
27              # think
```

Wait for both forks and then grab them both

# Simultaneous Acquisition in Harmony

```
5   mutex = synch.Lock()
6   forks = [False,] * N
7   conds = [synch.Condition(),] * N
8
9   def diner(which):
10      let left, right = (which, (which + 1) % N):
11          while choose({ False, True }):
12              synch.acquire(?mutex)
13
14              while forks[left]:
15                  synch.wait(?conds[left], ?mutex)
16              while forks[right]:
17                  synch.wait(?conds[right], ?mutex)
18              assert not (forks[left] or forks[right])
19              forks[left] = forks[right] = True
20              synch.release(?mutex)
21              # dine
22              synch.acquire(?mutex)
23              forks[left] = forks[right] = False
24              synch.notify(?conds[left])
25              synch.notify(?conds[right])
26              synch.release(?mutex)
27              # think
```

Wait for left fork, then wait for right fork. Wouldn't this be just as good?

# Simultaneous Acquisition in Harmony

```
5   mutex = synch.Lock()
6   forks = [False,] * N
7   conds = [synch.Condition(),] * N
8
9   def diner(which):
10      let left, right = (which, (which + 1) % N):
11          while choose({ False, True }):
12              synch.acquire(?mutex)
```

$$\mathbf{while}\ forks[left]:$$
$$\quad synch.\mathrm{wait}(?conds[left], ?mutex)$$
$$\mathbf{while}\ forks[right]:$$
$$\quad synch.\mathrm{wait}(?conds[right], ?mutex)$$

```
18              assert not (forks[left] or forks[right])
19              forks[left] = forks[right] = True
20              synch.release(?mutex)
21              # dine
22              synch.acquire(?mutex)
23              forks[left] = forks[right] = False
24              synch.notify(?conds[left])
25              synch.notify(?conds[right])
26              synch.release(?mutex)
27              # think
```

Wait for left fork, then wait for right fork. Wouldn't this be just as good?

NO!

(run through harmony if you don't believe me)

17

# Deadlock Prevention: Negate 3

## #3: Allow preemption

Requires mechanism to save / restore resource state:

multiplexing     vs     undo/redo

- Examples of multiplexing:
  - processor registers (contexts)
  - regions of memory (pages)

- Examples of undo/redo
  - database transaction processing

# Deadlock Prevention: Negate 4

**#4: Eliminate circular waits.**

Let R = {R1, R2, … Rn} be the set of resource types.
Let ( R , $<$ ) be a non-symmetric relation:
  for every $r, s, t$:
  - $\neg (r < r)$   [irreflexive]
  - $(r < s \wedge s < t) \Rightarrow r < t$   [transitive]
  - $\neg (r < s \wedge s < r)$   [non-symmetric]
  - $r \neq s \Rightarrow (r < s \vee s < r)$ [total order]

**Rule**:   Request resources in increasing order by $<$
(All resources from type Ri must be requested together)

**Rule**:   To request resources of type Ri, first release all resources
from type Rj where Ri $<$ Rj.

# Why < Rule Works

**Thm**:  Total order resource allocation avoids circular waits

Proof:  By contradiction.  Assume a circular wait exists
$$P1 \rightarrow P2 \rightarrow P3 \rightarrow \ldots \rightarrow Pn \rightarrow P1.$$
P1 requesting R1 held by P2.
P2 requesting R2 held by P3.  (So R1 < R2 holds)
…

Conclude:  R1 < R2,  R2 < R3,  …,   Rn < R1
By transitivity:  R1 < R1.  Violates irreflexivity.
A contradiction!

# Dining Philosophers (Again)

Pi: **do forever**
    acquire( F(i) );
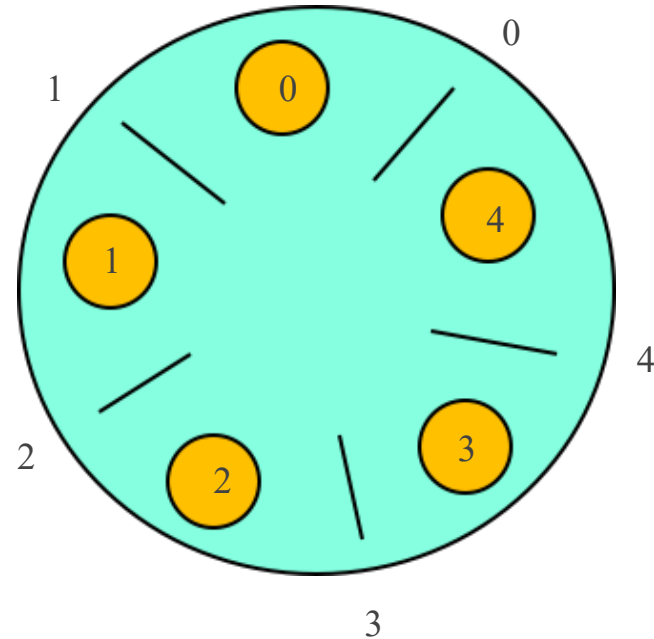    acquire( G(i) );
    eat
    release( F(i) );
    release( G(i) );
  **end**

F(i):   min(i, i+1 mod 5)
G(i):   max(i, i+1 mod 5)

# Ordering Resources in Harmony

```
1   if left < right:
2       synch.acquire(?forks[left])
3       synch.acquire(?forks[right])
4   else:
5       synch.acquire(?forks[right])
6       synch.acquire(?forks[left])
```

*or*

```
1   synch.acquire(?forks[min(left, right)])
2   synch.acquire(?forks[max(left, right)])
```
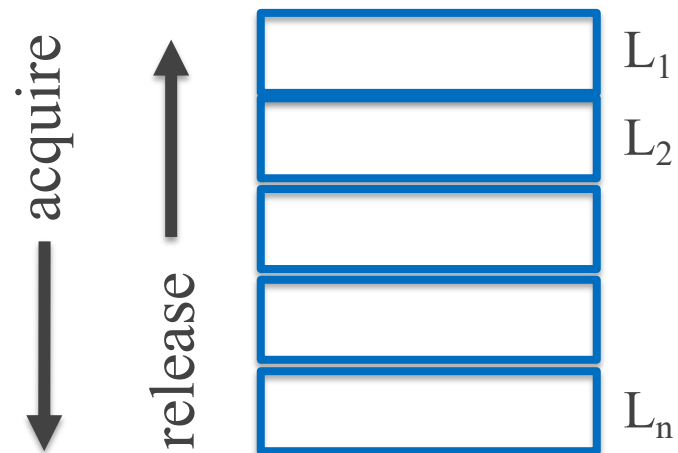
# Havender's Scheme (OS/360)

## Hierarchical Resource Allocation

Every resource is associated with a level.

- **Rule H1**: All resources from a given level must be acquired using a single request.
- **Rule H2**: After acquiring from level $L_j$ must not acquire from $L_i$ where $i < j$
- **Rule H3**: May not acquire from $L_i$ unless already released from $L_j$ where $j > i$.

Example of allowed sequence:
1. acquire(W@L1, X@L1)
2. acquire(Y@L3)
3. release(Y@L3)
4. acquire(Z@L2)

acquire

release

$L_1$
$L_2$

$L_n$

# Deadlock Detection

Create a <u>Wait-For Graph</u>
- 1 Node per Process
- 1 Outgoing Edge per Waiting Process, P
(from P to the process it's waiting for)

Note: graph holds for a single instant in time

**Cycle** in graph indicates deadlock

# Testing for cycles ( = deadlock)

**Reduction Algorithm**:
Find a node with no outgoing edges
- Erase node
- Erase any edges coming into it
- Repeat until no such node

**Intuition**: Deleted node is for process that is not waiting.  It will eventually finish and release its resources, so any process waiting for those resources will longer be waiting.

Erase whole graph ↔ graph has no cycles
Graph remains ↔ deadlock

# Graph Reduction: Example 1

Find node w/o outgoing edges
Erase node
Erase edges coming into it



Graph can be fully reduced, hence there was no deadlock at the time the graph was drawn. (Obviously, things could change later!)

# Graph Reduction: Example 2



*No node* with no outgoing edges…
Irreducible graph, contains a cycle
    (only some processes are in the cycle)
➤ deadlock

# Question:

Does choice of node for reduction matter?

**Answer: No.**
**Explanation:** an unchosen candidate at one step remains a candidate for later steps. Eventually—regardless of order—every node will be reduced (if there's no deadlock).

# Question:

Suppose no deadlock detected at time T.
Can we infer about a later time T+x?

**Answer:** Nothing.
**Explanation:**  The very next step could be to run some process that will request a resource…
        … establishing a cyclic wait
        … and causing deadlock

# Implementing Deadlock Detection

- Track resource allocation (who has what)
- Track pending requests (who's waiting for what)

Maintain a wait-for graph.

When to run graph reduction?

- Whenever a request is blocked?
- Periodically?
- Once CPU utilization drops below a threshold?

# Deadlock Recovery

Blue screen & reboot?

Kill one/all deadlocked processes
- Pick a victim
- Terminate
- Repeat if needed

Preempt resource/processes till deadlock broken
- Pick a victim (# resources held, execution time)
- Rollback (partial or total, not always possible)

# Deadlock in traffic

# Deadlock Avoidance

How do cars do it?
- Try not to block an intersection
- Don't drive into the intersection if you see that you might get stuck there

Why does this work?
- Prevents a wait-for relationship
- Cars won't take up a resource if they see they won't be able to acquire the next one…

# Deadlock Avoidance

**state**:  allocation to each process
**safe state**:  a state from which some execution is possible that does not cause deadlock

- Requires knowing max allocation for each process and who holds what resources
- Check that
  - Exists sequence P1 P2 … Pn of processes where:

    Forall i where $1 \leq i \leq n$:

    Pi can be satisfied by Avail + resources held by P1 … Pi-1.

Assumes no synchronization between processes, except for resource requests

# Safe State Example

Suppose: 12 tape drives and 3 processes: p0, p1, and p2

|     | max need | current usage | could still ask for |
| --- | --- | --- | --- |
| p0  | 10  | 5   | 5   |
| p1  | 4   | 2   | 2   |
| p2  | 9   | 2   | 7   |

3 drives remain

*Is this a safe state (i.e, is there a sequence of granting requests that will work without deadlock)?*

# Safe State Example

Suppose: 12 tape drives and 3 processes: p0, p1, and p2

|     | max need | current usage | could still ask for |
|-----|----------|---------------|---------------------|
| p0  | 10       | 5             | 5                   |
| p1  | 4        | 2             | 2                   |
| p2  | 9        | 2             | 7                   |

**3 drives remain**

Current state is *safe* because a safe sequence exists: [p1, p0, p2]
- p1 can complete with remaining resources
- p0 can complete with remaining+p1
- p2 can complete with remaining+p1+p0

What if p2 requests 1 drive? Grant or not?

# Safe State Example

Suppose: 12 tape drives and 3 processes: p0, p1, and p2

|      | max need | current usage | could still ask for |
|------|----------|---------------|---------------------|
| p0   | 10       | 5             | 5                   |
| p1   | 4        | 2             | 2                   |
| p2   | 9        | 3             | 6                   |

2 drives remain

Is this state safe?  (Is there a sequence of requests that works?)

# Safe State Example

Suppose: 12 tape drives and 3 processes: p0, p1, and p2

| | max<br>need | current<br>usage | could still<br>ask for |
|---|---|---|---|
| p0 | 10 | 5 | 5 |
| ~~p1~~ | ~~0~~ | ~~0~~ | ~~0~~ |
| p2 | 9 | 3 | 6 |

4 drives remain

Is this state safe?  (Is there a sequence of requests that works?)

# Safe State Example

Suppose: 12 tape drives and 3 processes: p0, p1, and p2

| | max need | current usage | could still ask for |
|---|---|---|---|
| p0 | 10 | 5 | 5 |
| ~~p1~~ | ~~0~~ | ~~0~~ | ~~0~~ |
| p2 | 9 | 3 | 6 |

4 drives remain

Is this state safe? (Is there a sequence of requests that works?)

*(potentially) STUCK…*
**(non-terminating state)**

# Safe State Example

Suppose: 12 tape drives and 3 processes: p0, p1, and p2

|    | max need | current usage | could still ask for |
|----|----------|---------------|---------------------|
| p0 | 10       | 5             | 5                   |
| p1 | 4        | 2             | 2                   |
| p2 | 9        | 2             | 7                   |

**3 drives remain**

Current state is *safe* because a safe sequence exists: [p1, p0, p2]
- p1 can complete with remaining resources
- p0 can complete with remaining+p1
- p2 can complete with remaining+p1+p0

What if p2 requests 1 drive? Grant or not?

# Safe State Example

Suppose: 12 tape drives and 3 processes: p0, p1, and p2

|     | max need | current usage | could still ask for |
| --- | --- | --- | --- |
| p0  | 10  | 5   | 5   |
| p1  | 4   | 2   | 2   |
| p2  | 9   | 2   | 7   |

3 drives remain

Current state is *safe* because a safe sequence exists: [p1, p0, p2]
- p1 can complete with remaining resources
- p0 can complete with remaining+p1
- p2 can complete with remaining+p1+p0

What if p2 requests 1 drive? Grant or not?   NO (block or deny)

# Banker's Algorithm

- from 10,000 feet:
  - Process declares its worst-case needs, asks for what it "really" needs, a little at a time
  - Algorithm decides when to grant requests
    - Build a graph assuming request granted
    - Reducible?  yes: grant request, no: wait

Problems:
- Fixed number of processes
- Need worst-case needs ahead of time
- Expensive

→ not used much in practice