

SALVATORE DI GIROLAMO <DIGIROLS@INF.ETHZ.CH>

# DPHPC: Balance Principles & Scheduling

*Recitation session*

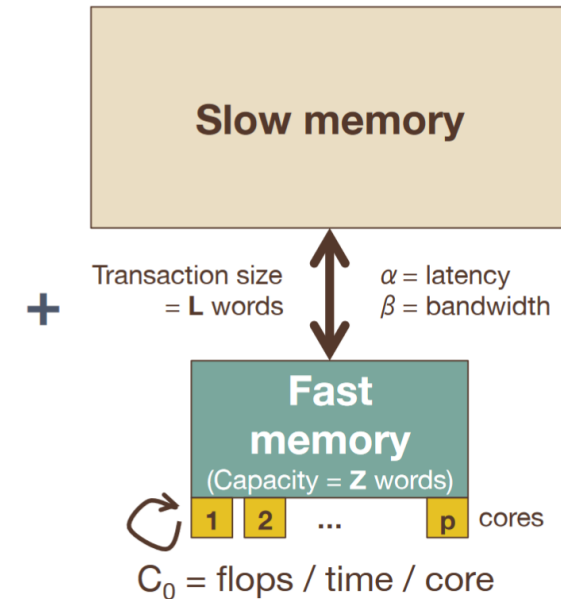
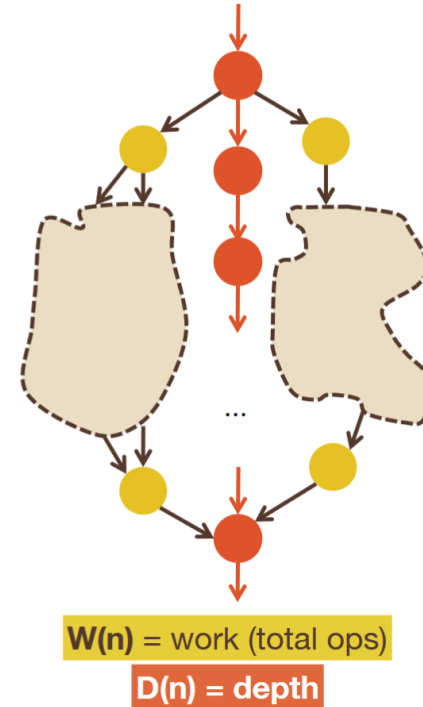


# Deriving a Balance Principle

- **Concept of balance:** a computation running on some machine is efficient if the compute-time dominates the I/O time. [Kung, 1986]

- **Deriving a balance principle:**

- Algorithmically analyze the parallelism
- Algorithmically analyze the I/O behavior (i.e., number of memory transfers)
- Combine these two analyses with a cost model for an abstract machine.



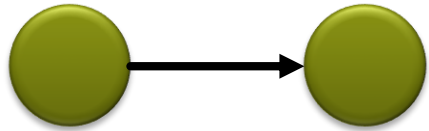
- **Goal: say precisely and analytically how**

- Changes to the architecture might affect the scaling of a computation
- Identify what classes of computation might execute efficiently on a given architecture

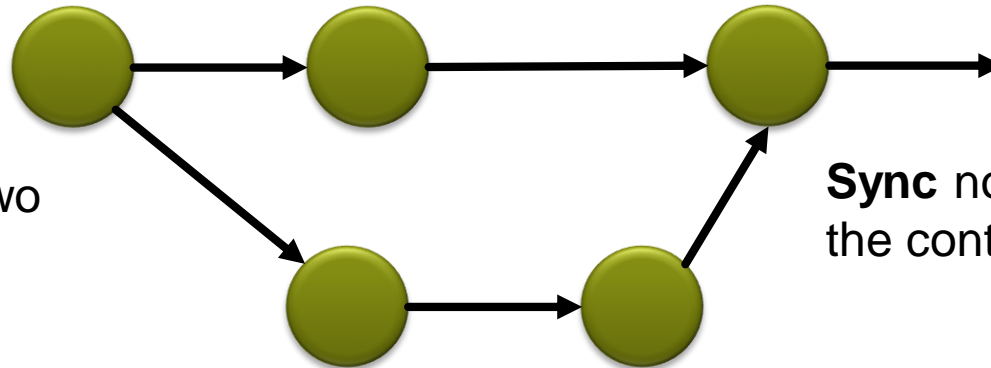
# The DAG Model



**Strand:** chain of serially executed instructions.



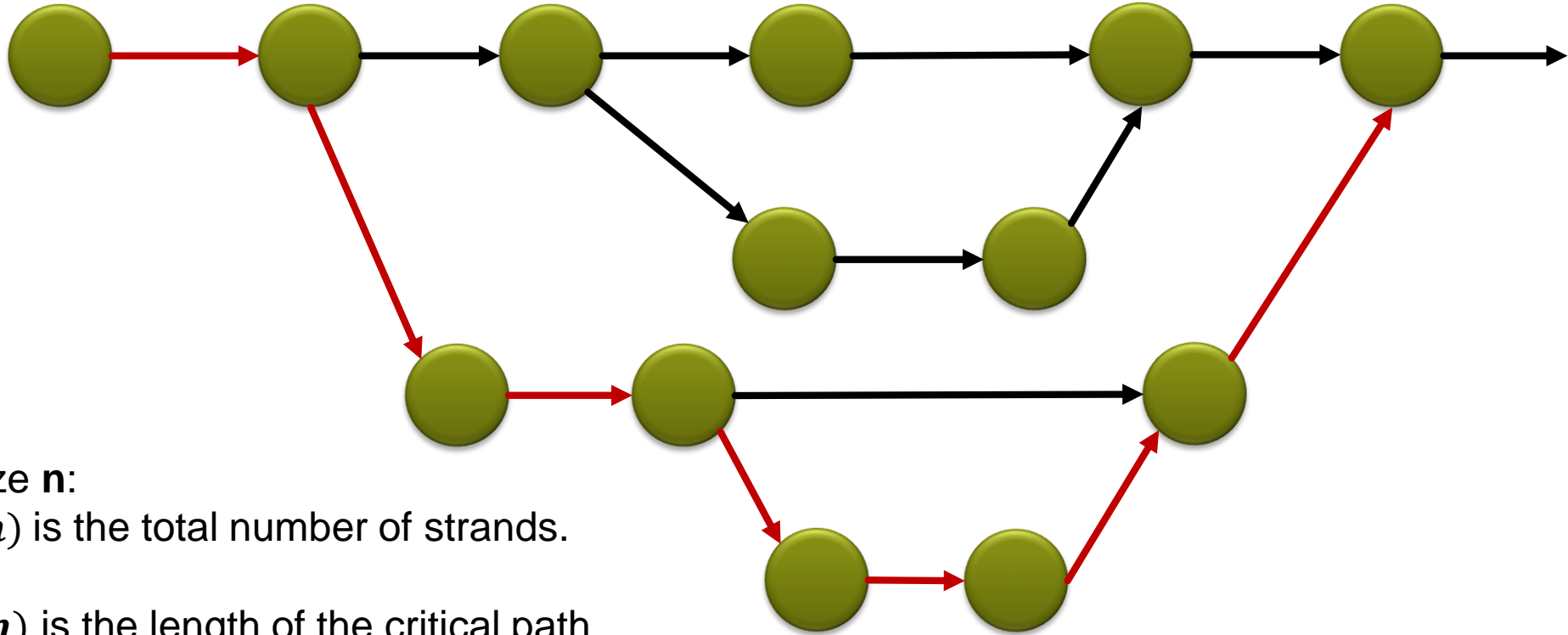
Strands are partially ordered with **dependencies**



**Spawn** nodes have two successors

**Sync** nodes are where the control flow merges

# The DAG Model



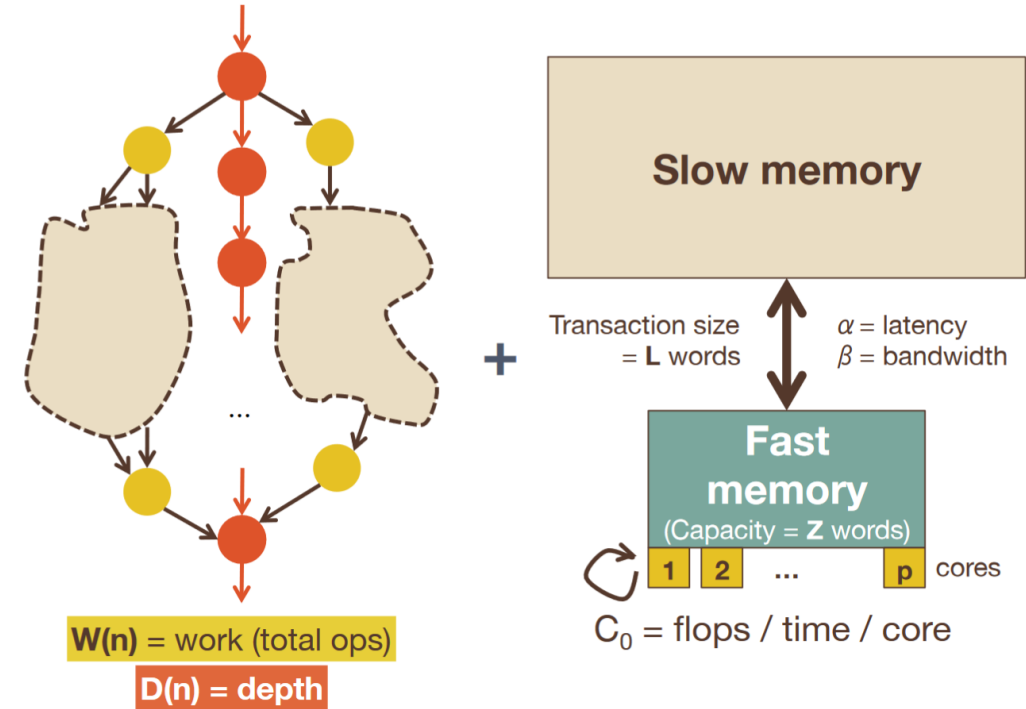
Given an input size  $n$ :

- The **work**  $W(n)$  is the total number of strands.
  - $W(n)=13$
- The **depth**  $D(n)$  is the length of the critical path (measured in number of strands).
  - Defines the minimum execution time of the computation
  - $D(n)=8$

The ratio  $\frac{W(n)}{D(n)}$  measures the average available parallelism

# Analyzing I/Os

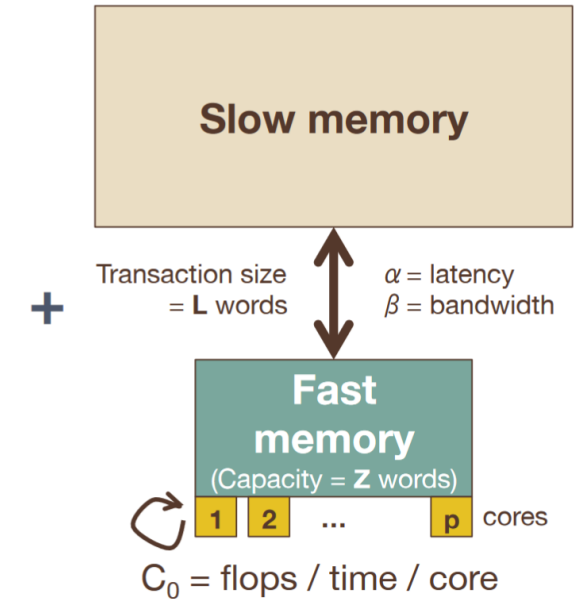
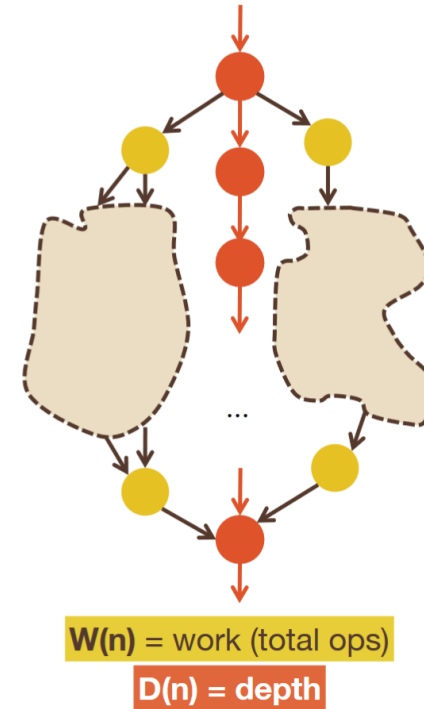
- We use the classical external memory model
- Two level memory
  - One large&slow
  - The other small&fast (capacity:  $Z$  words)
    - It can be an automatic cache or a software-controlled scratchpad*
- Work operations can be performed only on data in fast memory
- Slow $\leftrightarrow$ Fast memory transfers occur in blocks of  $L$  words
- $Q_{Z,L}(n)$  is the number of  $L$ -sized transfers between slow and fast memory for an input of size  $n$



Goal is to optimize the computational intensity:  $\frac{W(n)}{Q_{Z,L}(n) \cdot L}$

# Architecture-Specific Cost Model

- We need to introduce the time
  - This depends on the specific architecture
- $p$  cores
- Each core can deliver  $C_0$  operations per unit time
- The time to transfer  $m \cdot L$  words is:
  - $\alpha + m \cdot L / \beta$
  - $\alpha$  is the latency
  - $\beta$  is the bandwidth in units of words per time

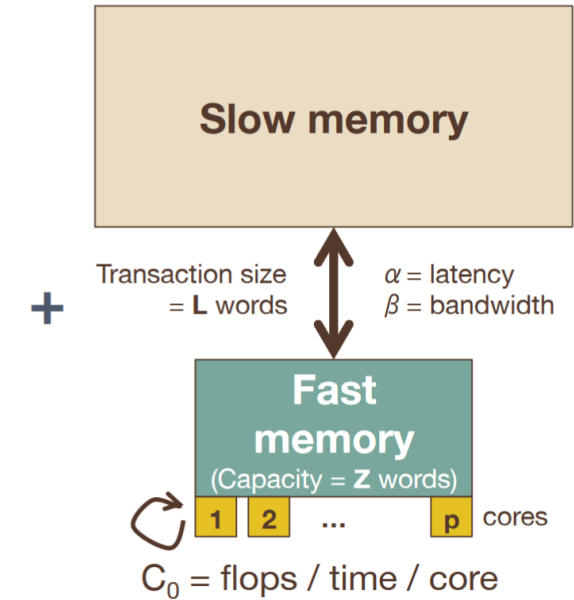
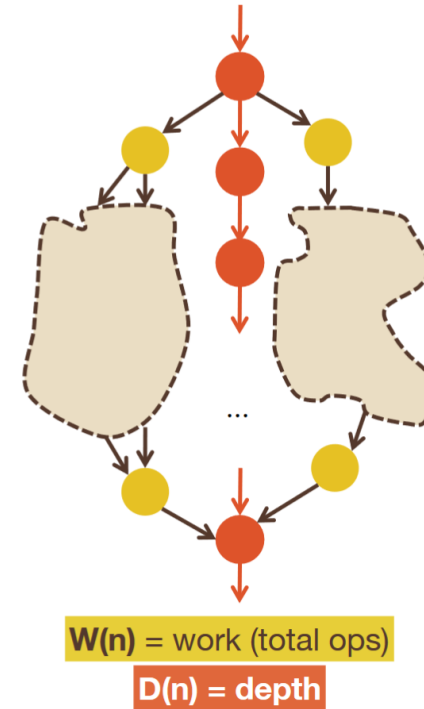


- The best possible compute time is (Brent's theorem):

$$T_{\text{comp}}(n; p, C_0) = \left( D(n) + \frac{W(n)}{p} \right) \cdot \frac{1}{C_0}$$

# Architecture-Specific Cost Model

- $Q_{Z,L}(n)$  is for the sequential case
- We need to move to the parallel case  $Q_{p;Z,L}(n)$ 
  - We can bound  $Q_{p;Z,L}(n)$  in terms of  $Q_{Z,L}(n)$   
*Blelloch et al, 2009, need to select a specific scheduler*
  - Compute it directly
- **Assumptions:**
  - the latency is accounted for each node in the critical path
  - all the  $Q_{p;Z,L}(n)$  are aggregated and pipelined by the memory system  
*Hence they are delivered at the peak bandwidth*
- We can estimate the memory cost as:



$$T_{\text{mem}}(n; p, Z, L, \alpha, \beta) = \alpha \cdot D(n) + \frac{Q_{p;Z,L}(n) \cdot L}{\beta}$$

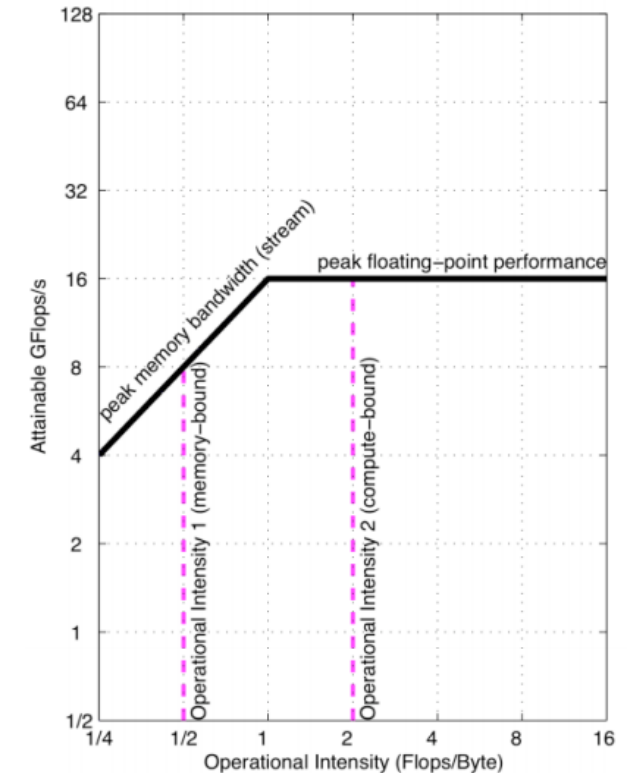
# The Balance Principle

- The balance principle follows by imposing  $T_{mem} \leq T_{comp}$

$$T_{mem}(n; p, Z, L, \alpha, \beta) = \alpha \cdot D(n) + \frac{Q_{p;Z,L}(n) \cdot L}{\beta}$$

$$T_{comp}(n; p, C_0) = \left( D(n) + \frac{W(n)}{p} \right) \cdot \frac{1}{C_0}$$

$$\underbrace{\frac{pC_0}{\beta}}_{\text{balance}} \left( 1 + \underbrace{\frac{\alpha\beta/L}{Q/D}}_{\text{Little's}} \right) \leq \underbrace{\frac{W}{QL}}_{\text{intensity}} \left( 1 + \underbrace{\frac{p}{W/D}}_{\text{Amdahl's}} \right)$$

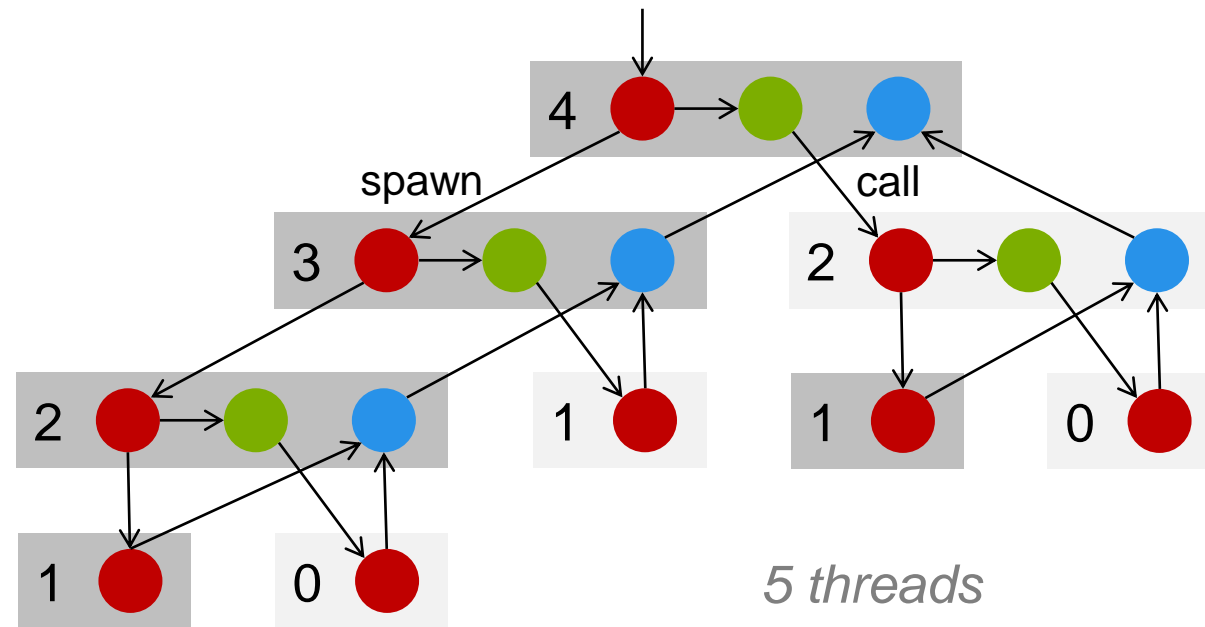




# Scheduling

The DAG unfolds dynamically:

```
int fib (int n) {
  if (n<2) return (n);
  else {
    int x,y;
    x = spawn fib(n-1);
    y = fib(n-2);
    sync;
    return (x+y);
  }
}
```

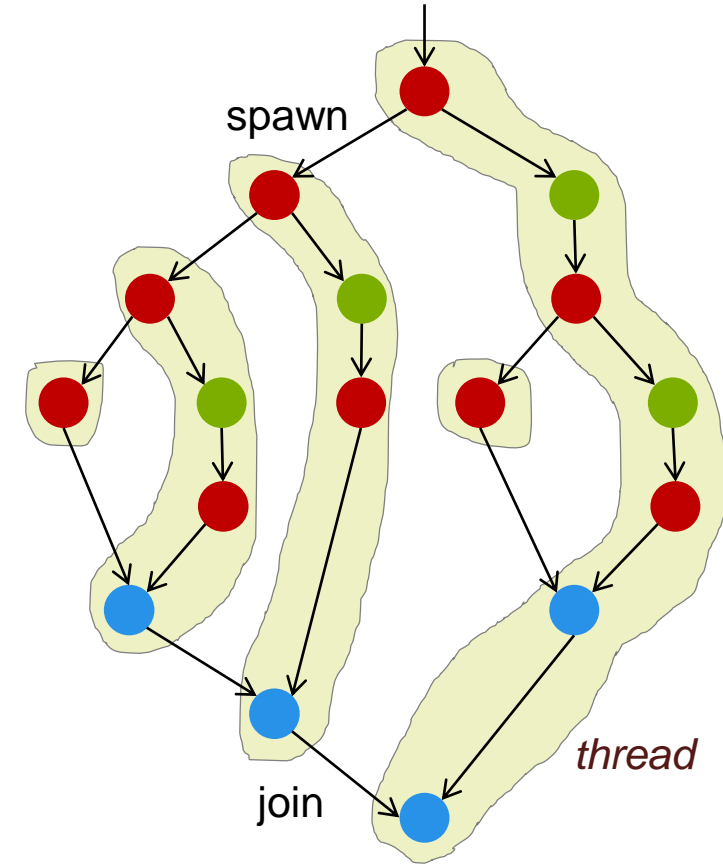
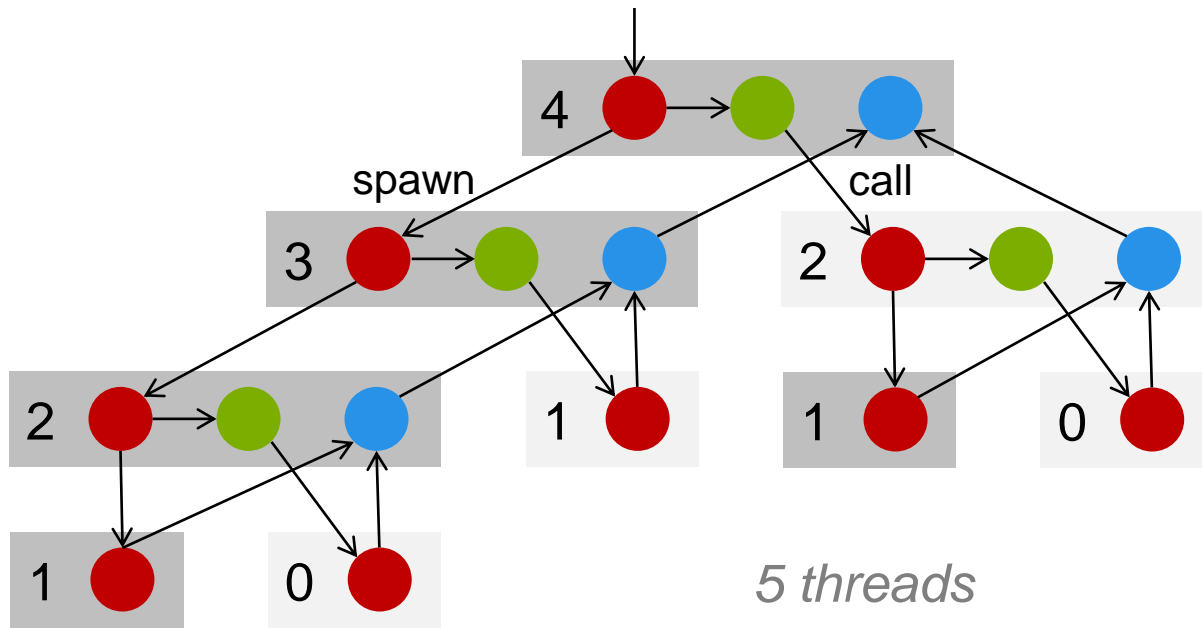


**Node:** Sequence of instructions without call, spawn, sync, return

**Edge:** Dependency

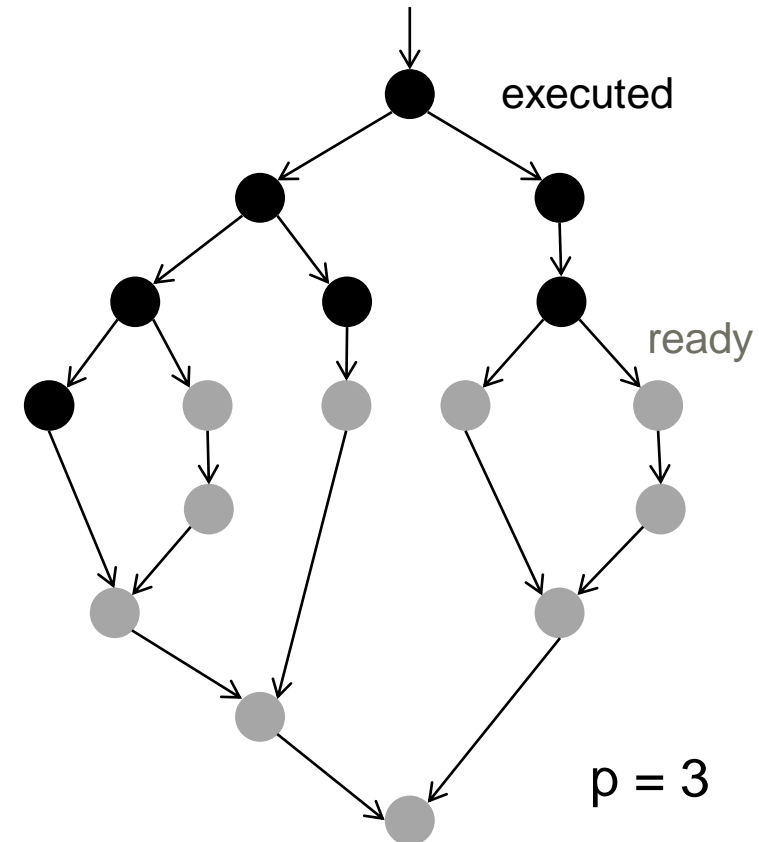
# Scheduling

*The DAG unfolds dynamically:*



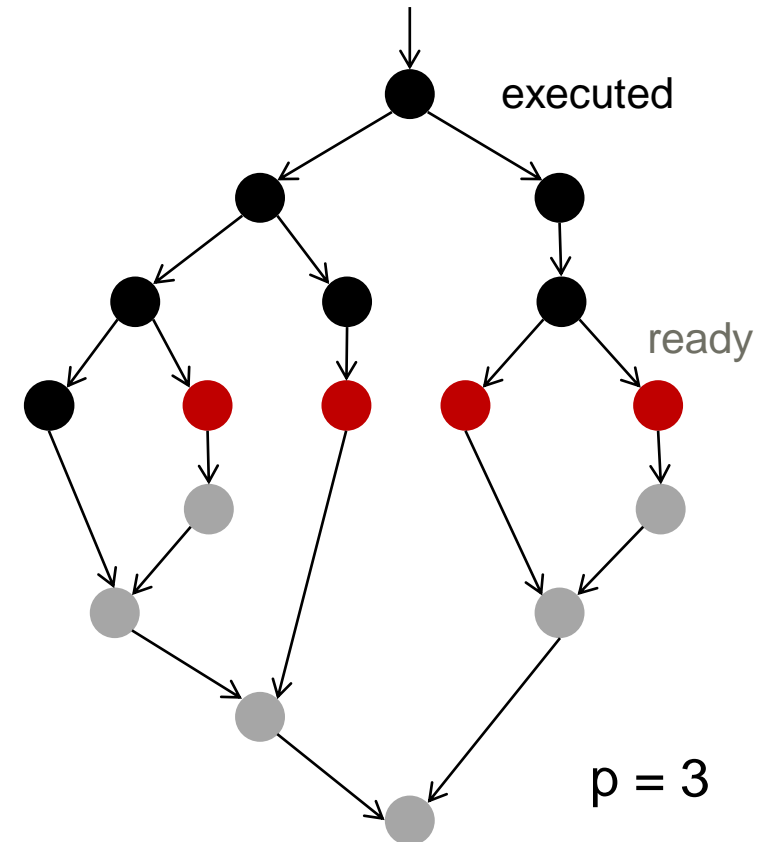
# Greedy Scheduler

- **Idea:** Do as much as possible in every step
- **Definition:** A node is ready if all predecessors have been executed



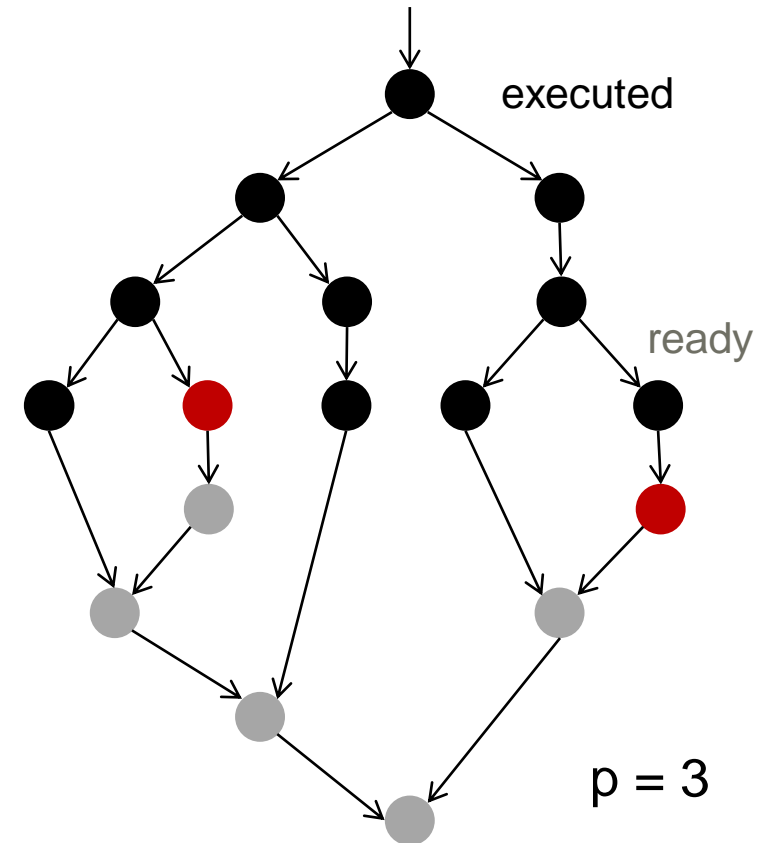
# Greedy Scheduler

- **Idea:** Do as much as possible in every step
- **Definition:** A node is ready if all predecessors have been executed
- **Complete step:**
  - $\geq p$  nodes are ready
  - run any  $p$



# Greedy Scheduler

- **Idea:** Do as much as possible in every step
- **Definition:** A node is ready if all predecessors have been executed
- **Complete step:**
  - $\geq p$  nodes are ready
  - run any  $p$
- **Incomplete step:**
  - $< p$  nodes ready
  - run all



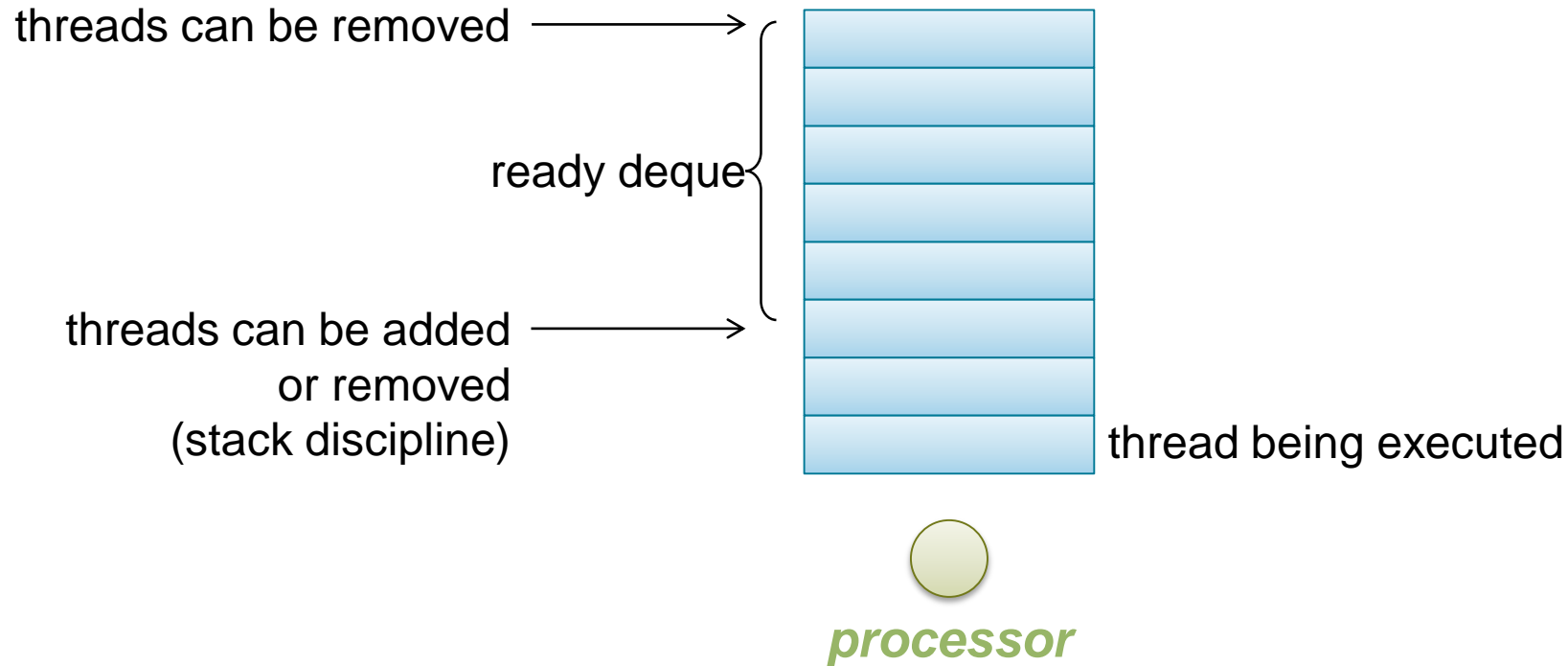
# Greedy Scheduler

**Maintain thread pool of live threads, each is ready or not**

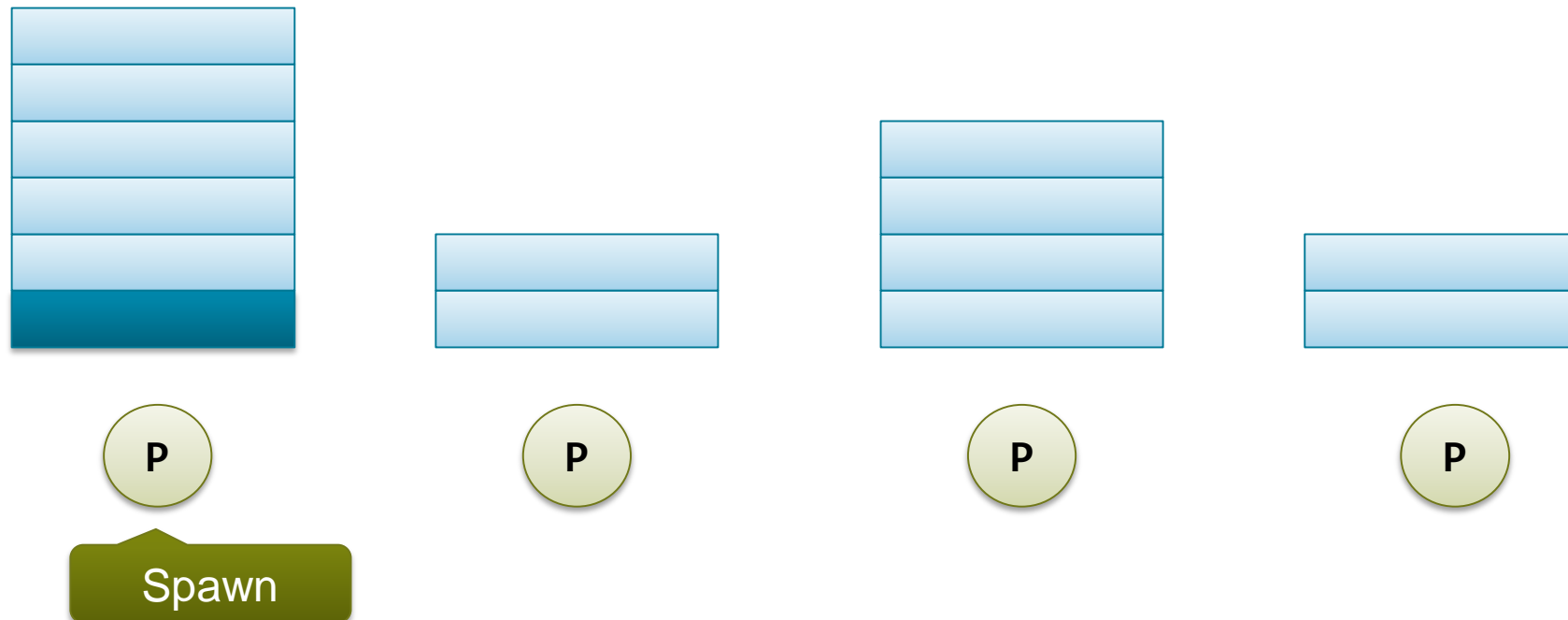
- **Initial:** Root thread in thread pool, all processors idle
- **At the beginning of each step each processor is idle or has a thread T to work on**
- **If idle**
  - *Get ready thread from pool*
- **If has thread T**
  - Case 0: T has another instruction to execute  
*execute it*
  - Case 1: thread T spawns thread S  
*return T to pool, continue with S*
  - Case 2: T stalls  
*return T to pool, then idle*
  - Case 3: T dies  
*if parent of T has no living children, continue with the parent, otherwise idle*

# Work Stealing Scheduler

- Each processor maintains a “ready deque:” deque of threads ready for execution; bottom is manipulated as a stack

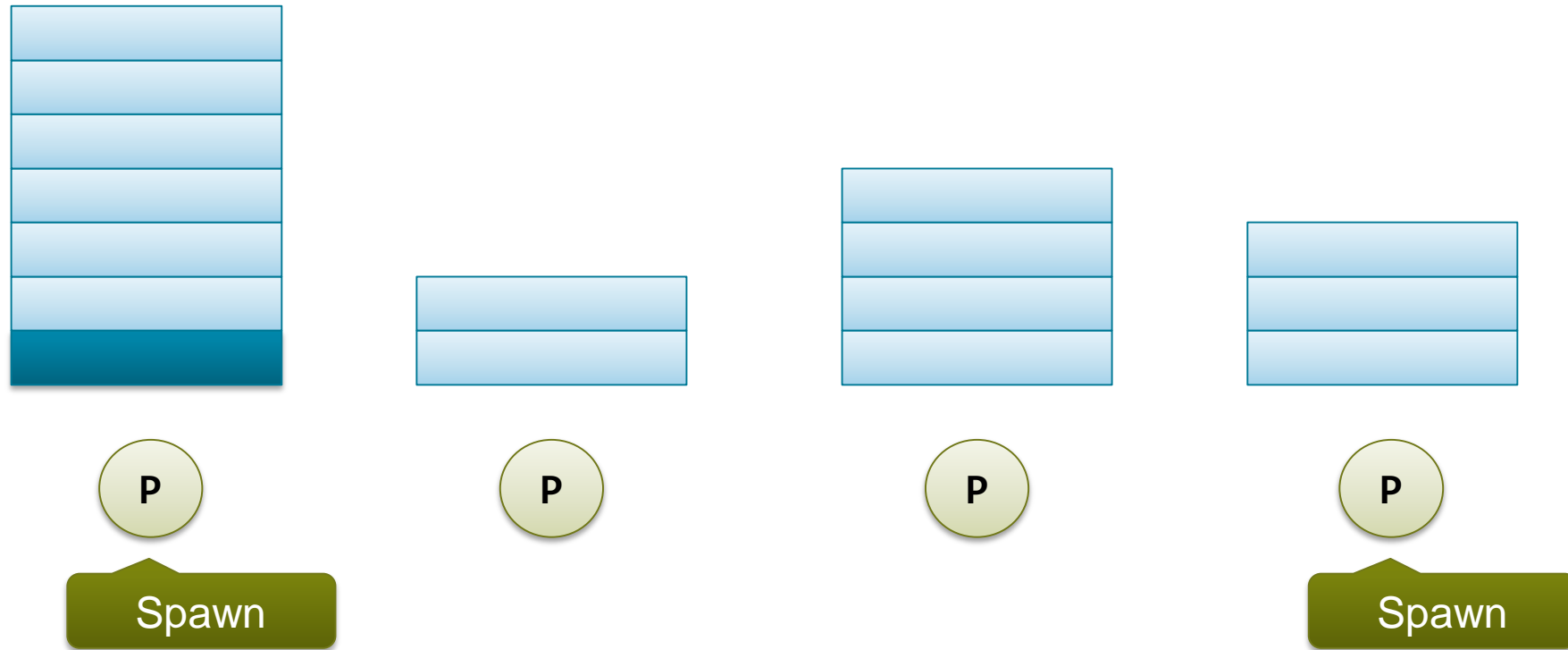


# Work Stealing Scheduler

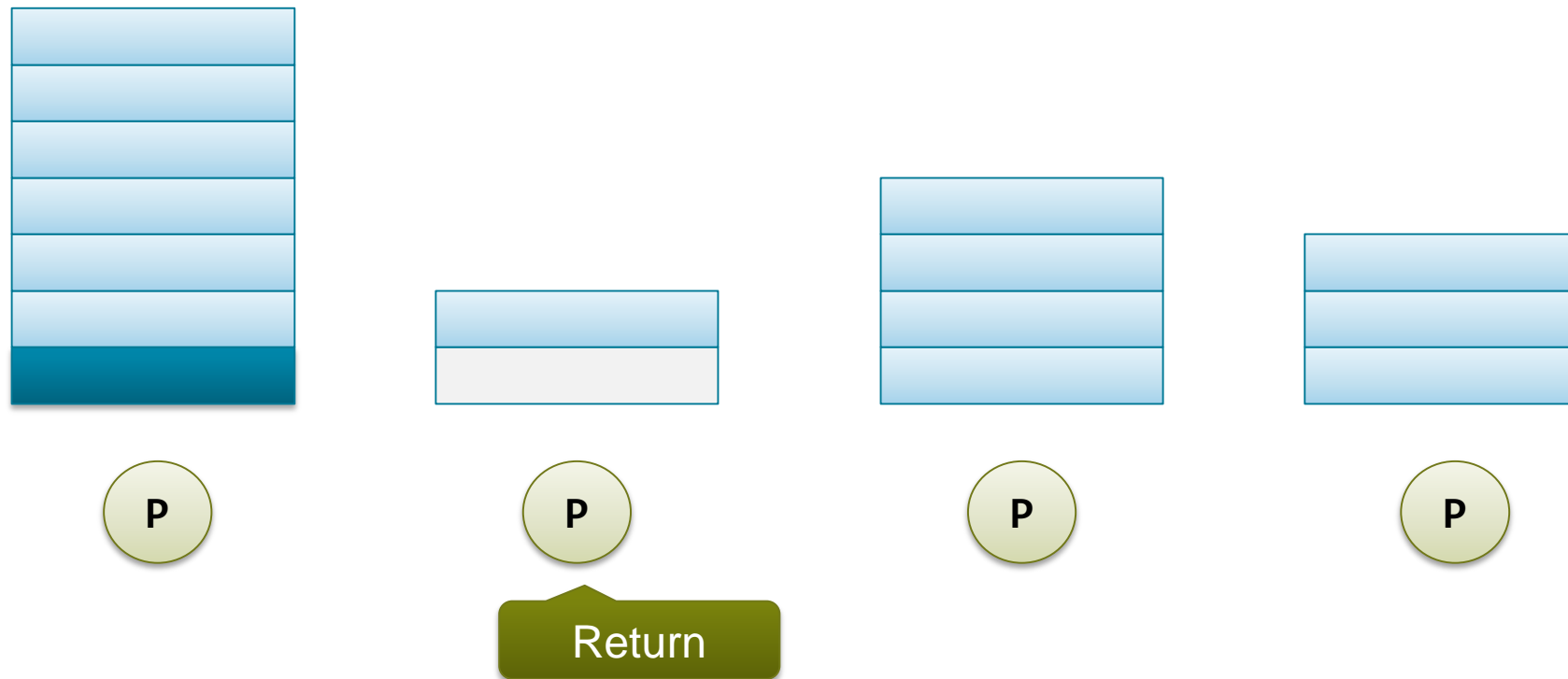




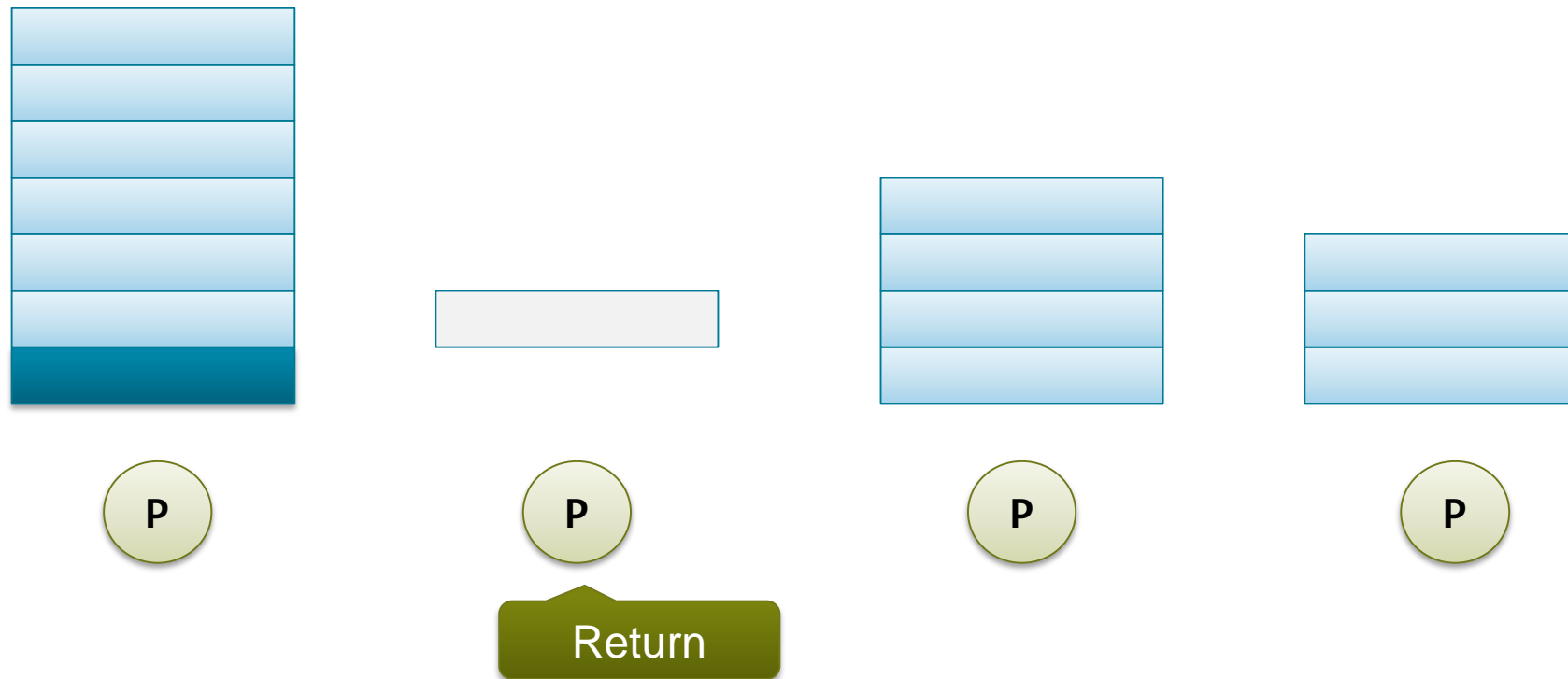
# Work Stealing Scheduler



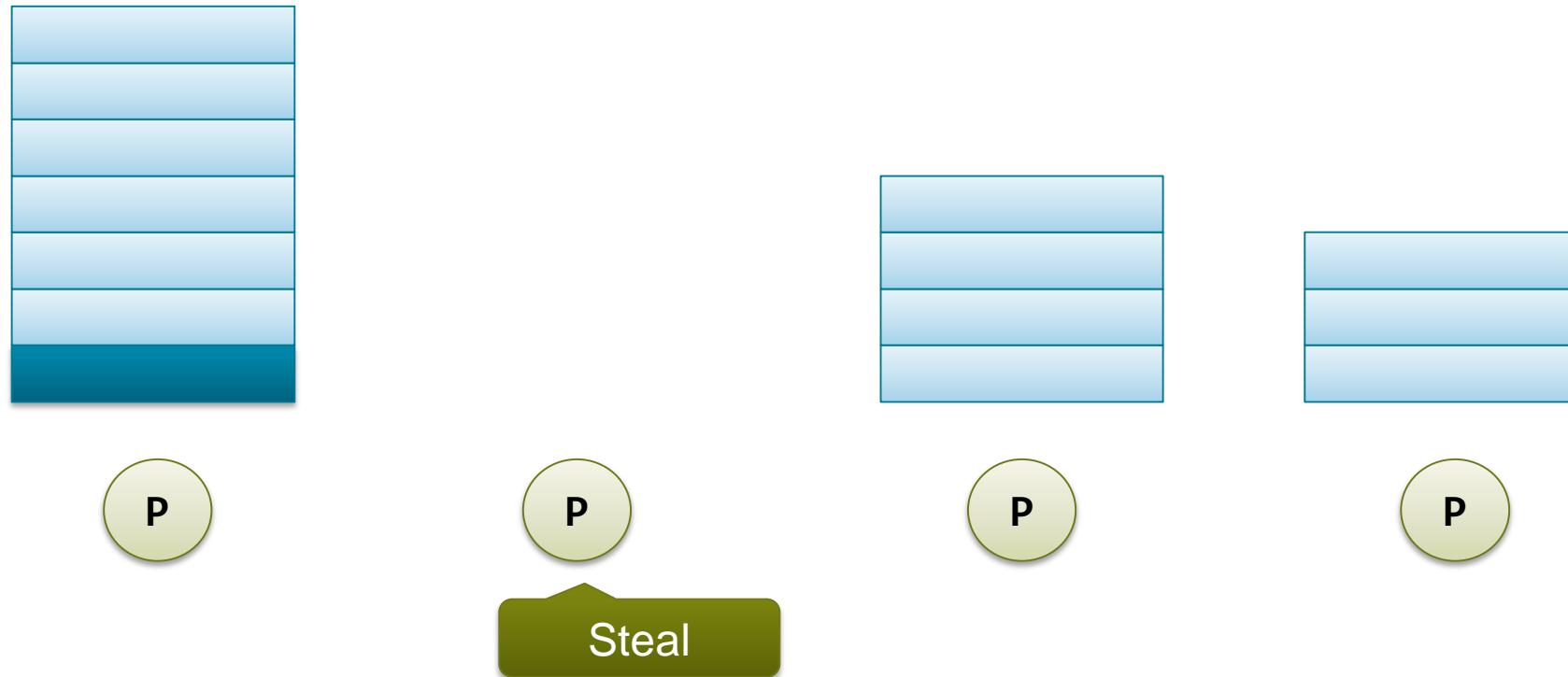
# Work Stealing Scheduler



# Work Stealing Scheduler

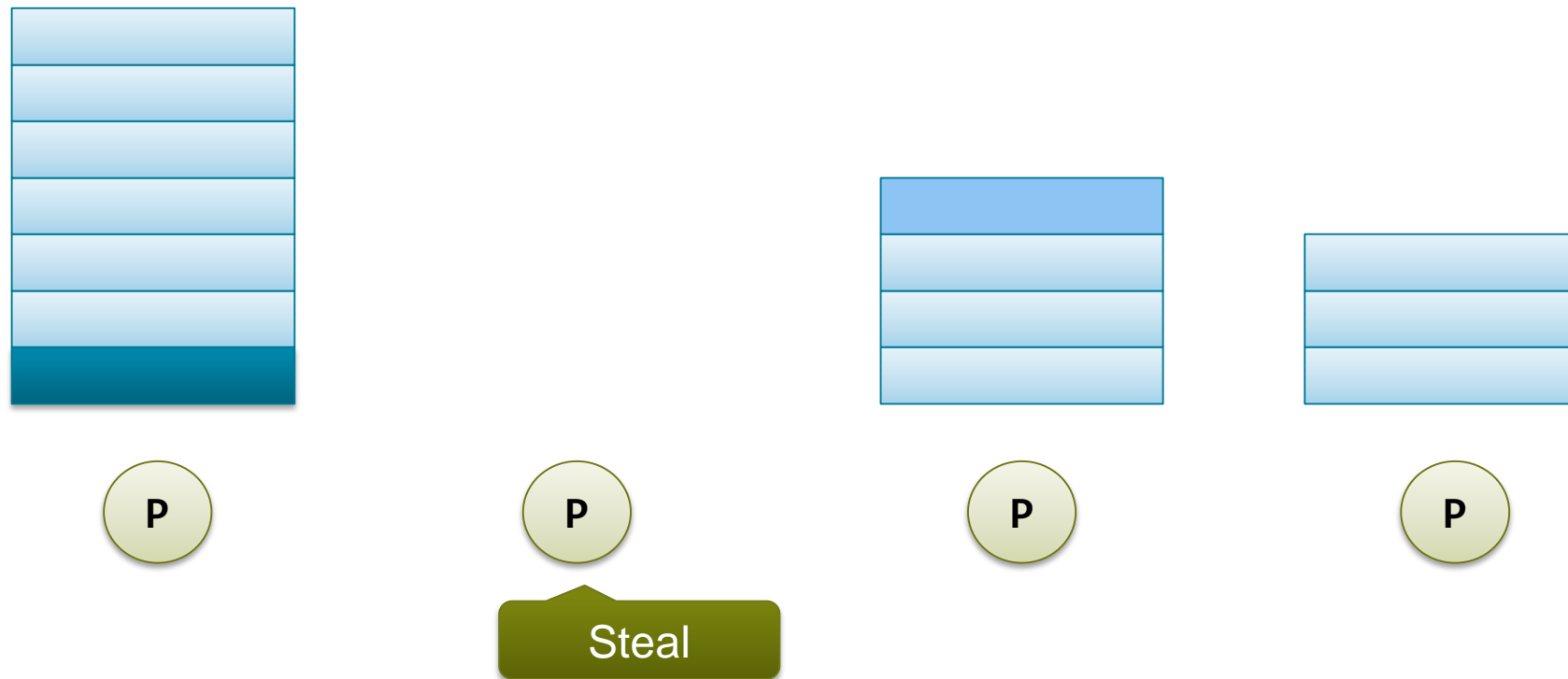


# Work Stealing Scheduler

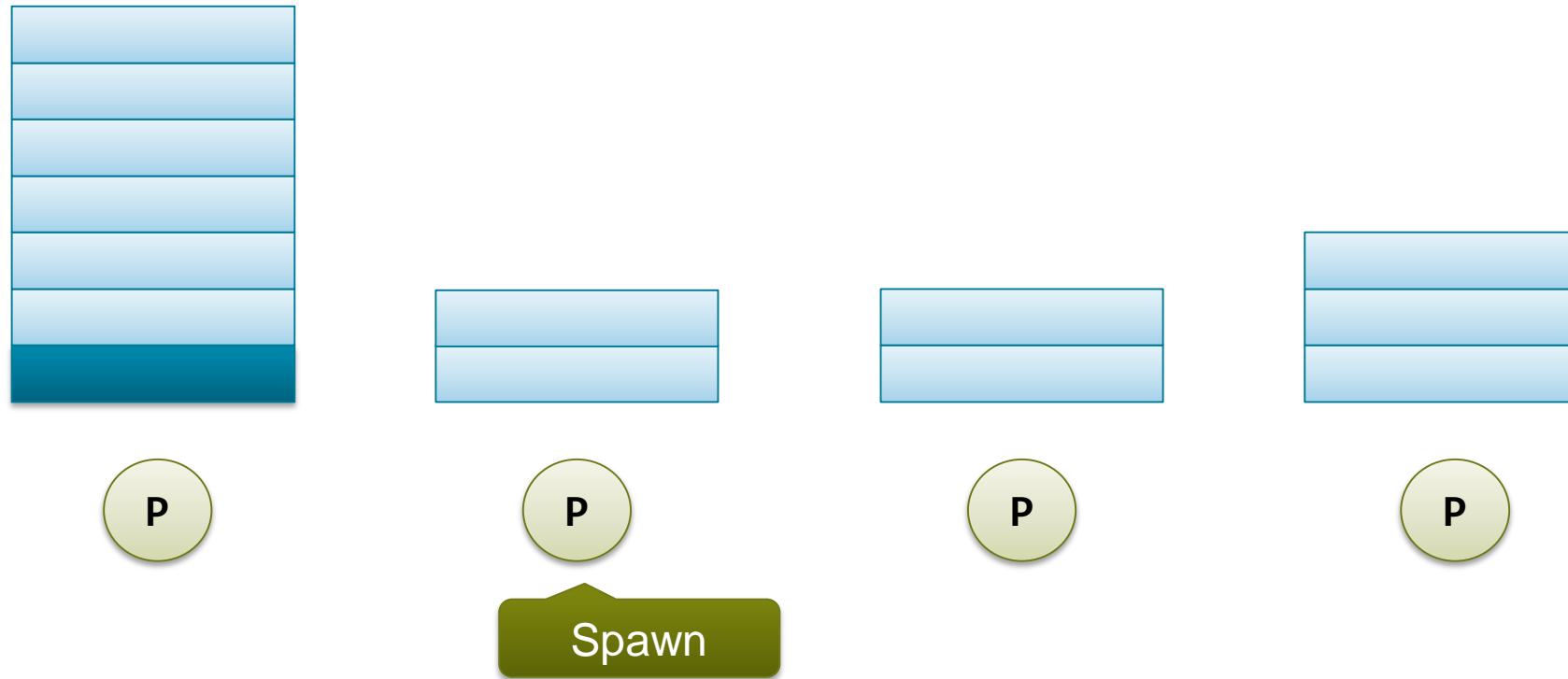


- **When a processor runs out of work, it steals a task from the top of a random victim's deque.**

# Work Stealing Scheduler



# Work Stealing Scheduler



# Work Stealing Scheduler

Each processor maintains a ready deque, bottom treated as stack

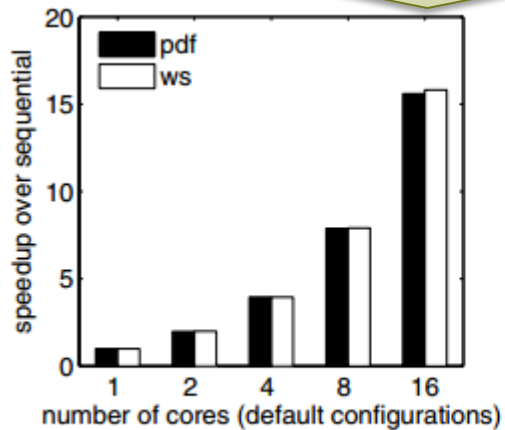
- **Initial:** Root thread in deque of a random processor
- **Deque not empty:**
  - Processor takes thread T from bottom and starts working
  - T spawns S: Put T on stack, continue with S
  - T stalls: Take next thread from stack
  - T dies: Take next thread from stack
  - If T enables a stalled thread S, S is put on the stack of T's processor
- **Deque empty:**
  - Steal thread from the top of a random (uniformly) processor's deque

# Parallel Depth First Scheduler

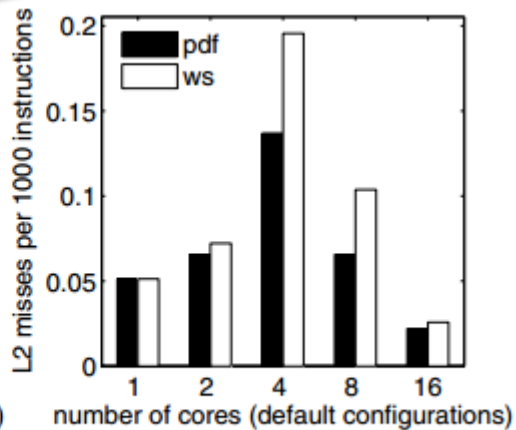
- **Based on the following insight:**
  - *Important (sequential) programs have already been highly tuned to get a good cache performance on a single score*
  - *Small working set*
  - *Good spatial and temporal reuse*

■ **Why the speedup is not that different? →** **Why the ready-to-execute task that the sequential**

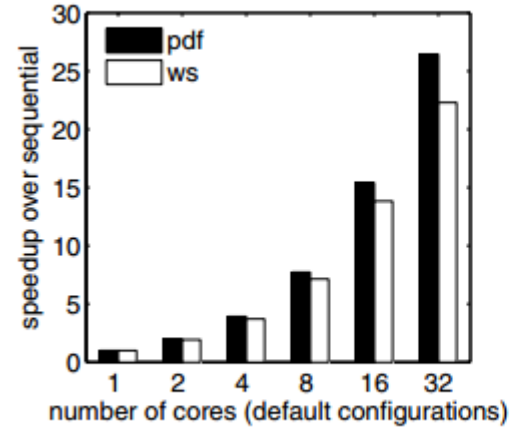
Why the speedup is not that different?  
Low miss/instruction ratio =>  
High Operational Intensity



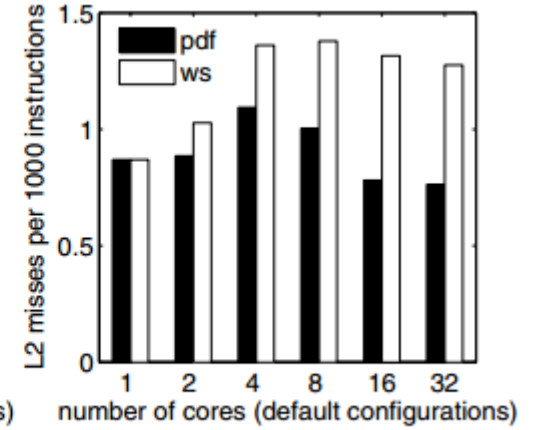
(a) LU



(b) LU



(e) Mergesort



(f) Mergesort