

Scheduling Multithreaded Computations By Work-Stealing

Robert D. Blumofe
The University of Texas, Austin

Charles E. Leiserson,
MIT Laboratory for Computer Science

Motivation

- Strict binding of multi-threaded computations on parallel computers.
- To find a parallel execution process by creating a directed acyclic graph and to traverse the instructions accordingly.

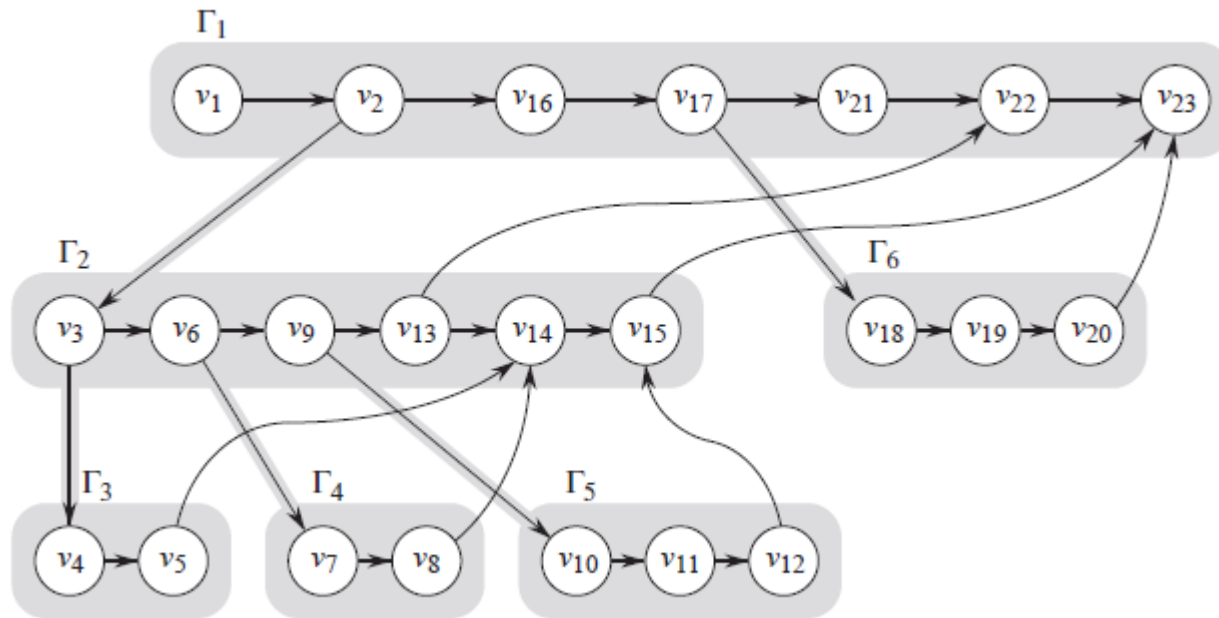
MIMD

- To maintain efficiency, enough threads must be active.
- Number of active threads should be within the hardware requirements.
- Related threads should be placed in the same processor.

Work Sharing vs Work Stealing

- Work Sharing:
 - Scheduler attempts to migrate threads to other processors hoping to distribute work to underutilized processors.
- Work Stealing:
 - Underutilized processors take the initiative by ‘steal’ing threads from other processors.

Multithreaded Computation



Blocks -> Threads

Circles -> Instructions

Right arrows -> Continue Threads

Curved Arrows -> Join Edges

Vertical/Slant arrows -> Spawn Edges

Strict vs Fully-Strict computation

- In a strict computation, all join edges from a thread go to ancestor of the thread in an activation tree.
- In a fully-strict computation all join edges from a thread go to the thread's parent.

Busy-Leaves

From the time thread T, is spawned until the time T dies, there is atleast one thread from the sub-computation that is ready.

Disadvantages:

- Not efficient in large environment of multiprocessors.
- Scalability

Busy-Leaves conditions

Spawn:

if T_a spawns T_b , T_a returns to thread pool. The processor starts next step with T_b .

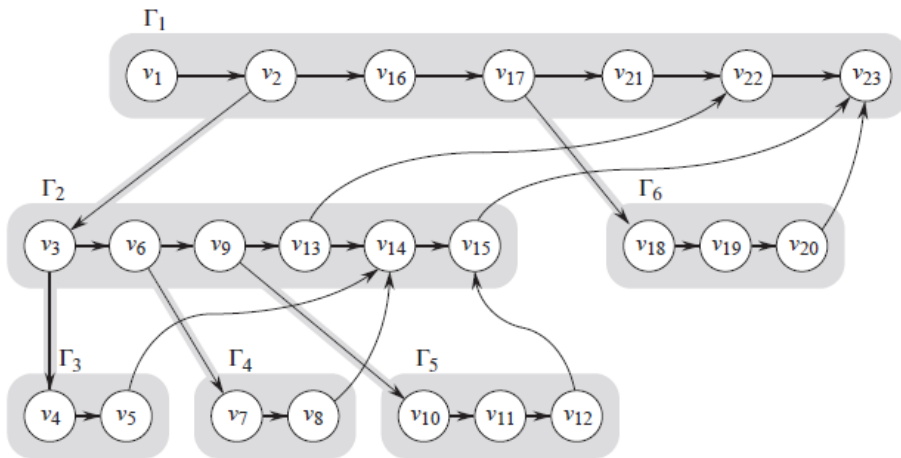
Stall:

if T_a stalls, it is returned to thread pool. Next step by the processor is idle.

Dies:

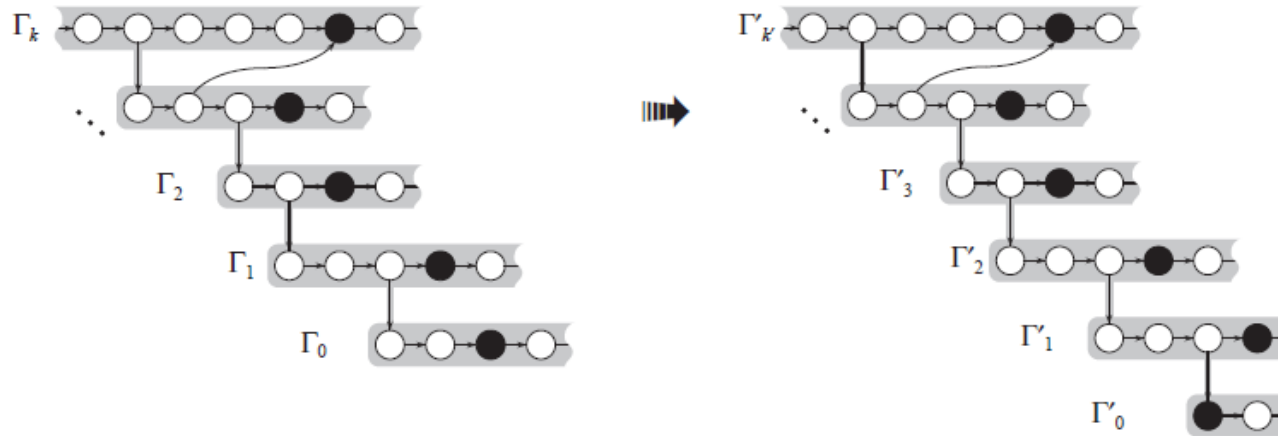
if T_a dies, processor checks if its parent T_b has any living children. If no children and no processor is working on T_b , it is taken from thread pool, else next step idle.

Busy-Leaves conditions



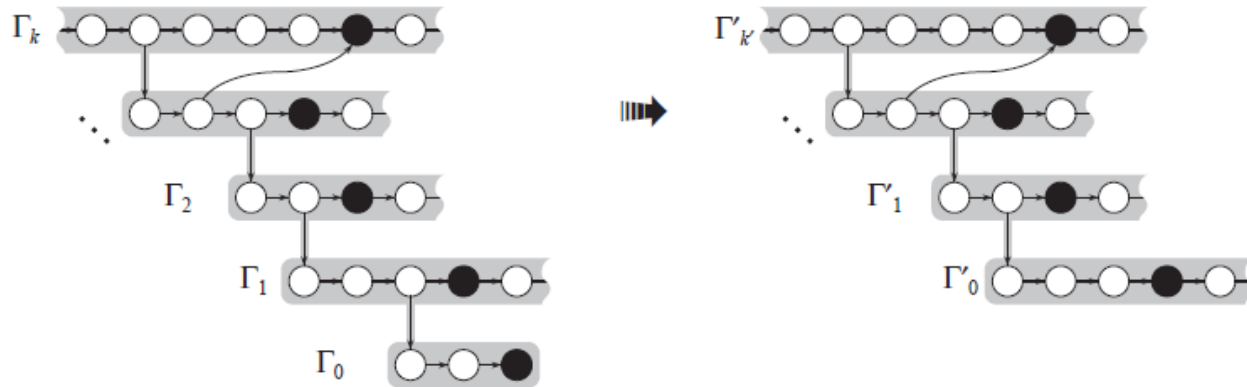
step	thread pool	processor activity	
		p_1	p_2
1		Γ_1 : v_1	
2		v_2	
3		Γ_2 : v_3	Γ_1 : v_{16}
4	Γ_2	Γ_3 : v_4	v_{17}
5	Γ_1 Γ_2	v_5	Γ_6 : v_{18}
6	Γ_1	Γ_2 : v_6	v_{19}
7	Γ_1	Γ_4 : v_7	v_{20}
8		v_8	Γ_1 : v_{21}
9	Γ_1	Γ_2 : v_9	
10	Γ_1	Γ_5 : v_{10}	Γ_2 : v_{13}
11	Γ_1	v_{11}	v_{14}
12		v_{12}	Γ_1 : v_{22}
13	Γ_1	Γ_2 : v_{15}	
14		Γ_1 : v_{23}	

Thread Spawn/Death



(a) Before spawn.

(b) After spawn.



(a) Before death.

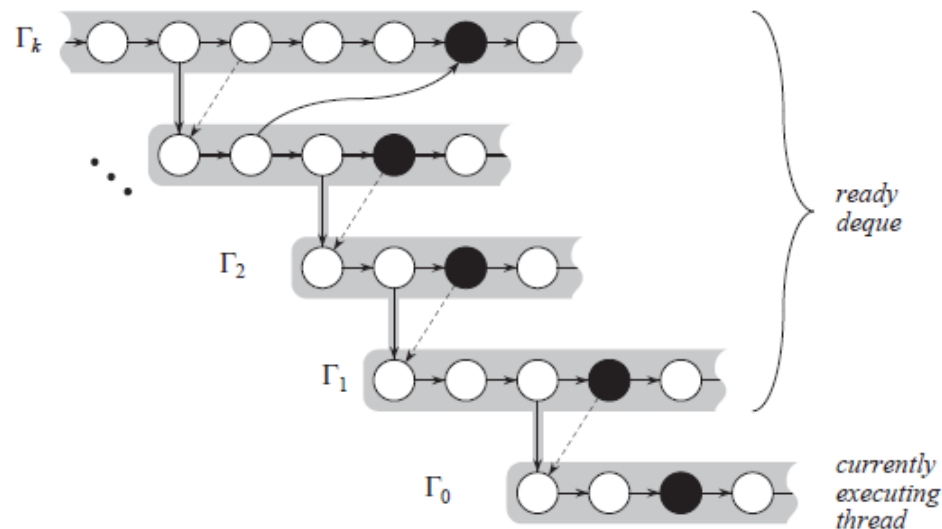
(b) After death.

Randomized Work-Stealing Algorithm

- The centralized pool of Busy-Leaves algorithm is distributed across the processors.
- Each processor maintains a ready deque.

Cases:

- Spawns
- Stalls
- Dies
- Enables



Randomized Work-Stealing Properties

- For a processor 'p', if there are 'k' threads in the deque, and if $k > 0$, the following properties are satisfied.
 - (1) For $i=1,2,3,\dots,k$, thread T_i is parent of T_{i-1}
 - (2) If $k > 1$, for $i=1,2,\dots,k-1$, thread T_i has not been worked on since it spawned T_{i-1}

Work-Stealing

- For a fully-strict computation with work ‘ T_1 ’ and critical path length T_∞ the expected running time with P processors is $T_1/P + O(T_\infty)$
- Execution time on P processors is
$$T_1/P + O(T_\infty + \lg P + \lg(1/P^2))$$
- Expected total communication is
$$O(PT_\infty(1+n_d)S_{\max})$$

Recycling Game for Atomic Access

- (P,M) recycling game:
 - P = number of balls = number of bins
 - M = total number of ball tosses
 - Adversary chooses some of the balls from reservoir selects one of the P bins randomly.
 - Adversary inspects each of the P bins that contains at least 1 ball and removes it from the bin.

Work Stealing

- Total delay incurred by M random requests made by P processors is $O(M + P \lg P + P \log(1/2))$

Conclusion

- The proposed work-stealing scheduling algorithm is efficient in terms of time, space and communication.
- ‘C’ base language called ‘Cilk’ being developed for programming multithreaded computations based on work-stealing.

<http://supertech.lcs.mit.edu/cilk>