

## **Lace: Non-Blocking Split Deque for Work-Stealing**

Tom van Dijk & Jaco van de Pol

FMV/Parallel Computing, a sunny morning in 2017

## Background

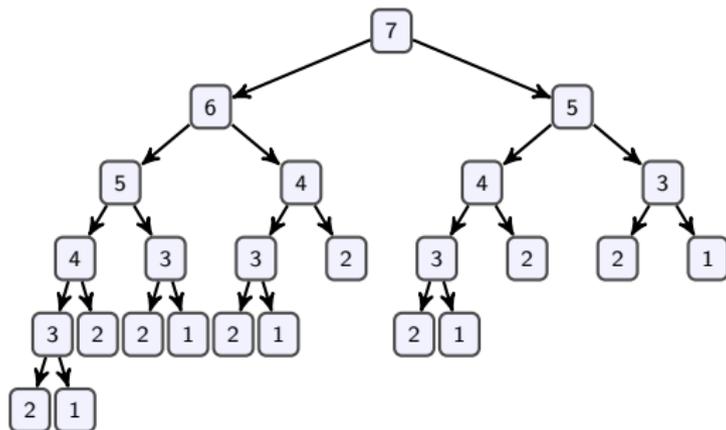
- ▶ PhD at Formal Methods & Tools, University of Twente
- ▶ PhD Research: Parallel Binary Decision Diagrams
  - ▶ Using work-stealing...
  - ▶ ...and lock-free hash tables
  - ▶ to implement **Sylvan** and **Lace**.
- ▶ Current research interests
  - ▶ Parallel Satisfiability
  - ▶ Using ZBDDs to store clause sets for Satisfiability
  - ▶ Solving Parity Games via Priority Promotion

## What to do as a **good** student?

- ▶ I want you to understand each slide.
- ▶ Ask me why I made certain choices.
- ▶ Ask me how to find performance problems.
- ▶ Ask me how to fine-tune the implementation.
- ▶ Ask me about the relation between shared-memory and message passing.
- ▶ Ask me why I think we cannot go much faster than this.

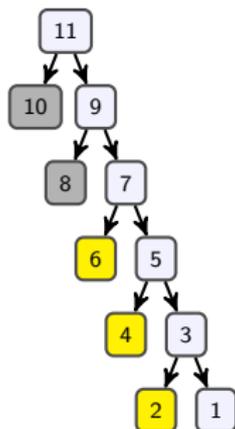
# Task parallelism

```
1 def fib(k):  
2   if k < 2 : return k  
3   spawn fib(k-2)  
4   spawn fib(k-1)  
5   n ← sync  
6   m ← sync  
7   return n + m
```



# Example: calculate `fib(11)`

Task graph:



Task deque (of first worker):



# Work-stealing related to its deque

Work-stealing operations	Deque operations
<code>spawn(task)</code>	<code>push(task)</code>
<code>sync</code>	<code>peek, pop</code>
<code>steal-and-run(victim)</code>	<code>steal</code>

- ▶ Each worker has 1 deque.
- ▶ Worker uses `push/peek/pop` on its own deque.
- ▶ Worker uses `steal` on other deques.
- ▶ Policy: steal from the thief.

## Implementations (blue = non-blocking)

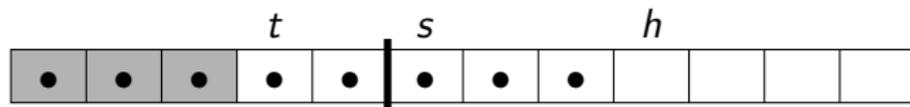
- ▶ Fully shared deque: Frigo ea ('Cilk' 1998), ABP (1998), Chase and Lev (2005), Hendler ea (2006)
- ▶ Private deque: Acar ea (2013)
- ▶ Split deque: Faxén ('Wool' 2008, 2010), Dinan ea (2009)
- ▶ Non-blocking split deque: Van Dijk & Van de Pol (2013)

## Challenges

- ▶ Avoid hidden and unnecessary communication
  - ▶ false sharing (variables accessed by thieves / owner)
  - ▶ unnecessary memory writes *and* reads
- ▶ Avoid using locks/mutexes
  - ▶ (solved using lock-free operations)
- ▶ Avoid expensive memory fences, e.g., Cilk-THE
  - ▶ (mostly solved using split principle)
- ▶ Avoid overhead, especially since most tasks are never stolen
  - ▶ (solved with “direct task stack”)

# Deque in Lace

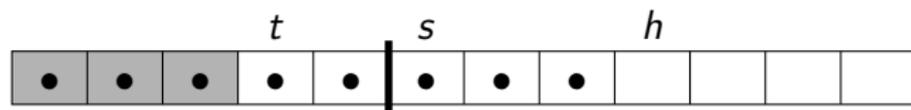
Deque is described by variables tail ( $t$ ), split ( $s$ ), head ( $h$ ).



- ▶ Tasks are **shared** or **private**.
- ▶ The first  $t$  tasks are **stolen**.
- ▶ Tasks steal by **atomic cas** on  $t$  and  $s$  together.
- ▶ Owner modifies  $h$  and  $s$  with normal memory operations.
- ▶ Extra flag: *movesplit*.

# Deque in Lace

Deque is described by variables tail ( $t$ ), split ( $s$ ), head ( $h$ ).

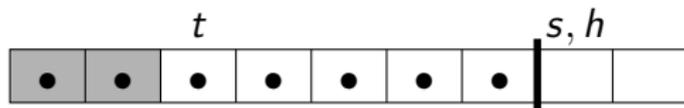


## Communication is key!!

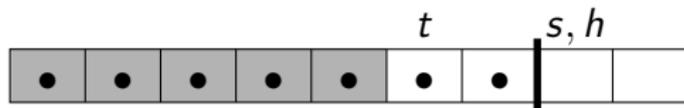
Cacheline	Contents	Thief access	Owner access
Shared 1	tail, split	Often	Sometimes
Shared 2	flag movesplit	Sometimes	Often
Private	head, osplit	–	Often

# Deque in Lace

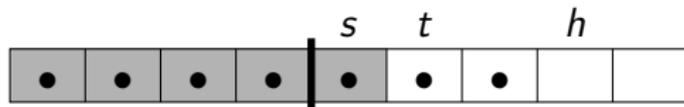
## Moving the split point back



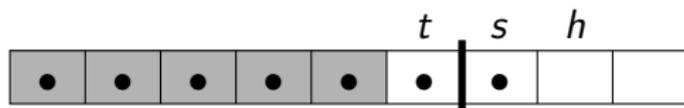
owner reads  $t, s$



thieves steal



owner sets  $s$



owner repairs  $s$

## Benchmarks

- ▶ fib(50) – 20,365,011,073 tasks
- ▶ uts(T3L) – Unbalanced Tree Search, 111,345,630 tasks
- ▶ queens(15) – 171,129,071 tasks
- ▶ matmul(4096) – 3,595,117 tasks
- ▶ No cut-off point, fine-grained, very small tasks.

## Measurements

- ▶ 48-core AMD machine (4 sockets, 12 cores per socket)
- ▶ Wallclock time around parallel part, 1, 48 workers.

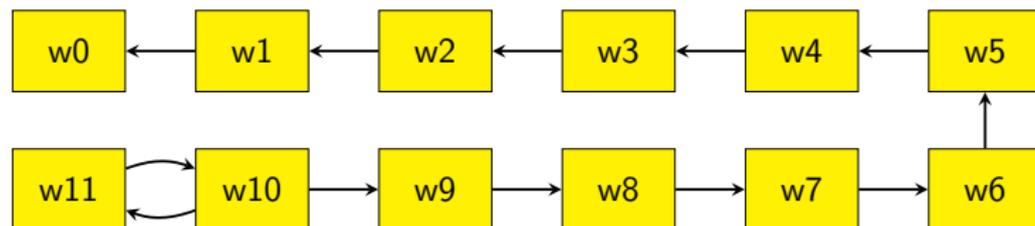
# Experimental results

Results	Benchmark time			Speedup	
	$T_S$	$T_1$	$T_{48}$	$T_S/T_{48}$	$T_1/T_{48}$
fib 50	149.2	144	4.13	34.5	34.9
uts T2L	84.5	86.0	1.81	46.1	47.4
uts T3L	43.11	44.2	2.23	18.7	19.9
uts T3L *	43.11	44.26	1.154	37.4	38.3
queens 15	533	602	12.63	42.2	47.7
matmul 4096	773	781	16.46	47.0	47.5

\* = with extension to fix issues with leapfrogging (next slides)

## Leapfrogging

- ▶ Waiting for stolen work? Steal from thief!
- ▶ Advantage: gives nice upper bound on deque size!
- ▶ Disadvantage: steal chaining...



- ▶ Work does not trickle down fast enough!

## Conclusions

- ▶ Non-blocking split deque has low overhead and good speedup
- ▶ Leapfrogging plus random stealing solves steal chaining
- ▶ Only require memory fence to shrink the shared portion
- ▶ Lace can be found at:
  - ▶ <http://github.com/trolando/lace>
  - ▶ Feel free to reproduce results (bench.py)
- ▶ Lace is used in our parallel BDD implementation Sylvan

# Algorithm outline

```
1 def steal():
2   if allstolen : return None
3   t, s  $\leftarrow$  (tail, split)
4   if t < s :
5     if cas ((tail,split), (t,s), (t+1,s)) : return Task(t)
6     else: return None
7   elif  $\neg$  movesplit : movesplit  $\leftarrow$  true
8   return None
```

# Algorithm outline

```
9 def push (task) :
10   if head = size : raise QueueFull
11   write task data at head
12   head  $\leftarrow$  head + 1
13   if oallstolen :
14     (tail,split)  $\leftarrow$  (head-1,head)
15     osplit  $\leftarrow$  head
16     allstolen  $\leftarrow$  false
17     oallstolen  $\leftarrow$  false
18     if movesplit : movesplit  $\leftarrow$  false
19   elif movesplit :
20     // Grow shared portion
21     new_split  $\leftarrow$  (osplit + head + 1) / 2
22     split  $\leftarrow$  new_split
23     osplit  $\leftarrow$  new_split
24     movesplit  $\leftarrow$  false
```

# Algorithm outline

```
24 def pop():  
25     head ← head - 1  
  
26 def pop-stolen():  
27     head ← head - 1  
28     if ¬ oallstolen :  
29         allstolen ← true  
30         oallstolen ← true
```

# Algorithm outline

```
31 def peek():
32     if head=0 : raise QueueEmpty
33     if oallstolen : return Stolen(head-1)
34     if osplit = head :
35         if  $\neg$  shrink-shared() :
36             allstolen  $\leftarrow$  true
37             oallstolen  $\leftarrow$  true
38             return Stolen(head-1)
39     if movesplit :
40         // Grow public section (excluding head-1)
41         new_split  $\leftarrow$  (osplit + head) / 2
42         split  $\leftarrow$  new_split
43         osplit  $\leftarrow$  new_split
44         movesplit  $\leftarrow$  false
45     return Work(head-1)
```

# Algorithm outline

```
45 def shrink-shared():
46     t, s ← (tail, split)
47     if t = s : return false
48     new_s ← (t + s) / 2
49     split ← new_s
50     osplit ← new_s
51     memory fence
52     t ← tail
53     if t = s : return false
54     if t > new_s :
55         new_s ← (t + s) / 2
56         split ← new_s
57         osplit ← new_s
58     return true
```

# Informal proof

