



Roll up and Conquer

John Feo
Cray Inc

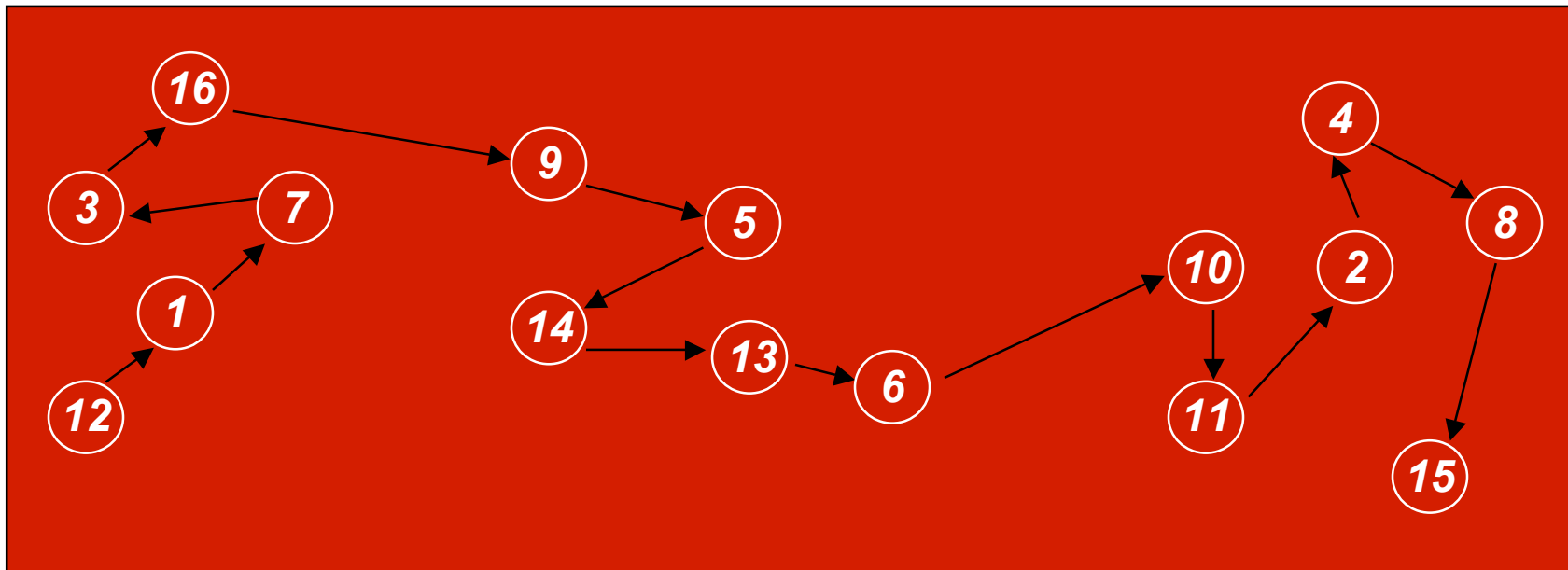
CUG 2006
Lugano, Switzerland

- 0 Divide-and-conquer is a well known technique that divides a large problem into smaller problems, solves the small problems, and then reduces their solutions
- 0 An opposite technique appears to work well for many graph problems
 - 0 roll up the graph into super nodes
 - 0 solve the problem for the graph of super nodes
 - 0 extend the solution of each super node to its included nodes

Prefix operations on lists

- 0 $P(1)$ = initial value
 $P(i) = P(i - 1) \oplus V(i), \quad i > 1$
- 0 If initial value is 1, operator is +, and $V(i)$ is 1 for all i , then $P(i)$ is the rank of node i
- 0 **List ranking** is a common procedure that occurs in many graph algorithms and is considered the HPL of graph algorithms

A linked list



	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
List	7	4	16	8	14	10	3	15	5	11	2	1	6	13	0	9
Rank	2	13	4	14	7	10	3	15	6	11	12	1	9	8	16	5

Sequential algorithm



0 Find start of list

$$\frac{n^2 + n}{2} - \sum_{i=1}^n List(i)$$

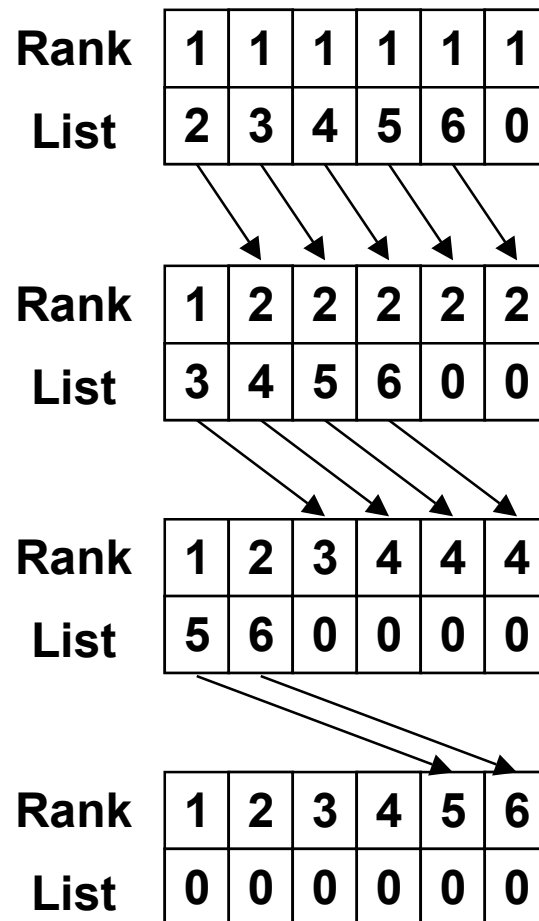
0 Walk list from start to end

0 $O(n)$ in number of instructions and time

Parallel algorithm (cyclic reduction)

- 0 **Rank**(i) = 1, $1 \leq i \leq n$
- 0 **Rank**(**List**(i)) += **Rank**(i), $1 \leq i \leq n$, **List**(i) \neq 0
- 0 **List**(i) = **List**(**List**(i)), $1 \leq i \leq n$, **List**(i) \neq 0
- 0 Repeat steps 2 and 3 until **List**(i) = 0, $\forall i$

Steps 2 and 3



$O(n \log n)$ in number of instructions

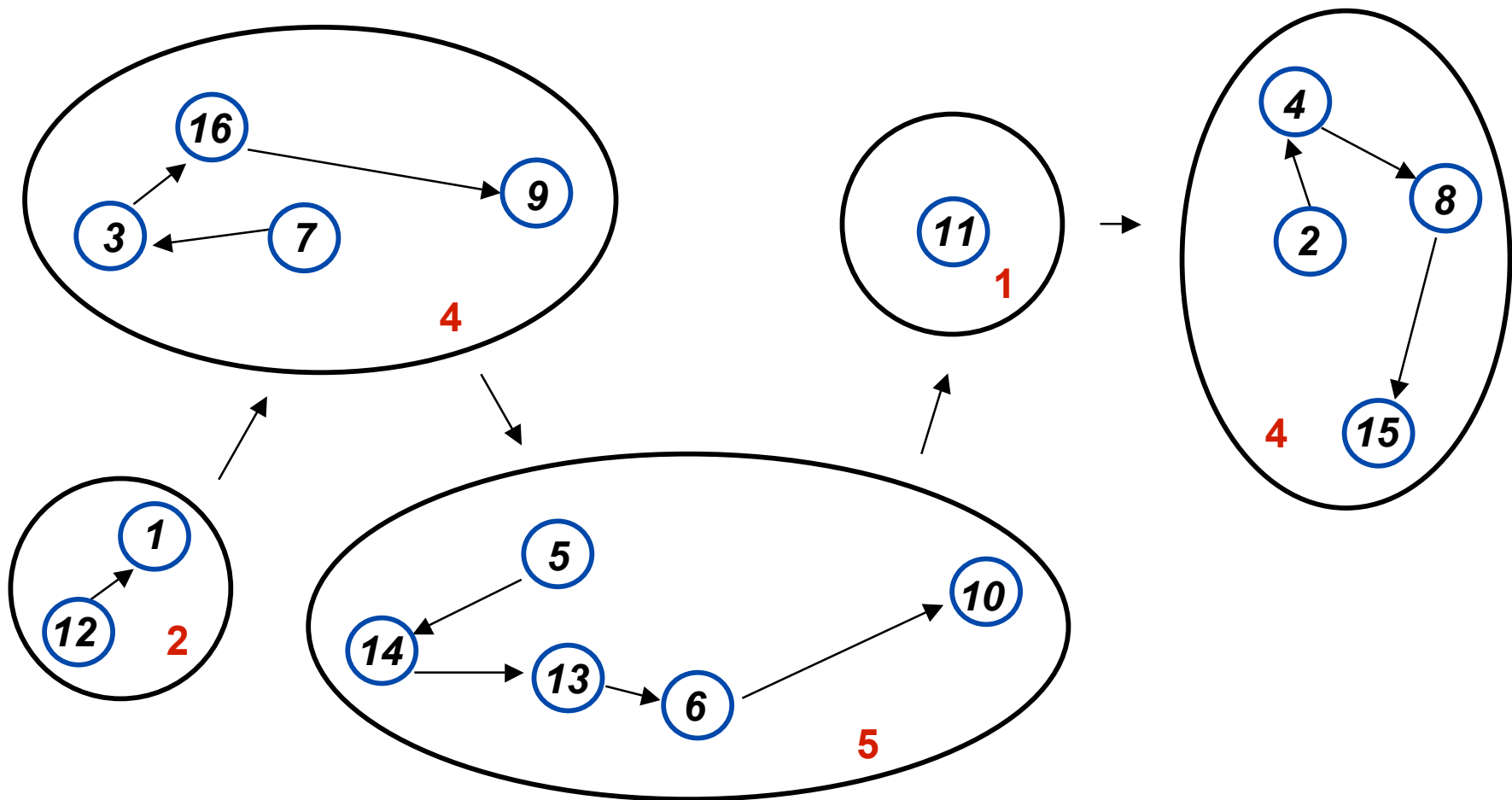
$O\left(\frac{n \log n}{P}\right)$ in time

We can do better ...

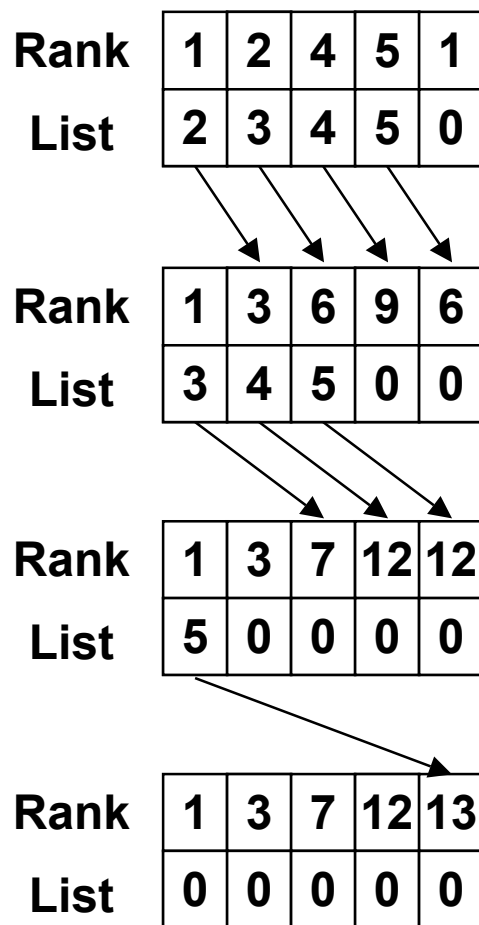


- 0 By **rolling up** the graph into super nodes we can reduce both the number of operations and time to $O(n)$
- 0 First, run the sequential algorithm to roll up the graph, and then run the parallel algorithm on the rolled up graph; finally, run the sequential algorithm on each super node to compute the rank of each original node
- 0 $O(n) + O(S \log S) + O(n)$, where S is the number of super nodes
- 0 A variation of the Hellman & JaJa algorithm for list ranking and the algorithm used by the MTA compiler to parallelize recurrences

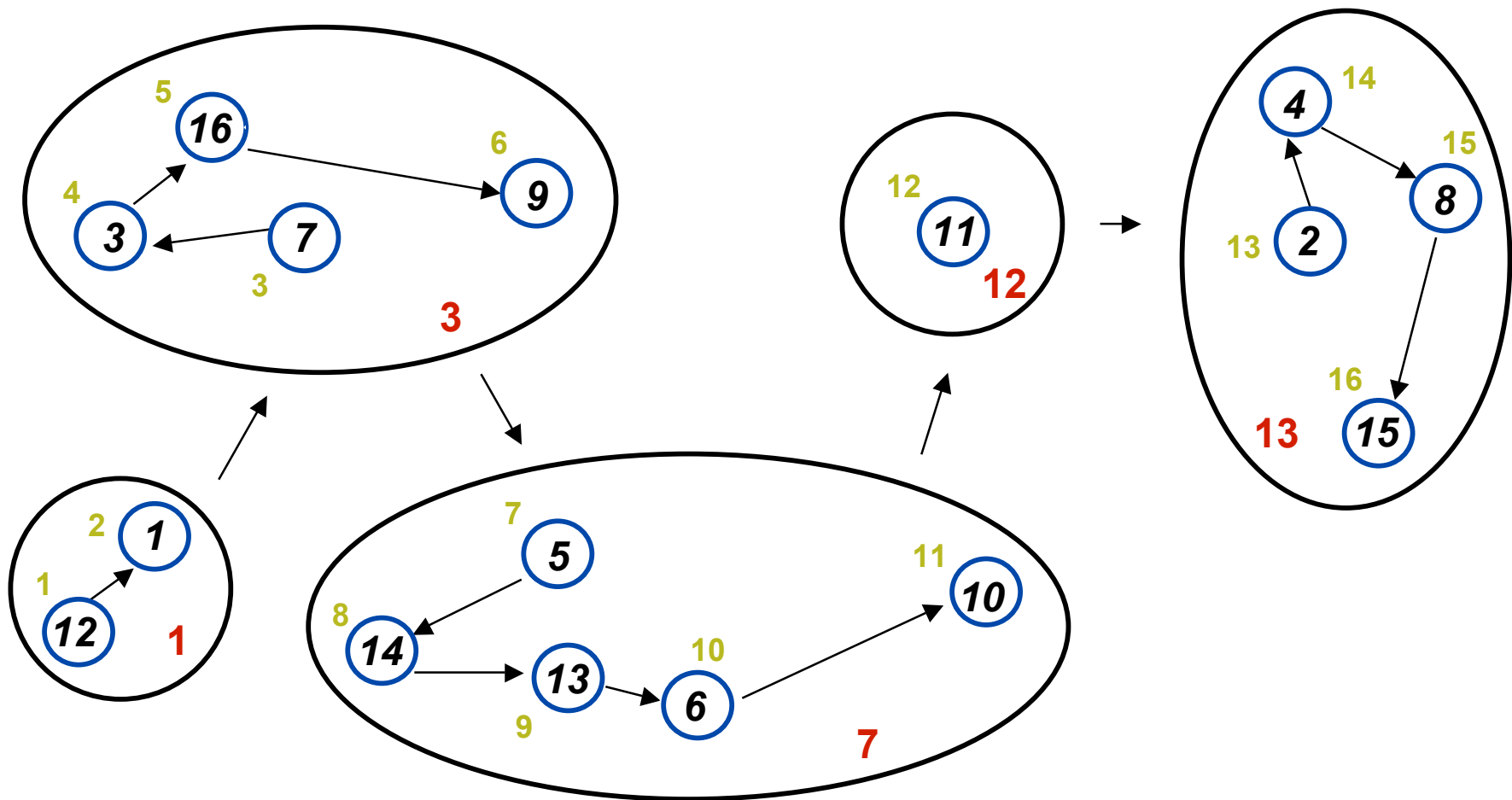
First step



Second step



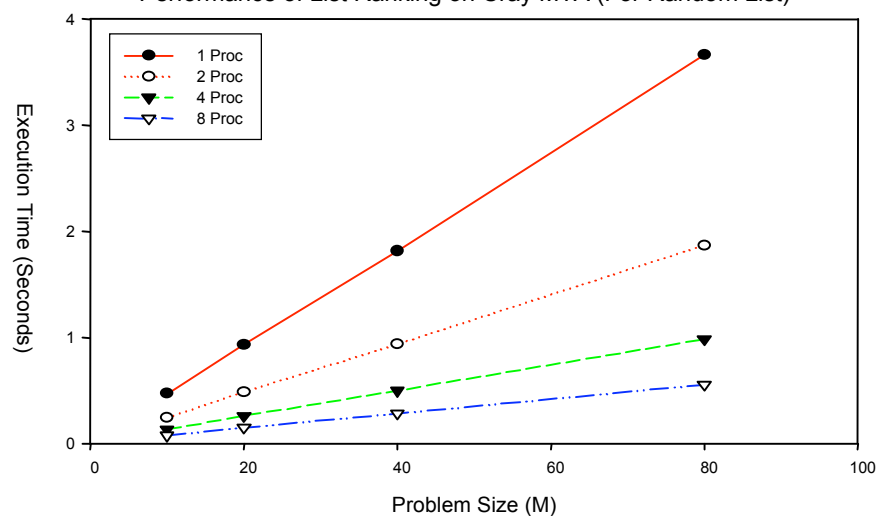
Third step



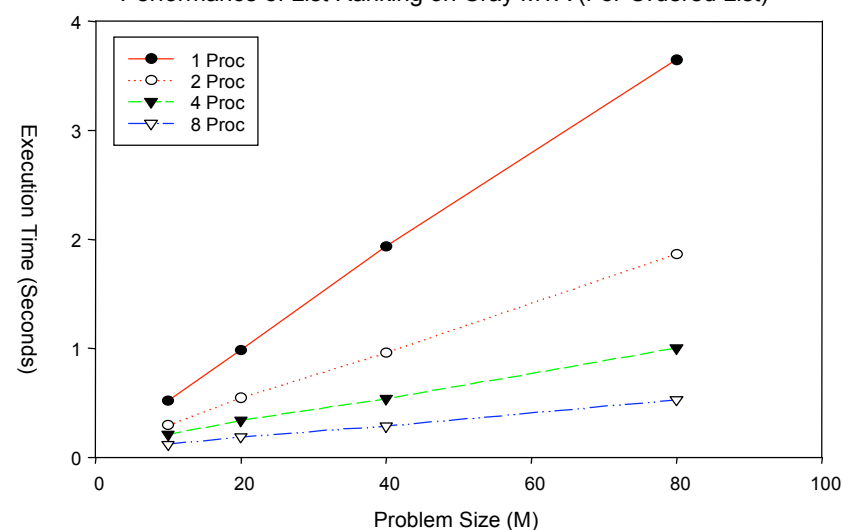
Performance comparison



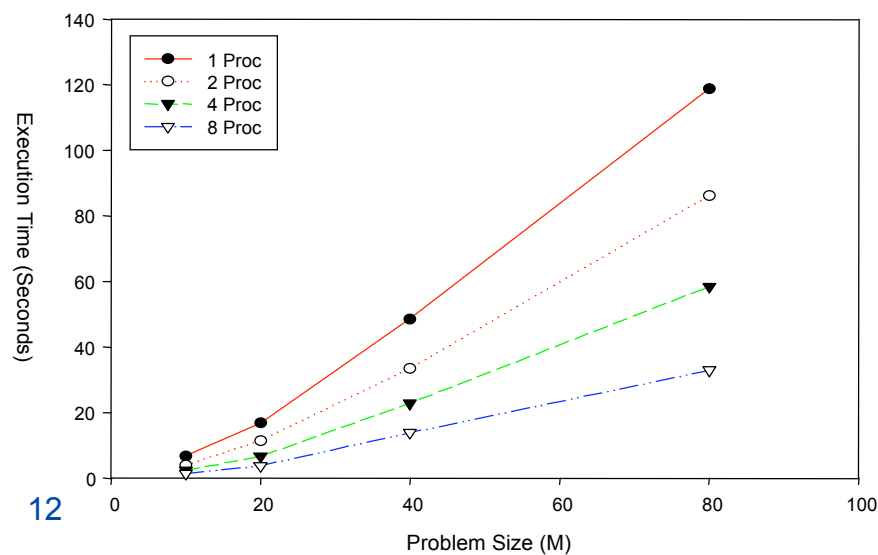
Performance of List Ranking on Cray MTA (For Random List)



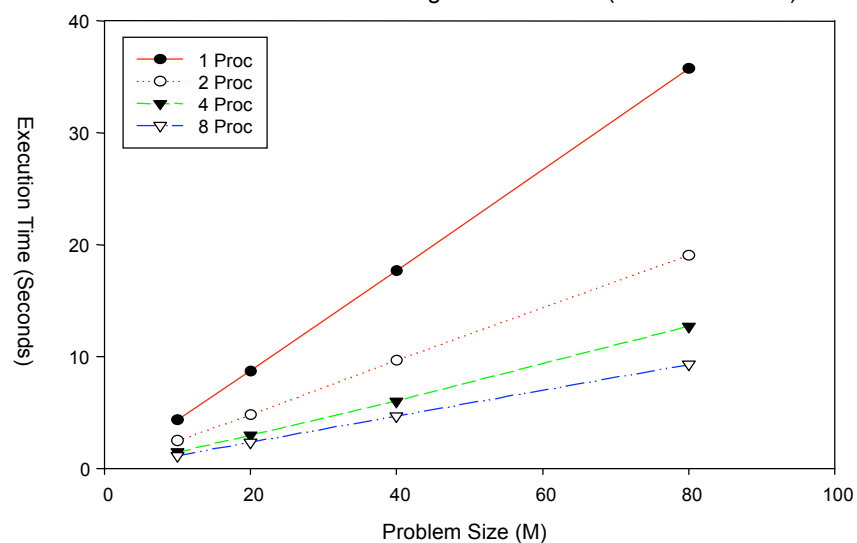
Performance of List Ranking on Cray MTA (For Ordered List)



Performance of List Ranking on SUN E4500 (For Random List)

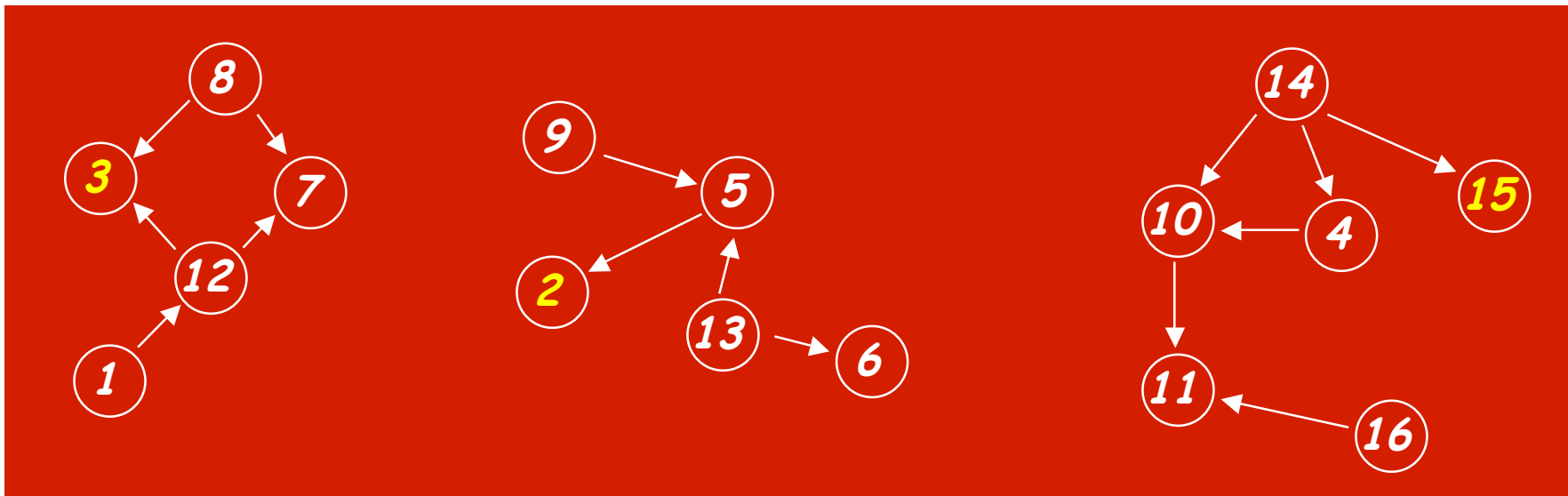


Performance of List Ranking on Sun E4500 (For Ordered List)



Connected Components

- 0 Label all nodes in a graph such that $\text{Label}[v] = \text{Label}[w]$ if and only if there is an undirected path between v and w



	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
Label	3	2	3	15	2	2	3	3	2	15	15	3	2	15	15	15

Sequential algorithms

- 0 Typically, based on either depth-first or breath-first search

```
for all nodes v  
    Label[v] = v  
for all nodes v  
    if v is unvisited, DFS(v)
```

where DFS(v) is

```
for all unvisited neighbors w of v  
    Label[w] = Label[v]  
    DFS(w)
```

- 0 If we parallelize the loop in the main program, multiple labels move through a component \Rightarrow to merge or clean up is expensive
- 0 If we parallelize the loop in DFS, parallelism is limited for nodes with small degree

Parallel algorithm

0 CRCW PRAM algorithm (Shiloach-Vishkin)

until no label changes

```
for all edges (v, w)           // hook components together  
  if Label[v] < Label[w] && Label[v] == Label[Label[v]]  
    Label[Label[v]] = Label[w]
```

```
for all nodes v               // compress components  
  if Label[v] != Label[Label[v]]  
    Label[v] = Label[Label[v]]
```

0 Parallel and performance is good, but

0 Outer loop may iterate many times

0 if $\text{Label}[v] = C$ for many v , then $\text{Label}[\text{Label}[v]]$ gets hot

Hybrid approach

// STEP 1, DFS w/o cleanup, rolls up the graph

```
for all nodes v
    if v is unvisited, mark v as a root and DFS(v)
```

← parallel loop

where DFS(v) is

```
for all unvisited neighbors w of v
    Label[w] = Label[v]
    DFS(w)
for all visited neighbors w of v
    store (Label[v], Label[w]) uniquely in a hash table
```

← parallel loop

← recursive call

← parallel loop

// STEP 2, run PRAM on rolled up graph

```
repeat until no label changes
    for all edges (v, w) in the hash table
        if Label[v] < Label[w] && Label[v] = Label[Label[v]]
            Label[Label[v]] = Label[w]
    for all roots v
        if Label[v] != Label[Label[v]]
            Label[v] = Label[Label[v]]
```

← parallel loop

← parallel loop

// STEP 3, relabel

```
for all roots v
    if v is unvisited, mark v as a root and relabel(v)
```

← parallel loop

where relabel(v) is

```
for all unvisited neighbors w of v,
    Label[w] = Label[v] and relabel(w)
```

← parallel loop

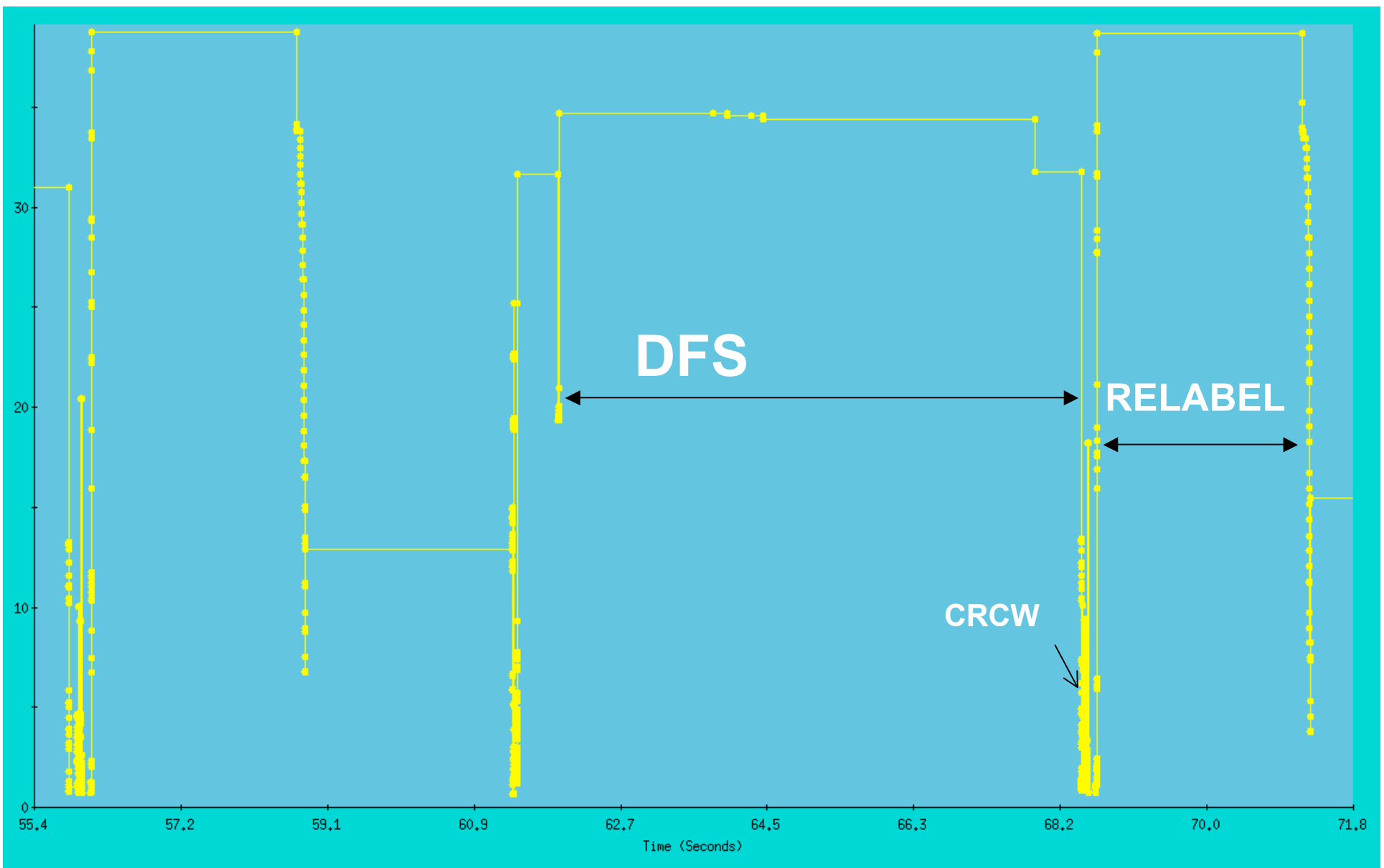
← recursive call

- 0 Large, pseudo-random graph
 - 0 67M nodes, 420M edges, 412 components
 - 0 64 nodes have degree $\sim 2^{20}$
 - 0 64K nodes have degree $\sim 2^{10}$
 - 0 the rest have degree ~ 12

P	Time
20	19.4
30	13.1
40	9.8

**That's about 1.25M edges per second per processor
!!!**

Traceview



Conclusions

- 0 “Roll up and conquer” is a programming technique well-suited for parallel graph algorithms
- 0 We have used the technique to develop high-performance, scalable parallel algorithms for several graph problems
- 0 The MTA’s **shared-memory, latency-tolerant processors, and full-and-empty bits** are critical for good performance