

# Parallel Breadth First Search

Rezaul Alam Chowdhury  
Dept. of Computer Science  
Stony Brook University  
Stony Brook, NY 11794-4400  
rezaul@cs.stonybrook.edu

Yonatan Fogel  
Dept. of Computer Science  
Stony Brook University  
Stony Brook, NY 11794-4400  
yfogel@cs.stonybrook.edu

Jesmin Jahan Tithi  
Dept. of Computer Science  
Stony Brook University  
Stony Brook, NY 11794-4400  
jtithi@cs.stonybrook.edu

**Abstract**—We have developed a multithreaded implementation of breadth-first search (BFS) of arbitrary graphs using the Cilk++ extensions to C++. We analyze our FBFS algorithm under the CILK model, where synchronizing  $n$  tasks takes  $\mathcal{O}(\log n)$  time. We introduce a lower bound for  $T_P$  in level-synchronous BFS. FBFS achieves high work-efficiency by applying parallel prefix sum to split up all work at a given distance from a source vertex.

FBFS algorithm deterministically matches the lower bound for  $T_P$ , is asymptotically optimal for  $T_1$ , and is scale-free. FBFS saves energy by letting processors become idle whenever possible: total non-idle work among all  $P$  processors is  $\mathcal{O}(T_1)$ . Our algorithm is non-deterministic in that it contains benign races which may affect performance (by up to a constant factor) but not its correctness. These races can be fixed with mutual-exclusion locks which will slow down FBFS empirically.

In particular, we show that for a graph  $G = (V, E)$  with diameter  $D$ , FBFS runs in time  $\mathcal{O}((|V| + |E|)/P + D \log P)$  on  $P$  processors, which means that it attains near-perfect linear speedup if  $P \ll (|V| + |E|)/(D \log(|V| + |E|))$ .

## I. INTRODUCTION

### A. Subsection Heading Here

Subsection text here.

1) Subsubsection Heading Here: Subsubsection text here.

## II. RELATED WORK

### A. Terms

- $G = (V, E)$  is the graph.
- $V$  is the set of nodes.
- $n = |V|$ .
- $m = |E|$ .
- $\Gamma(u)$  is the set of nodes adjacent to node  $u$ .
- $T_s$  is the running time for the most efficient serial algorithm.
- $T_1 \equiv \text{WORK}$  is the running time for a parallel algorithm running on one processing elements.
- $T_P$  is the running time for a parallel algorithm running on  $P$  processing elements.
- $T_\infty \equiv \text{SPAN}$  is the critical path for a parallel algorithm or the time it takes given infinite processing elements.

- $W_P \leq pT_P$  is the total amount of work for  $P$  processing elements. Idle processing elements do not count towards  $W_P$ .

We consider only *level-synchronous* BFS algorithms, which find all nodes at distance  $0 \leq d \leq |V|$  before any nodes at distance  $d' > d$ .

The standard serial approach for level-synchronous BFS can be seen in Fig. 1.

The bottlenecks for parallelization are the FIFO queue and the DIST array. The FIFO is fast but is inherently serial. When running in parallel, there is a benign race condition on lines 11–14 of Fig. 1. Multiple threads can enqueue the same vertex. The race is benign and rare: it can create extra work but does not affect correctness. This race is sometimes dealt with mutual exclusion, atomic instructions, or by simply performing the extra work when the race occurs.

### Existing algorithms

**MIT-Bag** [1] uses reducer hyperobjects and a new *bag* data structure to replace the FIFO. This algorithm relies on a work-stealing scheduler and runs in  $T_P = \mathcal{O}(n/P + m/P + D \lg^3(n/D))$  w.h.p in  $n$ . MIT-Bag uses mutual exclusion to deal with the race conditions on DIST for analysis, but in practice recommends simply doing the extra work.

**CRAY-BFS** [2] uses special hardware available on the Cray MTA-2 platform. It writes to the FIFO in parallel by using atomic increments, and uses special hardware mutexes for every 64-bit word.

**Block-BFS** [3] uses a block-accessed queue to replace the FIFO. It uses atomic primitives to allocate blocks from the FIFO that are small, but not so small so that they do not use atomics too often.

**Distinguished-BFS** [4] is a non level-synchronous BFS algorithm for arbitrary sparse graphs. Distinguished-BFS runs in  $T_P = \mathcal{O}(n^\epsilon)$  time with  $\mathcal{O}(mn^{1-2\epsilon})$  processors, if  $m \geq n^{2-3\epsilon}$ . The basic idea is to contract the graph to distinguished and then superdistinguished vertexes, at which point the graph will be dense.

## III. NEW PARALLEL BFS

### A. Lower Bound

We establish a lower bound for the parallel running time of level-synchronous BFS. Let  $D$  be the diameter of a graph  $G = (V, E)$ ,  $n = |V|$ ,  $m = |E|$ .

For level-synchronous BFS,  $T_s = \mathcal{O}(n + m + D) = \mathcal{O}(n + m)$  is a lower bound for both  $T_1$  and  $W_P$ . Let  $n_l$  be the number of nodes at distance  $l$  from the source vertex, and  $m_l$  be the sum of out-degrees of the  $n_l$  nodes at distance  $l$ . To establish the lower bound for  $T_P$ , consider a graph where for every level  $l$ ,  $(n_l + m_l) = \mathcal{O}(P)$ . There is therefore  $\mathcal{O}(P)$  work to be done per level. If we use  $P_l$  processors we have to sync on  $P_l$  tasks, which takes  $T_P = \mathcal{O}(P/P_l + \log P_l)$  which cannot be asymptotically better than  $\mathcal{O}(\log P)$ . Additionally  $T_s/P = \mathcal{O}(n/P + m/P)$  is a lower bound for  $T_P$ . Therefore the lower bound for  $T_P = \Omega(n/P + m/P + D \log P)$ .

## B. FBFS

The main innovation of FBFS is using prefix sums to efficiently split and combine the FIFO queues.

See Fig. 2.

## IV. PARALLEL PREFIX SUM

Parallel prefix sum is a common primitive used in many parallel algorithms. In the CILK model, the standard ( $\mathcal{O}(\log n)$  parallel fors) parallel prefix sum on  $n$  elements runs in  $T_\infty = \mathcal{O}(\log^2 n)$  or  $T_P = \mathcal{O}(n/P + \log^2 P)$ . In the PRAM model it runs in  $T_\infty = \mathcal{O}(\log n)$ .

We introduce a divide-and-conquer version of parallel prefix sum that runs in  $T_P = \mathcal{O}(n/P + \log P)$ .

## V. ANALYSIS

## VI. CONCLUSION

## VII. FUTURE WORK

TODO "WRITE" these instead of what we have now.

Add labels and links for each of these

- Optimize for false sharing
  - From parallel Prefix Sum :- pretty easy, just set appropriate GRAIN-Size and align properly.
  - For Dist array.
    - can use radix sort for large levels where num degrees =  $\Omega(n^c)$
    - Each proc can create an unsorted list of all cache lines its degrees touch, then each proc takes  $\#cachelines/p$  cache lines
  - For find sublist
- Optimize for cache-misses
- Make processor-oblivious
- Make cache-oblivious? Is it already?
- Optimize  $T_p, T_1, W_p$  for PRAM model

## REFERENCES

- [1] C. E. Leiserson and T. B. Schardl, "A work-efficient parallel breadth-first search algorithm (or how to cope with the nondeterminism of reducers)," in *Proceedings of the 22nd ACM symposium on Parallelism in algorithms and architectures*, ser. SPAA '10. New York, NY, USA: ACM, 2010, pp. 303–314. [Online]. Available: <http://doi.acm.org/10.1145/1810479.1810534>

- [2] D. A. Bader and K. Madduri, "Designing multithreaded algorithms for breadth-first search and st-connectivity on the cray mta-2," in *Proceedings of the 2006 International Conference on Parallel Processing*, ser. ICPP '06. Washington, DC, USA: IEEE Computer Society, 2006, pp. 523–530. [Online]. Available: <http://dx.doi.org/10.1109/ICPP.2006.34>
- [3] E. Saule and U. V. Catalyurek, "An early evaluation of the scalability of graph algorithms on the intel mic architecture," in *Proceedings of the 2012 IEEE 26th International Parallel and Distributed Processing Symposium Workshops & PhD Forum*, ser. IPDPSW '12. Washington, DC, USA: IEEE Computer Society, 2012, pp. 1629–1639. [Online]. Available: <http://dx.doi.org/10.1109/IPDPSW.2012.204>
- [4] J. Ullman and M. Yannakakis, "High-probability parallel transitive closure algorithms," in *Proceedings of the second annual ACM symposium on Parallel algorithms and architectures*, ser. SPAA '90. New York, NY, USA: ACM, 1990, pp. 200–209. [Online]. Available: <http://doi.acm.org/10.1145/97444.97686>

```

SERIAL-BFS(  $V, \Gamma, s$  )
1)  for each vertex  $u \in V$ 
2)     $\text{DIST}[u] \leftarrow \infty$ 
3)   $\text{DIST}[s] \leftarrow 0$ 
4)   $Q_{in} \leftarrow \emptyset$  {New empty input queue.}
5)   $\text{ENQUEUE}( Q_{in}, v )$ 
6)  while  $Q_{in} \neq \emptyset$  do {New empty output queue.}
7)     $Q_{out} \leftarrow \emptyset$ 
8)    while  $Q_{in} \neq \emptyset$  do
9)       $u \leftarrow \text{DEQUEUE}( Q_{in} )$ 
10)     for each vertex  $v \in \Gamma(u)$  do
11)       if  $\text{DIST}[v] = \infty$  then
12)          $\text{DIST}[v] \leftarrow \text{DIST}[u] + 1$ 
13)          $\text{ENQUEUE}( Q_{out}, v )$ 
14)      $Q_{in} \leftarrow Q_{out}$  {Output queue becomes input queue.}

```

Fig. 1. Standard serial level-synchronous breadth-first search implementation on a graph  $G = (V, E)$  with source vertex  $s \in V$ .  $\Gamma(u)$  is the adjacency list for node  $u$ .

```

PARALLEL-BFS(  $V, \Gamma, \gamma, s, P_{max}$  )
 $V[0 : n - 1]$  are the  $n$  nodes in the graph.  $\Gamma[u]$  is the sequence of adjacent nodes to node  $u$ .  $\gamma[u] = |\Gamma[u]|$ .  $s$  is the source vertex from which distance is
calculated.  $P_{max}$  is the maximum number of processors to use. Returns  $\text{DIST}[0 : n - 1]$  which represents the distance from  $s$  to each vertex.
1)  parallel for  $u \leftarrow 0$  to  $n - 1$  do
2)     $\text{DIST}[u] \leftarrow \infty$ 
3)     $\text{OWNER}[u] \leftarrow \infty$ 
4)   $\text{DIST}[s] \leftarrow 0$ 
5)   $\text{OWNER}[s] \leftarrow 0$ 
6)  if  $\gamma[s] = 0$  then
7)    return  $\text{DIST}$ 
8)   $\text{INPUT} \leftarrow \text{array}[0 : 0]$ 
9)   $\text{INPUT}[0] \leftarrow s$ 
10)  $\text{LEVEL} \leftarrow 0$ 
11)  $P \leftarrow 1$ 
12) while  $|\text{INPUT}| \neq 0$  do
13)    $\text{LEVEL} \leftarrow \text{LEVEL} + 1$ 
14)    $N \leftarrow |\text{INPUT}|$ 
15)    $\text{WORK} \leftarrow \text{array}[0 : N - 1]$ 
16)   parallel for  $u \leftarrow 0$  to  $N - 1$  do
17)      $\text{WORK}[u] \leftarrow \gamma[\text{INPUT}[u]]$ 
18)    $\text{PARALLEL-PREFIX-SUM}( \text{WORK}, \lfloor \frac{N}{P} \rfloor )$ 
19)    $W \leftarrow \text{WORK}[N - 1]$ 
20)    $P \leftarrow \text{MIN}( P_{max}, W )$  {Reduce #processors to save energy when  $W < P_{max}$ .}
21)    $\text{SUBLIST} \leftarrow \text{FIND-SUBLIST}( \text{WORK}, W, P )$ 
22)    $Q \leftarrow \text{LEVEL-TO-QUEUES}( \text{INPUT}, \text{WORK}, \text{SUBLIST}, \text{DIST}, \Gamma, \gamma, W, P, \text{LEVEL} )$ 
23)    $\text{SIZES} \leftarrow \text{array}[0 : P - 1]$ 
24)   parallel for  $i \leftarrow 0$  to  $P - 1$  do
25)      $Q\text{-NEW} \leftarrow \text{queue}$ 
26)     for  $v$  in  $Q[i]$  do
27)       if  $\text{OWNER}[v] = i$  then
28)          $Q\text{-NEW}.\text{ENQUEUE}( v )$ 
29)      $Q[i] \leftarrow Q\text{-NEW}$ 
30)      $\text{SIZES}[i] \leftarrow |Q[i]|$ 
31)    $\text{PARALLEL-PREFIX-SUM}( \text{SIZES}, 1 )$ 
32)    $\text{INPUT} \leftarrow \text{array}[0 : \text{SIZES}[P - 1]]$ 
33)   parallel for  $i \leftarrow 0$  to  $P - 1$  do
34)     if  $i = 0$  then
35)        $\text{OFFSET} \leftarrow 0$ 
36)     else
37)        $\text{OFFSET} \leftarrow \text{SIZES}[i - 1]$ 
38)     for  $j \leftarrow \text{OFFSET}$  to  $\text{OFFSET} + |Q[i]|$  do
39)        $\text{INPUT}[\text{OFFSET}] \leftarrow Q[i].\text{DEQUEUE}( )$ 

```

Fig. 2.

FIND-SUBLIST( WORK, W, P )

```

1)  N ← |WORK|
2)  Pn ← MIN( P, N )
3)  SUBLIST ← array[0 : P - 1]
4)  RANGESSTART ← array[0 : Pn - 1]
5)  RANGESEND ← array[0 : Pn - 1]
6)  parallel for i ← 0 to Pn - 1 do
7)    if i = 0 then
8)      FIRSTDEGREE ← 0
9)    else
10)     FIRSTDEGREE ← WORK[⌊ $\frac{iN}{P_n}$ ⌋ - 1]
11)     FIRSTDEGREE NEXT ← WORK[⌊ $\frac{(i+1)N}{P_n}$ ⌋ - 1]
12)     RANGESSTART[i] ← ⌊ $\frac{P \cdot \text{FIRSTDEGREE}}{W}$ ⌋
13)     RANGESEND[i] ← ⌊ $\frac{P \cdot \text{FIRSTDEGREE NEXT}}{W}$ ⌋ - 1
14)  parallel for i ← 0 to Pn - 1 do
15)    if RANGESSTART[i] ≤ RANGESEND[i] then
16)      parallel for j ← RANGESSTART[i] to RANGESEND[i] do
17)        SUBLIST[j] ← i
18)  return SUBLIST

```

{Use cores RANGESSTART[i] to RANGESEND[i]}

Fig. 3.

LEVEL-TO-QUEUES( INPUT, WORK, SUBLIST, DIST,  $\Gamma$ ,  $\gamma$ , W, P, LEVEL )

```

1)  N ← |WORK|
2)  Pn ← MIN( P, N )
3)  Q ← array[0 : P - 1]
4)  parallel for i ← 0 to P - 1 do
5)    Q[i].CLEAR( )
6)    FIRSTDEGREE ← ⌊ $\frac{iW}{P}$ ⌋
7)    WORKITEMS ← ⌊ $\frac{(i+1)W}{P}$ ⌋ - FIRSTDEGREE
8)    VERTEX ← BINARY-SEARCH-FOR-INDEX( WORK, ⌊ $\frac{N \cdot \text{SUBLIST}[i]}{P_n}$ ⌋, ⌊ $\frac{N \cdot (\text{SUBLIST}[i]+1)}{P_n}$ ⌋, FIRSTDEGREE + 1 )
9)    DEGREE ← FIRSTDEGREE - WORK[VERTEX - 1]
10)   while WORKITEMS > 0 do
11)     u ← INPUT[VERTEX]
12)     LIMIT ← MIN( WORKITEMS + DEGREE,  $\gamma[u]$  )
13)     for j ← DEGREE to LIMIT do
14)       v ←  $\Gamma[u][j]$ 
15)       if DIST[v] =  $\infty$  then
16)         DIST[v] ← LEVEL
17)         OWNER[v] ← i
18)         if  $\gamma[v] > 0$  then
19)           Q[i].ENQUEUE( v )
20)         WORKITEMS ← WORKITEMS - DEGREE
21)         DEGREE ← 0
22)         VERTEX ← VERTEX + 1
23)   return Q

```

{Benign race condition. All threads write the same value.}  
{Benign race condition. One thread's value will win.}

Fig. 4.

PARALLEL-PREFIX-SUM( V, GRAIN-SIZE )

$V[0 : n - 1]$  is a sequence of  $n$  integers. This function replaces  $V[i]$  with  $\sum_{0 \leq j \leq i} V[j]$

```

1)  if |V| > 1 then
2)    PARALLEL-PREFIX-SUM-UP( V, GRAIN-SIZE, 0, n )
3)    PARALLEL-PREFIX-SUM-DOWN( V, GRAIN-SIZE, 0, n, false, 0 )

```

Fig. 5.

PARALLEL-PREFIX-SUM-UP( V, GRAIN-SIZE, start, limit )

```

1)  size ← limit - start
2)  if size ≤ GRAIN-SIZE then
3)    return SERIAL-PREFIX-SUM( V, start, limit )
4)  else
5)    mid ← ⌊ $\frac{start+limit}{2}$ ⌋
6)    x ← spawn PARALLEL-PREFIX-SUM-UP( V, GRAIN-SIZE, start, mid )
7)    y ← PARALLEL-PREFIX-SUM-UP( V, GRAIN-SIZE, mid, limit )
8)    sync
9)    V[limit - 1] ← x + y
10)   return x + y

```

Fig. 6.

```

PARALLEL-PREFIX-SUM-DOWN( V, start, limit, rightmost_excluded, partial_sum )

1)  size  $\leftarrow$  limit - start
2)  if size  $\leq$  GRAIN-SIZE then
3)      SERIAL-PREFIX-SUM-DOWN( V, start, limit, partial_sum, rightmost_excluded )
4)      return
5)  else
6)      mid  $\leftarrow$   $\lfloor \frac{start+limit}{2} \rfloor$ 
7)      sum_left  $\leftarrow$  V[mid - 1]
8)      spawn PARALLEL-PREFIX-SUM-DOWN( V, GRAIN-SIZE, start, mid, false, partial_sum )
9)      if  $\neg$ rightmost_excluded then
10)         V[limit - 1]  $\leftarrow$  V[limit - 1] + partial_sum
11)     if limit - mid > 1
12)         PARALLEL-PREFIX-SUM-DOWN( V, GRAIN-SIZE, mid, limit, true, partial_sum + sum_left )

```

Fig. 7.

```

SERIAL-PREFIX-SUM( V, start, limit )

1)  V[start]  $\leftarrow$  V[start] + partial_sum
2)  for i  $\leftarrow$  start + 1 to limit - 1 do
3)      V[i]  $\leftarrow$  V[i] + V[i - 1]
4)  return V[limit - 1]

```

Fig. 8.

```

SERIAL-PREFIX-SUM-DOWN( V, start, limit, rightmost_excluded, partial_sum )

1)  for i  $\leftarrow$  start to limit - 2 do
2)      V[start]  $\leftarrow$  V[start] + partial_sum
3)  if  $\neg$ rightmost_excluded then
4)      V[limit - 1]  $\leftarrow$  V[limit - 1] + partial_sum

```

Fig. 9.