

What is a graph?

## Parallel Graph Algorithms

$$G = (V, E) \quad E: V \times V$$

Unweighted graphs

$$G = (V, E, W) \quad W: E \rightarrow \mathbb{R}$$

Directed:  $E$  is ordered

Undirected:  $E$  is unordered

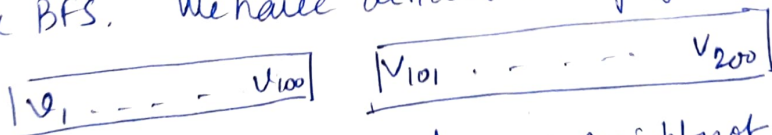
Very large graphs — Stanford Large Network Dataset Collection

Memory can be a bottleneck / compute can be a bottleneck

So both OpenMP parallelization for shared memory and MPI parallelization for distributed memory are important for graph algorithms.

GPU acceleration is hard for graph algorithms as typically graphs are sparse. Adjacency matrix for these large graphs are typically very sparse. Sparse matrix operations are not very optimal for GPU. Dense matrix operations work better — caches can be optimized for a block of contiguous access on CPU, on GPU the sparse access between DRAM and GPU memory is even worse.

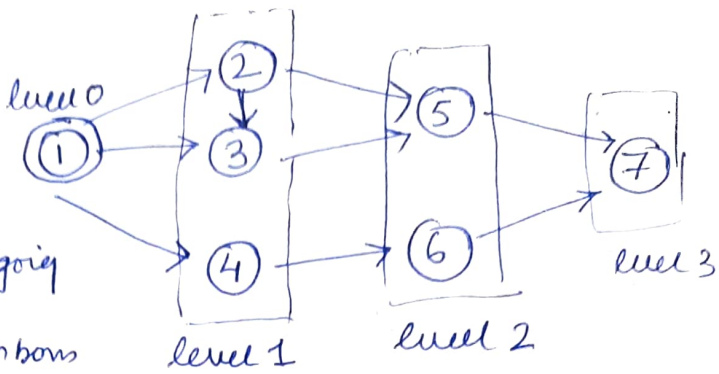
Let's take BFS. We have divided the graph nodes across processors.



The neighbors of a node on a processor might not be on the same processor. Thus to run BFS, communication will be needed. Communication might be much more costly than compute, and therefore the algorithms designed should also minimize communication.

# Breadth First Search

Root ①



First step: visit all nodes of ① reachable with outgoing edges.

Second step: visit all neighbors of ②, ③, ④.

Very large graph  $\rightarrow$  so all nodes are not fitting on one processor memory.

**Frontier**  $\rightarrow$  all nodes which got their levels marked in current iteration

level = 0

$$d = \begin{bmatrix} 0 & 2 & 3 & 4 & 5 & 6 & 7 \\ 0 & \infty & \infty & \infty & \infty & \infty & \infty \end{bmatrix}$$

distance from the root / level of the particular node in the BFS tree

$$\text{Frontier } f = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

while ( $f \neq 0$ )

{

$$u = A^T f$$

$f' = u$  & ( $d \neq \infty$ )  
update  $d$  with all elements  $f' = 1$

~~update~~

set  $d[i] = \text{level} + 1$

level = level + 1

$f = f'$

}

}

## Adjacency matrix

	1	2	3	4	5	6	7
1	0	1	1	0	0	0	0
2	0	0	1	0	1	0	0
3	0	0	0	0	1	0	0
4	0	0	0	0	0	1	0
5	0	0	0	0	0	0	1
6	0	0	0	0	0	0	1
7	0	0	0	0	0	0	0

$$a_{ij} = 1 \text{ if } (i, j) \in E, (i, j) \in V \times V$$

suppose a node  $j$  is in frontier i.e.  $f[j] = 1$ , then a node  $k$  will be candidate for next frontier

if  $a_{jk} = 1$

$$\begin{bmatrix} \vdots \\ j \\ \vdots \end{bmatrix}_{k^{\text{th}}} \begin{bmatrix} \vdots \\ j \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots \\ \text{---} \\ \vdots \end{bmatrix}_{k^{\text{th}}} u$$

$$u[k] = \sum_j (f[j] \wedge a[j, k])$$

$A^T f$  dot product  $\rightarrow$  addition  
multiplication replaced ( $\wedge, \wedge$ )

BFS with matrix algebra  $\rightarrow$  can be parallelized.

Non matrix based algorithm for BFS.

Sequential - BFS ( $G=(V,E)$ , vertex src)

```
{
   $\forall v \in V: \text{level}[v] = -1$ 
  level[src] = 0;
  curr_level = 0;
  frontier = { } // set of vertices at the current level.
  next-frontier = { } // set of vertices at the next level
  frontier.add(src);
  while (frontier is not empty)
  {
    for each vertex  $u \in \text{frontier}$ ;  $\leftarrow$  can be executed in parallel
    for each neighbor  $v$  of  $u$ :
      if (level[v] == -1) // not visited before.
      {
        next-frontier.add(v);
        level[v] = curr_level + 1;
      }
    }
    frontier = next-frontier;  $\leftarrow$  barrier
    next-frontier = { };
    curr_level++;
  }
}
```

Two threads looking at two different nodes of the frontier, but one of the neighbors of these two vertices are common. Both threads might see  $(\text{level}[v] = -1)$  and try to change it.



Parallel BFS using openMP

```
#pragma omp parallel
```

```
tid = omp_get_threadnum();
```

```
nthreads = omp_get_num_threads();
```

```
# traverse the frontier together
```

```
for (int i = tid; i < frontier.size(); i += nthreads)
```

```
{
  for each neighbor v of u
  if (level[v] == -1)
```

```
{ #pragma omp atomic capture
```

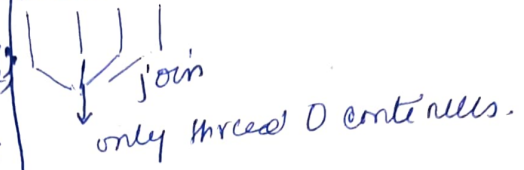
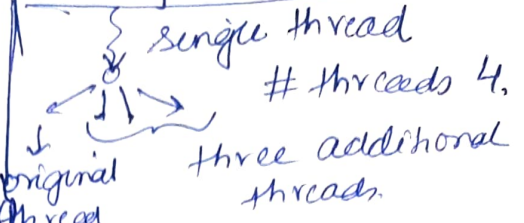
```
{ insert = level[v];
  level[v] = curr_level;
}
```

```
if (insert == -1)
  local_next_frontier.add(v);
}
```

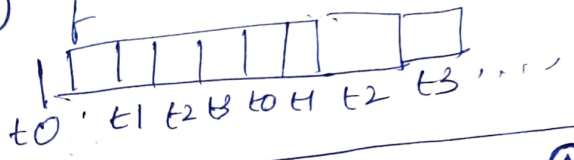
```
}
#pragma omp critical
{ nextfrontier.merge(local_next_frontier);
}
```

```
}
```

Fork join model



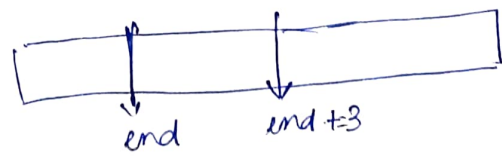
frontier: array of vertex ids



**(A)** critical section  
 #pragma omp critical  
 { enqueue vertex v in next-frontier  
 }  
 (-) lock at the code level  
 instructions cannot be executed by any other thread, even if that thread is working on a different piece of data. Instead we do lock on data.

**(B)** not a global new frontier

↑ how to efficiently do this?

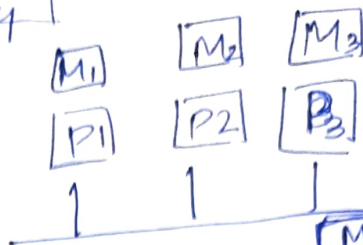


increment pointers only in critical section. A thread reserves space for its own entries, copy into three reserved space outside critical section.

# BFS in distributed memory setting

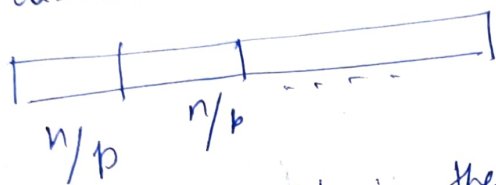
Cannot access the memory of another processor.

Instead use send/receiv



How do we distribute the data? (graph partitioning)

ID - data distribution



$n$  vertices  
 $p$  processors

Partition vertices

① randomize IDs then distribute for load balancing



cut/partition minimize edges across

for that vertex.

That processor will maintain the information  
 → adjacency list  
 → any other attribute

$owner(u) \rightarrow rank$

Distributed - BFS - ID (local  $G = (V, E)$ , vertices)

{

frontier = { }; // local

next-frontier = { }; // local

curr\_level = 0; // local

for all  $v \in V$ : level[v] = -1; // not visited yet.

if (owner(s) == me) // MPI\_COMM\_RANK used

{ level[s] = 0;

frontier.add(s);

}

while (true)

{

// cannot use frontier == NULL as not shared  
 my local frontier might be empty, but other processors might have nodes in their frontier. Those nodes might have neighbors, which I own - so I will be added in the next iteration. So a processor can be added in some iteration

// frontier contains vertices (local) in current frontier

for each vertex  $u \in \text{frontier}$ :

for  $v \in \text{neighbor}(u)$ : // send these vertices to their owners.

$w = \text{owner}(v)$

Buffer[w].add(v); // for each processor, I maintain a buffer.

p sends, p receives

for each  $p \neq \text{me}$

send buffer[p] to p;

recv RecvBuff[p]

from p;

// send buffers to respective processors.

Alltoallv(buffer, RecvBuff, ...); // each buffer to respective ranks, and

$\forall p = [0 \dots \text{numRanks} - 1]$

next-frontier.merge(RecvBuff[p]); // also receive from all other ranks

*duplicate entries.*

frontier = { };

for  $v \in \text{next-frontier}$ :

if level[v] == -1

level[v] = curr-level + 1;

frontier.add(v);

*v might be neighbor of multiple vertices at current level*

*only if level has changed, so frontier removes duplicate entries.*

next-frontier = { };

curr-level ++;

size = frontier.size();

AllReduce(size); // total global frontier

if (size == 0) // all frontiers are empty

break out of while loop

P1	P2	P3	...	P4
3	2	1		1
7	7	7		7

}

}