

Multicore Computing Lecture 19 – Parallel Graph Processing



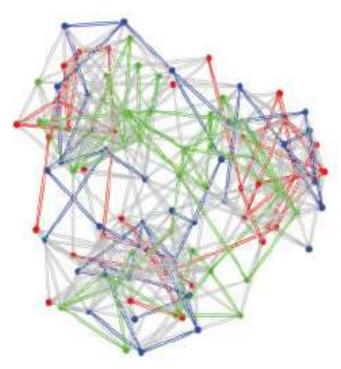
남 범 석 bnam@skku.edu



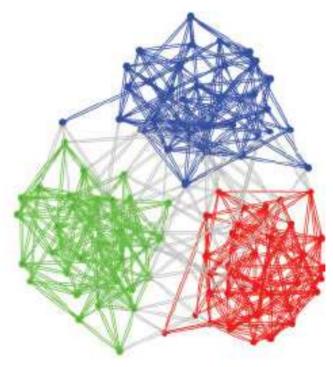
- Large Social Network Graph
 - Too large to store in a single machine

Parallel Graph Processing

Graph Partitioning and Clustering



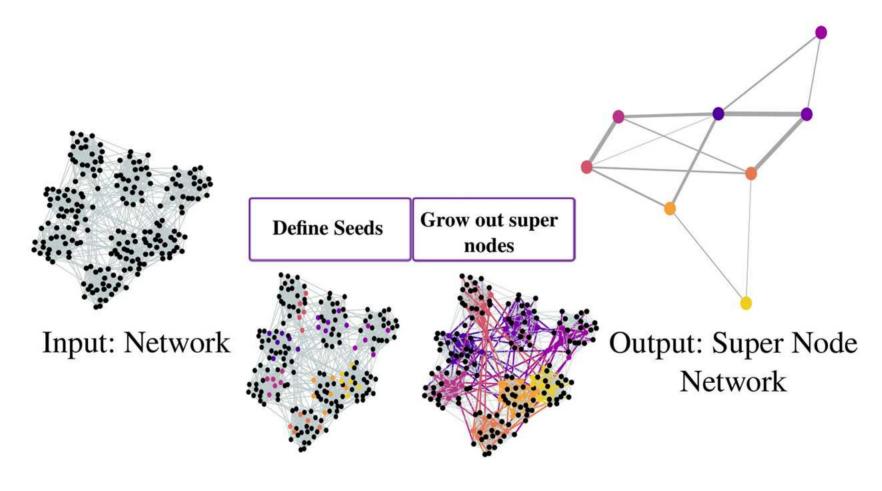
(a) A poor edge-cut partitioning. Vertices are assigned to partitions at random, thus, there are many inter-partition links.



(b) A good edge-cut of the same graph, where vertices that are highly connected are assigned to the same partition.

Parallel Graph Processing

Lots of interesting on-going research in this area



Back-to-the-Basic



- Minimum Spanning Tree
 - Prim's algorithm
- Single-Source Shortest Path
 - Dijkstra's algorithm
- All-pair Shortest Path
 - Floyd's algorithm

Minimum Spanning Tree

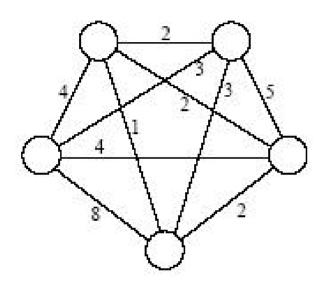


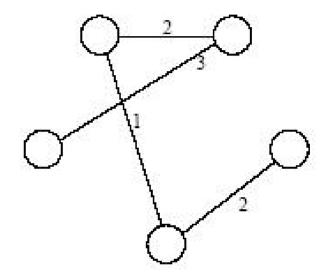
- A spanning tree of an undirected graph G is a subgraph of G that is a tree containing all the vertices of G.
- In a weighted graph, the weight of a subgraph is the sum of the weights of the edges in the subgraph.
- A minimum spanning tree (MST) for a weighted undirected graph is a spanning tree with minimum weight.

Minimum Spanning Tree



An undirected graph and its minimum spanning tree.



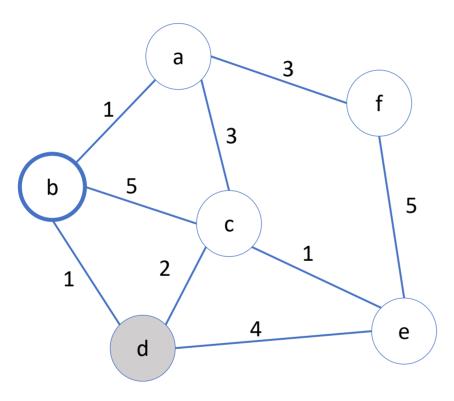




- Prim's algorithm for finding an MST is a greedy algorithm.
- Start by selecting an arbitrary vertex, include it into the current MST.
- Grow the current MST by inserting into the vertex closest to one of the vertices already in current MST.



(a) Start from b

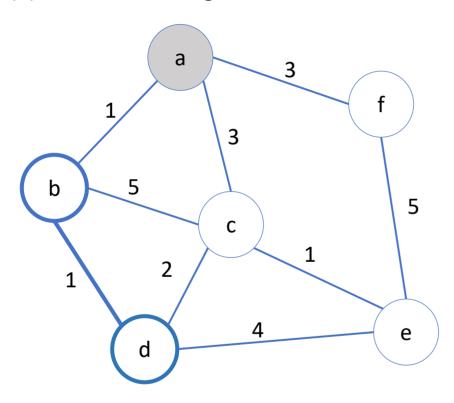


	a	b	С	d	е	f
d[]	1	0	5	1	∞	∞

a
$$\begin{bmatrix} 0 & 1 & 3 & \infty & \infty & 3 \\ 1 & 0 & 5 & 1 & \infty & \infty \\ c & 3 & 5 & 0 & 2 & 1 & \infty \\ d & \infty & 1 & 2 & 0 & 4 & \infty \\ e & \infty & \infty & 1 & 4 & 0 & 5 \\ f & 2 & \infty & \infty & \infty & 5 & 0 \end{bmatrix}$$



(b) After the first edge is selected



	а	b	С	d	е	f
d[]	1	0	2	1	4	∞

a
 0
 1
 3

$$\infty$$
 ∞
 3

 b
 1
 0
 5
 1
 ∞
 ∞

 c
 3
 5
 0
 2
 1
 ∞

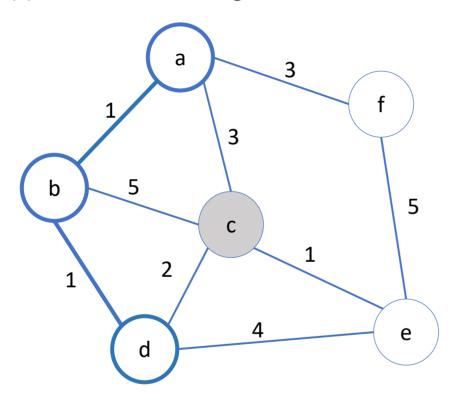
 d
 ∞
 1
 2
 0
 4
 ∞

 e
 ∞
 1
 4
 0
 5

 f
 2
 ∞
 ∞
 5
 0



(c) After the second edge is selected

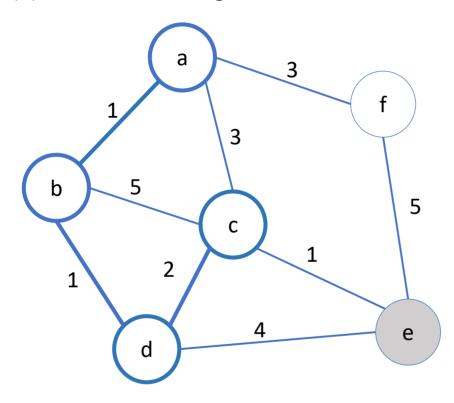


	а	b	С	d	е	f
d[]	1	0	2	1	4	3

a	0	1	3	∞	∞	3
b	1	0	5	1	∞	∞
С	3	5	0	2	1	∞
d	∞	1	2	0	4	∞
e	∞	∞	1	4	0	5
f	2	∞	∞	∞	5	0



(d) After the third edge is selected

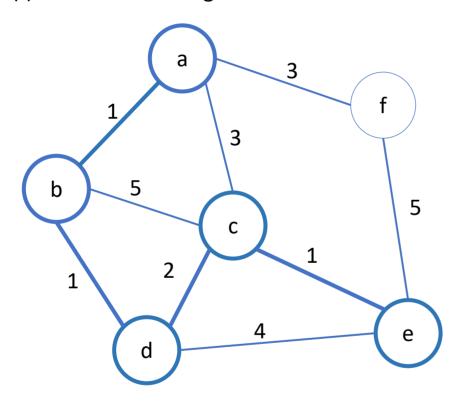


	а	b	С	d	е	f
d[]	1	0	2	1	1	3

a
$$\begin{bmatrix} 0 & 1 & 3 & \infty & \infty & 3 \\ 1 & 0 & 5 & 1 & \infty & \infty \\ c & 3 & 5 & 0 & 2 & 1 & \infty \\ d & \infty & 1 & 2 & 0 & 4 & \infty \\ e & \infty & \infty & 1 & 4 & 0 & 5 \\ f & 2 & \infty & \infty & \infty & 5 & 0 \end{bmatrix}$$



(f) After the fifth edge is selected

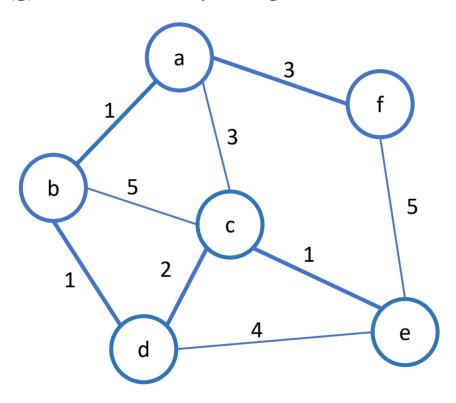


	а	b	С	d	е	f
d[]	1	0	2	1	1	3

a
$$\begin{bmatrix} 0 & 1 & 3 & \infty & \infty & 3 \\ 1 & 0 & 5 & 1 & \infty & \infty \\ c & 3 & 5 & 0 & 2 & 1 & \infty \\ d & \infty & 1 & 2 & 0 & 4 & \infty \\ e & \infty & \infty & 1 & 4 & 0 & 5 \\ f & 2 & \infty & \infty & \infty & 5 & 0 \end{bmatrix}$$



(g) Final minimum spanning tree



	а	b	С	d	е	f
d[]	1	0	2	1	1	3

a
 0
 1
 3

$$\infty$$
 ∞
 3

 b
 1
 0
 5
 1
 ∞
 ∞

 c
 3
 5
 0
 2
 1
 ∞

 d
 ∞
 1
 2
 0
 4
 ∞

 e
 ∞
 1
 4
 0
 5

 f
 2
 ∞
 ∞
 5
 0





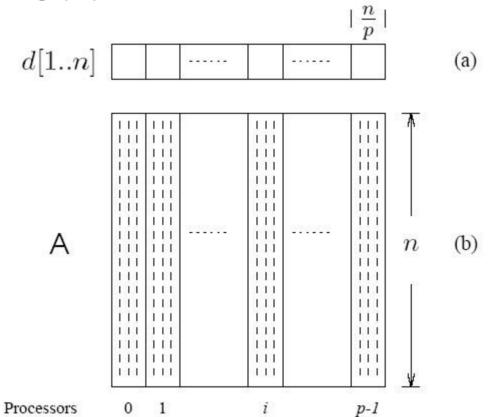
Prim's sequential minimum spanning tree algorithm.

```
1.
          procedure PRIM\_MST(V, E, w, r)
2.
          begin
3.
               V_T := \{r\};
4.
               d[r] := 0;
5.
               for all v \in (V - V_T) do
6.
                    if edge (r, v) exists set d[v] := w(r, v);
7.
                    else set d[v] := \infty;
8.
               while V_T \neq V do
9.
               begin
10.
                    find a vertex u such that d[u] := \min\{d[v] | v \in (V - V_T)\};
11.
                    V_T := V_T \cup \{u\};
                    for all v \in (V - V_T) do
12.
13.
                         d[v] := \min\{d[v], w(u, v)\};
14.
               endwhile
15.
          end PRIM_MST
```

- d[] change after every loop
- Greedy algorithm → Outer iteration is hard to parallelize
- The inner loop is relatively easy to parallelize.
 - The adjacency matrix is partitioned in a 1-D block fashion, with distance vector d partitioned accordingly.
 - In each step, a processor selects the locally closest node, followed by a global reduction to select globally closest node.
 - This node is inserted into MST, and the choice is broadcasted to all processors.
 - Each processor updates its part of the d vector locally.

```
1-D block partitioning: V_i per P_i. For each iteration: P_i \text{ computes a local min } d_i[u]. All-to-one reduction to P_0 to compute the global min. One-to-all broadcast of u. Local updates of d[v].
```

 The partitioning of the distance array d and the adjacency matrix A among p processes.



- The cost to select the minimum entry is O(n/p + log p).
- The cost of a broadcast is O(log p).
- The cost of local update of the d vector is O(n/p).
- The parallel time per iteration is 2(n/p + log p).
- The total parallel time is $2(n^2/p + n \log p)$
- In Big O notation, parallel cost is O(n²/p + n log p).
 - Constant coefficients are allowed to be dropped in Big O.
 - But it is often misleading in distributed and parallel algorithms since communication cost is not cheap

Single-Source Shortest Paths



- For a weighted graph G = (V,E,w), the single-source shortest paths problem is to find the shortest paths from a vertex v ∈ V to all other vertices in V.
- Dijkstra's algorithm is similar to Prim's algorithm.
 - maintains a set of nodes for which the shortest paths are known.
- It grows this set based on the node closest to source using one of the nodes in the current shortest path set.





Dijkstra's sequential single-source shortest paths algorithm.

```
procedure DIJKSTRA_SINGLE_SOURCE_SP(V, E, w, s)
1.
2.
          begin
3.
               V_T := \{s\};
               for all v \in (V - V_T) do
4.
5.
                    if (s, v) exists set l[v] := w(s, v);
6.
                    else set l[v] := \infty;
7.
               while V_T \neq V do
8.
               begin
9.
                    find a vertex u such that l[u] := \min\{l[v] | v \in (V - V_T)\};
10.
                    V_T := V_T \cup \{u\};
11.
                    for all v \in (V - V_T) do
12.
                         l[v] := \min\{l[v], l[u] + w(u, v)\};
13.
               endwhile
14.
          end DIJKSTRA SINGLE SOURCE SP
```

Dijkstra's Algorithm: Parallel Formulation

- Very similar to the parallel formulation of Prim's MST algorithm.
 - The weighted adjacency matrix is partitioned along columns
 - The node closest to the source in local node is globally reduced.
 - The node is broadcast to all processors and the I-vector updated.
- The parallel performance of Dijkstra's algorithm is identical to that of Prim's algorithm.

All-Pairs Shortest Paths



- Given a weighted graph G(V,E,w), the all-pairs shortest paths problem is to find the shortest paths between all pairs of vertices vi, vj \in V.
- A number of algorithms are known for solving this problem.

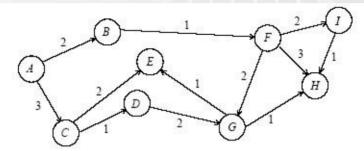
All-Pairs Shortest Paths: Matrix-Multiplication



- Multiplication of weighted adjacency matrix with itself → a matrix for shortest paths of length 2
 - Consider the multiplication of the weighted adjacency matrix with itself except, in this case, we replace the multiplication operation in matrix multiplication by addition, and the addition operation by minimization.
- Aⁿ contains all shortest paths.

All-Pairs Shortest Paths: Matrix-Multiplication





$$A^{1} = \begin{pmatrix} 0 & 2 & 3 & \infty & \infty & \infty & \infty & \infty \\ \infty & 0 & \infty & \infty & 1 & \infty & \infty & \infty \\ \infty & \infty & 0 & 1 & 2 & \infty & \infty & \infty \\ \infty & \infty & \infty & 0 & \infty & \infty & 2 & \infty & \infty \\ \infty & \infty & \infty & \infty & 0 & \infty & \infty & \infty & \infty \\ \infty & \infty & \infty & \infty & 0 & 2 & 3 & 2 \\ \infty & \infty & \infty & \infty & \infty & 0 & 1 & \infty \\ \infty & 0 & \infty \\ \infty & 1 & 0 \end{pmatrix}$$

$$A^{2} = \begin{pmatrix} 0 & 2 & 3 & 4 & 5 & 3 & \infty & \infty & \infty \\ \infty & 0 & \infty & \infty & \infty & 1 & 3 & 4 & 3 \\ \infty & \infty & 0 & 1 & 2 & \infty & 3 & \infty & \infty \\ \infty & \infty & \infty & 0 & 3 & \infty & 2 & 3 & \infty \\ \infty & \infty & \infty & \infty & 0 & \infty & \infty & \infty & \infty \\ \infty & \infty & \infty & \infty & 0 & \infty & \infty & \infty & \infty \\ \infty & \infty & \infty & \infty & 3 & 0 & 2 & 3 & 2 \\ \infty & \infty & \infty & \infty & \infty & 0 & 0 & \infty \\ \infty & 0 & \infty \\ \infty & 0 & \infty \end{pmatrix}$$

$$A^{4} = \begin{pmatrix} 0 & 2 & 3 & 4 & 5 & 3 & 5 & 6 & 5 \\ \infty & 0 & \infty & \infty & 4 & 1 & 3 & 4 & 3 \\ \infty & \infty & 0 & 1 & 2 & \infty & 3 & 4 & \infty \\ \infty & \infty & \infty & 0 & 3 & \infty & 2 & 3 & \infty \\ \infty & \infty & \infty & \infty & 0 & \infty & \infty & \infty & \infty \\ \infty & \infty & \infty & \infty & 0 & \infty & \infty & \infty & \infty \\ \infty & \infty & \infty & \infty & 3 & 0 & 2 & 3 & 2 \\ \infty & \infty & \infty & \infty & \infty & 1 & \infty & 0 & 1 & \infty \\ \infty & 0 & \infty \\ \infty & 1 & 0 \end{pmatrix}$$

$$A^{8} = \begin{pmatrix} 0 & 2 & 3 & 4 & 3 & 3 & 3 & 6 & 3 \\ \infty & 0 & \infty & \infty & 4 & 1 & 3 & 4 & 3 \\ \infty & \infty & 0 & 1 & 2 & \infty & 3 & 4 & \infty \\ \infty & \infty & \infty & 0 & 3 & \infty & 2 & 3 & \infty \\ \infty & \infty & \infty & \infty & 0 & \infty & \infty & \infty \\ \infty & \infty & \infty & \infty & 0 & \infty & \infty & \infty \\ \infty & \infty & \infty & \infty & 3 & 0 & 2 & 3 & 2 \\ \infty & \infty & \infty & \infty & 1 & \infty & 0 & 1 & \infty \\ \infty & 0 & \infty \\ \infty & 0 & \infty \\ \infty & 1 & 0 \end{pmatrix}$$

Matrix-Multiplication Based Algorithm: Parallel Formulation

- Aⁿ is computed by doubling powers
 - i.e., as A, A², A⁴, A⁸, and so on.
- We need log n matrix multiplications, each taking time O(n³).
- The serial complexity of this procedure is O(n³ log n).
- We can parallelize matrix multiplications and each of the log n matrix multiplications.

Dijkstra's Algorithm: Parallel Formulation

- Run single-source shortest path problem n times
 - For each source node
- Two parallelization strategies
 - Source Partitioned:
 - Execute each shortest path algorithm on a different processor
 - Source Parallel:
 - Parallelize each shortest path algorithm to further increase concurrency

Dijkstra's Algorithm: Source Partitioned Formulation

- Each processor P_i runs the shortest paths from vertex v_i
- Adjacency matrix is replicated at all processes.
 - Thus, no inter-process communication
- The parallel run time is: $\Theta(n^2)$ where p=n.
 - $n/p \times O(n^2)$
 - Complexity of Dijkstra's original algorithm = $O(n^2)$
- Fixed parallelism
 - It can use no more than n processors.

Dijkstra's Algorithm: Source Parallel Formulation

- Execute each of the shortest path problems in parallel
- We can use up to n² processors.
- Given p processors (p > n), each single source shortest path problem is executed by p/n processors.
- This takes time:
 - Parallel computation: n/n x O(n² / p/n)
 - Communication cost:
 - In each step, p processors need to find the global min distance → parallel reduction: O(log p), i.e., O(n log p)

$$T_P = \Theta\left(\frac{n^3}{p}\right) + \Theta(n\log p).$$

All-Pairs Shortest Paths: Floyd's Algorithm



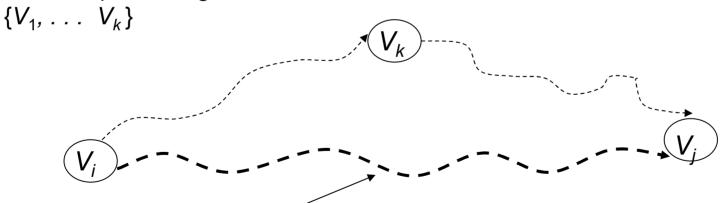
- Let D(k)[i,j]=weight of a shortest path from v_i to v_j using only vertices from $\{v_1, v_2, ..., v_k\}$ as intermediate vertices in the path
 - D(0)=W
 - D(n)=D which is the goal matrix
- How do we compute D(k) from D(k-1)?

The Recursive Definition:



- Case 1: A shortest path from v_i to v_j does not use v_k . Then D(k)[i,j] = D(k-1)[i,j].
- Case 2: A shortest path from v_i to v_j does use v_k . Then D(k)[i,j] = D(k-1)[i,k] + D(k-1)[k,j].

Shortest path using intermediate vertices



Shortest Path using intermediate vertices { $V_{1,...}$ V_{k-1} }

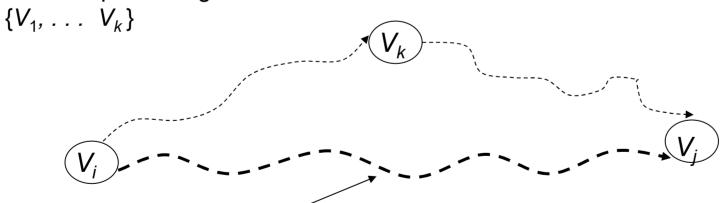
The Recursive Definition:



Since

$$D(k)[i,j] = D(k-1)[i,j]$$
 or $D(k)[i,j] = D(k-1)[i,k] + D(k-1)[k,j]$. We conclude: $D(k)[i,j] = min\{ D(k-1)[i,j], D(k-1)[i,k] + D(k-1)[k,j] \}$.

Shortest path using intermediate vertices



Shortest Path using intermediate vertices { $V_{1,...}$ V_{k-1} }

The pointer array P



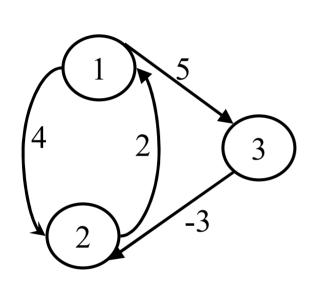
- Used to enable finding a shortest path
- Initially the array contains 0
- Each time that a shorter path from i to j is found → the k that provided the minimum is saved in P[i, j]
- P[i, j] = highest index node on the path from i to j
- To print the intermediate nodes on the shortest path, a recursive procedure that prints the shortest paths from i and k, and from k to j can be used

Floyd's Algorithm Using n+1 D matrices



```
// Computes shortest distance between all pairs of
// nodes, and saves P to enable finding shortest paths
// O(n^3)
1. D0 \leftarrow W // initialize D array to W []
2. P \leftarrow 0 // initialize P array to [0]
3. for k \leftarrow 1 to n
        do for i \leftarrow 1 to n
5.
             do for j \leftarrow 1 to n
                  if (D(k-1)[i, j] > D(k-1)[i, k] + D(k-1)[k, j]
6.
                        then D(k)[i, j] \leftarrow D(k-1)[i, k] + D(k-1)[k, j]
7.
                              P[i, j] \leftarrow k;
8.
9.
                        else D(k)[i, j] \leftarrow D(k-1)[i, j]
```

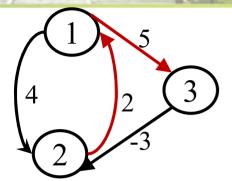




		1		3
$W = D^0 =$	1	0	4	5
$\mathbf{W} = \mathbf{D}_0 =$	2	2	0	8
	3	∞	-3	0

		1		3
	1	0	0	0
P =	2	0	0	0
	3	0	0	0





- 0	1	2	3
$D^0 = \frac{1}{1}$	0	4	5
2	2	0	8
3	8	-3	0

k = 1

Vertex 1 can be intermediate node

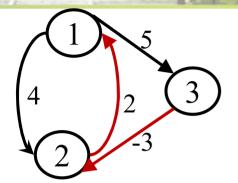
$$D^{1} = \begin{array}{c|cccc}
 & 1 & 2 & 3 \\
 & 0 & 4 & 5 \\
 & 2 & 0 & 7 \\
 & 3 & \infty & -3 & 0
\end{array}$$

D1[2,3] = min(D0[2,3], D0[2,1]+D0[1,3])
= min (
$$\infty$$
, 7)
= 7

$$P = \begin{array}{c|cccc}
 & 1 & 2 & 3 \\
1 & 0 & 0 & 0 \\
2 & 0 & 0 & 1 \\
3 & 0 & 0 & 0 \end{array}$$

D1[3,2] = min(D0[3,2], D0[3,1]+D0[1,2])
= min (-3,
$$\infty$$
)
= -3





_ 1	1	2	3
$D^1 = \frac{1}{1}$	0	4	5
2	2	0	7
3	∞	-3	0

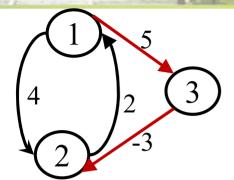
k = 2Vertices 1, 2 can be intermediate

$$D^{2} = \begin{array}{c|cccc}
 & 1 & 2 & 3 \\
 & 0 & 4 & 5 \\
 & 2 & 0 & 7 \\
 & 3 & -1 & -3 & 0
\end{array}$$

$$P = \begin{array}{c|cccc}
 & 1 & 2 & 3 \\
1 & 0 & 0 & 0 \\
2 & 0 & 0 & 1 \\
3 & 2 & 0 & 0
\end{array}$$

D2[3,1] = min(D1[3,1], D1[3,2]+D1[2,1])
= min (
$$\infty$$
, -3+2)
= -1





$$D^{2} = \begin{array}{c|cccc}
 & 1 & 2 & 3 \\
 & 0 & 4 & 5 \\
 & 2 & 2 & 0 & 7 \\
 & 3 & -1 & -3 & 0
\end{array}$$

k = 3Vertices 1, 2, 3 can be intermediate

$$D^{3} = \begin{array}{c|cccc}
 & 1 & 2 & 3 \\
 & 1 & 0 & 2 & 5 \\
 & 2 & 0 & 7 \\
 & 3 & -1 & -3 & 0
\end{array}$$

$$D3[1,2] = min(D2[1,2], D2[1,3]+D2[3,2])$$

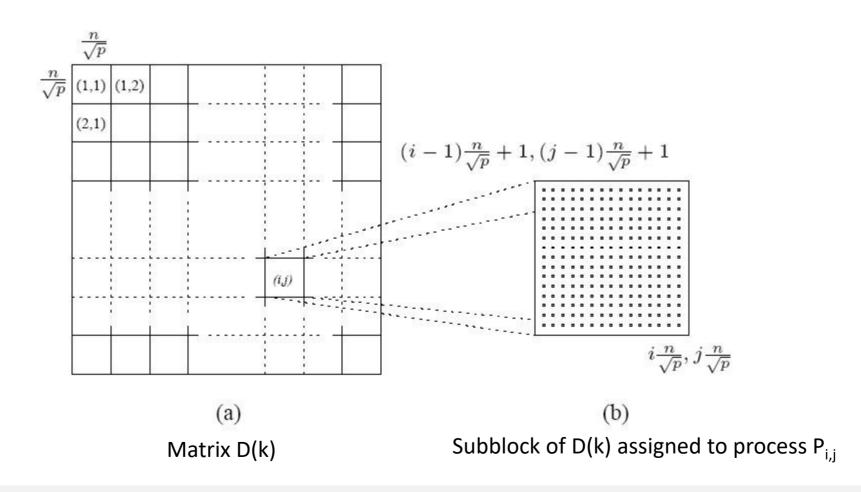
= min (4, 5+(-3))
= 2

$$P = \begin{array}{c|cccc}
 & 1 & 2 & 3 \\
 & 1 & 0 & 3 & 0 \\
 & 2 & 0 & 0 & 1 \\
 & 3 & 2 & 0 & 0 \\
\end{array}$$

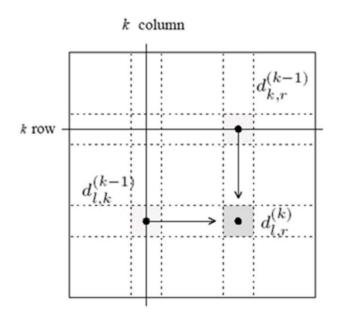
$$D3[2,1] = min(D2[2,1], D2[2,3]+D2[3,1])$$

= min (2, 7+ (-1))
= 2

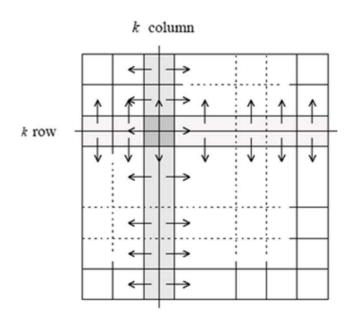
- Matrix D(k) is divided into p blocks of size $(n / \sqrt{p}) \times (n / \sqrt{p})$.
- Each processor updates its part of the matrix per iteration.



- To compute d(k)[l,r], $P_{i,j}$ must get d(k-1)[l,k] and d(k-1)[k,r].
 - d(k-1)[l,k] is on a process along the same row
 - d(k-1)[k,r] is on a process along the same column



- In general, during the kth iteration, each process containing the kth row sends it to the other processes in the same column.
- Communication patterns used in the 2-D block mapping.
 - Similarly, each process containing the kth column sends it to the other processes in the same row.



- Floyd's parallel formulation using the 2-D block mapping.
- $P_{*,j}$ denotes all the processes in the j^{th} column, and $P_{i,*}$ denotes all the processes in the i^{th} row.
- The matrix D(0) is the adjacency matrix.

```
procedure FLOYD_2DBLOCK(D^{(0)})
    begin
      for k = 1 to n do
      begin
4.
         each process P_{i,j} that has a segment of the k^{\text{th}} row of D^{(k-1)};
5.
              broadcasts it to the P_{*,j} processes;
         each process P_{i,i} that has a segment of the k<sup>th</sup> column of D^{(k-1)};
6.
              broadcasts it to the P_{i*} processes;
         each process waits to receive the needed segments;
7.
         each process P_{i,i} computes its part of the D^{(k)} matrix;
8.
9.
      end
10. end FLOYD_2DBLOCK
```

- During each iteration of the algorithm, the kth row and kth column of processors perform a one-to-all broadcast along their rows/columns.
- The size of this broadcast is n/\sqrt{p} elements, taking time $\Theta((n \log p)/\sqrt{p})$.
- The synchronization step takes time $\Theta(\log p)$.
- The computation time is $\Theta(n^2/p)$.
- The parallel run time of the 2-D block mapping formulation of Floyd's algorithm is

$$T_P = \Theta\left(\frac{n^3}{p}\right) + \Theta\left(\frac{n^2}{\sqrt{p}}\log p\right).$$

Floyd's Algorithm: Speeding Things Up by Pipelining

- The synchronization step in parallel Floyd's algorithm can be removed without affecting the correctness of the algorithm.
- A process starts working on the kth iteration as soon as it has computed the (k-1)th iteration and has the relevant parts of the D(k-1) matrix.

