

Parallelizing the Hungarian Method on GPU

Andrew Kope
Western University

The Assignment Problem

Problem Description

- ▶ There are a number of *agents* and *tasks*
- ▶ Any agent can be assigned to perform any task, incurring some *benefit* that may vary depending on the agent–task assignment
- ▶ It is required to perform all tasks by assigning exactly one agent to each task
- ▶ **Agents are assigned in such a way that the total benefit of the assignment is maximized**

Example

| | | Tasks | | |
|-------|---|--------|--------|--------|
| | | Task 1 | Task 2 | Task 3 |
| Agent | 1 | 250 | 300 | 150 |
| | 2 | 300 | 700 | 250 |
| | 3 | 100 | 200 | 150 |

Example

| | | Tasks | | |
|-------|---|--------|--------|--------|
| | | Task 1 | Task 2 | Task 3 |
| Agent | 1 | 250 | 300 | 150 |
| | 2 | 300 | 700 | 250 |
| | 3 | 100 | 200 | 150 |

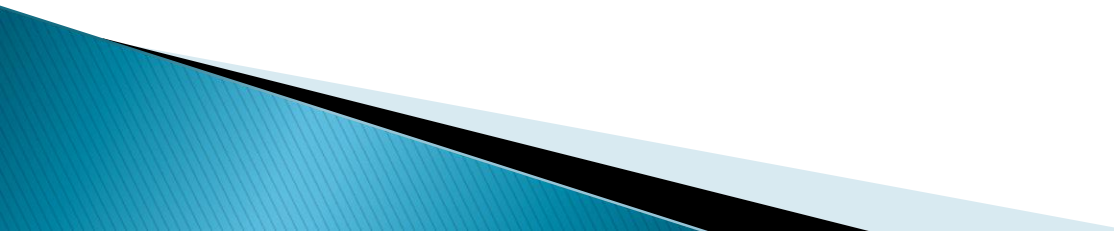
The optimal assignment is $250 + 700 + 150 = 1100$

Kuhn's Hungarian Method

History

- ▶ Published in 1955
- ▶ Based on the work of two Hungarian mathematicians

“one can always construct an optimal assignment by a succession of transfers followed by additional assignments until this is no longer available”

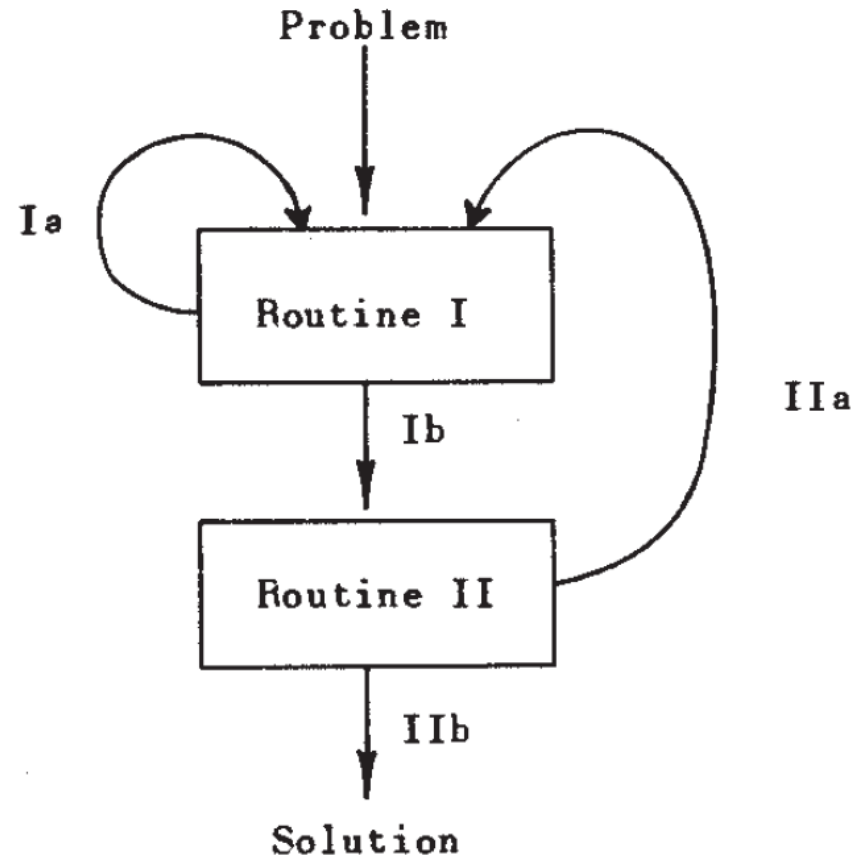


Description

- ▶ Step 1. Subtract the smallest entry in each row from all entries in its row
- ▶ Step 2. Subtract the smallest entry in each column from all entries in its column
- ▶ Step 3. Draw lines through appropriate rows and columns such that:
 - All the zero entries of the cost matrix are covered
 - The minimum number of such lines is used
- ▶ Step 4. Test for Optimality:
 - (i) If the minimum number of covering lines is n , an optimal assignment of zeros is possible and we are finished
 - (ii) If the minimum number of covering lines is less than n , an optimal assignment of zeros is not yet possible; proceed to Step 5
- ▶ Step 5. Determine the smallest entry not covered by any line. Subtract this entry from each uncovered row, and add it to each covered column; return to Step 3

Algorithm

- I: Routine 1; assign agents
- Ia: Increase number of assignments
- II: Routine 2; test for optimality
- IIa: Decrement current cover



The Cover

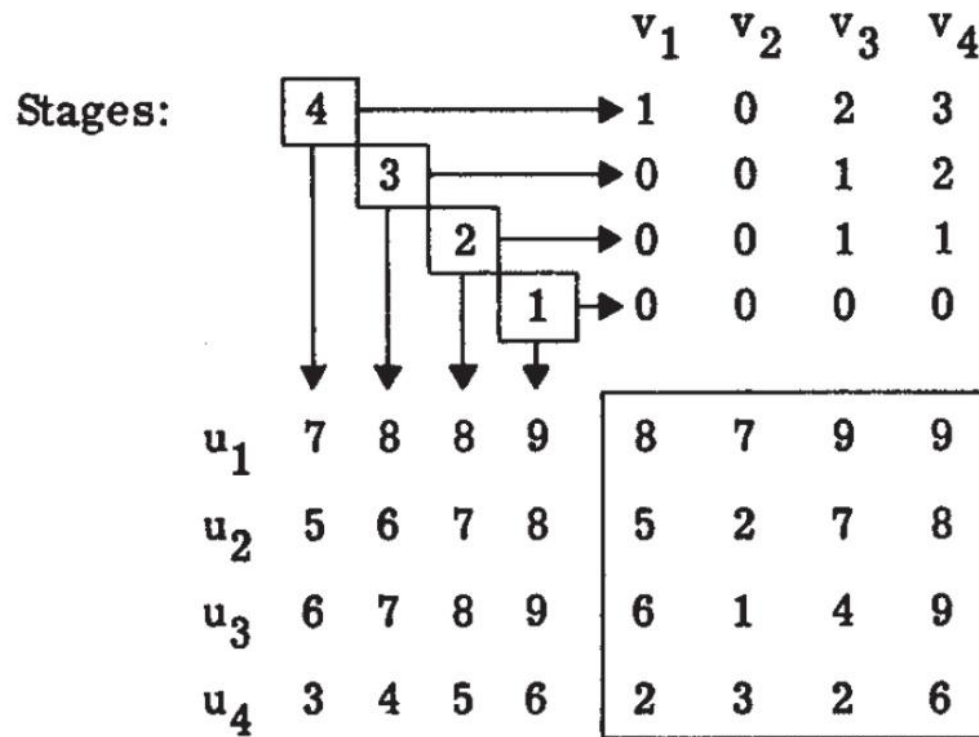
Description

- ▶ A set of lines is a cover if it contains all the zeroes in the cost/benefit matrix
- ▶ Composed of two vectors $\{U_i, V_j\}$
- ▶ With each iteration of the algorithm, the cover decreases in size
 - The minimal cover is an optimal solution

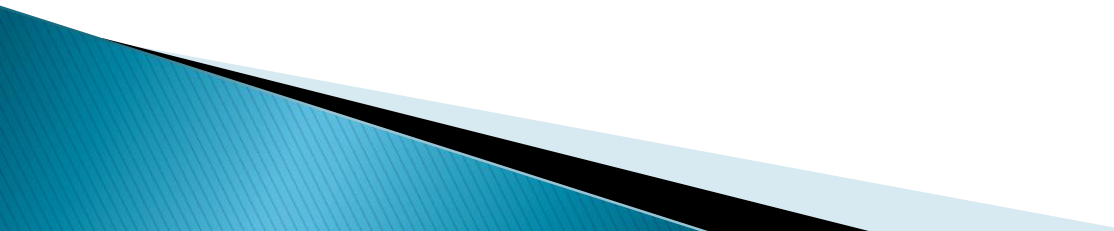
Example

Sum of row maxima = $9 + 8 + 9 + 6 = 32$

Sum of column maxima = $8 + 7 + 9 + 9 = 33$

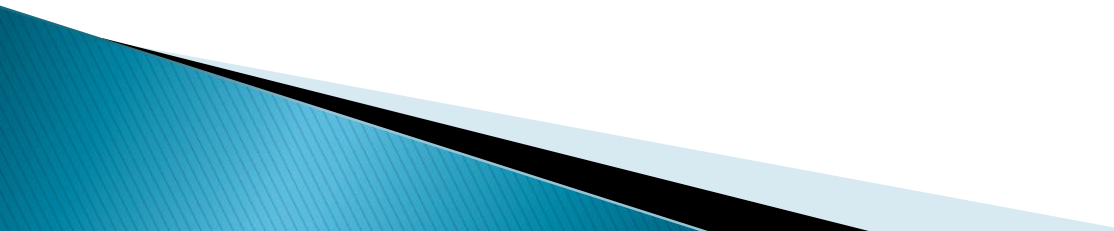


Calculating U and V

1. Find the maximum value of each row
 2. Find the maximum value of each column
 3. Sum the maximum values of every row
 4. Sum the maximum values of every column
 5. Compare the two sums
 - If the sum of the row maxima is least, the maximum values of each row are assigned to U
 - Conversely, if the sum of the column maxima is least, the maximum values of each column are assigned to V
- 

Parallelization

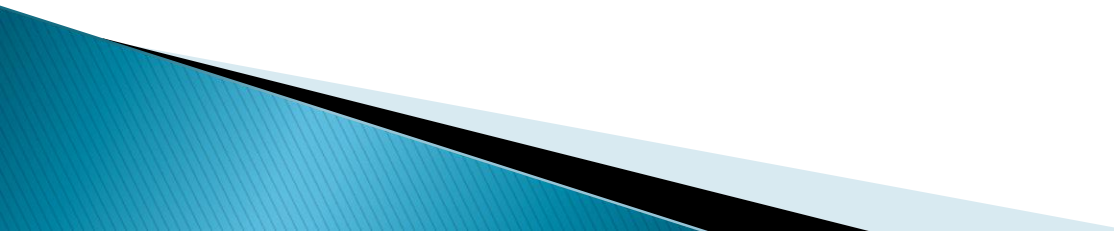
Calculating U and V

1. Find the maximum value of each row
 2. Find the maximum value of each column
 3. Sum the maximum values of every row
 4. Sum the maximum values of every column
 5. Compare the two sums
 - If the sum of the row maxima is least, the maximum values of each row are assigned to U
 - Conversely, if the sum of the column maxima is least, the maximum values of each column are assigned to V
- 

Finding Row Maxima

//Launch as many kernels as there are rows

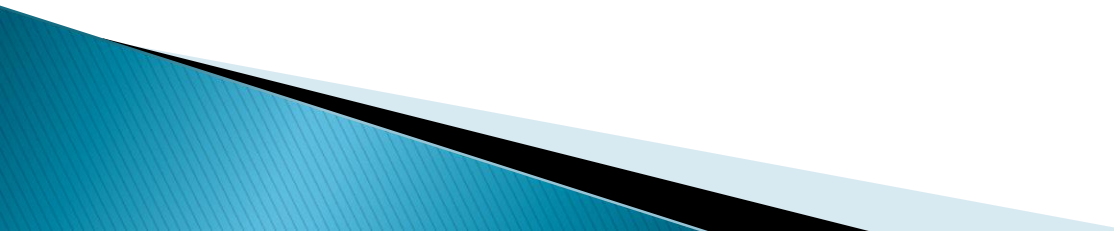
```
__global__ void FindMaxRow(int *d_r, int *d_row_max, int num_rows, int num_col) {  
    int currentRow = blockIdx.x*blockDim.x + threadIdx.x;  
    if (currentRow < num_rows) {  
        int currentCol;  
        //loop through columns  
        for (currentCol = 0; currentCol<num_col; currentCol++) {  
            //check to see if we have a new leader for the row  
            if (d_row_max[currentRow] < d_r[currentRow*num_col + currentCol]) {  
                d_row_max[currentRow] = d_r[currentRow*num_col + currentCol];  
            }  
        }  
    }  
}
```



Finding Column Maxima

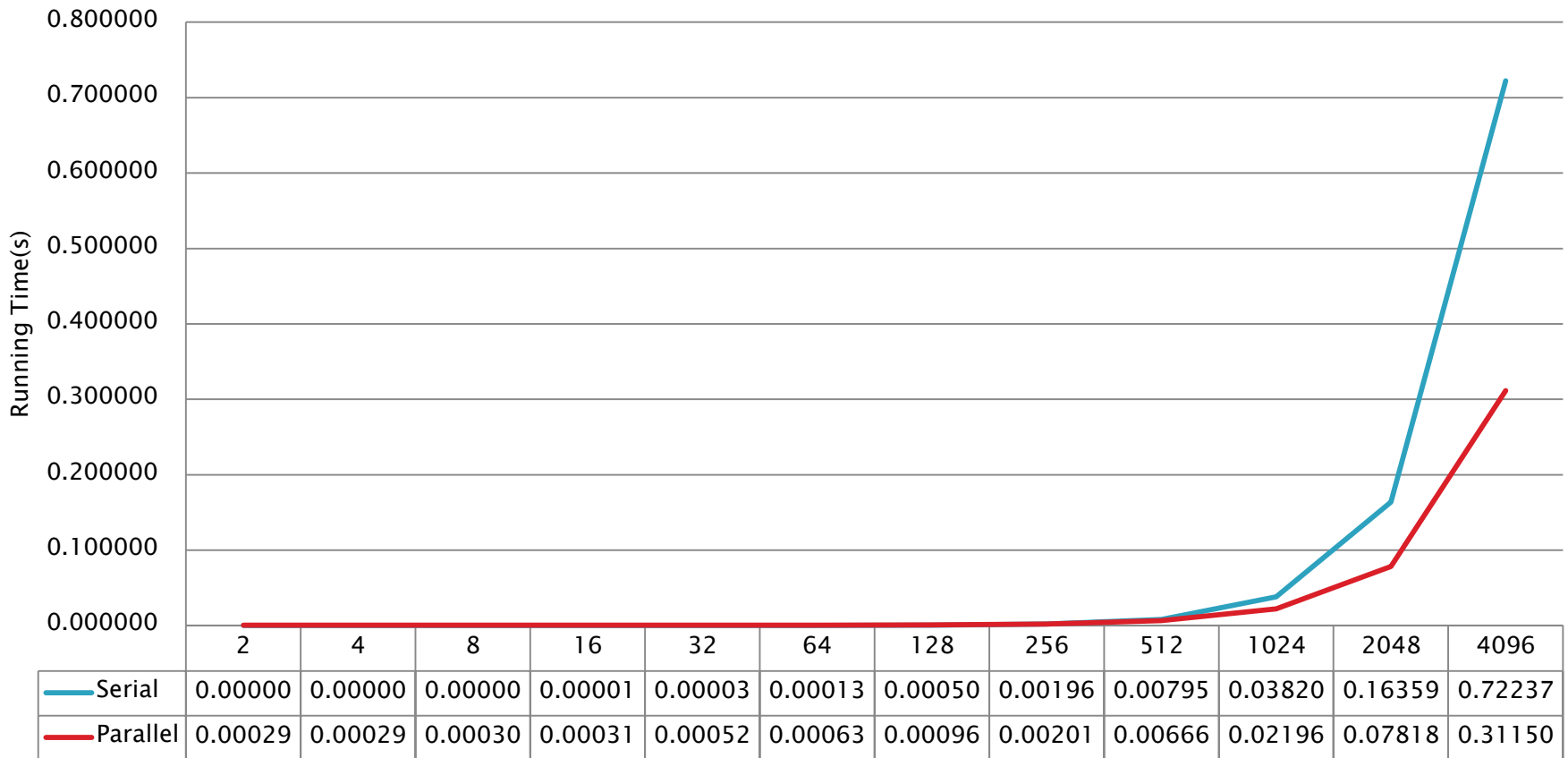
//Launch as many kernels as there are columns

```
__global__ void FindMaxCol(int *d_r, int *d_col_max, int num_rows, int num_col) {  
    int currentCol = blockIdx.x*blockDim.x + threadIdx.x;  
    if (currentCol < num_col) {  
        int currentRow;  
        //loop through rows  
        for (currentRow = 0; currentRow<num_row; currentRow++) {  
            //check to see if we have a new leader for the column  
            if (d_col_max[currentCol] < d_r[currentRow*num_col + currentCol]) {  
                d_col_max[currentCol] = d_r[currentRow*num_col + currentCol];  
            }  
        }  
    }  
}
```

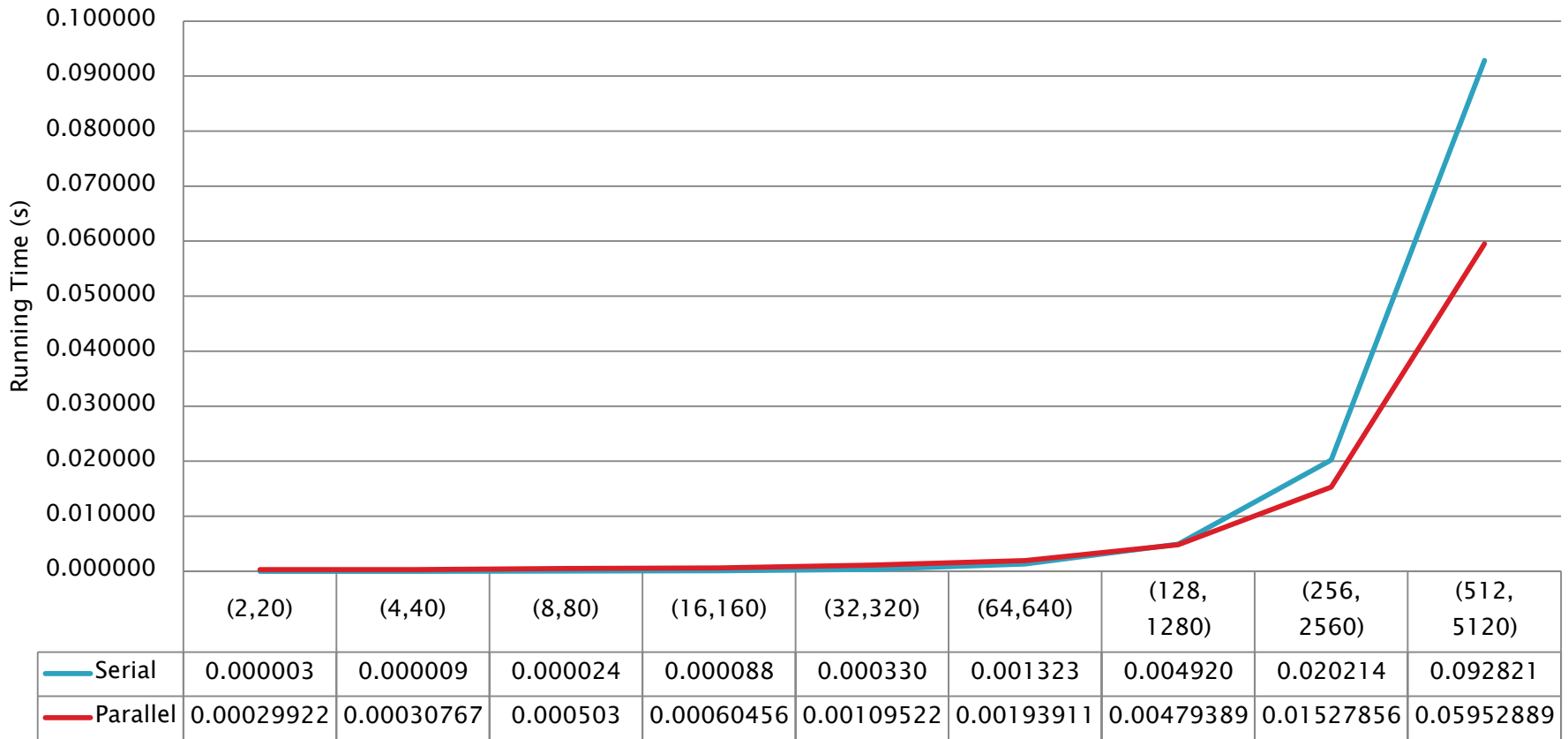


Results

Running Time vs. Input Size



Running Time vs. Input Size



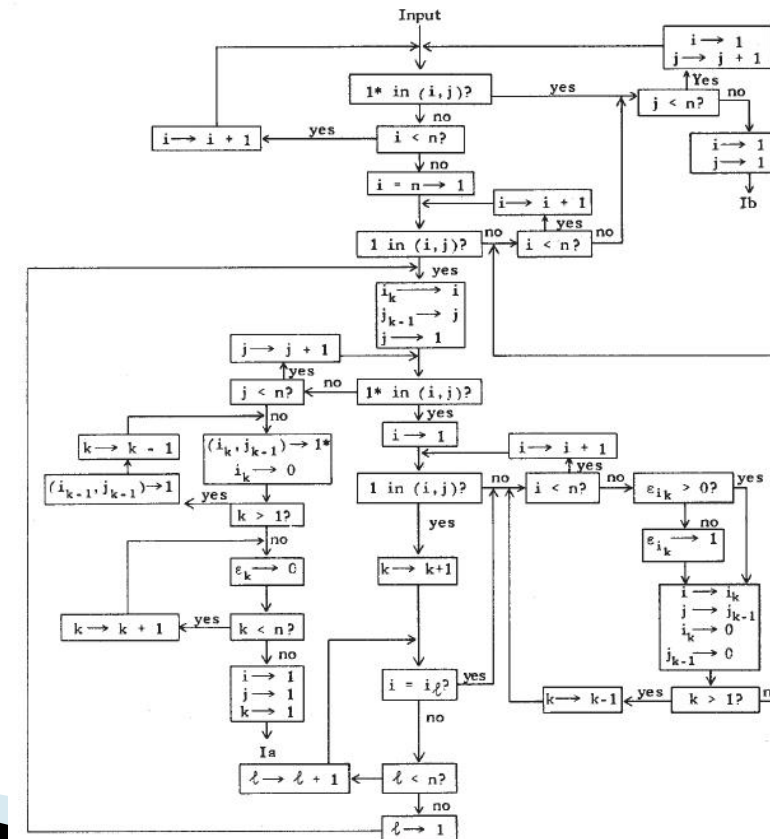
Limitations and Future Directions

Limitations

- ▶ No use of GPU shared memory
- ▶ Input size

Future Directions

- ▶ Use GPU to perform reductions
- ▶ Parallelize the transfer algorithm



References

- Bruff, D. (2005). The assignment problem and the hungarian method. Retrieved from:
http://www.math.harvard.edu/archive/20_spring_05/handouts/assignment_overheads.pdf
- Kuhn, H. W. (1955). The hungarian method for the assignment problem. 83–97.
- Kuhn, H. W. (1960). The hungarian method.