# Image processing techniques and their parallel performance comparisons using CUDA

Luke Unterman

# Problem:

- Identifying and conducting analysis on the various components within an image is the cornerstone of different computer vision applications
  - Optical character recognition
  - Medical imaging
  - Obstacle detection to improve automotive safety
- The dimensionality requirements for image analysis applications makes using the CPU infeasible
  - 4k image RGB image dimensions: 3840x2160x3=**24,883,200** pixels
- This project compares the GPU and CPU performance of several image processing techniques within a pipeline in order to produce the distinct segments.



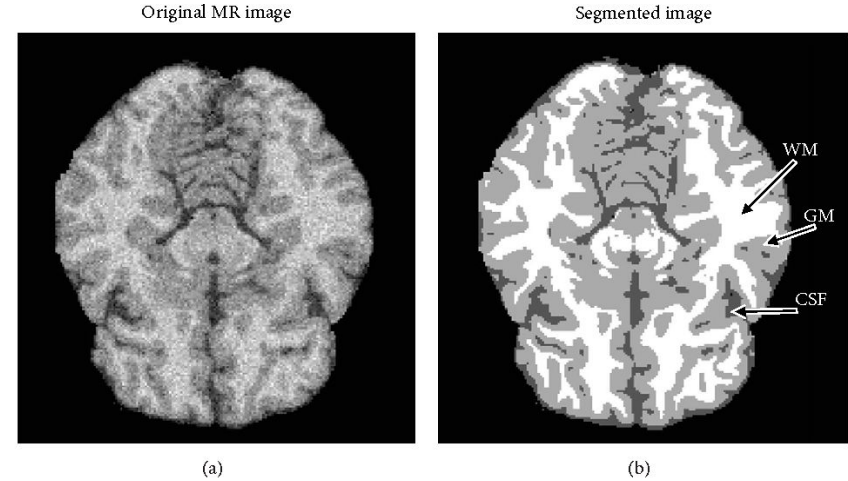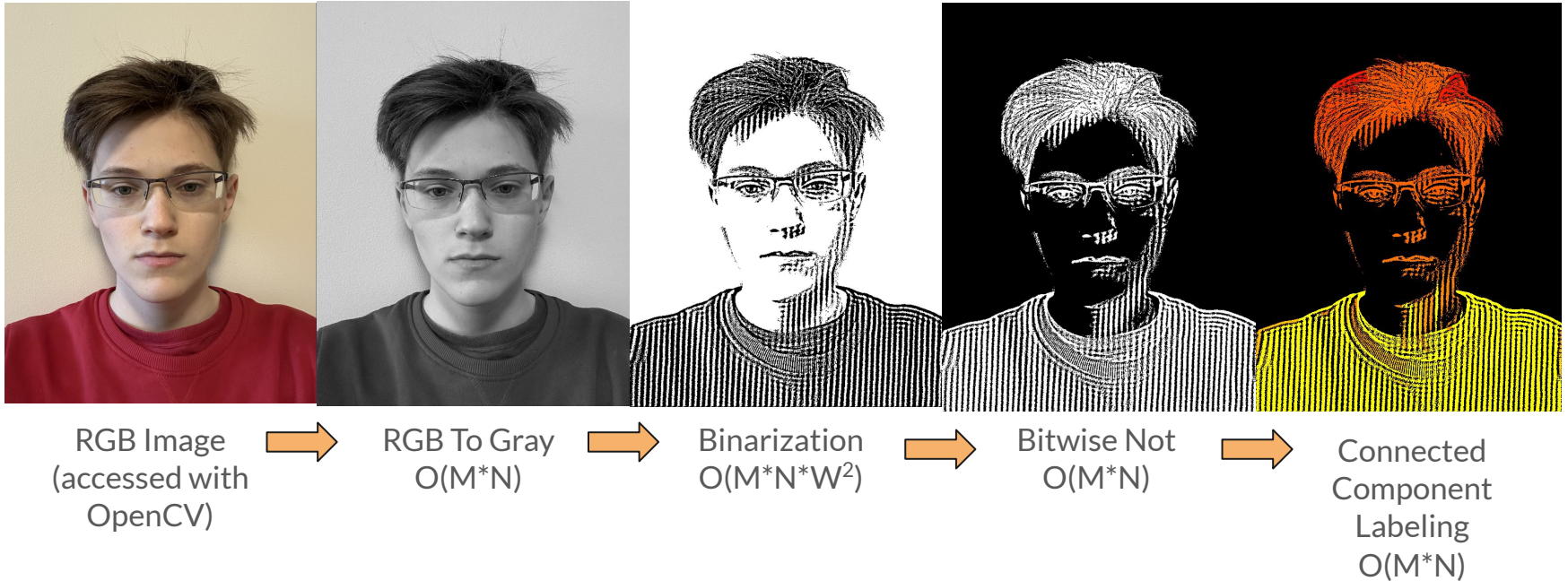Original MR image        Segmented image

(a)                (b)

Figure 1. Identifying white matter, grey matter, and cerebrospinal fluid in brain. [1]
- Greater WM volume associated with progression of Alzheimer's disease [2]

# Image Segmentation Algorithm Pipeline



RGB Image (accessed with OpenCV) → RGB To Gray $O(M*N)$ → Binarization $O(M*N*W^2)$ → Bitwise Not $O(M*N)$ → Connected Component Labeling $O(M*N)$

# Most of these img. processing algorithms are embarrassingly parallel

**Algorithm 1** SauvolaBinary

**Require:** Grayscale image, window size, parameter $k$, and dynamic range $R$

**Ensure:** Converts the image into a binary image using Sauvola's thresholding

1: **for** each pixel in the image **do**
2:     Define a window centered around the pixel
3:     Compute the mean intensity within the window
4:     Compute the standard deviation of intensities within the window
5:     Calculate the threshold:

$$\text{threshold} = \text{mean} \times \left(1 + k \times \left(\frac{\text{std\_dev}}{R} - 1\right)\right)$$

6:     **if** pixel intensity > threshold **then**
7:       Set pixel to white (255)
8:     **else**
9:       Set pixel to black (0)
10:     **end if**
11: **end for**

**Algorithm 2** BitwiseNot

**Require:** Grayscale image

**Ensure:** Negates each pixel value by performing a bitwise NOT operation

1: **for** each pixel in the image **do**
2:     Replace the pixel value with its bitwise NOT

$$\text{new\_value} = 255 - \text{current\_value}$$

3: **end for**

**Algorithm 3** PadImage

**Require:** Image, padding size, and padding value

**Ensure:** Creates a padded version of the image

1: **for** each pixel in the padded image **do**
2:     **if** pixel is in the padding region **then**
3:       Set pixel value to the padding value
4:     **else**
5:       Copy the corresponding pixel value from the input image, adjusted by the padding offset
6:     **end if**
7: **end for**

The nested for loops required for the M * N rows and columns can be removed via a 2-d CUDA block mapping

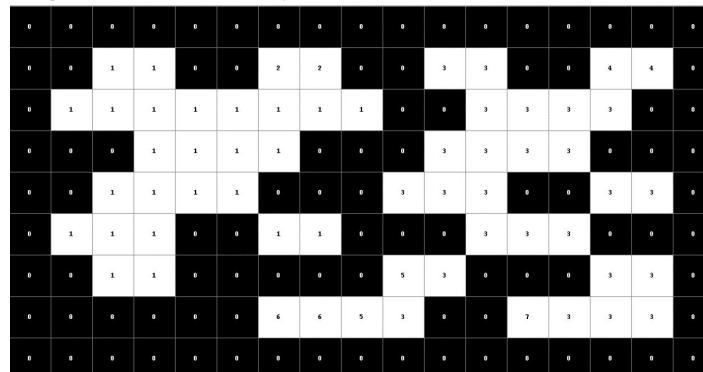# The Connected Component Labeling algorithm is much more complex...

**Algorithm 4** First Pass: Assign Provisional Labels and Record Equivalences

**Require:** Padded image `labelImage`, initial label index k = 1

**Ensure:** Assigns provisional labels and builds equivalence table

1: **for** each pixel $(row, col)$ in `labelImage` **do**
2:    **if** current pixel is black (0) **then**
3:       **continue**
4:    **end if**
5:    Get the labels of the 4-connected neighbors
6:    Filter out zero neighbors
7:    **if** no neighbors have labels **then**
8:       Assign a new label $k$ to the pixel and increment $k$
9:    **else**
10:      Assign the smallest neighbor label to the pixel
11:      Record equivalences between the assigned label and all neighbor labels
12:    **end if**
13: **end for**

Figure 2. First pass of two-pass connected component labeling algorithm. 4-connectivity.



4-connectivity

| Set ID | Equivalent Labels |
|--------|-------------------|
| 1 | 1,2 |
| 2 | 1,2 |
| 3 | 3,4,5,6,7 |
| 4 | 3,4,5,6,7 |
| 5 | 3,4,5,6,7 |
| 6 | 3,4,5,6,7 |
| 7 | 3,4,5,6,7 |

Note: CUDA implementation could require a global k value and global equivalence table

# The Connected Component Labeling algorithm is much more complex...

**Algorithm 4** Second Pass: Replace Provisional Labels with Resolved Labels

**Require:** Padded image `labelImage`, `flatEquivalenceMap`

**Ensure:** Updates labels in `labelImage` using resolved equivalences

1: **for** each pixel $(row, col)$ in `labelImage` **do**
2:     Get the label of the current pixel
3:     **if** label exists in `flatEquivalenceMap` **then**
4:         Replace label with resolved label
5:     **end if**
6: **end for**
7: Remove the padding from `labelImage`

There is an intermediate step of flattening the equivalence table. Note: the second pass is much more parallelizable.

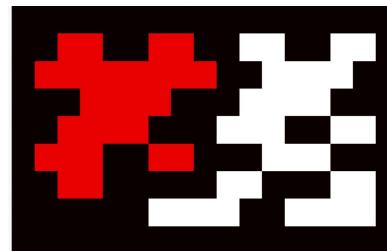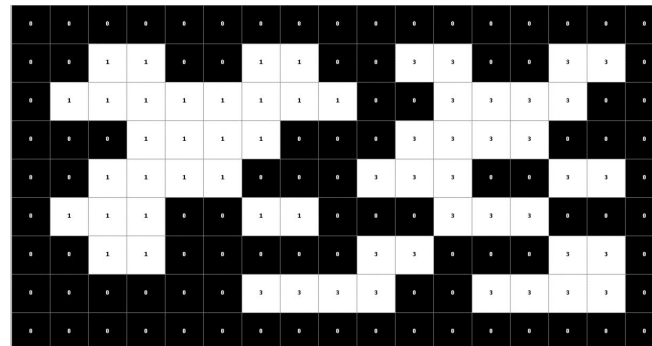Figure 3. Second pass of two-pass connected component labeling algorithm. 4-connectivity.





Figure 4. Image broken into distinct connected components

# Parallelization Approach

- To exploit the **GPU's** inherent speed when it comes to parallel **SIMD** computation, each of these algorithms will be parallelized with **CUDA**.
  - Higher instruction and memory bandwidth than the **CPU**
- The **same** parallelization procedure can be applied to almost every single one of the image processing algorithms in the pipeline

# General Parallelization Procedure:

1. Read in RGB image using OpenCV
2. Declare cv::Mat outputImage with same shape as input
3. Declare same variables on CUDA device, allocate memory, and copy inputImage memory from Host to Device
4. Devise blockDim and gridDim s.t. each thread is responsible for a pixel
5. Execute kernel to apply image transformation
6. Copy result from Device to Host
7. Save image to file

```
dim3 blockDim(threadsPerBlockPerDim, threadsPerBlockPerDim);
dim3 gridDim(gridDimx,gridDimy);

int column = blockIdx.x * blockDim.x + threadIdx.x;
int    row = blockIdx.y * blockDim.y + threadIdx.y;
```

Figure 5. Above: Designation of kernel block dimensions and grid dimensions. Below: assignment of unique index for each thread.

# Example CUDA kernels for GPU-based image processing

- All non-CCL algorithms apply the **same approach** to divide work
- Nested for loop to iterate over pixels removed
- **Highly parallel**



**Algorithm 1** Sauvola Binarization in CUDA
**Require:** $gray$: Grayscale image, $binary$: Output binary image, $width$, $height$: Dimensions of the image, $window\_size$: Window size, $k$: Sauvola parameter
1: $x = blockIdx.x \cdot blockDim.x + threadIdx.x$
2: $y = blockIdx.y \cdot blockDim.y + threadIdx.y$
3: **if** $x < width$ **and** $y < height$ **then**
4:   $half\_window = window\_size/2$
5:   $mean = 0.0,\ std\_dev = 0.0,\ count = 0$
6:   **for** $wy = -half\_window$ **to** $half\_window$ **do**
7:     **for** $wx = -half\_window$ **to** $half\_window$ **do**
8:       $nx = \min(\max(x + wx, 0), width - 1)$
9:       $ny = \min(\max(y + wy, 0), height - 1)$
10:       $pixel = gray[ny \cdot width + nx]$
11:       $mean = mean + pixel$
12:       $std\_dev = std\_dev + pixel \cdot pixel$
13:       $count = count + 1$
14:     **end for**
15:   **end for**
16:   $mean = mean/count$
17:   $std\_dev = \sqrt{\left(\frac{std\_dev}{count}\right) - (mean \cdot mean)}$
18:   $threshold = mean \cdot \left(1 + k \cdot \left(\frac{std\_dev}{R} - 1\right)\right)$
19:   **if** $gray[y \cdot width + x] > threshold$ **then**
20:     $binary[y \cdot width + x] = 255$
21:   **else**
22:     $binary[y \cdot width + x] = 0$
23:   **end if**
24: **end if**

**Algorithm 2** Connected Components Kernel (CUDA)
**Require:** $labelImage$: Image of labels, $equivalenceMap$: Map for label equivalences, $numberRows$: Number of rows, $numberColumns$: Number of columns
1: $column = blockIdx.x \cdot blockDim.x + threadIdx.x$
2: $row = blockIdx.y \cdot blockDim.y + threadIdx.y$
3: $idx = row \cdot numberColumns + column$
4: **if** $row < numberRows$ **and** $column < numberColumns$ **and** $column \geq 1$ **and** $row \geq 1$ **then**
5:   $label = labelImage[idx]$
6:   **if** $equivalenceMap[label] > 0$ **then**
7:     $labelImage[idx] = equivalenceMap[label]$
8:   **end if**
9: **end if**

**Algorithm 3** Connected Components Serial Code (Host)
1: $ccTuple = CCInit(inputImage, paddedImage)$
2: $mapMax = ccTuple[0]$
3: $equivalenceMap = ccTuple[1]$
4: $flatEquivalences = \text{vector}(mapMax + 1, 0)$
5: **for** each $(key, value)$ in equivalenceMap **do**
6:   $flatEquivalences[key] = value$
7: **end for**
8: Allocate memory to d_equivalence on device
9: Copy flatEquivalences memory on host to d_equivalen
10: Execute GPU ConnectedComponents kernel
11: Copy labelImage from device to outputImage on host

- Connected components labeling algorithm broken into **serial** and **parallel** components
- 1st-pass done on Host, 2nd-pass done on Device
- **Naive**, but does provide speedup

The nested for loops required for the M * N rows and columns can be removed via a 2-d CUDA block mapping`

# Performance Evaluation of CPU and GPU methods

- Runtime (ms) recorded for serial and parallel image processing methods
- Constant image "selfie.jpg" with original dimensions 3088 x 2316 x 3 reduced by different factors (1, 2, 4, 8, and 16)
- Speedup and Scalability calculated for each of the parallel image processing algorithms
- Speedup defined as $\frac{T_s}{T_p}$ where $T_s$ = time spent running algorithm serially and $T_p$ = time spent running algorithm in parallel.
- Scalability defined as $\frac{T_{p=1}}{T_{p=N}}$ where p=1 implies threadsPerDimPerBlock = 1, and p=N implies threadsPerDimPer Block=N.
- Constant threadsPerDimPerBlock value of 16 for the purposes of this experiment.
- All algorithmic executions performed on **cs.maple.vcu.edu** using NVIDIA GeForce 1080 Ti

# Results

- Time required for both CPU and GPU implementations of each image operation plotted
- Time(ms) for each operation ranked in increasing order:
    1. Bitwise Flip
    2. Padding
    3. Grayscale
    4. Sauvola Thresholding
    5. CCL
- Operations 1-3 clustered closely together (parallel, same time complexity)
- Operation 4 is 2nd highest due to increased time complexity (iteration over window size)
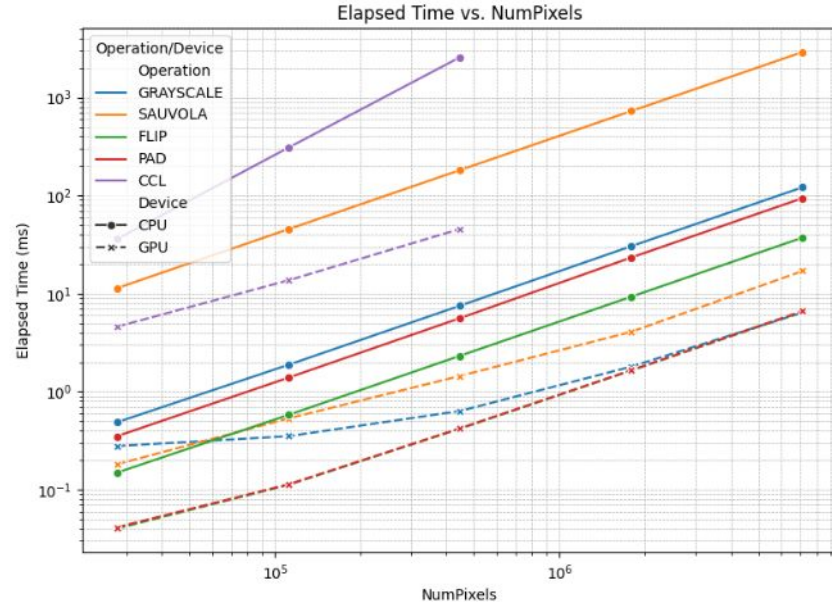- Operation 5: Same time complexity as 1-3, but reliant on serial pass.



Fig. 1. Plot of elapsed time spent running CPU and GPU algorithms over the number of pixels in inputImage used. Note: log x, y, scale.

# Results

- **Scalability** plotted over numPixels.

- Due to reliance on serial section, **CCL is least scalable** (no change despite increase in image size)

- Flip, and Padding algorithms mildly scalable , with peak at dimension reduction by scale of 8

- Sauvola thresholding algorithm highly scalable due to increased work

- Grayscale operation **exponentially scalable** (log plot) with increasing image size
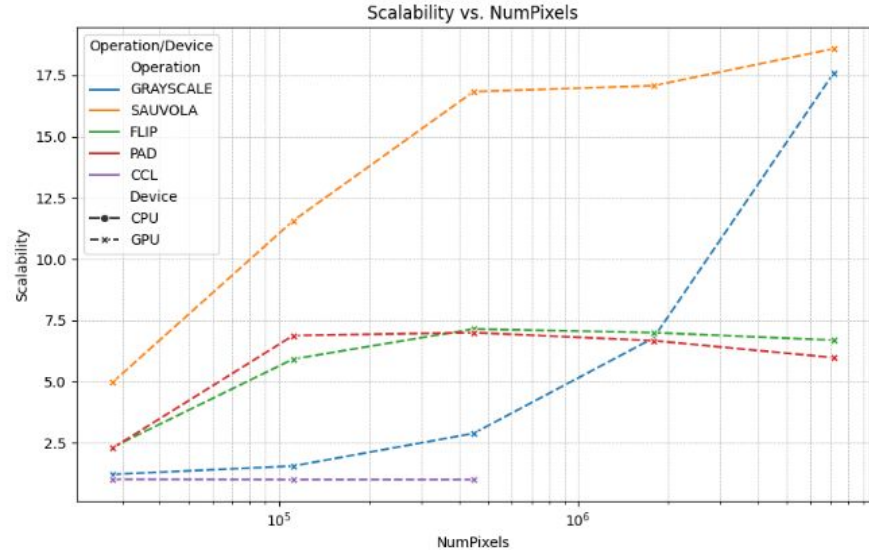


Fig. 2. Plot of the scalability of the parallel image processing operations at threadsPerNumPerDim=1,16 over number of pixels in inputImage. Note: log x scale.

# Results

- **Speedup** plotted over numPixels.

- **Sauvola** thresholding GPU implementation maintains highest speedup value due to **increased proportion** of GPU work over **overhead** ($O(W^2)$ vs $O(MNW^2)$).

- Grayscale algorithm's speedup is slightly higher than the Pad/Flip algorithms due to **additional work** dealing with 3 image channels.

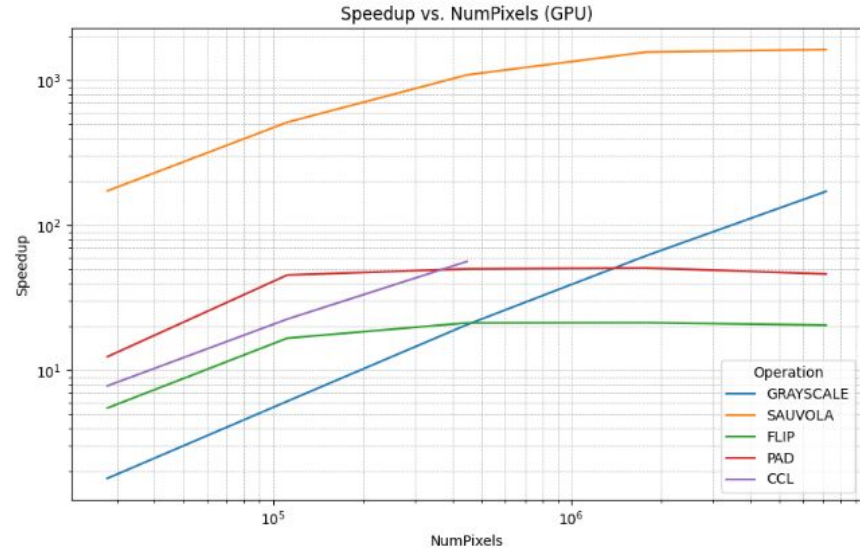- CCL, **despite serial sections**, achieves moderate speedup.



Fig. 3. Plot of the speedup of each of the parallel image-processing algorithms on the GPU over number of pixels in inputImage. Note: log x, y, scale.

# Conclusions

1. This project compared the **runtime**, **speedup**, and **scalability** values of several different image processing techniques on the GPU and CPU (**ranked in increasing order** of runtime in ms).
   a. Bitwise Not <
   b. Image Padding <
   c. Grayscale Conversion <
   d. Binarization/Thresholding <<
   e. Connected Component Labeling
2. In particular, the **Sauvola thresholding binarization** algorithms performs **quite well** on the GPU due to its large dimensionality reduction ($O(MNW^2)$ => $O(W^2)$).
3. Despite the serial sections in the parallel CCL implementation, it still achieves moderate speedup.
4. Image processing algorithms in general are fairly easy to implement on the GPU

# Citations

- Brain scan:
  - Heckemann, R. A., Keihaninejad, S., Aljabar, P., Rueckert, D., Hajnal, J. V., Hammers, A., May 2010. Improving intersubject image registration using tissue-class information benefits robustness and accuracy of multi-atlas based anatomical segmentation. NeuroImage 51 (1), 221-227. http://dx.doi.org/10.1016/j.neuroimage.2010.01.072
  - Garnier-Crussard, A., Bougacha, S., Wirth, M., Dautricourt, S., Sherif, S., Landeau, B., Gonneaud, J., De Flores, R., de la Sayette, V., Vivien, D., Krolak-Salmon, P., & Chételat, G. (2022). White matter hyperintensity topography in Alzheimer's disease and links to cognition. *Alzheimer's & dementia : the journal of the Alzheimer's Association*, *18*(3), 422–433. https://doi.org/10.1002/alz.12410
- Connected component labeling images:
  - Wikipedia contributors. (2024, November 5). Connected-component labeling. In Wikipedia, The Free Encyclopedia. Retrieved 02:41, December 10, 2024, from https://en.wikipedia.org/w/index.php?title=Connected-component_labeling&oldid=1255521016