



HAL
open science

Designing efficient SIMD algorithms for direct Connected Component Labeling

Arthur Hennequin, Ian Masliah, Lionel Lacassagne

► **To cite this version:**

Arthur Hennequin, Ian Masliah, Lionel Lacassagne. Designing efficient SIMD algorithms for direct Connected Component Labeling. WPMVP'19 - 5th Workshop on Programming Models for SIMD/Vector Processing, Feb 2019, Washington, United States. pp.4:1–4:8, 10.1145/3303117.3306164 . hal-02049029

HAL Id: hal-02049029

<https://hal.science/hal-02049029>

Submitted on 7 May 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Designing efficient SIMD algorithms for direct Connected Component Labeling

Arthur Hennequin^{1,2}, Ian Masliah¹, Lionel Lacassagne¹

firstname.name@lip6.fr

LIP6¹ – Sorbonne University, CNRS

LHCb experiment² – CERN

ABSTRACT

Connected Component Labeling (CCL) is a fundamental algorithm in computer vision, and is often required for real-time applications. It consists in assigning a unique number to each connected component of a binary image. In recent years, we have seen the emergence of direct parallel algorithms on multicore CPUs, GPUs and FPGAs whereas, there are only iterative algorithms for SIMD implementation. In this article, we introduce new direct SIMD algorithms for Connected Component Labeling. They are based on the new Scatter-Gather, Collision Detection (CD) and Vector Length (VL) instructions available in the recent Intel AVX512 instruction set. These algorithms have also been adapted for multicore CPU architectures and tested for each available SIMD vector length. These new algorithms based on SIMD Union-Find algorithms can be applied to other domains such as graphs algorithms manipulating Union-Find structures.

1 INTRODUCTION

Connected Component Labeling (CCL) is a central algorithm between low-level image processing (filtering and pre-processing) and high-level image processing (scene analysis, pattern recognition and decision). CCL consists in providing a unique number to each connected component of a binary image. There are several applications in computer vision that use CCL such as Optical Character Recognition, motion detection or tracking.

CCL algorithms are often required to be optimized to run in real-time and have been ported on a wide set of parallel machines [1][13]. After an era on single-core processors, where many sequential algorithms were developed [7] and codes were released [2], new parallel algorithms were developed on multicore processors [15][6].

The majority of CCL algorithms for CPUs are direct, by opposition to iterative ones – where the number of iterations (the number of image pass to process it) depends on the image structure – and require only 2 passes thanks to an equivalence table.

The first algorithms designed for GPUs [10][9][5] and SIMD [19][12] were iterative. Furthermore, the number of iterations required to reach a solution cannot be predicted beforehand and can reach a large number. In recent years, some direct algorithms for GPUs were designed and are better suited for real-time image processing in embedded systems [11][3][16][8].

In this article, we introduce a new direct CCL algorithm for SIMD processors which relies on *scatter-gather* operations. These algorithms are only available using Intel AVX512 extension and the

upcoming ARM/Fujitsu SVE.

Section 2 explains algorithmic design and the required SIMD instructions. Section 3 describes the classic Rosenfeld algorithm [18] and its derivatives, as well as concurrency issues that can appear using SIMD. Section 4 introduces two new SIMD algorithms with implementation details. The first one is based on a classic pixel-based approach, while the second one is based on sub-segments and uses bit manipulation and conflict detection to reduce memory accesses. Section 5 presents the benchmark protocol for reproducible experiments, the results and their analysis.

2 ALGORITHMIC DESIGNS TO EXPLOIT ARCHITECTURES

As clock frequencies of modern processors are expected to stay near their current levels, or even to get lower, the primary method to improve the computational capability of a chip is to increase either the number of processing units (cores) or the intrinsic parallelism of a core (SIMD). In this work, we target the most recent Intel architecture with the AVX512 SIMD extension set. The AVX512F extension has added the scatter instruction for each available vector size (128, 256, 512) which is required to develop an SIMD Rosenfeld algorithm. We have used the instructions *compress* and *expand* which are available in the AVX512VL extension (algo 8). From AVX512CD, we also used the *conflict detection* (CD) and *lzcnt* (count leading zeros) instructions (algo 11).

Each of these instructions are required to implement efficient SIMD CCL algorithms and are only available in the recent Intel AVX512 extension set. Developing new algorithms that use the architecture efficiently requires heavy algorithmic modifications and can only be done efficiently if the correct instructions are available. Due to this complex process, it is impossible for a compiler to provide any vectorization support. Domain Specific Languages (DSLs) such as Halide [17] would also require additional modifications from the DSL and the application programmer due to the new available instructions.

In this context, we developed a set of AVX512 C++ templates code that can be instantiated in 128-bit, 256-bit and 512-bit in order to evaluate the impact of SIMD width on the performance. The required sub-sets of AVX512 are AVX512F, AVX512CD and AVX512VL. We have further extended our portable template code to support the new ARM SVE instruction set but the code has not been tested in the scope of this paper.

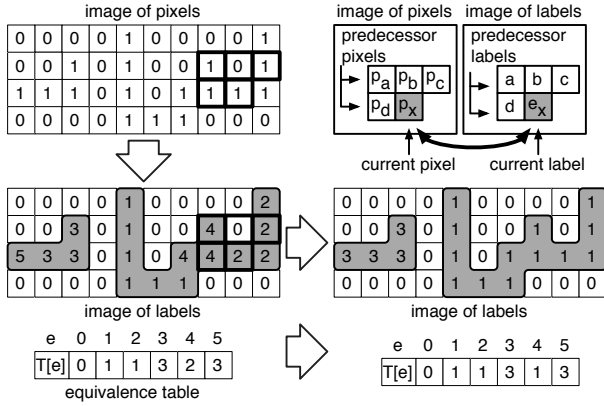


Figure 1: Example of 8-connected CCL with Rosenfeld algorithm: binary image (top), image of temporary labels (bottom left), image of final labels (bottom right) after the transitive closure of the equivalence table.

3 CLASSIC ALGORITHMS

Usually, CCL algorithms are split into three steps and perform two image scans (like the pioneer algorithm of Rosenfeld [18]). The first scan (or first labeling) assigns a temporary/provisional label to each connected component (CC) and some label equivalences are built if needed. The second step solves the equivalence table by computing the transitive closure (TC) of the graphs (in fact a forest) associated to the label equivalences. The third step performs a second scan (or second labeling) that replaces temporary labels of each CC by its final label.

Figure 1 defines some notations and gives an example of a classic Rosenfeld algorithm execution. Let p_x , e_x , the current pixel and its label. Let p_a, p_b, p_c, p_d , the neighbor pixels, and a, b, c, d , the associated labels. T is the equivalence table, e a label and r its root. The first scan of Rosenfeld is described in algorithm 3, the transitive closure in algorithm 4, while the classical union-find algorithms are provided in algorithms 1 & 2. In the example (Fig. 1), we can see that the rightmost CC requires three labels (1, 2 and 4). When the mask is in the position seen in figure 1 (in bold type), the equivalence between 2 and 4 is detected and stored in the equivalence table T . At the end of the first scan, the equivalence table is complete and applied to the image (like a Look-Up-Table).

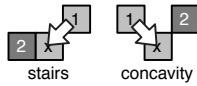


Figure 2: Basic patterns generating an equivalence: stairs & concavity

3.1 Algorithmic and SIMD optimisation

Decision Tree (DT) based algorithms for CCL [20] have been proved to be very efficient to enhance scalar implementations. The DT reduces the number of Union and Find function calls to the strict minimum, i.e. when there is an equivalence between two different

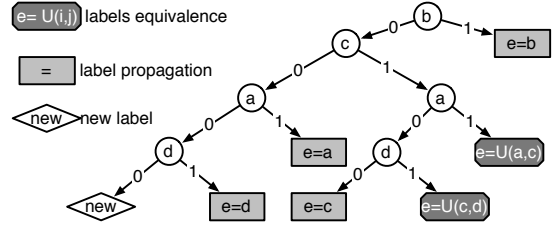


Figure 3: 8-connected Decision Tree for a 4-pixel mask. Labels equivalence (call to Union) in dark gray

Algorithm 1: Find(e, T)

Input: e a label, T an equivalence table
Result: r , the root of e

- 1 $r \leftarrow e$
- 2 **while** $T[r] \neq r$ **do**
- 3 $r \leftarrow T[r]$
- 4 **return** r

Algorithm 2: Union(e_1, e_2, T)

Input: e_1, e_2 two labels, T an equivalence table
Result: r , the least common ancestor of the e 's

- 1 $r_1 \leftarrow \text{Find}(e_1)$
- 2 $r_2 \leftarrow \text{Find}(e_2)$
- 3 **if** $r_1 < r_2$ **then**
- 4 $r \leftarrow r_1, T[r_2] \leftarrow r$
- 5 **else**
- 6 $r \leftarrow r_2, T[r_1] \leftarrow r$
- 7 **return** r

Algorithm 3: Rosenfeld algorithm – first labeling (step 1)

Input: a, b, c, d , four labels, p_x , the current pixel in (i, j)

- 1 **if** $p_x \neq 0$ **then**
- 2 $a \leftarrow E[i-1][j-1], \quad b \leftarrow E[i-1][j]$
- 3 $c \leftarrow E[i-1][j+1], \quad d \leftarrow E[i][j-1]$
- 4 **if** $(a = b = c = d = 0)$ **then**
- 5 $ne \leftarrow ne + 1, \quad e_x \leftarrow ne$
- 6 **else**
- 7 $r_a \leftarrow \text{Find}(a, T), \quad r_b \leftarrow \text{Find}(b, T)$
- 8 $r_c \leftarrow \text{Find}(c, T), \quad r_d \leftarrow \text{Find}(d, T)$
- 9 $e_x \leftarrow \min^+(r_a, r_b, r_c, r_d)$
- 10 **if** $(r_a \neq 0 \text{ and } r_a \neq e_x)$ **then** $\text{Union}(e_x, r_a, T)$
- 11 **if** $(r_b \neq 0 \text{ and } r_b \neq e_x)$ **then** $\text{Union}(e_x, r_b, T)$
- 12 **if** $(r_c \neq 0 \text{ and } r_c \neq e_x)$ **then** $\text{Union}(e_x, r_c, T)$
- 13 **if** $(r_d \neq 0 \text{ and } r_d \neq e_x)$ **then** $\text{Union}(e_x, r_d, T)$
- 14 **else**
- 15 $e_x \leftarrow 0$

labels. There are only two patterns generating an equivalence between labels: stairs and concavities which are depicted in figure 2. DT is more efficient than path-compression as it reduces the number of memory accesses. Considering the classical implementation of the Rosenfeld algorithm (algo. 3), there are four calls to Find and one to Union. The calls to find in the union can be omitted because the input labels are already equivalence trees' roots. Using a DT (algo. 5), there are at most one call to Union and two calls to Find as we have at most one equivalence between labels.

If we consider the extended topology of the mask in the SIMD Union Find (Fig. 4), there are, for a SIMD vector of cardinality $card$,

Algorithm 4: Sequential solve of equivalences (step 2)

```

1 for  $e \in [1 : ne]$  do
2    $T[e] \leftarrow T[T[e]]$ 

```

Algorithm 5: Rosenfeld with DT – optimized first labeling (step1)

Input: a, b, c, d , four labels, p_x , the current pixel in (i, j)

```

1 if  $p_x \neq 0$  then
2    $b \leftarrow E[i-1][j]$ 
3   if  $(b \neq 0)$  then
4      $e_x \leftarrow b$ 
5   else
6      $c \leftarrow E[i-1][j+1]$ 
7     if  $(c \neq 0)$  then
8        $a \leftarrow E[i-1][j-1]$ 
9       if  $(a \neq 0)$  then
10         $e_x \leftarrow U(a, c)$ 
11      else
12         $d \leftarrow E[i][j-1]$ 
13        if  $(a \neq 0)$  then
14           $e_x \leftarrow U(c, d)$ 
15        else
16           $e_x \leftarrow c$ 
17      else
18         $a \leftarrow E[i-1][j-1]$ 
19        if  $(a \neq 0)$  then
20           $e_x \leftarrow a$ 
21        else
22           $d \leftarrow E[i][j-1]$ 
23          if  $(a \neq 0)$  then
24             $e_x \leftarrow d$ 
25          else
26             $n_e \leftarrow n_e + 1$ 
27             $e_x \leftarrow n_e$ 
28 else
29    $e_x \leftarrow 0$ 

```

at most $card/2 + 1$ different labels, $card/2$ calls to Union, and thus, $card$ calls to Find, with $card/2 - 1$ which are redundant. If we take a vector with a cardinal of 8, we get $1/2$ Union per pixel and one call to Find. In the classical Union Find approach, we first perform a call to $Union(e_0, e_1)$ which calls $Find(e_0)$ and $Find(e_1)$. Then, we call $Union(e_1, e_2)$ which calls $Find(e_1)$ and $Find(e_2)$, and so on, until a call to $Union(e_3, e_4)$ which calls $Find(e_3)$ and $Find(e_4)$. As there can be more than one Union in an SIMD vector, concurrency issues can appear.

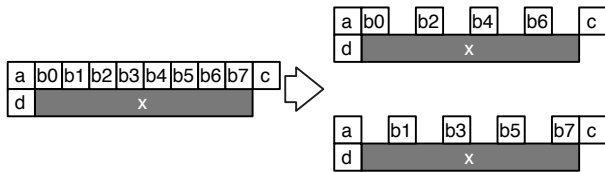


Figure 4: SIMD Union Find: number of different labels is equal to $card/2 + 1$ (here $card = 8$)

3.2 Algorithmic and SIMD concurrency issue

Depending on the neighboring labels of a pixel, concurrency issues can appear while performing simultaneous union operations. Figure 5 provides some examples (and the associated graphs) built with the basic patterns of figure 2. In the first row, there are no concurrency issues. There is one CC on the left ($T[3] \leftarrow 1, T[4] \leftarrow 1$)

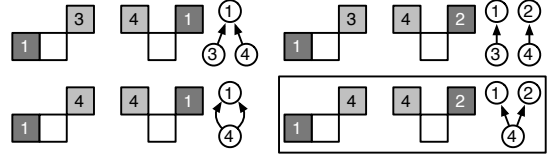


Figure 5: Composition of basic patterns leading to redundant accesses (bottom left) and concurrency issue (bottom right)

and two CCs on right ($T[3] \leftarrow 1, T[4] \leftarrow 2$).

In the second row, there is one CC on the left with twice the same equivalence ($T[4] \leftarrow 1, T[4] \leftarrow 1$) which is redundant but not an issue. However, on the right happens a concurrency, as we have $4 \equiv 1$ and $4 \equiv 2$. A possible correction could be $2 \equiv 1$ ($T[4] \leftarrow 1, T[2] \leftarrow 1$). Note that this problem can happen starting from a 3-pixel wide SIMD register if the two patterns are merged. The SIMD implementation of Union Find has to address these concurrency issues, which is – by far – non trivial to solve and to implement efficiently in SIMD.

4 NEW ALGORITHMS

4.1 SIMD Union-Find

The biggest challenge when designing an SIMD CCL algorithm is to design a fast and concurrency-free union-find algorithm to manage equivalences. The union algorithm must take into account the conflicts from multiple equivalence tree roots involved in simultaneous merge operations. Figure 6 shows an example of such complex merges. The top row shows the evolution of the root label pointers and the bottom row shows the pending merge operations in black and the finished merge operations in grey.

Init	Step 1	Step 2	Step 3
$3 \rightarrow 2$ $4 \rightarrow 1$ $4 \rightarrow 3$	$3 \rightarrow 3$ $4 \rightarrow 4$ $3 \rightarrow 1$	$3 \rightarrow 3$ $4 \rightarrow 4$ $2 \rightarrow 1$	$3 \rightarrow 3$ $4 \rightarrow 4$ $2 \rightarrow 2$

Figure 6: Execution of algorithm 7 *VecUnion* with arguments $\vec{e}_1 = [3, 1, 2]$, $\vec{e}_1 = [4, 4, 3]$ (example of simultaneous unions in 3 steps and their serialization).

Algorithms presented in this section and the followings are written for a cardinality of 8 (for the sake of clarity) but can easily be extended to 4 or 16 elements. Unmasked equalities and inequalities between vectors are tested using the intrinsics *cmpeq_epi32_mask* and *cmpneq_epi32_mask*, but are written as mathematical comparison to be more readable. Masked comparisons are expressed using their corresponding intrinsic.

Algorithm 6 uses gather loads to find the roots of all labels in \vec{e} . The loop must run until all roots have been found, so the number of iterations is equal to the maximum distance between involved

labels and their roots.

Algorithm 6: VecFind(\vec{e}, T, m)

Input: \vec{e} a vector of label, T an equivalence table, m a mask
Result: \vec{r} , the roots of \vec{e}

```

1  $\vec{r} \leftarrow \vec{e}$ 
2  $done \leftarrow 0$  // mask
3 while  $done \neq m$  do
4    $\vec{l} \leftarrow \text{mask\_i32gather\_epi32}(\vec{r}, \neg done \wedge m, \vec{r}, T, \text{sizeof}(\text{uint32}))$ 
5    $done \leftarrow \text{mask\_cmpeq\_epi32\_mask}(\vec{l}, \vec{r}, m)$  // mask
6    $\vec{r} \leftarrow \vec{l}$ 
7 return  $\vec{r}$ 

```

Algorithm 7 is inspired by the Playne-Equivalence reduce function designed for parallel CCL on GPU [16]. The main difference with a GPU algorithm is that we have to consider the parallelism within an SIMD vector instead of the parallelism between GPU threads in memory. To solve the concurrency issue the same way the GPU does, we introduce the *VecScatterMin* function that emulates the behavior of the CUDA *atomicMin* function. This function takes two vectors \vec{idx} and \vec{val} and tries to perform the operation: $T[\vec{idx}] \leftarrow \min(T[\vec{idx}], \vec{val})$. It then returns the old value of $T[\vec{idx}]$ if the operation succeeded or the current value if another vector element has written it first. Using this function, only one value is written to a memory address at a time, but it is always the minimum value of the concurrent store operations. This allows the *VecUnion* operation to retry until all involved equivalence trees have been merged. Because of the pixel topology, there can be at most $card/2$ simultaneous equivalence tree merges. In practice, the merge vectors are very sparse, allowing us to reduce the numbers of operations needed by compressing the vectors. The *VecScatterMin* function is described in Algorithm 8.

Algorithm 7: VecUnion($\vec{e}_1, \vec{e}_2, T, m$)

Input: \vec{e}_1, \vec{e}_2 two vectors of labels, T an equivalence table, m a mask

```

1  $\vec{r}_1 \leftarrow \text{VecFind}(\vec{e}_1, T, m)$ 
2  $\vec{r}_2 \leftarrow \text{VecFind}(\vec{e}_2, T, m)$ 
3  $m \leftarrow \text{mask\_cmpeq\_epi32\_mask}(m, \vec{e}_1, \vec{e}_2)$  // mask
4 while  $m$  do
5    $\vec{r}_{max}, \vec{r}_{min} \leftarrow \text{max\_epu32}(\vec{r}_1, \vec{r}_2), \text{min\_epu32}(\vec{r}_1, \vec{r}_2)$ 
6    $\vec{r}_1, \vec{r}_2 \leftarrow \vec{r}_{max}, \vec{r}_{min}$ 
7    $\vec{r}_3 \leftarrow \text{VecScatterMin}(\vec{r}_1, \vec{r}_2, T, m)$ 
8    $done \leftarrow \text{mask\_cmpeq\_epi32\_mask}(m, \vec{r}_1, \vec{r}_3)$  // mask
9    $\vec{r}_1 \leftarrow \vec{r}_3$ 
10   $m \leftarrow \neg done \wedge m$  // mask

```

4.2 SIMD Rosenfeld pixel algorithm

Like its scalar counterpart, the SIMD Rosenfeld pixel algorithm (v1) is a two pass direct CCL algorithm. In order to simplify the algorithm as well as improving the memory footprint and the performance, we embed the equivalence table into the image. The label creation can now easily be done in parallel as the new label is equal to the linear address of the pixel plus 1 to differentiate the background: $i \times w + j + 1$, where (i, j) are the pixel coordinates and w the width of the image. This bijection also allows for a faster relabeling as it can be done during the transitive closure step.

Algorithm 8: VecScatterMin($\vec{idx}, \vec{val}, T, m$)

Input: \vec{idx}, \vec{val} two vectors of labels, T an equivalence table, m a mask
Result: \vec{r} , the old values of $T[\vec{idx}]$

```

1  $\vec{rotate} \leftarrow \text{set\_epi32}(0, 7, 6, 5, 4, 3, 2, 1)$ 
2  $\vec{idx}_c \leftarrow \text{maskz\_compress\_epi32}(m, \vec{idx})$ 
3  $\vec{val}_c \leftarrow \text{maskz\_compress\_epi32}(m, \vec{val})$ 
4  $n \leftarrow \text{popcnt\_u32}(m)$ 
5  $\vec{rotate} \leftarrow \text{maskz\_move\_epi32}(0xFFFF >> (17 - n), \vec{rotate})$ 
6  $\vec{idx}_r, \vec{val}_r \leftarrow \vec{idx}_c, \vec{val}_c$ 
7 for  $i \leftarrow 0$  to  $n$  do
8    $\vec{idx}_r \leftarrow \text{permute\_epi32}(\vec{rotate}, \vec{idx}_r)$ 
9    $\vec{val}_r \leftarrow \text{permute\_epi32}(\vec{rotate}, \vec{val}_r)$ 
10   $\text{same\_addr} \leftarrow \vec{idx}_c = \vec{idx}_r$  // mask
11   $\vec{val} \leftarrow \text{mask\_min\_epu32}(\vec{val}, \text{same\_addr}, \vec{val}, \vec{old})$ 
12   $\vec{old} \leftarrow \text{mask\_i32gather\_epi32}(\vec{idx}, m, \vec{idx}, T, \text{sizeof}(\text{uint32}))$ 
13   $\vec{new} \leftarrow \text{maskz\_expand\_epi32}(m, \vec{val}_c)$ 
14   $\vec{new} \leftarrow \text{min\_epu32}(\vec{new}, \vec{old})$ 
15   $\text{mask\_i32scatter\_epi32}(T, m, \vec{idx}, \vec{new}, \text{sizeof}(\text{uint32}))$ 
16   $\vec{r} \leftarrow \text{mask\_mov\_epi32}(\vec{new}, \text{cmpeq\_epi32\_mask}(\vec{new}, \vec{val}), \vec{old})$ 
17 return  $\vec{r}$ 

```

Algorithm 9 describes the processing of a pixel vector during the first pass. The pixels are processed in a sequential natural reading order. Neighbor vectors $\vec{a}, \vec{b}, \vec{c}, \vec{d}$ and the current pixel \vec{x} can be obtained by doing aligned loads or by register rotation. Border and corners cases can be handled by setting the out of image pixel vectors to zero. We start by constructing unaligned vectors \vec{ab} and \vec{bc} from $\vec{a}, \vec{b}, \vec{c}$ by doing some element shifting (lines 2 and 3). We also compute the bitmask m corresponding to \vec{x} : a bit is set to 1 for a foreground pixel and to 0 for a background pixel (line 4). The next step is to initialize the labels in \vec{x} as shown in figure 8. Each pixel either points to itself or to its left neighbour (lines 6 to 9). We can now compute \vec{dx} from \vec{d} and \vec{x} and propagate the neighbor labels into \vec{x} (lines 10 and 11). The *vec_maskz_min_n⁺* function can be implemented using the property of unsigned integers overflow ($-1 = \text{MAX_UINT}$) and the *maskz_min_epu32* intrinsics. Finally, \vec{x} can be stored to memory (line 13) and we can call the *VecUnion* function described in subsection 4.1 (lines 15 to 18). As previously said, only the stairs and concavity patterns can lead to an union operation. The other configurations are handled with the label propagation step. In our implementation, the equivalence table pointer is moved by 1 pixel to account for the +1 in the labels and simplify memory accesses.

The second pass is the simultaneous transitive closure and relabeling. In this pass we don't need the neighbor information making the loading pattern straightforward. For each vector of pixel \vec{e} , we find the corresponding equivalence tree root and write it back to memory. Processing the pixels in the same order as in first pass allows the algorithm to capitalize on previous iterations to find the roots faster. We can compute the true number of label n , using the *popcnt_u32* (population count in a 32-bit integer) instruction and some vector comparisons. We use the property that by definition a root label points to itself. Algorithm 10 describes this process.

Figure 7 represents the execution of the SIMD Rosenfeld pixel algorithm on a 12×4 image and SIMD register of size of 4. The outlined area shows the steps of algorithm 9. The *VecUnion* operation is not detailed here but the pattern is similar to the one in figure 6. The modifications in the image from the scan (9 \rightarrow 4, 18 \rightarrow 4) in figure 7 are due to the equivalence table being embed in the image.

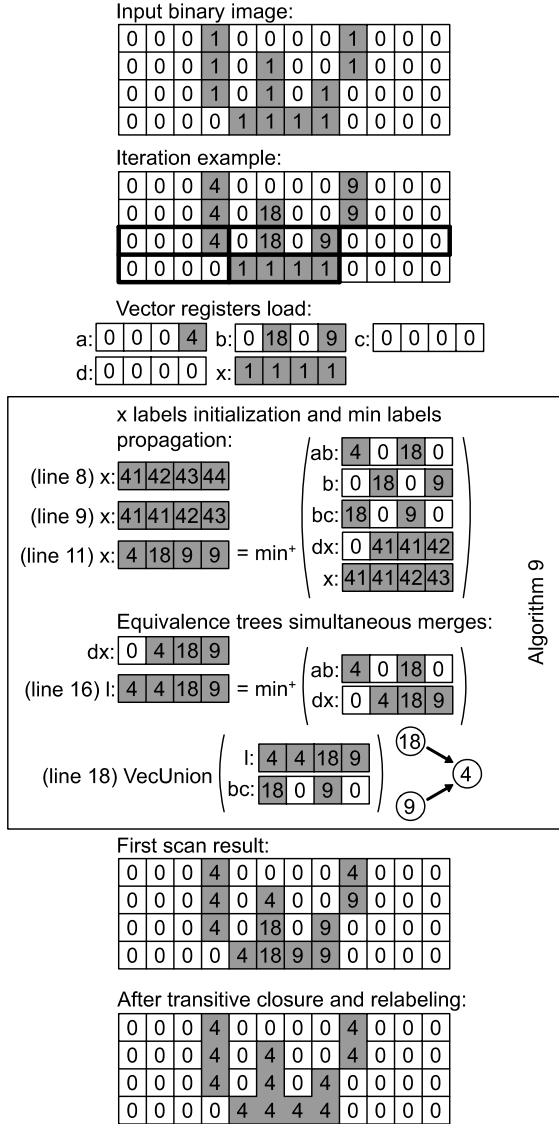


Figure 7: Example of an iteration of the SIMD Rosenfeld pixel algorithm.

4.3 SIMD Rosenfeld sub-segment algorithm

The SIMD rosenfeld pixel algorithm works well for low and medium image densities but for high image densities the performance collapses due to the pixel-recursive labels leading to more iterations in the *VecFind* while loop. To address this issue at the cost of a few more operations, we introduce the SIMD Rosenfeld sub-segment algorithm (v2). The only difference with the SIMD Rosenfeld pixel

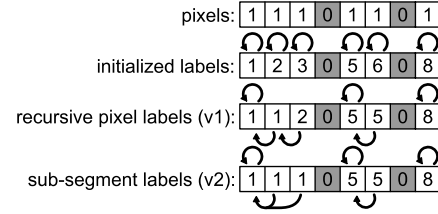


Figure 8: Creation of labels in SIMD Rosenfeld pixel algorithm (v1) and SIMD Rosenfeld sub-segment algorithm (v2).

Algorithm 9: SIMD Rosenfeld pixel (v1)

Input: T , the image / equivalence table, \vec{a} , \vec{b} , \vec{c} , \vec{d} , four vector of neighbor labels, \vec{x} , the current vector of pixels in (i, j) , w , the width of the image

- 1 // Shuffles :
- 2 $\vec{ab} \leftarrow \text{vec_shift_right_epi32}(\vec{a}, \vec{b})$
- 3 $\vec{bc} \leftarrow \text{vec_shift_left_epi32}(\vec{b}, \vec{c})$
- 4 $m \leftarrow \vec{x} \neq 0$ // mask
- 5 // x labels initialization and min labels propagation :
- 6 $\vec{inc} \leftarrow \text{set_epi32}(7, 6, 5, 4, 3, 2, 1, 0)$
- 7 $\vec{base} \leftarrow \text{set1_epi32}(i \times w + j + 1)$
- 8 $\vec{x} \leftarrow \text{maskz_add_epi32}(m, \vec{base}, \vec{inc})$
- 9 $\vec{x} \leftarrow \text{mask_sub_epi32}(\vec{x}, m \wedge (m < 1), \vec{x}, \vec{1})$
- 10 $\vec{dx} \leftarrow \text{vec_shift_right_epi32}(\vec{d}, \vec{x})$
- 11 $\vec{x} \leftarrow \text{vec_maskz_min}^+(m, \vec{ab}, \vec{b}, \vec{bc}, \vec{dx}, \vec{x})$
- 12 // Store x :
- 13 $\text{mask_store_epi32}(\&T[i][j], m, \vec{x})$
- 14 // Equivalence trees simultaneous merges :
- 15 $\vec{dx} \leftarrow \text{vec_shift_right_epi32}(\vec{d}, \vec{x})$
- 16 $\vec{l} \leftarrow \text{vec_maskz_min}^+(m, \vec{ab}, \vec{dx})$
- 17 $\text{merge} \leftarrow (\vec{bc} \neq 0) \wedge (\vec{b} = 0) \wedge (\vec{l} \neq 0) \wedge m$ // mask
- 18 $\text{VecUnion}(\vec{l}, \vec{bc}, \&T[0][-1], \text{merge})$

Algorithm 10: SIMD Transitive closure (solve equivalences)

Input: T , the image / equivalence table, w , the width of the image

Result: n , the true number of labels in the image (optional)

- 1 $n \leftarrow 0$
- 2 **for each** $\vec{e} \in T$ **do**
- 3 $m \leftarrow \vec{e} = 0$ // mask
- 4 $\vec{e} \leftarrow \text{VecFind}(\vec{e}, T, m)$
- 5 $\text{mask_store_epi32}(\&T[i][j], m, \vec{e})$
- 6 // (i, j) are the coordinates of \vec{e}
- 7 $\vec{inc} \leftarrow \text{set_epi32}(7, 6, 5, 4, 3, 2, 1, 0)$
- 8 $\vec{base} \leftarrow \text{set1_epi32}(i \times w + j + 1)$
- 9 $\vec{l} \leftarrow \text{maskz_add}(m, \vec{base}, \vec{inc})$
- 10 $n \leftarrow n + \text{popcnt_u32}(\text{mask_cmpeq_epi32_mask}(m, \vec{l}, \vec{e}))$

algorithm lies in the way we produce new labels (Fig. 8).

We define a segment as a sequence of same value pixels. A sub-segment is a segment bounded by the size of a vector. We use the conflict detection instructions from AVX512CD to compute a conflict-free subset ($c\vec{f}ss$) which allows us, with the count leading zero instruction ($lzcnt_epi32$) and some bit manipulation, to retrieve the index of the first element of a sub-segment. By making the pixel point directly to the sub-segment start, we can reduce the number of *VecFind* iterations to only 1 jump per sub-segment. Algorithm 11 describes these changes and figure 9 shows the key states of the sub-segment start computation code. Lines 14 to 16 of algorithm 11

are optional for the correctness of the algorithm but improve the performance.

	0	1	2	3	4	5	6	7
pixels:	1	1	1	0	1	1	0	1
cfss:	00000000	00000001	00000011	00000000	00000111	00010111	00001000	00110111
bitmask:	00000000	00000001	00000011	00000111	00001111	00011111	00111111	01111111
andnot:	00000000	00000000	00000000	00000111	00001000	00001000	00110111	01001000
lzct:	8	8	8	5	4	4	2	1
x:	0	0	0	3	4	4	6	7

Figure 9: Use of conflict detection to find the sub-segment’s start indices. In this example, elements have 8 bits and $lzct$ count from the 8th bit instead of the 32th.

Algorithm 11: SIMD Rosenfeld sub-segment (v2)

Input: T , the image / equivalence table, $\vec{a}, \vec{b}, \vec{c}, \vec{d}$, four vector of neighbor labels, \vec{x} , the current vector of pixels in (i, j) , w , the width of the image

- 1 // Shuffles :
- 2 $\vec{ab} \leftarrow \text{vec_shift_right_epi32}(\vec{a}, \vec{b})$
- 3 $\vec{bc} \leftarrow \text{vec_shift_left_epi32}(\vec{b}, \vec{c})$
- 4 $m \leftarrow \vec{x} \neq 0$
- 5 // x labels initialization and min labels propagation :
- 6 $\vec{bitmask} \leftarrow \text{set_epi32}(0x7F, 0x3F, 0x1F, 0xF, 7, 3, 1, 0)$
- 7 $\vec{cfss} \leftarrow \text{maskz_conflict}(m, x)$
- 8 $\vec{lzct} \leftarrow \text{lzcnt_epi32}(\text{andnot_epi32}(\vec{cfss}, \vec{bitmask}))$
- 9 $\vec{base} \leftarrow \text{set1_epi32}(i \times w + j + 1 + 32)$
- 10 $\vec{x} \leftarrow \text{maskz_sub_epi32}(m, \vec{base}, \vec{lzct})$
- 11 $\vec{dx} \leftarrow \text{vec_shift_right_epi32}(\vec{d}, \vec{x})$
- 12 $\vec{x} \leftarrow \text{vec_maskz_min}^+(m, \vec{ab}, \vec{b}, \vec{bc}, \vec{dx}, \vec{x})$
- 13 // Optional propagation:
- 14 $\vec{perm} \leftarrow \text{maskz_sub_epi32}(m, 32, \vec{lzct})$
- 15 $\vec{x}_p \leftarrow \text{permute_epi32}(\vec{perm}, \vec{x})$
- 16 $\vec{x} \leftarrow \text{vec_maskz_min}^+(m, \vec{x}, \vec{x}_p)$
- 17 // Store:
- 18 $\text{mask_store_epi32}(\&T[i][j], m, \vec{x})$
- 19 // Equivalence trees simultaneous merges :
- 20 $\vec{dx} \leftarrow \text{vec_shift_right_epi32}(\vec{d}, \vec{x})$
- 21 $\vec{l} \leftarrow \text{vec_maskz_min}^+(m, \vec{ab}, \vec{dx})$
- 22 $\text{merge} \leftarrow (\vec{bc} \neq 0) \wedge (\vec{b} = 0) \wedge (\vec{l} \neq 0) \wedge m$
- 23 $\text{VecUnion}(\vec{l}, \vec{bc}, \&T[0][j-1], \text{merge})$

4.4 Parallel SIMD algorithms

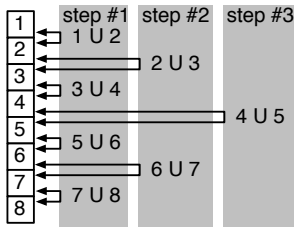


Figure 10: Pyramidal border merging of disjoint sets

We use OpenMP for the parallel implementation and make the assumption that the memory model is NUMA with shared memory between processors. SIMD algorithms have an increased pressure on memory bandwidth which tends to reduce the multicore-parallelism efficiency if the application is not compute bound.

The approach we follow in our parallel implementations of the SIMD Rosenfeld pixel and sub-segment algorithms is based on a divide-and-conquer method described in [4]. The image is split into p sub-images, with p the number of cores. This parallel algorithm minimizes the number of merges required by taking a pyramidal approach but diminishes the number of active cores at a given time. It needs $\log_2(p)$ steps to complete the merge. Each step is fully parallel and does not require *atomic* instructions to update the equivalence table as this scheme merges borders of disjoint sets (Fig. 10).

First, each processor core takes a sub-image horizontal strip and applies the first pass of either algorithm. Except for the modified loop indexes, there is no difference between the sequential and parallel code. Then, the next step consists in applying a pyramidal merge. As we divide the image by the number of cores available p , the total number of merges required is equal to $p - 1$. For each merge step, we divide the number of active cores doing a two way merge until we have only one core left. On the border line between two sub-images, we apply a merging pass on each column. For each SIMD vector of pixels, we have at most two calls to *VecUnion*. Algorithm 12 describes the body of the border merge loop. Finally, we apply the second pass which does not require any code modification compared to the sequential version.

Algorithm 12: SIMD Border merging

Input: T , the image / equivalence table, $\vec{a}, \vec{b}, \vec{c}, \vec{d}$, four vector of neighbor labels, \vec{x} , the current vector of pixels in (i, j) , w , the width of the image

- 1 // Shuffles :
- 2 $\vec{ab} \leftarrow \text{vec_shift_right_epi32}(\vec{a}, \vec{b})$
- 3 $\vec{bc} \leftarrow \text{vec_shift_left_epi32}(\vec{b}, \vec{c})$
- 4 $\vec{dx} \leftarrow \text{vec_shift_right_epi32}(\vec{d}, \vec{x})$
- 5 $m \leftarrow \vec{x} \neq 0$
- 6 // Equivalence trees simultaneous merges :
- 7 $\vec{l} \leftarrow \text{vec_maskz_min}^+(m, \vec{ab}, \vec{b})$
- 8 $\text{merge} \leftarrow (\vec{l} \neq 0) \wedge (\vec{dx} = 0) \wedge m$ // mask
- 9 $\text{VecUnion}(\vec{x}, \vec{l}, \&T[0][j-1], \text{merge})$
- 10 $\vec{l} \leftarrow \text{vec_maskz_min}^+(m, \vec{b}, \vec{bc})$
- 11 $\text{merge} \leftarrow (\vec{l} \neq 0) \wedge (\vec{b} = 0) \wedge m$ // mask
- 12 $\text{VecUnion}(\vec{x}, \vec{l}, \&T[0][j-1], \text{merge})$

5 BENCHMARK

5.1 Benchmark protocol

In our benchmark protocol, we aim to study the behavior of CCL algorithms in the widest range of configurations available, and to see how they compare to each other. For reproducible results, MT19937 [14] was used to generate images of varying density ($d \in [0\% - 100\%]$) and granularity ($g \in \{1 - 4\}$). The image density is the ratio of black and white pixels in the image: a black image (with zero pixel set) has a density of 0% and a white image (with all pixel set) has a density of 100%. An image of granularity g

is composed of $g \times g$ macro-pixels set to 1 or 0. Granularity increases the local homogeneity of an image. Natural and structured images lead to processing times roughly equal to those of random images with $g = 4$, while $g = 1$ provides difficult and unstructured configurations to stress the algorithms and especially those manipulating segments or run-length's. This benchmark methodology for CCL algorithms was previously used in [12], [3], [4], [8].

We tested the performance of the available SIMD vector length (128, 256, 512) in single and multicore on a dual socket Intel Xeon(R) Gold 6126 running at 2.6 Ghz (turbo-boost off). Before any testing, it was unclear whether the new 512 vector length would have a positive impact on the performance due to the additional stress on the memory bandwidth and the frequency throttle. This is especially true while fully exploiting SIMD and multicore due to bandwidth saturation. Our results are summarized in table 1, 2 and figures 11, 12 and are discussed in the following section. The performance of these new algorithms is compared to the classic pixel-based algorithms with DT (Rosenfeld and Rosenfeld+DT) and to the fastest run-length based segment labeling algorithm (LSL_{STD} and LSL_{RLE}) [4].

5.2 Results analysis

Figure 11 shows the execution time (in cycles per pixels) on images of varying densities. We observe that for both versions, there is a bump around 45% which corresponds to the percolation threshold in 8-connectivity. We can also see that the main difference in performance between the pixel and sub-segment algorithms happens for higher densities. The pixel version being slower because of the recursive pixel labels leading to a longer while loop in *VecFind* (as seen in Sec. 4.3). This performance difference grows with the number of cores due to the added cost of find operations in the border merging step.

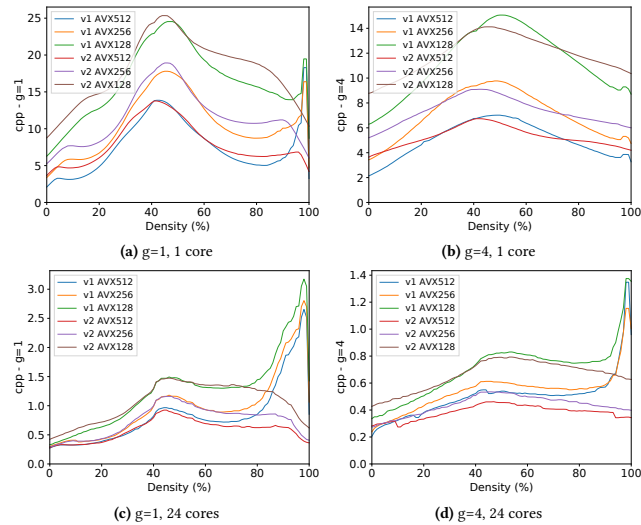


Figure 11: Cycles per pixels for SIMD Rosenfeld pixel (v1) and SIMD Rosenfeld sub-segment (v2), applied to 2048×2048 images

SIMD vs scalar: In the single-threaded case, SIMD versions are two times faster than the Rosenfeld scalar versions. They are also

threading	1-core mono thread			24-core multi thread		
	g=1	g=2	g=4	g=1	g=2	g=4
LSL _{STD}	10.47	6.98	5.91	0.71	0.36	0.28
LSL _{RLE}	16.96	9.31	6.05	0.95	0.45	0.24
Rosenfeld	30.61	19.01	12.49	2.56	1.69	1.27
Rosenfeld+DT	20.01	11.81	8.30	1.95	1.51	1.09
SIMD v1 ₅₁₂	7.61	6.07	5.06	0.84	0.52	0.50
SIMD v1 ₂₅₆	10.65	8.64	7.08	0.99	0.57	0.55
SIMD v1 ₁₂₈	16.40	13.23	11.33	1.26	0.76	0.71
SIMD v2 ₅₁₂	7.97	6.54	5.26	0.58	0.41	0.38
SIMD v2 ₂₅₆	11.54	9.22	7.35	0.74	0.53	0.44
SIMD v2 ₁₂₈	17.89	14.15	11.89	1.07	0.71	0.66

Table 1: Average cycles per pixels for 2048×2048 images - best SIMD in bold (lower values is better)

threading	1-core mono thread			24-core multi thread		
	g=1	g=2	g=4	g=1	g=2	g=4
LSL _{STD}	0.27	0.38	0.44	4.12	7.53	9.45
LSL _{RLE}	0.21	0.32	0.46	3.56	6.91	11.80
Rosenfeld	0.13	0.17	0.24	1.15	1.70	2.17
Rosenfeld+DT	0.17	0.25	0.33	1.47	2.02	2.46
SIMD v1 ₅₁₂	0.44	0.49	0.56	4.22	5.06	5.66
SIMD v1 ₂₅₆	0.29	0.34	0.40	3.57	4.21	5.18
SIMD v1 ₁₂₈	0.18	0.21	0.24	2.74	3.37	4.04
SIMD v2 ₅₁₂	0.37	0.42	0.51	4.96	5.84	6.92
SIMD v2 ₂₅₆	0.25	0.30	0.36	4.06	4.98	6.12
SIMD v2 ₁₂₈	0.16	0.19	0.22	2.77	3.48	4.05

Table 2: Average throughput (Gpx/s) for 2048×2048 images - best SIMD in bold (higher is better).

faster than the LSL versions. In the multi-threaded case, this ratio grows to $\times 3$. Compared to LSL, the SIMD versions are faster only for $g=1$, which is the worst case for full-segment labeling like LSL (the strategy to save memory accesses, especially with the RLE version becomes more profitable than in the mono-threaded case where the pressure on the memory is lower).

SIMD Scalability (128 / 256 / 512) In the single-threaded case, doubling the SIMD size provides a speedup around 1.5. In multi-threaded case, this ratio still exists but only for 128 / 256 registers. For 256 / 512 registers, the ratio drops to 1.2, due to bandwidth saturation.

Thread scalability: Depending on the SIMD size, the scalability of the algorithms varies in the interval $\times 10 - \times 13$ for 512-bit registers and up to $\times 15 - \times 18$ for 128-bit registers. That is an efficiency between 40% (for 512-bit registers) up to 75% (for 128-bit registers). Considering all the instructions for the control-flow and the label propagation within a register, for such an irregular algorithm, we consider the results acceptable.

Performance with image size: The best performance is achieved when the image fits in the cache for all algorithms. (see Tab. 3). For a granularity equal to 4, LSL_{RLE} is still the fastest algorithm. On

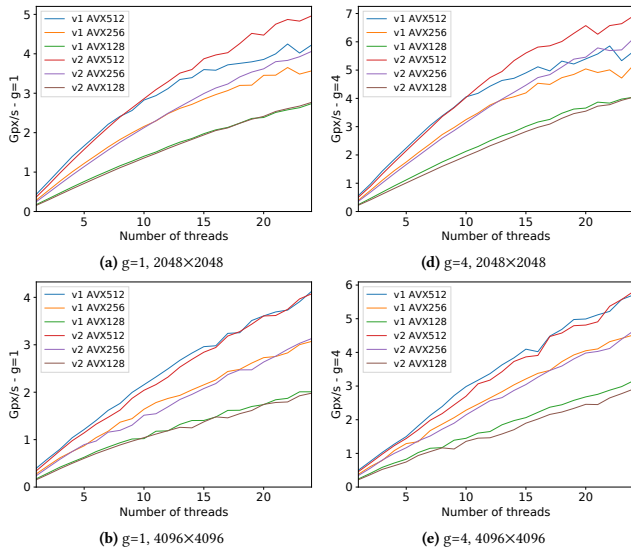


Figure 12: Average throughput (Gpx/s) for SIMD Rosenfeld pixel (v1) and SIMD Rosenfeld sub-segment (v2)

lower granularities, the new algorithms outperform Rosenfeld+DT by a factor $\times 2.3$ up to $\times 2.8$ for $g=4$, and from $\times 2.6$ up to $\times 3.4$ for $g=1$. We also observe that, for larger images, the pixel-based algorithm SIMD v1 becomes faster than the sub-segment-based algorithm SIMD v2. These two new algorithms are specially well-suited to very complex / un-structured / quasi-random images.

image size	2k images		4k images		8k images	
	g=1	g=4	g=1	g=4	g=1	g=4
LSL _{RLE}	3.56	11.80	3.45	9.16	3.25	7.55
Rosenfeld+DT	1.47	2.46	1.41	2.22	1.37	1.94
SIMD v1 ₅₁₂	4.22	5.66	4.13	5.71	3.83	4.89
SIMD v2 ₅₁₂	4.96	6.92	4.07	5.81	3.59	4.45

Table 3: Average throughput (Gpx/s) for 2k, 4k, 8k images on 24 cores (best performance in bold, for each column).

6 CONCLUSION

This paper presented two new CCL direct algorithms for SIMD multi-core architectures. The work is motivated by a large number of applications in computer vision and could also be of interest for the graph community with a new Union-Find SIMD algorithm. The use of SIMD instructions in CCL and more generally in computer vision is scarce due to the complexity of the code vectorization process. We have developed a performance efficient and portable SIMD algorithms for the Rosenfeld algorithm that outperform scalar pixel-based algorithms by a factor of $\times 2.3$ up to $\times 3.4$.

We found out that the new optimized SIMD algorithms are performing better than the State-of-the-Art algorithms on small images of fine granularity and on single core CPU. While it does not scale as well as run-length based algorithms, a lot of applications which do not require big images but instead search for high throughput could benefit from it. In future work, we plan to test our algorithms on new ARM/Fujitsu SVE architectures.

REFERENCES

- [1] D. A. Bader and J. Jaja. Parallel algorithms for image histogramming and connected components with an experimental study. *Parallel and Distributed Computing*, 35,2:173–190, 1995.
- [2] F. Bolelli, M. Cancilla, L. Baraldi, and C. Grana. Toward reliable experiments on the performance of connected components labeling algorithms. *Journal of Real-Time Image Processing (JRTIP)*, pages 1–16, 2018.
- [3] L. Cabaret, L. Lacassagne, and D. Etiemble. Distanceless label propagation: an efficient direct connected component labeling algorithm for GPUs. In *IEEE International Conference on Image Processing Theory, Tools and Applications (IPTA)*, pages 1–8, 2017.
- [4] L. Cabaret, L. Lacassagne, and D. Etiemble. Parallel Light Speed Labeling for connected component analysis on multi-core processors. *Journal of Real Time Image Processing*, 15(1):173–196, 2018.
- [5] M. Ceska. Computing strongly connected components in parallel on cuda. In Nvidia, editor. *GPU Technology Conference*, 2010.
- [6] S. Gupta, D. Palsetia, M. A. Patwary, A. Agrawal, and A. Choudhary. A new parallel algorithm for two-pass connected component labeling. In *Parallel & Distributed Processing Symposium Workshops (IPDPSW)*, pages 1355–1362. IEEE, 2014.
- [7] L. He, X. Ren, Q. Gao, X. Zhao, B. Yao, and Y. Chao. The connected-component labeling problem: a review of state-of-the-art algorithms. *Pattern Recognition*, 70:25–43, 2017.
- [8] A. Hennequin, Q. L. Meunier, L. Lacassagne, and L. Cabaret. A new direct connected component labeling and analysis algorithm for GPUs. In *IEEE International Conference on Design and Architectures for Signal and Image Processing (DASIP)*, pages 1–6, 2018.
- [9] W. W. Hwu, editor. *GPU Computing Gems*, chapter 35: Connected Component Labeling in CUDA. Morgan Kaufman, 2001.
- [10] A. Kalentev, A. Rai, S. Kemnitz, and R. Schneider. Connected component labeling on a 2d grid using CUDA. *Journal of Parallel and Distributed Computing*, 71:615–620, 2011.
- [11] Y. Komura. Gpu-based cluster-labeling algorithm without the use of conventional iteration: application to swendsen-wang multi-cluster spin flip algorithm. *Computer Physics Communications*, pages 54–58, 2015.
- [12] L. Lacassagne, L. Cabaret, F. Hebach, and A. Petreto. A new SIMD iterative connected component labeling algorithm. In *ACM Workshop on Programming Models for SIMD/Vector Processing (PPoPP)*, pages 1–8, 2016.
- [13] A. Lindner, A. Bieniek, and H. Burkhardt. PISA-Parallel image segmentation algorithms. pages 1–10. Springer, 1999.
- [14] M. Matsumoto and T. Nishimura. Mersenne twister: A 623-dimensionally equidistributed uniform pseudorandom number generator. *Transactions on Modeling and Computer simulation*, 8(1):3–30, 1998.
- [15] M. Niknam, P. Thulasiraman, and S. Camorlinga. A parallel algorithm for connected component labeling of gray-scale images on homogeneous multicore architectures. *Journal of Physics - High Performance Computing Symposium (HPCS)*, 2010.
- [16] D. P. Playne and K. Hawick. A new algorithm for parallel connected-component labelling on GPUs. *IEEE Transactions on Parallel and Distributed Systems*, 2018.
- [17] J. Ragan-Kelley, C. Barnes, A. Adams, S. Paris, F. Durand, and S. Amarasinghe. Halide: a language and compiler for optimizing parallelism, locality, and recomputation in image processing pipelines. *ACM SIGPLAN Notices*, 48(6):519–530, 2013.
- [18] A. Rosenfeld and J. Platz. Sequential operator in digital pictures processing. *Journal of ACM*, 13,4:471–494, 1966.
- [19] F. Wende and T. Steinke. Swendsen-wang multi-cluster algorithm for the 2d/3d Ising Model on Xeon Phi and GPU. In ACM, editor, *International Conference on High Performance Computing (SuperComputing)*, pages 1–12, 2013.
- [20] K. Wu, E. Otoo, and A. Shoshani. Optimizing connected component labeling algorithms. *Pattern Analysis and Applications*, 2008.