# EE216 Re-configurable Computing Project Report

Zhaojun Ni, 2022231102

January 19, 2024

## 1 Performance and Resource Utilization

After increasing the frequency to **300MHz**, the latency for processing a single 480p image from TUM RGB-D Dataset on the ZCU104 Evaluation board is only **1.78ms**. It uses 40578 LUTs, 73346 FFs, 149 BRAMs and 193 DSPs.
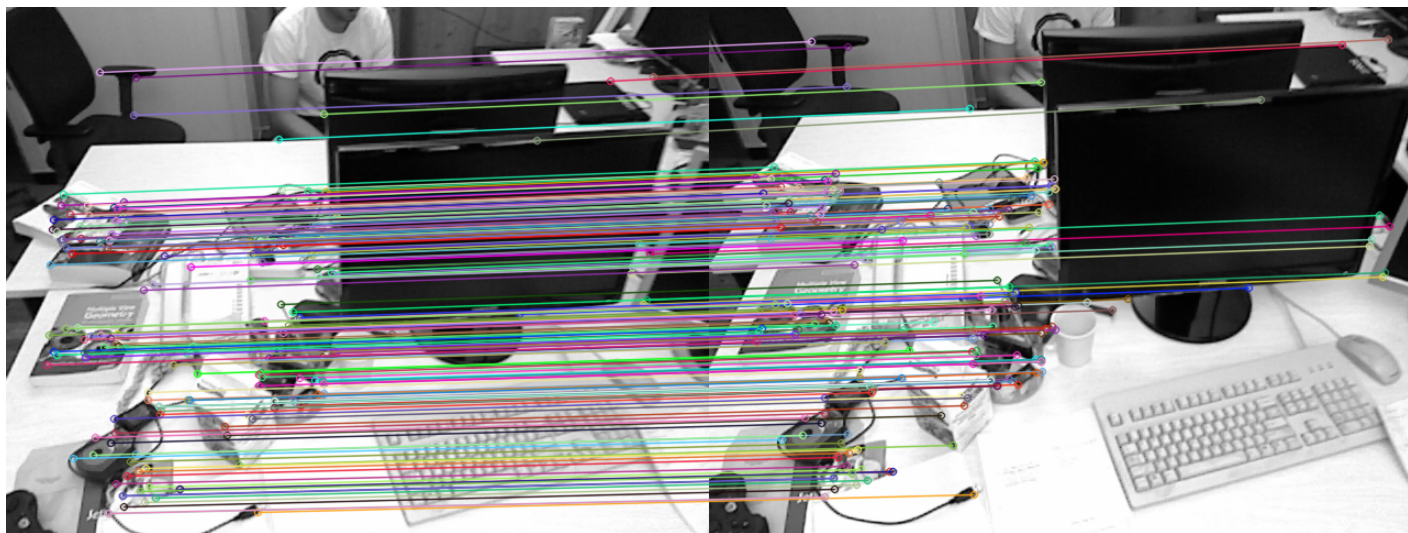


Figure 1: Matching result on ZCU104 Board @ 300MHz.



Figure 2: Pynq running result on ZCU104 Board @ 300MHz.

# 2 Experimental Records

## # 1. C Simulation



Figure 3: Csim Result

## # 2. Synthesis

The synthesis report shows negative slack which equals to the uncertainty while it works fine on the board. And the interval column shows that the system has balanced throughput. And it indicates that 300 MHz may be the peak frequency of the design.



Figure 4: Synthesis Report

# 3. Co-Simulation



Figure 5: Co-sim Report



Figure 6: Co-sim Result

# 3   Optimization Methods

### # 0. Loop Tripcount

In HLS, specifying the tripcount for a loop is the first step in optimization. This allows the synthesis report to show the expected number of cycles each function is likely to consume. Based on the order of magnitude of the cycle count in the synthesis report, along with the relationship between frequency and the proportion of processed images, one can roughly estimate the delay.

### # 1. Loop Unroll

Loop unrolling in HLS is an optimization technique that involves expanding the loop by replicating the loop body's iterations, aiming to reduce loop control overhead and enhance instruction-level parallelism. In certain situations, this can improve the performance of hardware descriptions. For example, in "fast.hpp", loop unrolling has been applied to the loop inside the "fast_score" function. Similarly, in the "process_fast" function, loop unrolling has been performed on the explicitly bounded "for" loop, with the additional specification of the unroll factor. The "unroll factor" is the number of times a loop is unfolded, indicating the repetition count of each original iteration in the unfolded copies.

### # 2. Array Partition

Array partitioning is used to divide an array, allowing the partitioned sections to be read within a single clock cycle. For instance, when dealing with a line buffer that performs row-wise read operations, and considering a two-dimensional array where the first dimension represents rows and the second dimension represents columns, performing a complete partition on the first dimension is sufficient. Similarly, for small windows like a 7x7 convolution kernel, complete partitioning can be done in both dimensions.

In the "loop_descriptor", I observed potential conflicts between accesses by "t0" and "t1" to the "window". To address this, I created two copies, "window_buf0" and "buf1", dedicated to "t0" and "t1" respectively. When attempting complete partitioning in both dimensions for these two buffers, the synthesis process became very slow. To mitigate this, I explicitly specified partitioning in only the second dimension for both buffers. Otherwise, synthesis would generate certain warnings.

### # 3. Pipeline

In this project, multiple functions starting with "process" are annotated with the "pragma HLS DATAFLOW". Consequently, from the synthesis report, it is observed that modules traversed by DATAFLOW seem to undergo automatic pipelining within their internals after adjusting loop boundaries. For certain pixel traversal operations within functions, such as iterating over columns, the report indicates that pipelining is enabled (pipeline: yes).

### # 4. Non-Blocking Optimization

From the initial synthesis report, it is evident that the main delay occurs in the "process_rBRIEF" operation, and it consumes the majority of hardware resources. From the perspective of the "process_rBRIEF" function, it begins by traversing all pixels through two nested loops. If the central pixel of the mask window is identified as a key point, computations for the functions "IC Angle 31" and "descriptor' are then performed. The former, "IC Angle 31", involves several loops with undefined loop boundaries, while the latter, "descriptor", contains an additional loop with 256 iterations.

### a. Unroll loops in IC Angle 31.

In this function, the outer loop's iteration count, denoted as variable "v," is stored in an array. I manually unrolled the loop by extracting the values from the 15-column array and individually set them as the loop's upper limit. This way, each sub-loop can be separately specified with an UNROLL FACTOR for parallel computation.

### b. Cache window columns[1].

Upon closer inspection, the initial version's slow speed can be attributed to its blocking nature. It requires the completion of a series of operations in "IC Angle 31" and "descriptor" before accessing the next pixel. Drawing inspiration from the referenced literature[1], to achieve non-blocking computation for rBRIEF, a FIFO is employed to cache the window columns. In specific operations, it is possible to modify the original "process_rBRIEF" and "process_output" functions by specifying the output of the former and the input of the latter as HLS stream types.
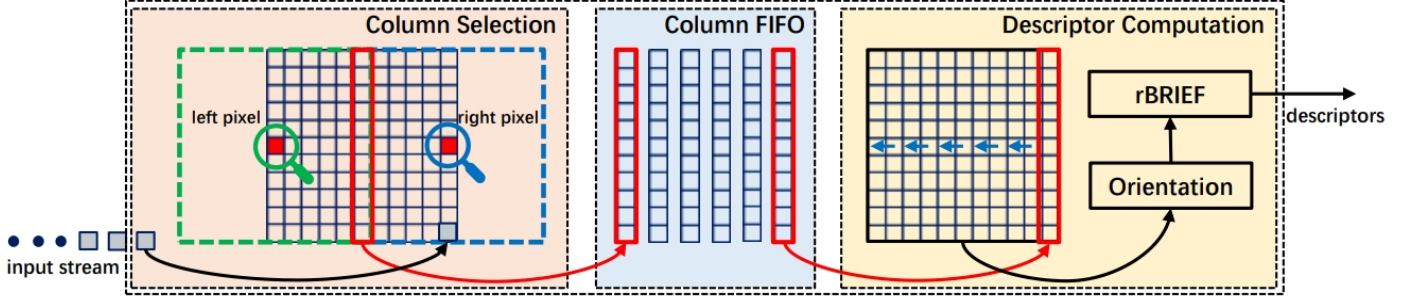
Figure 7: Non-blocking computation is achieved by caching the window columns.[1]

Figure 7 is an rBRIEF module proposed by the authors in Reference 1. I have implemented it using HLS, with the Column Selection part realized in the "process_rBRIEF_col_sel" function and the Descriptor Computation implemented in the "process_rBRIEF_des_compute" function. The data type of the column FIFO in the middle of the figure directly utilizes the "WindowCol" structure from the ORB.h file, and data flow is facilitated through the use of "hls::stream".

The author's idea is that, as a window of length 37 traverses from left to right, attention should be paid to whether there are key points in the leftmost and rightmost columns of the window. If a key point is detected in the rightmost column, the author starts outputting the middle column's window column through a FIFO to the subsequent descriptor computation module. After precisely 36 cycles, 37 columns will be written into the FIFO. At this point, the complete orientation and rBRIEF for the window can be calculated.

The author does not use a counter to determine whether the window is intact. This is because it is possible to encounter two key points in the same row with a separation of less than a window size. The author's idea of detecting the leftmost column of the window is **clever**. This is because the first key point encountered in the rightmost column will naturally move to the position detected by the leftmost column after exactly 36 cycles. At this moment, a column-end label signal can be written into the FIFO to initiate the computation of the key point's window.

This algorithm is advantageous as it utilizes a FIFO. In the window buffer of the descriptor computation, each column shifts to the left every other cycle, and the newest column is the one most recently written. In cases of adjacent key points, this allows for **efficient data reuse**, thereby reducing latency to a significant extent.

**c. Implement Column Selection.**

My "process_rBRIEF_col_sel" function is essentially implemented according to the algorithm in the author's paper. However, during simulation, I encountered some issues. If a 37x37 window directly slides from the top-left corner of the image to the right, the simulated output only has 2 feature points instead of 6. The reason is that if there is no horizontal padding, the leftmost column written into when the sliding window reaches the position of the 19th column (HALF WINDOW SIZE) – meaning the feature points at the leftmost and rightmost positions of the window – will not be detected.

To address this, I introduced padding in the column traversal of the algorithm. I filled the leftmost 18 columns and the rightmost 18 columns with zeros. After making this modification, the simulation results were error-free, and it correctly outputted all six feature points.

**d. Column FIFO data component.**

The first three bits of this structure serve as flags, representing the "last," "begin," and "end" signals. The "last" flag indicates whether the current column being written is the last one in the entire image. The "begin" flag indicates whether the current column is the first column in the key point window. The "end" flag indicates whether the current column is the last column in the key point window. Following the flags, there are two uint16 data fields representing the coordinates of the key point. Then, there are two uint data fields with a bit width of 8 * WIND_SIZE, representing the columns of the image and blured image. The final uint8 data is the value from the mask window, i.e., the response value of the key

point. The depth of FIFO is set to 200 according to the reference[1].

   **e. Implement Descriptor Computation.**

I replaced the original "output" function with my self-implemented "process_rBRIEF_des_compute" function to handle the output of descriptors. Similarly, I use a while loop to continuously process data from the FIFO. As long as the received data's "flag last" is not 1, the processing continues. In each cycle, if the data in the FIFO is valid, the function reads the structure data from the FIFO. Then, it performs a circular left shift on the "img window" and "blur window". Next, it inserts the "img col" and "blur col" from the structure into the last columns of their respective windows. When the read "flag_end" is 1, it indicates that the data in the current window is a complete window for a key point. Subsequently, the angle and descriptor calculations are carried out. The coordinate information, response value, descriptor calculation results, and other details are then written to the "dst_axi". Following the style of the original "process_output" function, the AXI signals "last" are set to 0, "keep" to -1. Additionally, for sorting convenience, the results are written into "idxTree" and "cntTree".

# 4   Reference

[1]Q. Zhang, H. Sun, Q. Deng, H. Yu and Y. Ha, "NORB: A Stream-Based and Non-Blocking FPGA Accelerator for ORB Feature Extraction," 2023 30th IEEE International Conference on Electronics, Circuits and Systems (ICECS), Istanbul, Turkiye, 2023, pp. 1-4, doi: 10.1109/ICECS58634.2023.10382726.