*Article*

# Template Matching. Optimizing performance and area in Vivado HLS

**Vlad-Eusebiu Baciu** (iD)

Student ID: 0577953
E-mail: vlad-eusebiu.baciu@vub.be
Course: Multiprocessors and Reconfigurable Architectures

1 **Abstract:** The task consists to implement of a simple image detector technique that attempts to
2 evaluate the matching of two images. A simple cross-correlation technique is proposed: sum
3 of absolute differences (SAD). It uses an image patch (template), tailored to a specific feature of
4 the search image. In order to apply this template matching method the image is normalized and
5 grayscaled. Two designs with optimizations targeting performance and area have been created in
6 Vivado HLS.

7 **Keywords:** template matching; image detector; VIVADO HLS

8 ## 0. Introduction

9 The algorithm is able to receive two stream of pixels from the original image and
10 template image in a BMP format (RGB, 3 bytes/pixel). On both input images two
11 algorithms are applied: grayscaling and normalization. The algorithm returns a new
12 image in which the original input template image is superimposed over the rest of
13 the grayscale image. The grayscale algorithm supports three standards: LUMA YUV,
14 ITU-R BT.709 and ITU-R BT.2100. The top function calls the grayscale function seting the
15 enconding type as LUMA YUV and the normalization function with a minumum value
16 of 0 and a maximum value of 255 as parameters. These parameters can be configured
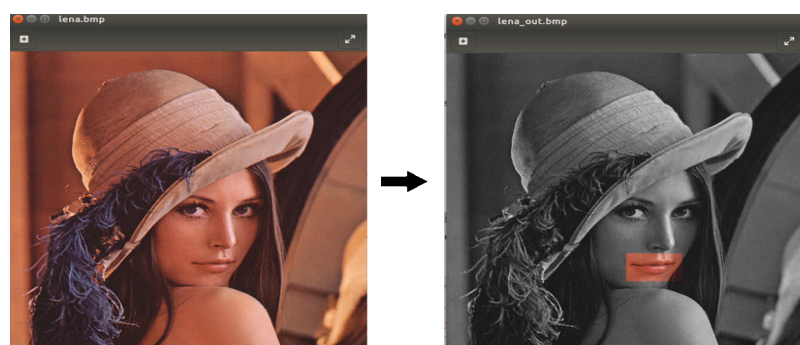17 by the user if access to top function is permitted and can not be modified from the test
18 bench.



**Figure 1.** Template matching result: (**a**) Input image (**b**) Output image

19 ## 1. Vivado HLS - optimization

20 Three Vivado HLS solution are presented in this report:

21 - Solution 1: default solution, no optimizations are applied except HLS loop trip-
22 count;
23 - Solution 2: few optimization directives applied in order to achive the highest
24 performance;
25 - Solution 3: few optimization directives applied in order to reduce area and resource
26 consumption.

27  Vivado HLS provides pragmas that can be used to optimize the design: reduce la-
28  tency, improve throughput performance, and reduce area and device resource utilization
29  of the resulting RTL code. These pragmas can be added directly to the source code for
30  the kernel. The following table contains the pragma directives used in this project.

**Table 1.** Vivado HLS pragmas used in the project.

| Pragma | Optimization type |
| --- | --- |
| LOOP TRIPCOUNT | performance |
| ARRAY PARTITION | performance |
| LOOP UNROLL | performance |
| DATAFLOW | performance |
| STREAM | performance |
| ARRAY MAP | area |
| INLINE | area |
| BIT-ACCURATE TYPES | area |
| ALLOCATION | area |

31  **LOOP TRIPCOUNT:** Can be applied to a loop to manually specify the total number of
32  iterations performed by a loop.
33  **ARRAY PARTITION:** Partitions an array into smaller arrays or individual elements
34  and increases the amount of read and write ports for the storage.
35  **LOOP UNROLL:** Transforms loops by creating multiples copies of the loop body in the
36  RTL design, which allows some or all loop iterations to occur in parallel.
37  **DATAFLOW:** Enables task-level pipelining, allowing functions and loops to overlap in
38  their operation, increasing the concurrency of the RTL implementation.
39  **STREAM:** Top-level function array parameters are implemented as a RAM interface
40  port.
41  **ARRAY MAP:** Combines multiple smaller arrays into a single large array to help
42  reduce block RAM resources.
43  **INLINE:** Allows operations within the function to be shared and optimized more
44  effectively with surrounding operations.
45  **BIT-ACCURATE TYPES:** Allow any arbitrary bit width to be specified.
46  **ALLOCATION:** Limits resource allocation in the implemented kernel.

47  **2. Results**
48  *2.1. Default solution*

49  The default solution is presented below in which no optimization is used. Vivado
50  HLS reports the total latency of each loop, which is the number of clock cycles to execute
51  all iterations of the loop. The loop latency is therefore a function of the number of loop
52  iterations, or tripcount. Since the tripcount is a constant value in my case the **pragma**
53  **loop tripcount** has been used for each loop.

**Figure 2.** Performance estimates for default solution. No optimization included.

Without optimization the solution is implemented using a total of 15.6 % of resources.



**Figure 3.** Utilization estimates for default solution. No optimization included.

In the following picture the operations and control steps are presented. The algorithm "executes" first the imGrayScale instance for the input image then in control step C2 starts executing the same instance for the template input image. At the same time the imGreyNormalization instance is "executed" for the previous image. In the control step C4 the template image, now in grayscale format, gets normalized. Starting with C6 the algorithm executes the template matching algorithm. Between C8 and C11 the output image is constructed as seen in Figure 1.

**Figure 4.** Performance analysis. No optimization included.

*2.2. Performance optimization*

The first design targets the performance. In order to achive the highest performance the following methods have been used:

2.2.1. Loop unrolling

The UNROLL pragma transforms loops by creating multiples copies of the loop body in the RTL design, which allows some or all loop iterations to occur in parallel. I have chosen an unrolling factor of 4. Partially unrolling a loop lets you specify a factor N, to create N copies of the loop body and reduce the loop iterations accordingly. The impact of this method can be seen in Figure 6. Doing this the performance have been increased with 14 % and the resource utilization reached 20 %.



**Figure 5.** Loop labeled L99 after loop unrolling optimization.



**Figure 6.** Loop labeled L99 before loop unrolling optimization. Unrolling factor 4.

**Figure 7.** Loop unrolling impact on performance and utilization.

### 2.2.2. Array partitioning

Partitioning the arrays results in an increase of read and write ports for the storage. The RTL design will include multiple small memories and registers instead of one large memory. This method requires more memory instances or registers. Resource utilization increase up to 24.3 % but the performance improves with 51 %. Regarding the performance increase one has to keep in mind that cumulates with the later optimization method which was 14 %. It can be seen from Figure 8 and Figure 9 the impact of array partitioning with a factor of 3.



**Figure 8.** Array partitioning impact on performance and utilization. Memory port read.



**Figure 9.** Array partitioning impact on performance and utilization. Memory port read.

Figure 10 shows the performance and control steps of one particular loop from imGreyNormalization function after pragma unroll has benn used. Further optimization can be done if one utilize array partitioning. The differences between Figure 10 and Figure 11 are visible, the number of control steps decrease when array partitioning is used.
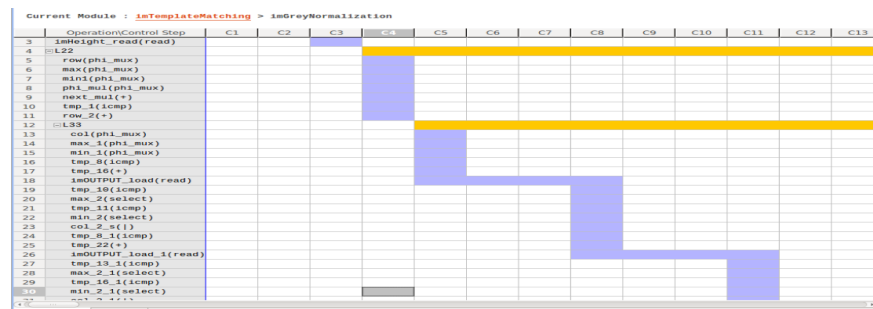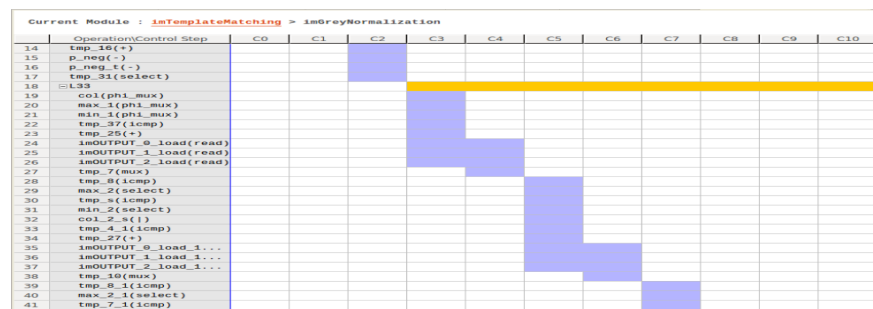
**Figure 10.** Loop unrolling impact on performance.



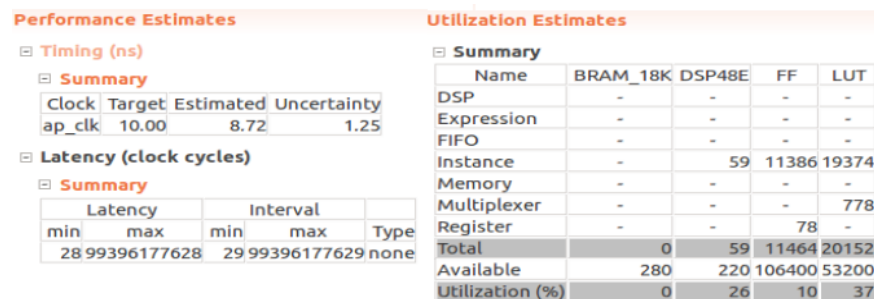**Figure 11.** Loop unrolling and array partitioning impact on performance.



**Figure 12.** Array partitioning impact on performance and utilization.

### 2.2.3. Dataflow

Vivado HLS analyzes the dataflow between sequential functions or loops and create channels (based on pingpong RAMs or FIFOs) that allow consumer functions or loops to start operation before the producer functions or loops have completed. This allows functions or loops to operate in parallel, which decreases latency and improves the throughput of the RTL.



**Figure 13.** Loop unrolling impact on performance and utilization.

**Figure 14.** Loop unrolling impact on performance and utilization.

### 2.2.4. Stream

Top-level function array parameters are implemented as a RAM interface port when STREAM directive is used. I have used this directive inside imTemplateMatching function which is the top function. If the data stored in the array is consumed or produced in a sequential manner, a more efficient communication mechanism is to use streaming data as specified by the STREAM pragma, where FIFOs are used instead of RAMs. To see the difference one has to compare Figure 4 with Figure 15. Now the gray scale transformation for both input images (original, template) occurs at the same time in C0 and C1 steps. This applies also for the normalization algorithm. The final performance improves with (cumulated and compared with no optimization solution) is 99 % but the resource utilization is 44.6 %.



**Figure 15.** Streaming data - impact on performance.



**Figure 16.** Streaming data - impact on performance and utilization.

### 2.3. *Area optimization*

As seen before, one has to find the balance between performance and resource utilization. To find the sweet spot between the two some method for area optimization have been used.

### 2.3.1. Allocation

Allocation limits the number of RTL instances and hardware resources used to implement specific. Because I have function instances that are duplicate, to reduce area I

110 can limit the number of instances to 1 utilizing allocation. This reduces resources utilized
111 by the function, but negatively impacts performance. This method reduces dramatically
112 the resource usage to 18 % and decrease the performance just with 1 % (all the other
113 optimization directives are still applied).



**Figure 17.** Allocation - impact on performance and utilization.

114 2.3.2. Array mapping

115     Resource consumption can be reduced up to 16 % if array mapping is used. Array
116 mapping directive combines multiple smaller arrays into a single larger array. This
117 larger array can then be targeted to a single larger memory (RAM or FIFO) resource. The
118 performance decreases but not too much compared to last method.



**Figure 18.** Array mapping - impact on performance and utilization.

119 2.3.3. Bit-accurate types

120     Another two area optimization methods have been used but starting from the
121 default solution. In the previous methods the area optimization started from the solu-
122 tion that targets the highest performance and the purpose was to find the sweet spot
123 between performance and area. The default area consumption is around 17 %. Using
124 bit-accurate types for some variables declared inside each function such as indexes the
125 area consumption decrease to 14.5 % and the performance is the same.



**Figure 19.** Bit-accurate types - impact on performance and utilization.

### 2.3.4. Inline

Other method proposed is to inline the imGrayScale and imNormalization instances. The resource utilization drops to 14.3 %.

Unfortunately no other methods were found to decrease considerably the resource consumption starting from the default optimization. As presented, the decrease in the resource consumption was higher when the performance optimization was used as a starting point.

**Performance Estimates**

**Timing (ns)**

**Summary**

| Clock | Target | Estimated | Uncertainty |
|-------|--------|-----------|-------------|
| ap_clk | 10.00 | 9.40 | 1.25 |

**Latency (clock cycles)**

**Summary**

| Latency | | Interval | | |
|---------|-----|----------|-----|------|
| min | max | min | max | Type |
| 25 | 200638376825 | 26 | 200638376826 | none |

**Detail**

**Instance**

**Utilization Estimates**

**Summary**

| Name | BRAM_18K | DSP48E | FF | LUT |
|------|----------|--------|------|------|
| DSP | - | 1 | - | - |
| Expression | - | - | 0 | 3408 |
| FIFO | - | - | - | - |
| Instance | - | 26 | 5598 | 7872 |
| Memory | - | - | - | - |
| Multiplexer | - | - | - | 1640 |
| Register | - | - | 2451 | - |
| Total | 0 | 27 | 8049 | 12920 |
| Available | 280 | 220 | 106400 | 53200 |
| Utilization (%) | 0 | 12 | 7 | 24 |

**Figure 20.** Inlining - impact on performance and utilization.

133 **References**

134 1. Vivado HLS Optimization Methodology Guide, UG1270 (v2017.4) December 20, 2017
135 2. Vivado HLS – Tips and Tricks, Xilinx
136 3. Vivado Design Suite User Guide, UG902 (v2019.1) July 12, 2019
137 4. Multiprocessors and Reconfigurable Architectures course and lab materials