

DIGITAL SYSTEM DESIGN & METHODOLOGIES II A.Y. 2020/2021

Optional project report

High Level Synthesis – Optimization of a hardware accelerator

Gioele Mombelli

Personal Code: 10520552 gioele.mombelli@mail.polimi.it

0 – Introduction and goals

I have developed this project as a part of the course "Digital System Design & Methodologies" at Politecnico di Milano. The goal of the project is to learn the fundamentals of High Level Synthesis and its tools, with a particular focus on Vitis HLS by Xilinx (the new version of the well-known Vivado HLS) and the features that it provides.

This work is not meant to be more than an academic exercise. I will briefly explain in this report the way of reasoning behind the project, the motivations that driven each solution and the issues faced during the development.

Sources and study material

- Lectures of the course "Digital System Design & Methodologies" by Professor Pilato and Professor Soldavini (AY 2020/2011);
- Vivado Design Suite User Guide High Level Synthesis UG902 (v2017.4) February 2, 2018 by Xilinx;
- Lecture "Introduction to HLS" by Simone Bologna University of Bristol;
- Vivado HLS: debug guide for investigating C/RTL co-simulation issues by Xilinx;
- Designing High-Quality Hardware on a Development Effort Budget: A Study of the Current State of High-Level Synthesis by Sun, Zuo, Campbell, Rupnow, Gurumani, Doucet, Chen.

During the whole project I referred multiple times to the lectures provided by Professor Pilato and Professor Soldavini during the course. While the lectures had a fundamental role in my understanding of the dynamics behind HLS, I also had to search for additional material in order to have a deeper understanding of the processes involved. Last but not the least, some advice and guidance by the teachers themselves have been crucial for the development of the project.

1 – Project overview

The focus of the project was to synthesize and optimize a hardware accelerator. The original code was provided by Professor Pilato and professor Soldavini (file: kernel_gemm.c).

High Level Synthesis allows a designer to obtain a hardware implementation of the functionalities described by a program written in a high level language, such as C or C++, without having to directly write it in a hardware description language, such as VHDL or Verilog. While this makes the process easier for the designer, at least for the part related to bridging the gap between high level algorithms and their hardware implementation, it also requires some skills and knowledge in order fine-tune the synthesis tool and to obtain better performances. One of the biggest issues I faced during the development was indeed predicting the effects of a modification or a directive: some solutions that I thought wouldn't perform well got instead good results, while other solutions that I believed to be optimized worsened the performances of the component.

First analysis of the source code

I started by performing a manual inspection of the code, in order to check if the program was correctly synthesizable. The original source code didn't include any system calls (such as memory allocation functions like malloc or calloc), any recursive functions (which are not synthesizable) and its constructs had a finite or bounded dimension. This first inspection confirmed that the functionality could be implemented in hardware. A guidance and more information on the synthesizability of a generic C code and the properties that must be satisfied can be found in the Xilinx HLS guide.

The second analysis consisted in reading the code to get an understanding of the operations performed. Since the code relied heavily on operations between matrices and tensors, it was particularly prone to a hardware implementation in order to be able to exploit the implicit parallelism of the technology. The various loops executed in the code were the parts that could be exploited in order to obtain improvements, so one of the first step I did was to label each of them. I will briefly report here the original source code with the labels:

[file: src/vanilla/kernel_gemm.c]

```
/*inverse helmholtz gemm.c*/
    /* Copyright (C) 2017-2020 TU Dresden */
3
    /* All Rights Reserved
   /* Authors: Gerald Hempel
               Karl Friebel
   #define X
8
               11
11
9
   #define Y
10
   #define Z
               11
12
   /* general 2D matrix multiplication */
13 void gemm(
     unsigned m1,
     unsigned n1,
    unsigned m2,
16
     unsigned n2,
17
18
    const double in1[m1][n1],
19
     const double in2[m2][n2],
20
     double out[m1][n2]
21 )
22
    GEMM_LOOP1: for (unsigned i = 0; i < m1; i++)</pre>
23
24
     GEMM LOOP2: for (unsigned j = 0; j < n2; j++) {
2.5
         out[i][i] = 0;
2.6
          GEMM_LOOP3: for (unsigned k = 0; k < m2; k++) {
           out[i][j] += in1[i][k] * in2[k][j];
29
        }
30
31
    /** Transposes a matrix. */
33
   void matrix transpose(const double in[X][Y], double out[Y][X])
34
     MTRANS_LOOP1: for (unsigned i = 0; i < X; ++i)
3.5
36
       MTRANS_LOOP2: for (unsigned j = 0; j < Y; ++j)
37
          \operatorname{out}[\overline{j}][i] = \operatorname{in}[i][j];
38
    /** Reshapes and transposes a tensor [X, Y, Z] into [Z, X*Y]. */
39
40
    void reshape transpose(const double in[X][Y][Z], double out[Z][X*Y])
42
     RES LOOP1: for (unsigned i = 0; i < X; ++i)
       RES LOOP2: for (unsigned j = 0; j < Y; ++j)
43
          RES_LOOP3: for (unsigned k = 0; k < Z; ++k)
44
            out[k][i*Y + j] = in[i][j][k];
45
    /** Transposes and reshapes a tensor [Z, X*Y] into [X, Y, Z]. */
47
48
    void transpose reshape(const double in[Z][X*Y], double out[X][Y][Z])
49
      TRANSRES_LOOP1: for (unsigned i = 0; i < X; ++i)
```

```
51
         TRANSRES_LOOP2: for (unsigned j = 0; j < Y; ++j)
           TRANSRES LOOP3: for (unsigned k = 0; k < Y; ++k)
 53
             out[i][j][k] = in[k][i*Y + j];
 54
 5.5
     /** Transposes a tensor [X, Y*Z] into [Z, X*Y]. */
    void tensor transpose(const double in[X][Y*Z], double out[Z][X*Y])
 57
 58
       TTRANS LOOP1: for (unsigned i = 0; i < X; ++i)
         TTRANS_LOOP2: for (unsigned j = 0; j < Y; ++j)
 59
 60
           TTRANS_LOOP3: for (unsigned k = 0; k < Z; ++k)
             out[\bar{k}][i*Y + j] = in[i][j*Y + k];
 61
 62
 63
 64 void kernel body (
      const double S[X][Y],
      const double D[X][Y][Z],
 66
      const double u[X][Y][Z],
 67
 68
      double r[X][Y][Z],
     double S T[Y][X],
 69
 70
      double t\overline{1}[Z][X*Y],
      double t2[Z][X*Y]
 71
 72 )
 73
 74
      reshape transpose(u, t1);
 75
       gemm(X, Y, Z, X*Y, S, t1, t2);
 76
 77
      tensor_transpose(t2, t1);
 78
       gemm(X, Y, Z, X*Y, S, t1, t2);
 79
      tensor_transpose(t2, t1);
 80
      gemm(X, Y, Z, X*Y, S, t1, t2);
 81
 82
       matrix_transpose(S, S_T);
 83
       MULT LOOP1: for (unsigned i = 0; i < X; ++i)
 84
         MULT_LOOP2: for (unsigned j = 0; j < Y; ++j)
 8.5
 86
           MULT_LOOP3: for (unsigned k = 0; k < Z; ++k)
             t2[i][j*Y + k] *= D[i][j][k];
 88
       gemm(X, Y, X, Y*Z, S_T, t2, t1);
 89
 90
      tensor_transpose(t1, t2);
       gemm(X, Y, X, Y*Z, S_T, t2, t1);
tensor_transpose(t1, t2);
 91
 92
 93
      gemm(X, Y, X, Y*Z, S T, t2, t1);
 94
 95
       transpose reshape(t1, r);
 96
 97
 98
 99 void kernel (
     const double S[11][11],
101
      const double D[11][11][11],
      const double u[11][11][11],
102
103
       double r[11][11][11]
104
105 {
      double S_T[Y][X];
106
107
       double t\overline{1}[Z][X*Y];
108
      double t2[Z][X*Y];
109
110
       kernel body(S, D, u, r, S T, t1, t2);
111 }
112
```

Once the code has been analysed and the environment has been set up, the following step was about obtaining a test bench for the code.

2 – Test bench generation

Some time has been dedicated to the development and the generation of a test bench that could be used to test the correctness of the code through the process. The goal of the test bench was *not to test the correctness of the original code*, but to have a reliable input/output association obtained *from the original code* that could be used to check that the behaviour of the application remained constant through the whole process of synthesis and optimization. The test files, called "input" and "output", were generated by the following code, which will be explained later:

[file: test/dataset_generator.c]

```
1 #include <stdio.h>
   #include <stdlib.h>
   #include <math.h>
4 #include <time.h>
5 #include "kernel gemm.c"
   #define RANGE MIN 0.0
   #define RANGE MAX 10.0
10 #define X 11
11 #define Y 11
12 #define Z 11
13
14 double random double(){
16 double range = (RANGE MAX - RANGE MIN);
   double div = RAND_MAX / range;
17
18
   return RANGE MIN + (rand() / div);
19
20
21
22 void write input (const char* file) {
23
24 FILE* input;
25 input = fopen(file, "w");
26 int i = 0;
27 int j = 0;
   int k = 0;
28
   srand(17);
30
31
   //Section 1: writing of S matrix
32 for(i=0; i<X; i++){
       for (j=0; j<Y; j++) {
33
       fprintf(input, "%.15lf\n", random double());
34
35
36
38
39 //Section 2: writing of D tensor:
40 for(i=0; i<X; i++){
    for (j=0; j<Y; j++) {</pre>
               for (k=0; k<Z; k++) {
42
                    fprintf(input, "%.15lf\n", random_double());
43
44
45
46 }
47
48 //Section 2: writing of u tensor:
49 for(i=0; i<X; i++) {
50
    for (j=0; j<Y; j++) {</pre>
               for (k=0; k<Z; k++) {
    fprintf(input, "%.15lf\n", random_double());</pre>
51
52
54
        }
55 }
```

```
fclose(input);
 58
 59
 60
     void load input(const char* file, double S[X][Y], double D[X][Y][Z], double u[X][Y][Z]){
 62
 63 FILE* input;
 64 input = fopen(file, "r");
 65
     double number = 1.0;
 66
    int i = 0;
     int j = 0;
 67
 68
     int k = 0;
 69
 70
     //Section 1: reading of S matrix
     printf("S matrix:\n");
 71
 72
     for (i=0; i<X; i++) {</pre>
          for(j=0; j<Y; j++) {
    fscanf(input, "%lf[^\n]", &number);
    printf("S[%d][%d] = %.15lf\n", i, j, number);</pre>
 7.3
 74
 75
 76
               S[i][j] = number;
 77
 78
 79
     //Section 2: reading of D tensor
printf("D tensor:\n");
 8.0
 81
 82
     for(i=0; i<X; i++) {</pre>
 83
          for (j=0; j<Y; j++) {</pre>
 84
                   for(k=0; k<Z; k++) {
                        fscanf(input, "%lf[^\n]", &number);
printf("D[%d][%d][%d] = %.15lf\n", i, j, k, number);
 8.5
 86
 87
                        D[i][j][k] = number;
 88
                   }
 89
 90
 91
     //Section 3: reading of u tensor
 93
     printf("u tensor:\n");
     for (i=0; i<X; i++) {
 94
 9.5
          for (j=0; j<Y; j++) {</pre>
 96
                   for (k=0; k<Z; k++) {
                        fscanf(input, "%lf[^\n]", &number);
 98
                        printf("u[%d][%d] [%d] = %.15lf\n", i, j, k, number);
 99
                        u[i][j][k] = number;
100
101
          }
102
103
104
105
106
107
108 void write output (const char* file, double r[X][Y][Z]) {
109
110 FILE* output;
111 output = fopen(file, "w");
112
113 printf("r tensor:\n");
     for (int i=0; i<X; i++) {</pre>
114
          for (int j=0; j<Y; j++) {
115
                   for(int k=0; k<Z; k++){
    printf("r[%d][%d] [%d] = %.15lf\n", i, j, k, r[i][j][k]);</pre>
116
117
118
                        fprintf(output, "%.15lf\n", r[i][j][k]);
119
120
121
122
123 fclose (output);
124
125
126
127
     int main(){
128
129 double S[X][Y];
130 double D[X][Y][Z];
131 double u[X][Y][Z];
     double r[X][Y][Z];
132
133
```

```
134 write input("C:/Users/gio m/Desktop/ProgettoDSDM/src/input");
135 printf("The following input test case has been generated:\n");
load input("C:/Users/gio m/Desktop/ProgettoDSDM/src/input", S, D, u);
137
138
    printf("Start execution of original source code kernel gemm...\n");
139
140
    kernel(S, D, u, r);
141
142 printf("The following output has been generated:\n");
143
    write output("C:/Users/gio m/Desktop/ProgettoDSDM/src/output", r);
144
145
146 printf("Test generation complete. You can find the results in the input file and in the output
file.\n");
147
148 return 0;
149
```

This code simply generates random floating point values between 0 and 10 and writes them to a file. There are just enough values to fill the inputs of the original source code: S, D, u and r. The values are then read from the file and used as inputs for an execution of the "kernel" function, the one to be synthesized. The output of the function is finally written to another file, called "output". In this way I have obtained an input file and its associated output from the original source code. They can be found in the attached folder "test".

The next step was to write the actual test bench that the HLS tool would have used. This is the code of the test bench:

[file: test/kernel_gemm_tb.c]

```
#include <stdio.h>
   #include <stdlib.h>
   #include <math.h>
   #define X 11
 6
   #define Y 11
 7
    #define Z 11
 8
   #define DOUBLE EPSILON 1E-10
10
   #define INPUT "C:/Users/gio m/Desktop/ProgettoDSDM/src/input"
11
   #define OUTPUT "C:/Users/gio_m/Desktop/ProgettoDSDM/src/output"
12
14 int almost_equal(double a, double b, double epsilon) {
15     printf("Epsilon: %.15lf\n", fabs(a-b));
16
        return fabs(a-b) <= epsilon;</pre>
17
19 void load input(const char* file, double S[X][Y], double D[X][Y][Z], double u[X][Y][Z]) {
20
21 FILE* input;
   input = fopen(file, "r");
22
23 double number = 1.0;
24 int i = 0;
25 int j = 0;
26 int k = 0;
   //Loading 1: reading of S matrix
2.8
29 for(i=0; i<X; i++){
30
    for (j=0; j<Y; j++) {</pre>
         fscanf(input, "%lf[^\n]", &number);
32
            S[i][j] = number;
33
        }
34
36
   //Loading 2: reading of D tensor
37 for(i=0; i<X; i++){
38
     for(j=0; j<Y; j++){
39
                for(k=0; k<Z; k++) {
                     fscanf(input, "%lf[^\n]", &number);
40
                     D[i][j][k] = number;
41
```

```
42
 44
 4.5
 46
     //Loading 3: reading of u tensor
 47 for(i=0; i<X; i++) {
     for(j=0; j<Y; j++){
 48
                for (k=0; k<Z; k++) {
 49
                    fscanf(input, "%lf[^\n]", &number);
 5.0
 51
                     u[i][j][k] = number;
 52
 53
    }
 54
 55
 57
 5.8
 59
 60 int compare(const char* file, double r[X][Y][Z]){
 62 FILE* output;
 63 output = fopen(file, "r");
 64 double current value = -1.0;
 66 for(int i=0; i<X; i++) {
       for (int j=0; j<Y; j++) {</pre>
 67
 68
                 for(int k=0; k<Z; k++) {</pre>
 69
                     fscanf(output, "%lf[^\n]", &current value);
                     printf("Current value: %.15lf - Test value: %.15lf\n", current value,
r[i][j][k]);
 71
 72
                     if(!almost_equal(current_value, r[i][j][k], DOUBLE_EPSILON)){
 73
                         fclose (output);
 74
                         return (-1);
 7.5
 76
                }
 77
 78 }
 79
 80 fclose (output);
 81
    return(0);
 83
 84
 85 int main(){
 86
 87 double S[X][Y];
 88 double D[X][Y][Z];
 89
    double u[X][Y][Z];
 90 double r[X][Y][Z];
 91
 92 load input (INPUT, S, D, u);
 93
 94 kernel(S, D, u, r);
 96 int test_result = compare(OUTPUT, r);
 97
 98 if(test result == 0){
       printf("The test was successful!\n");
 99
100
        return 0;
101 }
102
103 if(test result == -1){
       printf("Test failed!\n");
104
105
        return -1;
106 }
108 printf("Unexpected error.");
109 return 1;
110
111
```

The most useful resource in writing the test bench was the Xilinx HLS user guide, which provided me a complete explanation of the best practices that should be adopted when writing a test bench, such as the required return values. One of the issues I faced was the comparison of two floating

point numbers: the function at line 14 performs the comparison by stating that two values are equal if their difference is less than a chosen epsilon, in this case 10⁻¹⁰. This seemed to me a reasonably small value in order to consider two numbers as equal, since in real application such a small difference could be neglectable most of the times.

3 – Solutions evaluation

Between all the solutions I have tried, some of them turned out to be improvements while others didn't affect at all the overall performance, or even worsened it. Others brought an improvement but failed the Co-Simulation. Here I will report them in chronological order, providing explanations and the reasoning behind each one. The first solution I implemented was the trivial one, in order to have a reference value to evaluate the improvement.

Note: all the solutions provided here are obtained by running, in order, C simulation, C synthesis and C/RTL co-simulation. Every solution provided passes C simulation correctly. The target clock is set as 10.0ns and the target device is the same for each solution presented here.

Solution 1: Trivial (with default pipeline)

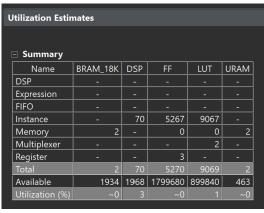
Source file: src/vanilla/kernel_gemm.c

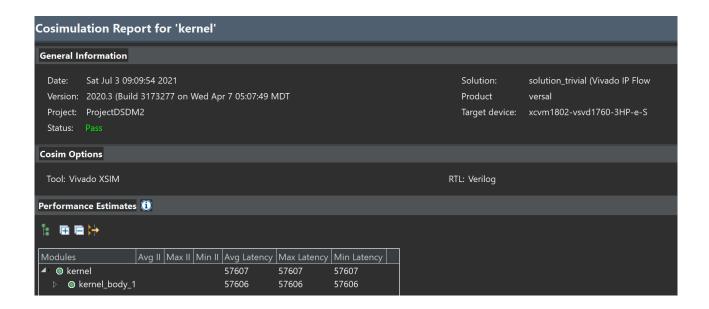
User defined directives: none

Latency: 57607 cycles Co-simulation result: pass

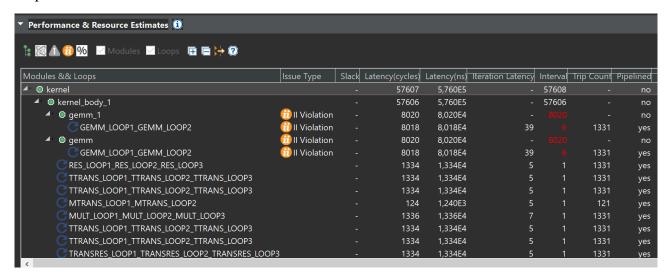
This solution is obtained by running HLS on the vanilla version of the code, without any optimization or directive applied. Actually, Vitis HLS tried to apply some default optimization like loop pipelining, but I decided to consider this as the starting point anyway. The synthesis required a few minutes and the following results were obtained:







As we can see from the utilization estimates, the resource utilization is practically zero, but the number of clock cycles needed to complete the computation is quite high. This is an exhibition of the classic area/latency trade-off, so I decided to focus my project on improving the latency of the component since there is plenty of resource available. In order to do this, I relied on the insight given to me by Vitis, in particular the section that highlighted the various cycles needed to each loop:



As we can see, Vitis is trying to pipeline all the loops of the code with a II value of 1, but some of them result in a II Violation, and the initiation interval is longer. This gave me the idea of trying to partition the variables involved in the heaviest computations. But first, I tried out of curiosity to give manually the Pipeline directive, as seen in the next solution.

Solution 2: Pipelined

Source file: src/vanilla/kernel_gemm.c

User defined directives: pipeline

Latency: 57607 cycles Co-simulation result: pass

This solution will only be reported briefly since it didn't bring better performances. The only thing done here has been to manually apply the PIPELINE directive on the innermost loop of "gemm" function (GEMM_LOOP3). The fact that nothing changed in the results confirmed that Vitis already tried to pipeline the loops.

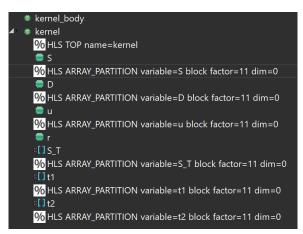
Solution 3: Pipelined and Partitioned

Source file: src/vanilla/kernel_gemm.c

User defined directives: pipeline and partition

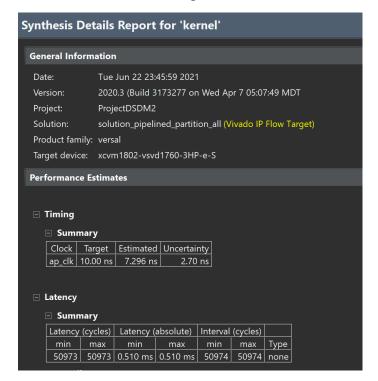
Latency: 50974 cycles Co-simulation result: fail

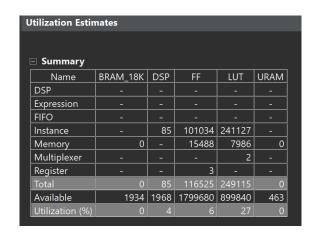
This solution was obtained from the former one by applying also the ARRAY_PARTITION directive to all the variables involved in the code, with a block factor of 11 on each dimension (this was a brute-force approach that basically decomposed each element of the array into a single value assigned to a register). The reasoning behind this has been the following: since Vitis notified that there weren't enough memory ports to manage the data during the computation given the optimizations applied, the placement of the data into register structures should have led to an improvement.

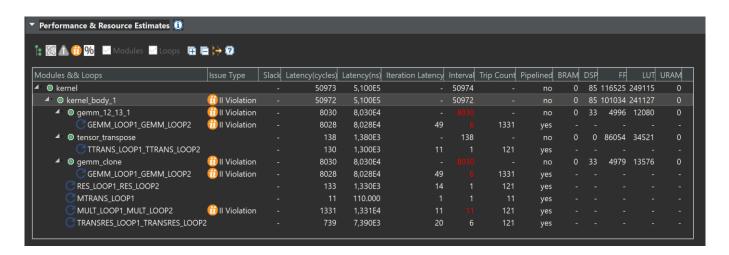


Directives implemented in this solution.

Here the results of this implementation:







As we can see from the report, there has been indeed an improvement of almost 7000 cycles, paired obviously with an increase of the resources exploited. The synthesis of this kind of implementation took a huge amount of time, probably because the tool had to build a large multiplexer to manage the access to each single data, and the results aren't as good as I was expecting. Moreover, this solution failed the C Co-Simulation. I decided to include it anyway in order to provide feedback on the work done, but I have noticed that a lot of solutions that used Partitioning ended up failing Co-Simulation. The debugging of this issue has been tough and even after some work I couldn't figure out what was going wrong, so I decided to try other approaches, such as unrolling loops.

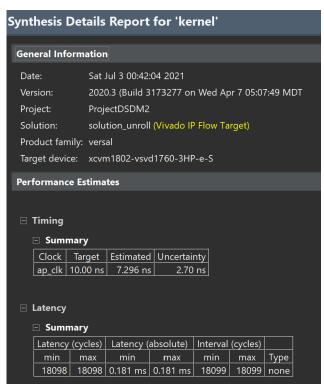
Solution 4: Loop Unrolling

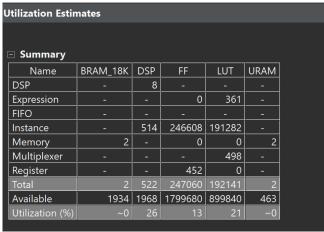
Source file: src/vanilla/kernel_gemm.c

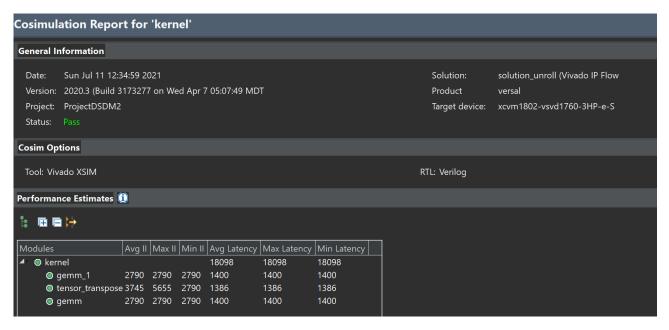
User defined directives: unrolling of GEMM_LOOP3 and TTRANS_LOOP2

Latency: 18098 cycles Co-simulation result: pass

By analyzing the code and the report given by Vitis, it can be seen that the loop executed by the functions "gemm" and "tensor_transpose" are executed multiple times in the component, contributing for a large part in the overall count of clock cyles. I decided to start applying unrolling only to those two loops in order to check if the improvement was greater than before, and eventually progressing with the unrolling of every loop. Here the results obtained:

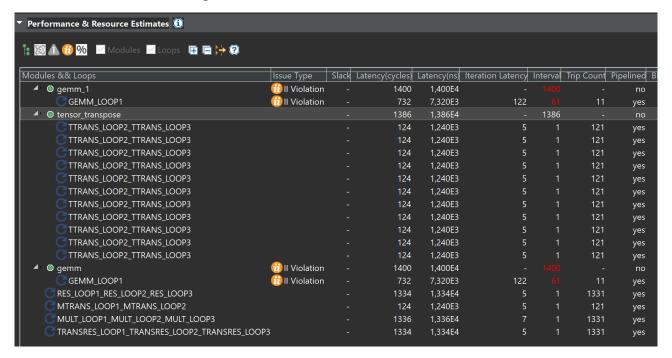






As the metrics of Vitis show, the unrolling of the loops allowed to obtain the first significant improvement in the performance of the accelerator, by saving 39.509 clock cycles. This is an increase in performance of **68,5%** with respect to the trivial solution. The overall resource exploitation increased; however, the DSP consumption had an increase of just +22% and the flip flops of +7%, and the LUT utilization even decreased from 27% to 21%.

Here is the metric of each loop:



The next solution naturally followed from the results of this one: a full unrolling of each loop.

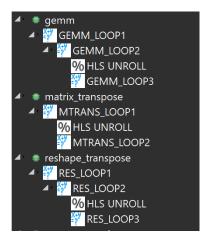
Solution 5: Every Loop Unrolling

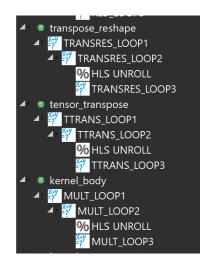
Source file: src/vanilla/kernel_gemm.c

User defined directives: unrolling of every loop executed

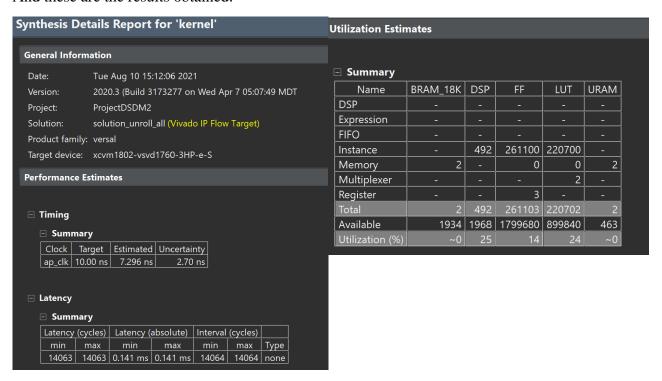
Latency: 14063 cycles Co-simulation result: pass

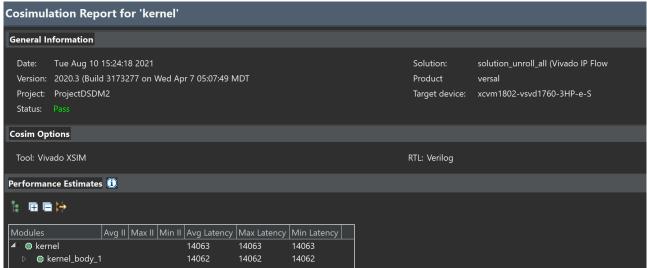
Very similar to the previous one but with slightly better performances. Here the optimizations applied:



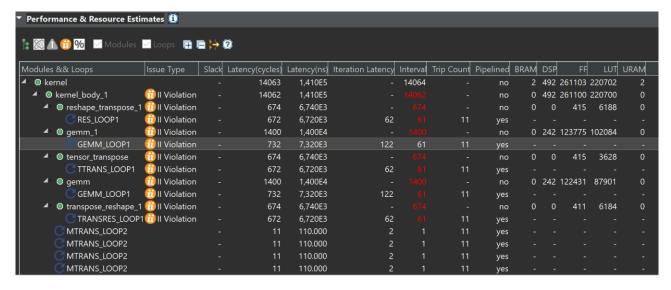


And these are the results obtained:





This time the improvement was of more than 4000 cycles on the previous one value. This is far less than the amount of the previous solution but still a better result, especially if compared to the trivial approach. The utilization increased but each resource is occupied for less than 25%. It is important to notice that Vitis is still trying to automatically apply pipelining where possible:



Following this solution, several approaches were tried in order to obtain improvements. I tried to apply another time the partitioning, but even if applied only to one or two variables, the code **always failed co-simulation.** This issue remained during the whole process and in the end no solution was able to work if the partitioning was enabled.

After various attempts, no further directives (or combination of them) brought improvements to the project. I decided that it was time to start working on the source code, keeping the loops unrolled, in order to find a version of it that could perform better. The unsuccessful solutions are briefly reported, while the working ones are described in detail.

Solution 6: Gemm Modules Duplication

This solution will be briefly described just for keeping track of the work done. I tried to duplicate the various "gemm" functions by making copies of them and calling one and only one of them each time in the kernel function. I thought that this approach could force Vitis to implement a single module for each function and speed up the computation, but this didn't bring any improvement.

Solution 7: Removing Ping-Pong Buffer

The kernel function uses two arrays, t1 and t2, as a ping pong buffer: the output of a computation is saved in one of them, in order to be used as an input on the next step and vice-versa. This solution tried to use dedicated arrays for each operation in order to have arrays always used as input and others always dedicated to the output. The solution however didn't present improvements or different behaviours from the previous ones, so it is not reported here in detail. I will just mention how the different arrays have been implemented:

Original code:

```
74 reshape_transpose(u, t1);
75
76 gemm(X, Y, Z, X*Y, S, t1, t2);
77 tensor transpose(t2, t1);
```

Edited code:

```
74 reshape_transpose(u, t1);
75 gemm(X, Y, Z, X*Y, S, t1, t2);
76 double t3[Z][X*Y];
77 tensor_transpose(t2, t3);
78 double t4[Z][X*Y];
79 gemm(X, Y, Z, X*Y, S, t3, t4);
and so on...
```

Solution 8: Transpose Deleted

Source file: src/deleted_transpose/kernel_gemm.c

User defined directives: unrolling of every loop executed

Latency: 11359 cycles Co-simulation result: pass

One of the best solutions has been found thanks to a modification of the original source code. This variation isn't actually hardware related, since it does not expose more parallelism, but is a kind of programming trick that allows to save some clock cycles. Analyzing the behavior of the code, it was noticed that some cycles were lost to carry out the transposition of the matrix by executing the function "tensor_transpose". In particular, the transposition happened always after an execution of the "gemm" function on the same matrix. To speed up the operation, I decided to condense the two functions "tensor_transpose" and "gemm" into a single "gemm_transpose" function, which performs the multiplication on the matrix without transposing it but using instead the right sequence of values. In practice, instead of transposing the matrix, it is left as it is but read as if it were transposed. To be clearer, I will briefly report here the modification applied to the source code.

Each occurrence of the following sequence of function calls:

```
gemm(X, Y, Z, X*Y, S, t1, t2);
tensor_transpose(t2, t1);
```

Has been substituted with the single function call:

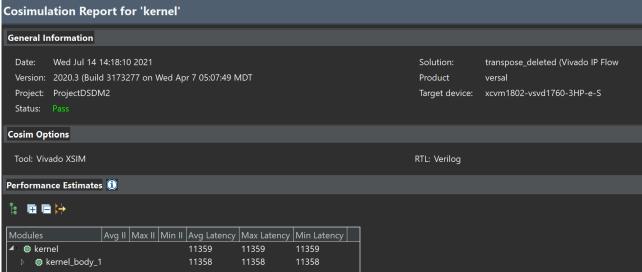
```
gemm transpose (X, Y, Z, X*Y, S, t2, t1);
```

And as it can be seen from the code, the latter function performs the matrix multiplication by reading the values of in2 as if it were transposed:

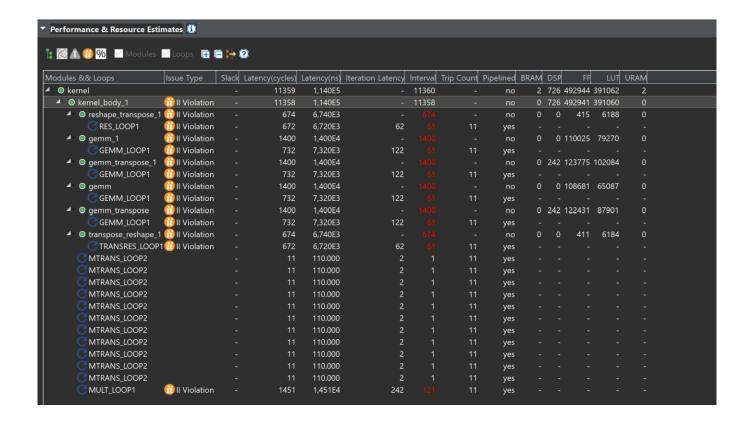
```
34 void gemm transpose (
35 unsigned m1,
36
      unsigned n1,
 37
      unsigned m2.
 38
    unsigned n2,
 39
      const double in1[m1][n1],
     const double in2[m2][n2],
      double out[m1][n2]
 41
 42
 43 {
 44
       int h = 0;
        int 1 = 0;
 46
 47
     GEMMTR LOOP1: for (unsigned i = 0; i < m1; i++) {
 48
         GEMMTR LOOP2: for (unsigned j = 0; j < n2; j++) {
 50
 51
          out[i][j] = 0;
 52
           GEMMTR LOOP3: for (unsigned k = 0; k < m2; k++) {
            out[\bar{i}][j] += in1[i][k] * in2[h][l];
 54
            1++;
 55
          }
        }
 56
 57
        h++;
 58
        1=0;
 59
      }
 60
```

This modification allowed to skip the call of the "tensor_transpose" function, allowing us to remove it from the code. The result can be seen from the report of the synthesis tool:





This version improved the former one of almost 3000 cycles, reaching a performance of 11359 cycles, which is an improvement of **80,3%** of the original performance. As the testing results and the Co-Simulation report show, the code still behaves as intended and the results are the same. On the next page is reported the complete performance estimate of each loop. It is worth to remember that the same directives of the previous working solution ("Every Loop Unrolling") are applied here, so there is a combination of automatic pipelining and unrolling on each loop. Some II Violations are still present, but the results are much better than the first version.



Quite satisfied with this result, I decided to try to solve the II Violations to make pipelining work better. So, I tried again to partition the variables. This can be seen on the next solution.

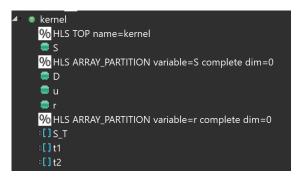
Solution 9: Transpose Deleted + Double Partition

Source file: src/deleted_transpose/kernel_gemm.c

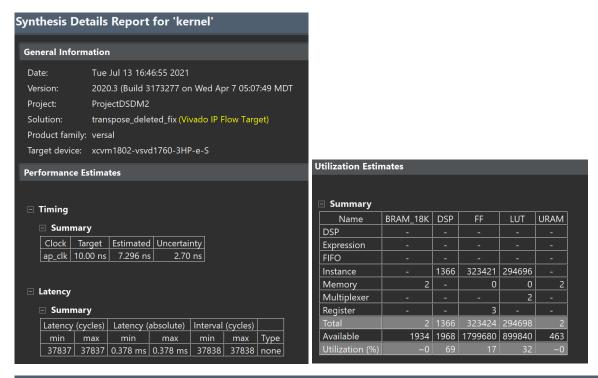
User defined directives: unrolling of every loop executed, complete partitioning on S and r

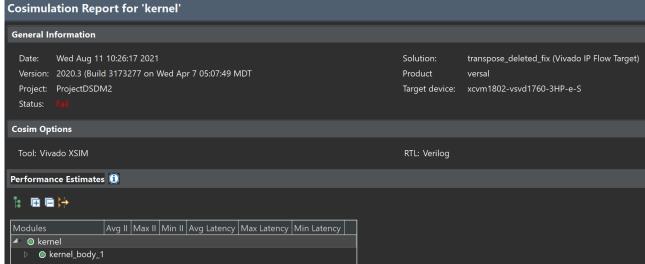
Latency: 37837 cycles Co-simulation result: fail

In order to fix the II Violations given by the previous solution, I tried to apply partitioning on S and r with the "complete" option on every dimension. For clarity, here the directives related to the partitioning:

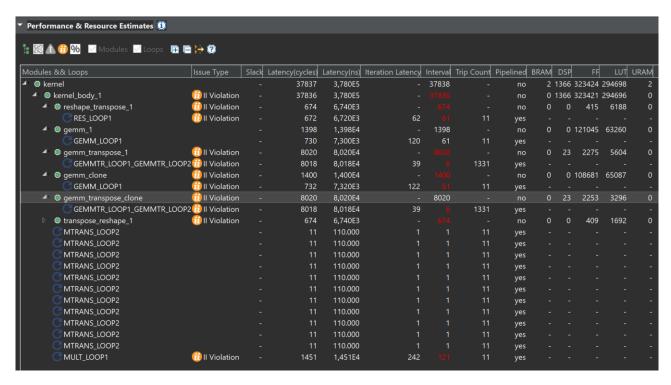


I was confident that, despite the usual long and boring synthesis time, this solution would have brought a considerable improvement on the performance of the component. As the report of Vitis shows, I was wrong.





The first thing to be noticed is that the number of clock cycles considerably increased in this solution: more than 26.000 cycles with respect to the previous solution. The DSP utilization also increased, but the remarkable thing is that Co-Simulation failed.



At this point I decided that, despite the results, I would have kept trying with partitioning. The next solution shows some promising but not satisfying results.

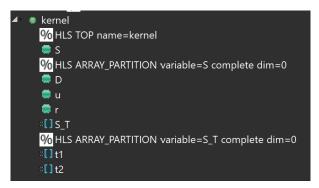
Solution 10: Transpose Deleted + S Partition

Source file: src/deleted_transpose/kernel_gemm.c

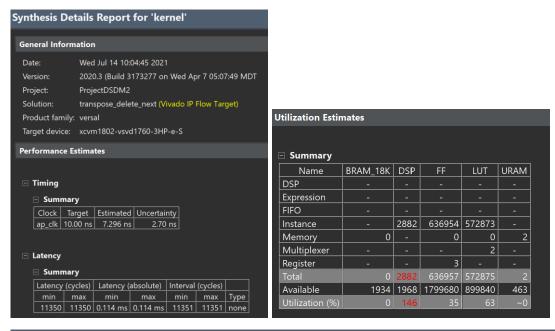
User defined directives: unrolling of every loop executed, complete partitioning on S and S_T

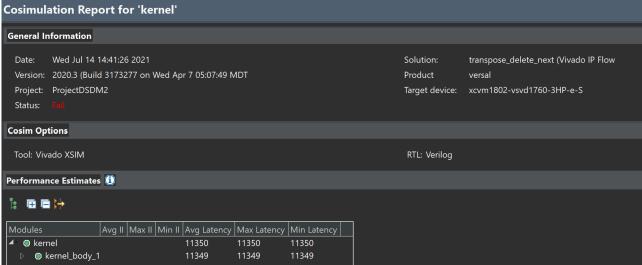
Latency: 11350 cycles Co-simulation result: fail

The idea behind this solution was to apply partitioning only on the smaller matrix S (and the S_T variable used to store it), being it an 11x11 matrix, smaller than the 11x11x11 tensors. Here the directives applied:



At first, I was surprised that this result actually obtained the best cycles estimation of all the solutions tried so far, even if it was an improvement of only 9 cycles with respect to the best solution without partitioning:





Inspecting the reports given by the tool we can however find some problems: the DSP utilization is 146%, and the Co-Simulation presents a failure. Since the improvement on the speed of the component is neglectable, I decided to stop trying to apply partitioning to this version of the code, since the utilization of this directive brought more disadvantages than improvements and the debugging of Co-Simulation results was not possible. I decided to revert the changes and try other modifications to be applied to the code.

Solution 11: Mult Extracted

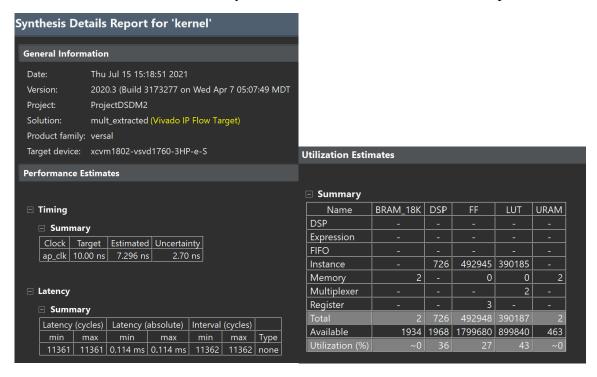
Source file: src/mult_extracted/kernel_gemm.c

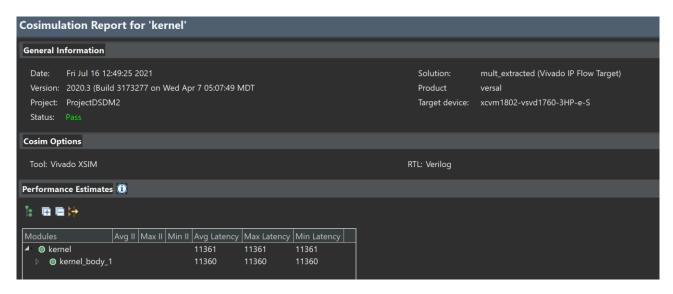
User defined directives: unrolling of every loop executed

Latency: 11361 cycles Co-simulation result: pass

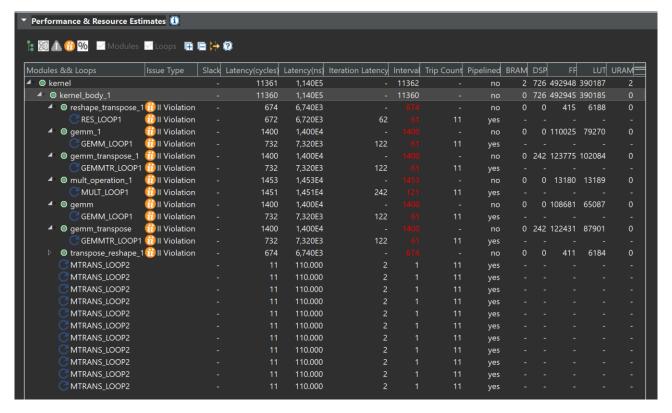
The rationale behind this solution was trying to extract the multiplication between t2 and D performed in the "kernel_body" function and substituting it with a function call to a dedicated function. This is the definition of the function:

This modification led to a 11361 cycles solution, as can be seen from the reports:





This result is surely better than the last one, but still could not outperform the results obtained in Solution 8. For completeness, in the next page the details of each loop:



Once again, I tried to apply partitioning to some variables of this code in order to improve performances.

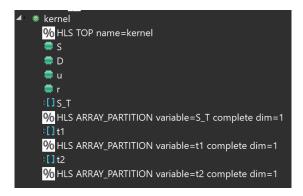
Solution 12: Mult Extracted Partition

Source file: src/mult_extracted/kernel_gemm.c

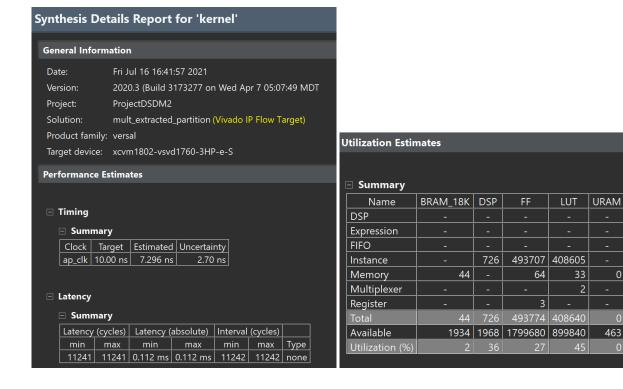
User defined directives: unrolling of every loop executed + partitioning on S_T, t1, t2

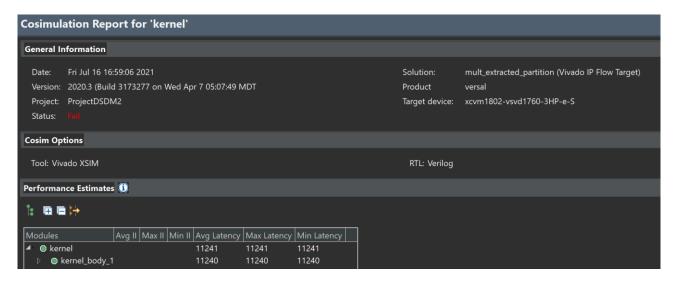
Latency: 11241 cycles Co-simulation result: fail

This solution is the same as number 11 with the partitioning directive applied to S_T, t1 and t2. This time the partitioning has been applied only to the first dimension of the variables and the arrays involved.



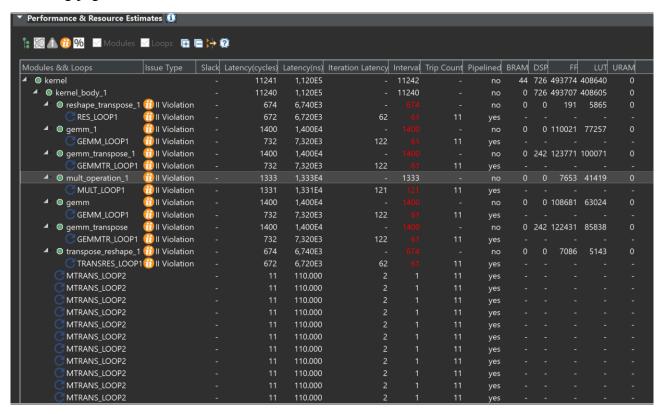
The result seems good with respect to the number of clock cycles, with the code obtaining a value of 11241, which makes this the best performing solution so far. Unfortunately, as the following reports show, the solution fails Co-Simulation.





The utilization percentages are all still under 50%, so the component has still resources available to further improvement of the performance.

The II Violation are still present on the loops, as can be seen from the detailed report in the following page:



This solution seems to be a good starting point for further work, since the estimation of clock cycles is very good, but unfortunately inspecting the failure in Co-Simulation is a difficult task.

Solution 13: Mult Extracted Loop Merging

This solution will be just mentioned because it didn't bring any improvement with respect to solution 11. The solution is the same as solution 11 with the LOOP_MERGE directive set on the "gemm" loops. It is reported here just for completeness.

4 - Final considerations

Among all the solutions considered, the number 8 is the one that obtained the best results. It can perform the computation in 11359 cycles on the target architecture without using too much of the space available.

This is surely not the best solution in the absolute way, since I believe that better results could be obtained: the target hardware has still resources available and not all the warnings have been solved by the solutions presented here. Moreover, even if the original source code includes a certain degree of serial computations, which cannot really be optimized for the execution in hardware beside some pipelining (for example, every time an operation requires the results of a previous one), other things could be really sped up. I am referring in particular to the operations between elements of the matrices: with the correct number of multipliers they could be done in parallel by saving a considerable amount of time, spending just the cycles needed to read data, perform a multiplication and propagating the output. While this is quite easy to think at an intuitive level, it is not immediate to instruct the HLS tool to do exactly that.

In the end, what I learned from this project is that HLS tools offer a great bridge between high level code and hardware implementation, making the process of obtaining an architecture much faster and, most importantly, allowing the designer to reason at a higher level, closer to human thinking, without having to deal with hardware description languages. This is very powerful since the architecture can be tested and eventually modified according to the desired behaviour by editing a more understandable code. The drawback is that the synthesis tools sometimes act like a "black box" where the output is hard to predict. Some results surprised me since they performed way worse or way better than my expectations. I believe that this kind of tools require great experience and an extensive knowledge of the hardware architectures in order to be exploited at full power.