

High Level Synthesis of Digital Systems

2023-2024

Prof.ssa PERRI

Prof. FRUSTACI

FIR Filter Analysis

Giorgio Ubbriaco

247284

bbrgrg00h11d086x@studenti.unical.it

Maggio 2024

Index

1 Introduction	5
1.1 LTI Systems	5
1.2 Causal Systems	5
1.3 FIR Filter	5
1.4 High Level Synthesis	6
2 Tasks to be performed	7
3 Definitions	8
4 C Simulations	11
5 Solutions	16
5.1 Unoptimized Solution	16
5.2 Operation Chaining Solution	23
5.3 Code Hoisting Solution	30
5.4 Loop Fission Solution	34
5.5 Loop Unrolling Solution	38
5.5.1 Loop Unrolling Factor=2	40
5.5.2 Loop Unrolling Factor=4	47
5.6 Loop Pipelining Solution	54
5.7 Bitwidth Optimization Solution	58
5.8 AXI Solution	59
6 Conclusions	60

Listings

definitions/definitions.h	8
c_simulations/fir_tb.cpp	11
c_simulations/fir_tb_output.txt	13
solutions/unoptimized/fir_unoptimized.cpp	16
solutions/unoptimized/firConvolutionUnoptimized_top.vhd	19
solutions/unoptimized/firConvolutionUnoptimized_tb.vhd	20
solutions/unoptimized/clk_constraint.xdc	21
solutions/operation_chaining/10ns/clk_constraint.xdc	25
solutions/operation_chaining/9ns/clk_constraint.xdc	25
solutions/operation_chaining/8ns/clk_constraint.xdc	25
solutions/operation_chaining/7ns/clk_constraint.xdc	25
solutions/operation_chaining/6ns/clk_constraint.xdc	25
solutions/operation_chaining/5ns/clk_constraint.xdc	25
solutions/operation_chaining/4ns/clk_constraint.xdc	25
solutions/operation_chaining/3ns/clk_constraint.xdc	25
solutions/code_hoisting/fir_code_hoisting.cpp	30
solutions/code_hoisting/clk_constraint.xdc	32
solutions/loop_fission/fir_loop_fission.cpp	34
solutions/loop_fission/clk_constraint.xdc	36
solutions/loop_unrolling/factor2/fir_loop_unrolling_manual_factor2.cpp	40
solutions/loop_unrolling/factor2/fir_loop_unrolling_automatic_factor2.cpp	40
solutions/loop_unrolling/factor2/fir_loop_unrolling_automatic_partitioning_factor2.cpp	40
solutions/loop_unrolling/factor2/clk_constraint.xdc	43
solutions/loop_unrolling/factor4/fir_loop_unrolling_manual_factor4.cpp	47
solutions/loop_unrolling/factor4/fir_loop_unrolling_automatic_factor4.cpp	47
solutions/loop_unrolling/factor4/fir_loop_unrolling_automatic_partitioning_factor4.cpp	47
solutions/loop_unrolling/factor2/clk_constraint.xdc	50
solutions/loop_pipelining/fir_loop_pipelining.cpp	55
solutions/loop_pipelining/clk_constraint.xdc	57

List of Figures

1 HLS Design Flow	6
2 HLS Clock Period from Synthesis Report	17
3 HLS Unoptimized Solution If Statement Analysis	17
4 HLS Unoptimized Solution Else Statement Analysis	18
5 HLS Unoptimized Solution Read Operations Analysis	18
6 HLS Unoptimized Solution Read Operations Analysis	18
7 HLS Unoptimized Solution Write Operations Analysis	18
8 HLS Unoptimized Solution Write Operations Analysis	18
9 Vivado Operation Chaining Solution clk=3ns Timing Constraint Violated	26
10 Vivado Operation Chaining Solution clk=3ns Timing Constraint Violated	26
11 Vivado Operation Chaining Utilization Plot	26
12 Vivado Operation Chaining Timing Plot	27
13 Vivado Operation Chaining Power Plot	29
14 HLS Code Hoisting Solution Trip Count Analysis	31
15 HLS Loop Fission Solution Analysis	35
16 HLS Array Partitioning	38
17 HLS Loop Unrolling	39
18 HLS Loop Unrolling Manual Factor=2 Analysis	42
19 HLS Loop Unrolling Automatic Factor=2 Analysis	42
20 Vivado Loop Unrolling Factor=2 Utilization Plot	44

21	Vivado Loop Unrolling Factor=2 Timing Plot	45
22	Vivado Loop Unrolling Factor=2 Dynamic Power Plot	46
23	HLS Loop Unrolling Manual Factor=4 Analysis	49
24	HLS Loop Unrolling Automatic Factor=4 Analysis	49
25	Vivado Loop Unrolling Factor=4 Utilization Plot	51
26	Vivado Loop Unrolling Factor=4 Timing Plot	52
27	HLS Loop Pipelining	54
28	HLS Synthesis Report	61
29	HLS Synthesis Utilization Estimates Report Plot	62
30	HLS Synthesis Trip Count Report Plot	62
31	HLS Synthesis Latency Reportv Plot	62
32	HLS Synthesis Interval Report Plot	63
33	HLS C/RTL Cosimulation Report	64
34	HLS C/RTL Cosimulation Report Report Plot	65
35	HLS Export RTL Report	66
36	HLS Export RTL Report Report Plot	67

List of Tables

1	HLS Unoptimized Solution Timing Summary (ns)	16
2	HLS Unoptimized Solution Latency Summary (clock cycles)	16
3	HLS Unoptimized Solution Latency Loops Summary	17
4	HLS Unoptimized Solution Utilization Estimates Summary	19
5	HLS Unoptimized Solution C/RTL Cosimulation Summary	19
6	HLS Unoptimized Solution Export RTL Resource Usage	19
7	HLS Unoptimized Solution Export RTL Final Timing	19
8	Vivado Unoptimized Solution Utilization Report [#]	21
9	Vivado Unoptimized Solution Timing Report	22
10	Vivado Unoptimized Solution Dynamic Power Report [mW]	22
11	Vivado Unoptimized Solution Dynamic Power Report [mW]	22
12	Vivado Unoptimized Solution Energy Single Operation Report [pJ]	22
13	HLS Operation Chaining Solution Timing Summary (ns)	23
14	HLS Operation Chaining Solution Latency Summary (clock cycles)	23
15	HLS Operation Chaining Solution Latency Loops Summary	23
16	HLS Operation Chaining Solution Utilization Estimates [#]	24
17	HLS Operation Chaining Solution C/RTL Cosimulation Report	24
18	HLS Operation Chaining Solution Export RTL Report	25
19	Vivado Operation Chaining Solution Utilization Report [#]	26
20	Vivado Operation Chaining Solution Timing Report	27
21	Vivado Operation Chaining Solution Dynamic Power Report [mW]	28
22	Vivado Operation Chaining Solution Dynamic Power Report [mW]	28
23	Vivado Operation Chaining Solution Energy Single Operation Report [pJ]	28
24	HLS Code Hoisting Solution Timing Summary (ns)	30
25	HLS Code Hoisting Solution Latency Summary (clock cycles)	30
26	HLS Code Hoisting Solution Latency Loops Summary	30
27	HLS Code Hoisting Solution Utilization Estimates Summary	31
28	HLS Code Hoisting Solution C/RTL Cosimulation Summary	32
29	HLS Code Hoisting Solution Export RTL Resource Usage	32
30	HLS Code Hoisting Solution Export RTL Final Timing	32
31	Vivado Code Hoisting Solution Utilization Report [#]	32
32	Vivado Code Hoisting Solution Timing Report	32
33	Vivado Code Hoisting Solution Dynamic Power Report [mW]	32
34	Vivado Code Hoisting Solution Dynamic Power Report [mW]	33
35	Vivado Code Hoisting Solution Energy Single Operation Report [pJ]	33

36	HLS Loop Fission Solution Timing Summary (ns)	34
37	HLS Loop Fission Solution Latency Summary (clock cycles)	34
38	HLS Loop Fission Solution Latency Loops Summary	34
39	HLS Loop Fission Solution Utilization Estimates Summary	35
40	HLS Loop Fissiong Solution C/RTL Cosimulation Summary	36
41	HLS Loop Fissiong Solution Export RTL Resource Usage	36
42	HLS Loop Fissiong Solution Export RTL Final Timing	36
43	Vivado Loop Fission Solution Utilization Report [#]	36
44	Vivado Loop Fission Solution Timing Report	36
45	Vivado Loop Fission Solution Dynamic Power Report [mW]	36
46	Vivado Loop Fission Solution Dynamic Power Report [mW]	36
47	Vivado Loop Fission Solution Energy Single Operation Report [pJ]	36
48	HLS Loop Unrolling Factor=2 Solution Timing Summary (ns)	41
49	HLS Loop Unrolling Factor=2 Solution Latency Summary (clock cycles)	41
50	HLS Loop Unrolling Factor=2 Solution Latency Loops Summary	41
51	HLS Loop Unrolling Factor=2 Solution Utilization Estimates [#]	43
52	HLS Loop Unrolling Factor=2 Solution C/RTL Cosimulation Report	43
53	HLS Loop Unrolling Factor=2 Solution Export RTL Report	43
54	Vivado Loop Unrolling Factor=2 Solution Utilization Report [#]	44
55	Vivado Loop Unrolling Factor=2 Solution Timing Report	44
56	Vivado Loop Unrolling Factor=2 Solution Dynamic Power Report [mW]	45
57	Vivado Loop Unrolling Factor=2 Solution Dynamic Power Report [mW]	46
58	Vivado Loop Unrolling Factor=2 Solution Energy Single Operation Report [pJ]	46
59	HLS Loop Unrolling Factor=4 Solution Timing Summary (ns)	48
60	HLS Loop Unrolling Factor=4 Solution Latency Summary (clock cycles)	48
61	HLS Loop Unrolling Factor=4 Solution Latency Loops Summary	48
62	HLS Loop Unrolling Factor=4 Solution Utilization Estimates [#]	50
63	HLS Loop Unrolling Factor=4 Solution C/RTL Cosimulation Report	50
64	HLS Loop Unrolling Factor=4 Solution Export RTL Report	50
65	Vivado Loop Unrolling Factor=4 Solution Utilization Report [#]	50
66	Vivado Loop Unrolling Factor=4 Solution Timing Report	51
67	Vivado Loop Unrolling Factor=4 Solution Dynamic Power Report [mW]	52
68	Vivado Loop Unrolling Factor=4 Solution Dynamic Power Report [mW]	52
69	Vivado Loop Unrolling Factor=4 Solution Energy Single Operation Report [pJ]	52
70	HLS Loop Pipelining Solution Timing Summary (ns)	55
71	HLS Loop Pipelining Solution Latency Summary (clock cycles)	55
72	HLS Loop Pipelining Solution Latency Loops Summary	55
73	HLS Loop Pipelining Solution Utilization Estimates Summary	56
74	HLS Loop Pipelining Solution Register Detail Estimates Summary	56
75	HLS Loop Pipelining Solution C/RTL Cosimulation Summary	57
76	HLS Loop Pipelining Solution Export RTL Resource Usage	57
77	HLS Loop Pipelining Solution Export RTL Final Timing	57
78	Vivado Loop Pipelining Solution Utilization Report [#]	57
79	Vivado Loop Pipelining Solution Timing Report	57
80	Vivado Loop Pipelining Solution Dynamic Power Report [mW]	57
81	Vivado Loop Pipelining Solution Dynamic Power Report [mW]	57
82	Vivado Loop Pipelining Solution Energy Single Operation Report [pJ]	57

1 Introduction

1.1 LTI Systems

Un **sistema a tempo discreto (TD) SLS o LTI** è un sistema che risulta essere contemporaneamente lineare e stazionario. In particolare, un sistema si dice **lineare** se ad una combinazione lineare degli ingressi corrisponde la combinazione lineare delle uscite (*Principio di Sovrapposizione degli Effetti*):

$$\alpha x_1 + \beta x_2 \xrightarrow{\tau} \alpha y_1 + \beta y_2 \quad (1)$$

Invece, un sistema TD si dice **stazionario** (o tempo-invariante) se ritardando/anticipando l'ingresso di una quantità n_0 l'uscita viene ritardata/anticipata della stessa quantità:

$$x[n - n_0] \xrightarrow{\tau} y[n - n_0] \quad (2)$$

Un sistema SLS-TD è caratterizzato dalla sua risposta impulsiva $h[n]$. Essa rappresenta l'uscita del sistema quando all'ingresso viene applicato un impulso discreto.



Conoscendo la risposta impulsiva $h[n]$ e applicando le proprietà di linearità e stazionarietà, è possibile calcolare l'uscita $y[n]$ corrispondente ad un qualunque ingresso $x[n]$:

$$y[n] = \sum_{k=-\infty}^{+\infty} x[k]h[n-k] = x[n] * h[n] \quad (3)$$

Pertanto, l'uscita si ottiene effettuando la convoluzione tra l'ingresso e la risposta impulsiva.

1.2 Causal Systems

Un sistema (non necessariamente SLS) si dice causale se l'uscita all'istante n dipende solo dai valori precedenti dell'ingresso:

$$y[n] = \sum_{k=-\infty}^{+\infty} x[k]h[n-k] = \sum_{k=-\infty}^{+\infty} h[k]x[n-k] \quad (4)$$

Affinché un sistema SLS-TD sia causale è necessario e sufficiente che:

$$h[k] = 0 \quad \text{per ogni } k < 0 \quad (5)$$

1.3 FIR Filter

I principali filtri digitali sono i filtri FIR (Finite Impulse Response) e IIR (Infinite Impulse Response). In particolare, un filtro FIR è un sistema LTI causale con risposta finita all'impulso, la cui funzione di trasferimento risulta essere un polinomio in z^{-1} . In particolare, considerando $h[n]$, cioè la sequenza di coefficienti del filtro in questione, e la sequenza di ingresso $x[n]$, l'obiettivo è calcolare l'uscita $y[n]$ mediante convoluzione:

$$y[n] = \sum_{k=0}^{N-1} h[k]x[n-k] \quad (6)$$

1.4 High Level Synthesis

Il tool **Xilinx® Vivado® High-Level Synthesis (HLS)** permette di trasformare le specifiche C in un'implementazione RTL (Register Transfer Level) che può essere sintetizzata e implementata in un FPGA. Pertanto, mediante linguaggi di programmazione ad alto livello come C, C++ e SystemC, è possibile descrivere in maniera più semplice e veloce un'architettura hardware complessa rispetto ad un linguaggio descrittivo dell'hardware come il VHDL o il Verilog. Ovviamente, tale vantaggio si traduce in una minore precisione ed efficienza della descrizione hardware in questione rispetto a quella che si otterrebbe mediante linguaggi HDL. Pertanto, la sintesi ad alto livello dei sistemi digitali risulta essere una congiunzione tra il mondo hardware e quello software. In particolare, tale approccio comporta i seguenti **vantaggi**:

- **Maggiore produttività per i progettisti hardware**

I progettisti di hardware possono lavorare a un livello di astrazione superiore mentre creano hardware ad alte prestazioni.

- **Migliori prestazioni del sistema per i progettisti di software**

Gli sviluppatori di software possono accelerare le parti computazionalmente intensive dei loro algoritmi su un nuovo target di compilazione, cioè l'FPGA.

Inoltre, l'utilizzo di una metodologia di progettazione di sintesi di alto livello consente di:

- **Sviluppare algoritmi tramite linguaggi ad alto livello come C e C++**

Lavorare a un livello astratto dai dettagli di implementazione, che consumano tempo di sviluppo.

- **Verificare il codice mediante compilatori ad alto livello**

Convalidare la correttezza funzionale del progetto in modo più rapido rispetto ai linguaggi di descrizione hardware tradizionali come VHDL e Verilog.

- **Controllare il processo di sintesi C attraverso le direttive di ottimizzazione**

Creare implementazioni hardware specifiche ad alte prestazioni.

- **Creare più implementazioni dal codice sorgente C utilizzando le direttive di ottimizzazione**

Esplorare diversi design per lo stesso progetto permette di trovare implementazioni più efficienti e ottimali.

- **Sfruttare tutti i vantaggi di portabilità e ri-usabilità del codice ad alto livello**

Utilizzare tutte le metodologie ad alto livello basate sull'ingegneria del software e, quindi, risparmiare tempo e costi di progettazione software ulteriori.

Per quanto riguarda l'**HLS Design Flow**, questo si articola nei seguenti passi:

1. Compilare, eseguire (simulare) ed eseguire il debug del codice C o C++.
2. Sintetizzare il codice C in un'implementazione RTL, eventualmente utilizzando direttive di ottimizzazione (pragma).
3. Generare e analizzare i report corrispondenti
4. Packaging dell'implementazione RTL per effettuare l'esportazione dell'IP corrispondente.

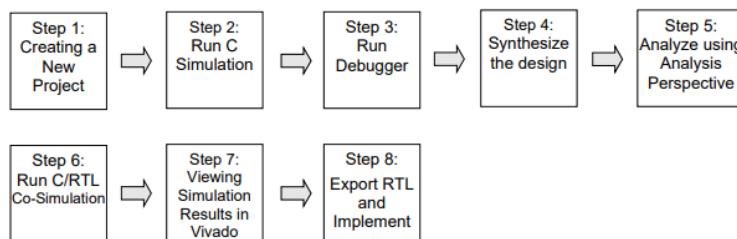


Figure 1: HLS Design Flow

2 Tasks to be performed

Prendendo come riferimento la struttura del filtro FIR, citato precedentemente, e considerando il tool di sintesi ad alto livello per sistemi digitali, fornito da Xilinx® Vivado®, cercare di ottimizzare secondo le caratteristiche principali il circuito in questione. In particolare, utilizzare diverse tecniche di progettazione ad alto livello e sfruttare, se possibile, le direttive di ottimizzazione (pragma) per ottenere valori ottimali in termini di latenza, dissipazione di potenza e utilizzazione delle risorse. Fondamentalmente, implementare le seguenti tecniche nella maniera opportuna:

- **Unoptimized Solution**

Rappresenta la soluzione iniziale (non ottimizzata) del filtro FIR in questione.

- **Operation Chaining Solution**

Rappresenta una soluzione basata sul cambiamento del periodo del clock.

- **Code Hoisting Solution**

Rappresenta una soluzione basata su un'alternativa gestione delle condizioni di controllo.

- **Loop Fission Solution**

Rappresenta una soluzione basata sulla scissione di cicli di operazioni complesse in più loop semplici.

- **Loop Unrolling Solution**

Rappresenta una soluzione basata sul parallelismo delle operazioni.

- **Loop Pipelining Solution**

Rappresenta una soluzione basata sulla scissione di operazioni combinatorie complesse in più procedure semplici.

- **Bitwidth Optimization Solution**

Rappresenta una soluzione basata sull'ottimizzazione del numero di bit utilizzati per ogni tipo di dato.

- **AXI Solution**

Rappresenta una soluzione basata sull'utilizzo dell'interfaccia AXI-Stream.

3 Definitions

Qui di seguito vengono riportate le MACRO e le intestazioni dei metodi corrispondenti alle soluzioni implementate per il filtro FIR in questione. In particolare, ogni definizione presenta la documentazione associata.

```
1 #include "ap_int.h"
2 #include <ap_axi_sdata.h>
3 #include <hls_stream.h>
4 #include <stdio.h>
5
6
7
8 #ifndef DEFINITIONS_H
9
10 /**
11  * Represents FIR filter size
12  * @param SIZE FIR filter size
13  */
14
15 #define SIZE 11
16
17 /**
18  * Represents the data type used to store the coefficients of the FIR filter
19  * @param coeffsType data type for FIR filter coefficients
20  */
21 typedef int coeffsType;
22
23 /**
24  * Represents the data type used to store the filter input samples
25  * @param samplesType data type for FIR filter input samples
26  */
27 typedef int samplesType;
28
29 /**
30  * Represents the data type used to store the intermediate accumulator when calculating the
31  * filter output.
32  * @param accType data type for FIR filter intermediate accumulator
33  */
34 typedef int accType;
35
36 /**
37  * Unoptimized Design.
38  * @param inputFilter
39  * @param outputFilter
40  */
41 void firConvolutionUnoptimized(samplesType inputFilter, samplesType* outputFilter);
42
43 /**
44  * Operation Chaining Design.
45  * @param inputFilter
46  * @param outputFilter
47  */
48 void firConvolutionOperationChaining(samplesType inputFilter, samplesType* outputFilter);
49
50 /**
51  * Code Hoisting Design.
52  * @param inputFilter
53  * @param outputFilter
54  */
55 void firConvolutionCodeHoisting(samplesType inputFilter, samplesType* outputFilter);
56
57 /**
58  * Loop Fission Design.
59  * @param inputFilter
60  * @param outputFilter
61  */
62 void firConvolutionLoopFission(samplesType inputFilter, samplesType* outputFilter);
```

```

63 /**
64 * Loop Unrolling Factor=2 Design.
65 * @param inputFilter
66 * @param outputFilter
67 */
68 void firConvolutionLoopUnrollingFactor2(samplesType inputFilter, samplesType* outputFilter);
69
70 /**
71 * Loop Unrolling Factor=2 with Unrolling Pragma Design.
72 * @param inputFilter
73 * @param outputFilter
74 */
75 void firConvolutionLoopUnrollingFactor2Pragma(samplesType inputFilter, samplesType*
    outputFilter);
76
77 /**
78 * Loop Unrolling Factor=2 with Unrolling Pragma and Partitioning Pragma Design.
79 * @param inputFilter
80 * @param outputFilter
81 */
82 void firConvolutionLoopUnrollingFactor2PP(samplesType inputFilter, samplesType* outputFilter
    );
83
84 /**
85 * Loop Unrolling Factor=4 Design.
86 * @param inputFilter
87 * @param outputFilter
88 */
89 void firConvolutionLoopUnrollingFactor4(samplesType inputFilter, samplesType* outputFilter);
90
91 /**
92 * Loop Unrolling Factor=4 with Unrolling Pragma Design.
93 * @param inputFilter
94 * @param outputFilter
95 */
96 void firConvolutionLoopUnrollingFactor4Pragma(samplesType inputFilter, samplesType*
    outputFilter);
97
98 /**
99 * Loop Unrolling Factor=4 with Unrolling Pragma and Partitioning Pragma Design.
100 * @param inputFilter
101 * @param outputFilter
102 */
103 void firConvolutionLoopUnrollingFactor4PP(samplesType inputFilter, samplesType* outputFilter
    );
104
105 /**
106 * Loop Pipelining Design.
107 * @param inputFilter
108 * @param outputFilter
109 */
110 void firConvolutionLoopPipelining(samplesType inputFilter, samplesType* outputFilter);
111
112 /**
113 * Bitwidth Optimization Design.
114 * @param inputFilter
115 * @param outputFilter
116 */
117 void firConvolutionBitwidthOptimization(ap_int<8> inputFilter, ap_int<18+SIZE>* outputFilter
    );
118
119 /**
120 * Represents the AXI type for AXI-stream interface.
121 * @param AXI_TYPE AXI type
122 */
123 typedef ap_axiu<32, 2, 5, 6> AXI_TYPE;
124
125 /**

```

```

126 * AXI Design.
127 * @param inputFilter
128 * @param outputFilter
129 */
130 void firConvolutionAXI(hls::stream<AXI_TYPE> &inputStreamFilter, hls::stream<AXI_TYPE> &
131   outputStreamFilter);
132
133 #endif

```

Bisogna, inoltre, specificare che ogni solution (tranne per alcune come verrà appositamente specificato) è stata progettata inizialmente in Xilinx® Vivado® High-Level Synthesis (HLS) e poi, successivamente, dopo aver esportato l'IP, è stata prevista sintesi e implementazione in Xilinx® Vivado®.

Oltre a ciò, si specifica che per ogni implementazione, effettuata in Xilinx® Vivado®, è stato previsto un constraint di clock associato e corrispondente alla solution progettata in Xilinx® Vivado® High-Level Synthesis (HLS) e, in aggiunta, ogni design implementato è stato simulato considerando input random così da ottenere una dati di potenza accurati mediante generazione di file .saif.

4 C Simulations

Qui di seguito viene riportato il file testbench per la *C Simulation* in HLS e il corrispondente output ottenuto. Bisogna precisare che ogni risultato ottenuto è stato verificato dal punto di vista numerico.

```
1 #include "definitions.h"
2 #include <stdio.h>
3 #include "ap_int.h"
4 #include <iostream>
5
6
7 void firConvolutionUnoptimized(samplesType inputFilter, samplesType* outputFilter);
8 void firConvolutionOperationChaining(samplesType inputFilter, samplesType* outputFilter);
9 void firConvolutionLoopFission(samplesType inputFilter, samplesType* outputFilter);
10 void firConvolutionCodeHoisting(samplesType inputFilter, samplesType* outputFilter);
11 void firConvolutionLoopUnrollingFactor2(samplesType inputFilter, samplesType* outputFilter);
12 void firConvolutionLoopUnrollingFactor2Pragma(samplesType inputFilter, samplesType*
13     outputFilter);
14 void firConvolutionLoopUnrollingFactor2PP(samplesType inputFilter, samplesType* outputFilter
15 );
16 void firConvolutionLoopUnrollingFactor4(samplesType inputFilter, samplesType* outputFilter);
17 void firConvolutionLoopUnrollingFactor4Pragma(samplesType inputFilter, samplesType*
18     outputFilter);
19 void firConvolutionLoopUnrollingFactor4PP(samplesType inputFilter, samplesType* outputFilter
20 );
21 void firConvolutionLoopPipelining(samplesType inputFilter, samplesType* outputFilter);
22 void firConvolutionBitwidthOptimization(ap_int<8> inputFilter, ap_int<18+SIZE>* outputFilter
23 );
24 void firConvolutionAXI(hls::stream<AXI_TYPE> &inputStreamFilter, hls::stream<AXI_TYPE> &
25     outputStreamFilter);
26
27 int main() {
28
29     samplesType inputFilter[SIZE] = {1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1};
30     samplesType outputFilter;
31
32     printf("\n*****firConvolutionUnoptimized*****\n");
33     printf("%-20s%-20s%-20s\n", "iteration", "inputFilter[i]", "outputFilter[i]");
34     for( int i=0; i<SIZE; ++i ) {
35         firConvolutionUnoptimized( inputFilter[i], &outputFilter );
36         printf("%-20d%-20d%-20d\n", i, inputFilter[i], outputFilter);
37     }
38
39     printf("\n*****firConvolutionOperationChaining*****\n");
40     for( int i=0; i<SIZE; ++i ) {
41         firConvolutionOperationChaining( inputFilter[i], &outputFilter );
42         printf("%-20d%-20d%-20d\n", i, inputFilter[i], outputFilter);
43     }
44
45     printf("\n*****firConvolutionCodeHoisting*****\n");
46     for( int i=0; i<SIZE; ++i ) {
47         firConvolutionCodeHoisting( inputFilter[i], &outputFilter );
48         printf("%-20d%-20d%-20d\n", i, inputFilter[i], outputFilter);
49     }
50
51     printf("\n*****firConvolutionLoopFission*****\n");
52     for( int i=0; i<SIZE; ++i ) {
53         firConvolutionLoopFission( inputFilter[i], &outputFilter );
54         printf("%-20d%-20d%-20d\n", i, inputFilter[i], outputFilter);
55     }
56
57     printf("\n*****firConvolutionLoopUnrollingFactor2*****\n");
58     for( int i=0; i<SIZE; ++i ) {
59         firConvolutionLoopUnrollingFactor2( inputFilter[i], &outputFilter );
60         printf("%-20d%-20d%-20d\n", i, inputFilter[i], outputFilter);
61     }
```

```

58     printf("\n*****firConvolutionLoopUnrollingFactor2Pragma*****\n");
59     for( int i=0; i<SIZE; ++i ) {
60         firConvolutionLoopUnrollingFactor2Pragma( inputFilter[i], &outputFilter );
61         printf("%-20d%-20d%-20d\n", i, inputFilter[i], outputFilter);
62     }
63
64     printf("\n*****firConvolutionLoopUnrollingFactor2PragmaPartitioning*****\n");
65     for( int i=0; i<SIZE; ++i ) {
66         firConvolutionLoopUnrollingFactor2PP( inputFilter[i], &outputFilter );
67         printf("%-20d%-20d%-20d\n", i, inputFilter[i], outputFilter);
68     }
69
70     printf("\n*****firConvolutionLoopUnrollingFactor4*****\n");
71     for( int i=0; i<SIZE; ++i ) {
72         firConvolutionLoopUnrollingFactor4( inputFilter[i], &outputFilter );
73         printf("%-20d%-20d%-20d\n", i, inputFilter[i], outputFilter);
74     }
75
76     printf("\n*****firConvolutionLoopUnrollingFactor4Pragma*****\n");
77     for( int i=0; i<SIZE; ++i ) {
78         firConvolutionLoopUnrollingFactor4Pragma( inputFilter[i], &outputFilter );
79         printf("%-20d%-20d%-20d\n", i, inputFilter[i], outputFilter);
80     }
81
82     printf("\n*****firConvolutionLoopUnrollingFactor4PragmaPartitioning*****\n");
83     for( int i=0; i<SIZE; ++i ) {
84         firConvolutionLoopUnrollingFactor4PP( inputFilter[i], &outputFilter );
85         printf("%-20d%-20d%-20d\n", i, inputFilter[i], outputFilter);
86     }
87
88     printf("\n*****firConvolutionLoopPipelining*****\n");
89     for( int i=0; i<SIZE; ++i ) {
90         firConvolutionLoopPipelining( inputFilter[i], &outputFilter );
91         printf("%-20d%-20d%-20d\n", i, inputFilter[i], outputFilter);
92     }
93
94     ap_int<8> inputFilterAp[SIZE] = {42, 7, 58, 120, -64, 15, -89, -75, 43, 43, -14};
95     ap_int<18+SIZE> outputFilterAp;
96     std::cout << "\n*****firConvolutionBitwidthOptimization*****\n";
97     for( int i=0; i<SIZE; ++i ) {
98         firConvolutionBitwidthOptimization( inputFilterAp[i], &outputFilterAp );
99         std::cout << std::left << std::setw(20) << i
100            << std::left << std::setw(20) << inputFilterAp[i]
101            << std::left << std::setw(20) << outputFilterAp << "\n";
102    }
103
104    hls::stream<AXI_TYPE> inputStreamFilter;
105    hls::stream<AXI_TYPE> outputStreamFilter;
106    AXI_TYPE inFilter, outFilter;
107    for (int i=0; i<SIZE; i++) {
108        inFilter.data=inputFilter[i];
109        inFilter.strb = -1;
110        inFilter.keep = 15;
111        inFilter.user = 0;
112        inFilter.last = (i==SIZE-1) ? 1 : 0;
113        inFilter.id = 0;
114        inFilter.dest = 0;
115        inputStreamFilter << inFilter;
116    }
117    printf("\n*****firConvolutionAXI*****\n");
118    firConvolutionAXI( inputStreamFilter, outputStreamFilter );
119    for( int i=0; i<SIZE; ++i ) {
120        outFilter = outputStreamFilter.read();
121        printf("%-20d%-20d%-20d\n", i, inFilter.data.to_int(), outFilter.data.to_int());
122    }
123
124    return 0;
125
```

126 }

```
1 INFO: [SIM 2] **** CSIM start ****
2 INFO: [SIM 4] CSIM will launch GCC as the compiler.
3 Compiling ../../fir_tb.cpp in debug mode
4 Compiling ../../fir_axi.cpp in debug mode
5 Compiling ../../fir_bitwidth_optimization.cpp in debug mode
6 Compiling ../../fir_code_hoisting.cpp in debug mode
7 Compiling ../../fir_loop_fission.cpp in debug mode
8 Compiling ../../fir_loop_pipelineing.cpp in debug mode
9 Compiling ../../fir_loop_unrolling_factor2.cpp in debug mode
10 Compiling ../../fir_loop_unrolling_factor2_pp.cpp in debug mode
11 Compiling ../../fir_loop_unrolling_factor2_pragma.cpp in debug mode
12 Compiling ../../fir_loop_unrolling_factor4.cpp in debug mode
13 Compiling ../../fir_loop_unrolling_factor4_pp.cpp in debug mode
14 Compiling ../../fir_loop_unrolling_factor4_pragma.cpp in debug mode
15 Compiling ../../fir_operation_chaining.cpp in debug mode
16 Compiling ../../fir_unoptimized.cpp in debug mode
17 Generating csim.exe
18
19 *****firConvolutionUnoptimized*****
20 iteration      inputFilter[i]      outputFilter[i]
21 0             1                  53
22 1             1                  53
23 2             1                  -38
24 3             1                  -38
25 4             1                  275
26 5             1                  775
27 6             1                  1088
28 7             1                  1088
29 8             1                  997
30 9             1                  997
31 10            1                  1050
32
33 *****firConvolutionOperationChaining*****
34 0             1                  53
35 1             1                  53
36 2             1                  -38
37 3             1                  -38
38 4             1                  275
39 5             1                  775
40 6             1                  1088
41 7             1                  1088
42 8             1                  997
43 9             1                  997
44 10            1                  1050
45
46 *****firConvolutionCodeHoisting*****
47 0             1                  53
48 1             1                  53
49 2             1                  -38
50 3             1                  -38
51 4             1                  275
52 5             1                  775
53 6             1                  1088
54 7             1                  1088
55 8             1                  997
56 9             1                  997
57 10            1                  1050
58
59 *****firConvolutionLoopFission*****
60 0             1                  53
61 1             1                  53
62 2             1                  -38
63 3             1                  -38
64 4             1                  275
65 5             1                  775
66 6             1                  1088
```

```

67 7           1           1088
68 8           1           997
69 9           1           997
70 10          1           1050
71
72 *****firConvolutionLoopUnrollingFactor2*****
73 0           1           53
74 1           1           53
75 2           1           -38
76 3           1           -38
77 4           1           275
78 5           1           775
79 6           1           1088
80 7           1           1088
81 8           1           997
82 9           1           997
83 10          1           1050
84
85 *****firConvolutionLoopUnrollingFactor2Pragma*****
86 0           1           53
87 1           1           53
88 2           1           -38
89 3           1           -38
90 4           1           275
91 5           1           775
92 6           1           1088
93 7           1           1088
94 8           1           997
95 9           1           997
96 10          1           1050
97
98 *****firConvolutionLoopUnrollingFactor2PragmaPartitioning*****
99 0           1           53
100 1          1           53
101 2          1           -38
102 3          1           -38
103 4          1           275
104 5          1           775
105 6          1           1088
106 7          1           1088
107 8          1           997
108 9          1           997
109 10         1           1050
110
111 *****firConvolutionLoopUnrollingFactor4*****
112 0           1           53
113 1           1           53
114 2           1           -38
115 3           1           -38
116 4           1           275
117 5           1           775
118 6           1           1088
119 7           1           1088
120 8           1           997
121 9           1           997
122 10          1           1050
123
124 *****firConvolutionLoopUnrollingFactor4Pragma*****
125 0           1           53
126 1           1           53
127 2           1           -38
128 3           1           -38
129 4           1           275
130 5           1           775
131 6           1           1088
132 7           1           1088
133 8           1           997
134 9           1           997

```

```

135 10          1          1050
136
137 *****firConvolutionLoopUnrollingFactor4PragmaPartitioning*****
138 0            1          53
139 1            1          53
140 2            1          -38
141 3            1          -38
142 4            1          275
143 5            1          775
144 6            1          1088
145 7            1          1088
146 8            1          997
147 9            1          997
148 10           1          1050
149
150 *****firConvolutionLoopPipelining*****
151 0            1          53
152 1            1          53
153 2            1          -38
154 3            1          -38
155 4            1          275
156 5            1          775
157 6            1          1088
158 7            1          1088
159 8            1          997
160 9            1          997
161 10           1          1050
162
163 *****firConvolutionBitwidthOptimization*****
164 0            42         2226
165 1            7          371
166 2            58         -748
167 3            120        5723
168 4            -64        4476
169 5            15          13066
170 6            -89        35907
171 7            -75        63411
172 8            43          64678
173 9            43          18722
174 10           -14        -48096
175
176 *****firConvolutionAXI*****
177 0            1          53
178 1            1          53
179 2            1          -38
180 3            1          -38
181 4            1          275
182 5            1          775
183 6            1          1088
184 7            1          1088
185 8            1          997
186 9            1          997
187 10           1          1050
188 INFO: [SIM 1] CSim done with 0 errors.
189 INFO: [SIM 3] ***** CSIM finish *****

```

5 Solutions

Di seguito verranno illustrate e analizzate diverse implementazioni del filtro FIR. In particolare, verranno mostrati i report generati da HLS e Vivado così da effettuare ulteriori considerazioni a riguardo. Bisogna precisare che la prima soluzione presentata, cioè quella iniziale non ottimizzata, verrà discussa approfonditamente così da affrontare e spiegare i parametri di confronto che verranno utilizzati nelle successive sezioni.

5.1 Unoptimized Solution

Qui di seguito viene riportato il codice relativo alla soluzione iniziale (non ottimizzata) del filtro FIR in questione.

```

1 #include "definitions.h"
2
3 void firConvolutionUnoptimized(samplesType inputFilter, samplesType* outputFilter) {
4     coeffsType coefficientsFilter[SIZE] = {53, 0, -91, 0, 313, 500, 313, 0, -91, 0, 53};
5     static samplesType shiftRegister[SIZE];
6     accType accumulator;
7     accumulator = 0;
8     int i;
9     loop: for( i=SIZE-1; i>=0; --i ) {
10         if( i==0 ) {
11             accumulator += inputFilter * coefficientsFilter[0];
12             shiftRegister[0] = inputFilter;
13         } else {
14             shiftRegister[i] = shiftRegister[i-1];
15             accumulator += shiftRegister[i] * coefficientsFilter[i];
16         }
17     }
18     *outputFilter = accumulator;
19 }
```

In particolare, si può notare come l'input, l'output e i coefficienti utilizzati risultano essere definiti rispettivamente tramite i tipi *samplesType*, *samplesType* e *coeffsType*, cioè tutti e tre definiti come tipi di dato *int* come descritto in *definitions.h*. Bisogna specificare che l'operazione di convoluzione è stata ottenuta mediante il loop qui sopra definito. Tale operazione è data dallo shifting dei valori di input al filtro e, successivamente, dall'accumulo. Quest'ultimo è ottenuto mediante la somma, tra la corrente calcolata e la precedente accumulata, e il prodotto, tra il valore shiftato nella precedente iterazione e il coefficiente *i-esimo* corrispondente.

Infine, il valore accumulato nella variabile *accumulator* viene assegnato all'output atteso dalla procedura. Effettuando la sintesi è possibile evidenziare il seguente report:

Clock	Target	Estimated	Uncertainty
ap_clk	10.00	8.510	1.25

Table 1: HLS Unoptimized Solution Timing Summary (ns)

Latency		Interval	
min	max	min	max
23	45	23	45

Table 2: HLS Unoptimized Solution Latency Summary (clock cycles)

Si può notare come il periodo di clock previsto è pari a 10 ns, come impostato durante la creazione della solution in questione, mentre quello stimato è minore. Questo è dovuto al fatto che il tool stima un'incertezza (assimilabile ad un margine) che permette di calcolare il periodo di clock effettivo utilizzato dal processo di sintesi.

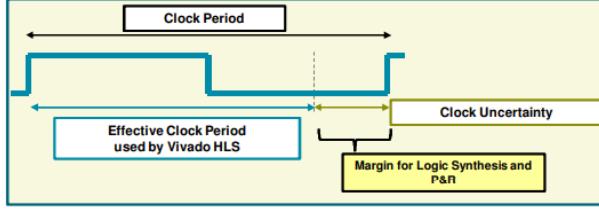


Figure 2: HLS Clock Period from Synthesis Report

In particolare, il periodo di clock effettivo, calcolato dal tool, non deve essere considerato come parametro di progettazione ma come un margine di sicurezza previsto da HLS per il successivo place and route. Nel caso in cui si volesse operare sul periodo di clock della solution come design parameter, si dovrebbe impostare un periodo di clock differente durante la creazione della solution e poi far calcolare automaticamente l'incertezza ed il periodo del clock effettivo dallo stesso tool.

Loop Name	Latency		Iteration Latency		Initiation Interval		Trip Count
	min	max	min	max	achieved	target	
- loop	22	44	2	4	-	-	11

Table 3: HLS Unoptimized Solution Latency Loops Summary

Considerando, invece, il loop descritto nel file sorgente, si può notare come esso presenti 11 iterazioni (corrispondenti alla capienza di processing del filtro corrispondente) e una iteration latency (IL), cioè latenza per ogni iterazione, compresa tra 2 e 4 cicli di clock. Si deduce che la latenza totale del loop in questione sia compresa tra 22 e 44 cicli di clock.

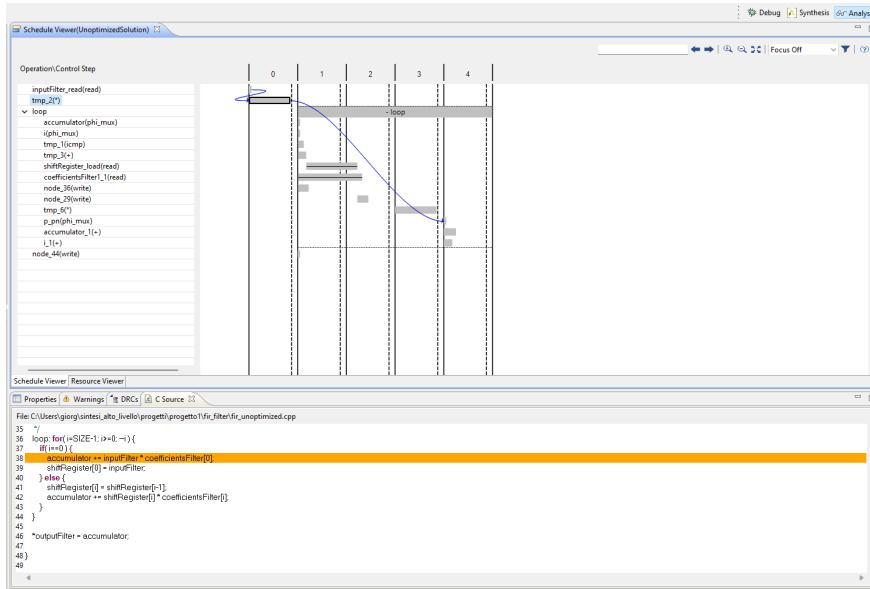


Figure 3: HLS Unoptimized Solution If Statement Analysis

In particolare, si può notare come il ciclo di clock in più, previsto nella latenza totale dell'architettura, sia dovuto alla verifica, al calcolo e all'assegnazione in corrispondenza dell'iterazione 0 all'interno del loop. È come se l'if all'interno del ciclo for venga trattato al di fuori dello stesso loop. Questo può essere notato facendo riferimento alle operazioni svolte all'interno del loop.

Inoltre, si può notare come le operazioni di lettura previste all'interno del loop siano svolte in parallelo (come mostrato nel primo gruppo di figure allegato), mentre quelle di scrittura scheduled in differente

maniera (come mostrato nel secondo gruppo di figure allegato).

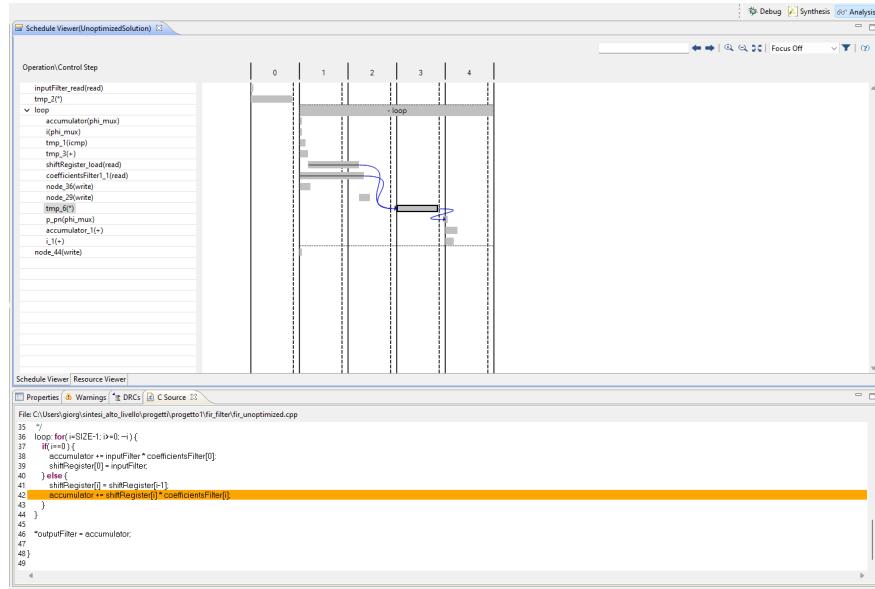


Figure 4: HLS Unoptimized Solution Else Statement Analysis

Infatti, si può notare come la latenza associata alle operazioni previste nelle rimanenti iterazioni del loop venga correlata alla latency associata al ciclo stesso.

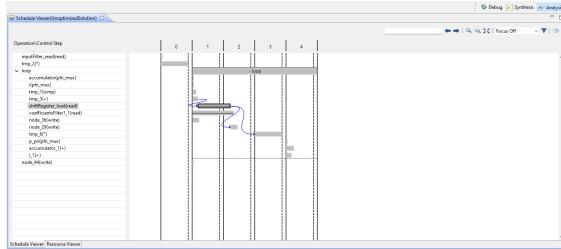


Figure 5: HLS Unoptimized Solution Read Operations Analysis

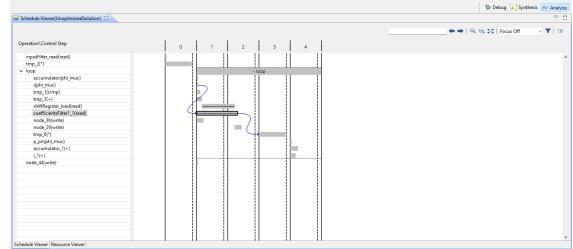


Figure 6: HLS Unoptimized Solution Read Operations Analysis

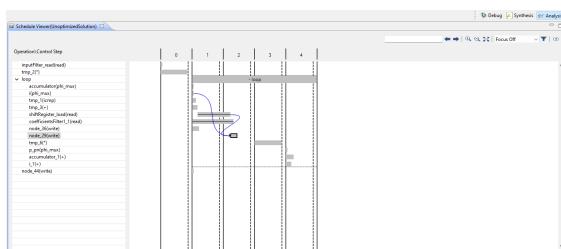


Figure 7: HLS Unoptimized Solution Write Operations Analysis

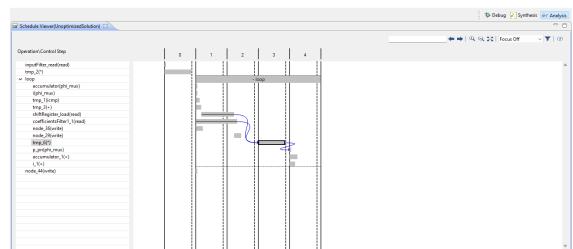


Figure 8: HLS Unoptimized Solution Write Operations Analysis

Qui di seguito, viene allegato l'utilizzazione delle risorse stimata dal processo di sintesi.

Name	BRAM_18K	DSP48E	FF	LUT
DSP	-	-	-	-
Expression	-	4	0	105
FIFO	-	-	-	-
Instance	-	-	-	-
Memory	0	-	74	8
Multiplexer	-	-	-	120
Register	-	-	213	-
Total	0	4	287	233
Available	280	220	106400	53200
Utilization (%)	0	1	~0	~0

Table 4: HLS Unoptimized Solution Utilization Estimates Summary

Successivamente effettuando la C/RTL Cosimulation è possibile evidenziare il seguente report:

RTL	Status	Latency			Interval		
		min	avg	max	min	avg	max
VHDL	Pass	43	43	44	43	43	44

Table 5: HLS Unoptimized Solution C/RTL Cosimulation Summary

Questa tabella mostra quanti cicli di clock sono necessari affinché venga fornito l'ingresso e l'uscita successivi.

Dopodiché si può procedere con l'Export RTL.

Resource	VHDL
SLICE	81
LUT	275
FF	160
DSP	2
BRAM	0
SRL	0

Timing	VHDL
CP required	10.000
CP achieved post-synthesis	5.745
CP achieved post-implementation	6.410

Table 7: HLS Unoptimized Solution Export RTL Final Timing

Table 6: HLS Unoptimized Solution Export RTL Resource Usage

Si può notare come il numero di Flip Flop e di DSP sia diminuito rispetto a quello stimato nel report di sintesi. Questo è dovuto al fatto che, l'Export RTL non è una vera e propria implementazione ma si tratta di un processo durante il quale il tool effettua dei miglioramenti.

Successivamente a tale fase si procede con l'importazione dell'IP in Vivado.

```

1 library IEEE;
2 use IEEE.STD_LOGIC_1164.ALL;
3
4 entity firConvolutionUnoptimized_top is
5   port (
6     ap_clk : IN STD_LOGIC;
7     ap_rst : IN STD_LOGIC;
8     ap_start : IN STD_LOGIC;
9     ap_done : OUT STD_LOGIC;
10    ap_idle : OUT STD_LOGIC;
11    ap_ready : OUT STD_LOGIC;
12    inputFilter : IN STD_LOGIC_VECTOR (31 downto 0);
13    outputFilter : OUT STD_LOGIC_VECTOR (31 downto 0);
14    outputFilter_ap_vld : OUT STD_LOGIC

```

```

15     );
16 end firConvolutionUnoptimized_top;
17
18 architecture Behavioral of firConvolutionUnoptimized_top is
19   component firConvolutionUnoptimized_0 is
20     port (
21       ap_clk : IN STD_LOGIC;
22       ap_rst : IN STD_LOGIC;
23       ap_start : IN STD_LOGIC;
24       ap_done : OUT STD_LOGIC;
25       ap_idle : OUT STD_LOGIC;
26       ap_ready : OUT STD_LOGIC;
27       inputFilter : IN STD_LOGIC_VECTOR (31 downto 0);
28       outputFilter : OUT STD_LOGIC_VECTOR (31 downto 0);
29       outputFilter_ap_vld : OUT STD_LOGIC
30     );
31   end component;
32
33 begin
34   firConvolutionUnoptimized_IP : firConvolutionUnoptimized_0 port map(
35     ap_clk => ap_clk,
36     ap_rst => ap_rst,
37     ap_start => ap_start,
38     ap_done => ap_done,
39     ap_idle => ap_idle,
40     ap_ready => ap_ready,
41     inputFilter => inputFilter,
42     outputFilter => outputFilter,
43     outputFilter_ap_vld => outputFilter_ap_vld
44   );
45 end Behavioral;
46
1 library IEEE;
2 use IEEE.STD_LOGIC_1164.ALL;
3
4 entity firConvolutionUnoptimized_tb is
5 end firConvolutionUnoptimized_tb;
6
7 architecture Behavioral of firConvolutionUnoptimized_tb is
8   component firConvolutionUnoptimized_top is
9     port (
10       ap_clk : IN STD_LOGIC;
11       ap_rst : IN STD_LOGIC;
12       ap_start : IN STD_LOGIC;
13       ap_done : OUT STD_LOGIC;
14       ap_idle : OUT STD_LOGIC;
15       ap_ready : OUT STD_LOGIC;
16       inputFilter : IN STD_LOGIC_VECTOR (31 downto 0);
17       outputFilter : OUT STD_LOGIC_VECTOR (31 downto 0);
18       outputFilter_ap_vld : OUT STD_LOGIC
19     );
20   end component;
21
22   constant clk_period : Time := 10 ns;
23   signal ap_clk : STD_LOGIC := '0';
24   signal ap_rst : STD_LOGIC;
25   signal ap_start : STD_LOGIC;
26   signal ap_done : STD_LOGIC;
27   signal ap_idle : STD_LOGIC;
28   signal ap_ready : STD_LOGIC;
29   signal inputFilter : STD_LOGIC_VECTOR (31 DOWNTO 0);
30   signal outputFilter : STD_LOGIC_VECTOR (31 DOWNTO 0);
31   signal outputFilter_ap_vld : STD_LOGIC;
32   signal cycles: integer := 44;
33
34 begin
35   uut : firConvolutionUnoptimized_top port map(
36     ap_clk,

```

```

37         ap_rst,
38         ap_start,
39         ap_done,
40         ap_idle,
41         ap_ready,
42         inputFilter,
43         outputFilter,
44         outputFilter_ap_vld
45     );
46
47     clk_process : process
48     begin
49         wait for clk_period/2;
50         ap_clk <= not ap_clk;
51     end process clk_process;
52
53     uut_process : process
54     begin
55         ap_rst <= '1';
56         ap_start <= '0';
57         wait for 250ns;
58         ap_rst <= '0';
59         wait for clk_period;
60         ap_start <= '1';
61             inputFilter <= "0110110001111000000011010001110";
62             wait for cycles*clk_period;
63             inputFilter <= "10101100001011101011101101100010";
64             wait for cycles*clk_period;
65             inputFilter <= "1111000001110001001011000111100";
66             wait for cycles*clk_period;
67             inputFilter <= "01101101000111100000010100000100";
68             wait for cycles*clk_period;
69             inputFilter <= "1100010001001111001010011011001";
70             wait for cycles*clk_period;
71             inputFilter <= "011010011001100101111010100100100";
72             wait for cycles*clk_period;
73             inputFilter <= "1011010111011110100000101110011";
74             wait for cycles*clk_period;
75             inputFilter <= "11110011001101000000001100001010";
76             wait for cycles*clk_period;
77             inputFilter <= "01010001001110100111010101001001";
78             wait for cycles*clk_period;
79             inputFilter <= "10100010011110110010000010011000";
80             wait for cycles*clk_period;
81             inputFilter <= "11101000110011101001011101010001";
82             wait for cycles*clk_period;
83             wait;
84     end process uut_process;
85 end Behavioral;

```

```
1 create_clock -period 10.000 -name myclk -waveform {0.000 5.000} [get_ports ap_clk]
```

In particolare, impostando un constraint di clock pari a 10 ns (come impostato anche in HLS) e considerando un numero di ciclo di clock pari a quelli ottenuti nel report di C/RTL Cosimulation, è possibile eseguire il processo di sintesi e implementazione.

LUT	LUTRAM	FF	BRAM	DSP	IO	BUFG
275	32	160	0	2	71	1

Table 8: Vivado Unoptimized Solution Utilization Report [#]

Cycles [#]	Clock Constraint [ns]	WNS [ns]	Maximum Clock Frequency [Mhz]
44	10	3.654	157.5795777

Table 9: Vivado Unoptimized Solution Timing Report

Pertanto, dopo aver effettuato la Post-Implementation Timing Simulation e aver inserito il file *.saif* all'interno dell'Implementation Power Report, è possibile analizzare con dettaglio la potenza dinamica e l'energia per singola operazione associata.

BRAM	Clock Enable	Clocks	DSP	Logic	Set/Reset [mW]	Data
0	0.454227469	1.215484925	0.335011806	0.921480532	3.57E-03	1.007059589

Table 10: Vivado Unoptimized Solution Dynamic Power Report [mW]

Dynamic Total	Energy Single Operation
3.936829417	39.36829417

Table 11: Vivado Unoptimized Solution Dynamic Power Report [mW]

Table 12: Vivado Unoptimized Solution Energy Single Operation Report [pJ]

Quindi, dopo aver generato i report associati alla soluzione iniziale non ottimizzata, si può procedere con le successive soluzioni che cercheranno di migliorare diversi aspetti dell'architettura appena proposta. In particolare, nelle successive sezioni verranno commentati i risultati in termini di potenza, timing e risorse facendo riferimento alla soluzione iniziale o a quella precedente su cui sarà basata la corrente.

5.2 Operation Chaining Solution

Considerando, ora, la stessa architettura software precedentemente esposta, si potrebbe fare un ragionamento riguardo il periodo di clock considerato. In particolare, si potrebbe pensare di utilizzare un periodo minore per analizzare come il tool gestisce la situazione associata. L'idea è quella di diminuire tale periodo senza andare oltre un margine. Questo vuol dire considerare lo slack associato alla soluzione iniziale e ridurlo fin quando è possibile, cioè fin quando i constraint di timing non sono violati. Ciò che dovremmo aspettarci è che il tool riesca a gestire la stessa architettura con un periodo di clock minore, rispetto a quello della soluzione iniziale non ottimizzata, garantendo in alcuni casi delle ottimizzazioni effettuate dallo stesso tool. L'ideale sarebbe ottenere una soluzione che mi permette di avere un giusto trade-off tra latenza (intesa come cicli di clock per ottenere un risultato), potenza e risorse utilizzate.

Qui di seguito verranno presentati i report associati alla sintesi, C/RTL Cosimulation ed Export RTL considerando un range di periodo di clock pari a [10ns, 3ns]. In particolare, le soluzioni analizzate non sono state considerate con un periodo del clock minore a 3ns dato che, come verrà mostrato e analizzato successivamente, in corrispondenza di un periodo del clock pari a 3ns, la soluzione hardware implementata all'interno di Vivado presenta un WNS negativo comportando, di conseguenza, la violazione dei constraint di timing.

Solution	Clock	Target	Estimated	Uncertainty
clk=10ns	ap_clk	10.00	8.510	1.25
clk=9ns	ap_clk	9.00	6.912	1.12
clk=8ns	ap_clk	8.00	6.912	1
clk=7ns	ap_clk	7.00	5.745	0.88
clk=6ns	ap_clk	6.00	4.644	0.75
clk=5ns	ap_clk	5.00	4.321	0.62
clk=4ns	ap_clk	4.00	3.254	0.5
clk=3ns	ap_clk	3.00	2.552	0.38

Table 13: HLS Operation Chaining Solution Timing Summary (ns)

Solution	Latency		Interval	
	min	max	min	max
clk=10ns	23	45	23	45
clk=9ns	24	57	24	57
clk=8ns	24	57	24	57
clk=7ns	25	69	25	69
clk=6ns	27	93	27	93
clk=5ns	27	93	27	93
clk=4ns	40	139	40	139
clk=3ns	40	139	40	139

Table 14: HLS Operation Chaining Solution Latency Summary (clock cycles)

Solution	Loop Name	Latency		Iteration Latency		Initiation Interval achieved	Trip Count
		min	max	min	max		
clk=10ns	- loop	22	44	2	4	-	11
clk=9ns	- loop	22	55	2	5	-	11
clk=8ns	- loop	22	55	2	5	-	11
clk=7ns	- loop	22	66	2	6	-	11
clk=6ns	- loop	22	88	2	8	-	11
clk=5ns	- loop	22	88	2	8	-	11
clk=4ns	- loop	33	132	3	12	-	11
clk=3ns	- loop	33	132	3	12	-	11

Table 15: HLS Operation Chaining Solution Latency Loops Summary

Si può notare come al diminuire del periodo di clock, la latenza totale e la latenza per ogni iterazione aumentano rispettivamente. Questo è dovuto al fatto che, considerando un periodo di clock sempre minore, il tool cerca di garantire il minor tempo di latenza affinché le operazioni presenti all'interno del loop vengano eseguite correttamente. Pertanto, se il periodo del clock diminuisce, questo si traduce in un aumento dei cicli di clock richiesti per ogni iterazione. Bisogna notare che, in corrispondenza di alcune soluzioni, alcuni parametri assumono valori uguali. In particolare, tale situazione si verifica in corrispondenza delle soluzioni aventi periodo del clock pari a 9ns e 8ns e anche in corrispondenza di 6ns e 5ns e, infine, 4ns e 3ns. Questa peculiarità potrebbe essere sfruttata successivamente, considerando ulteriori design parameters di confronto,

per la scelta di un trade-off opportuno. Dal momento che il tool cerca di garantire il soddisfacimento di tali operazioni complesse con un minore periodo di clock, quello che ci si aspetta è una utilizzazione delle risorse stimate e, in particolare, un aumento di quest'ultime al diminuire del periodo di clock.

Solution	BRAM	DSP	FF	LUT
clk=10ns	0	4	287	233
clk=9ns	0	4	619	301
clk=8ns	0	4	619	301
clk=7ns	0	4	623	305
clk=6ns	0	4	725	221
clk=5ns	0	4	725	221
clk=4ns	0	4	832	286
clk=3ns	0	4	832	286

Table 16: HLS Operation Chaining Solution Utilization Estimates [#]

Si evidenzia quanto citato precedentemente e, inoltre, si può notare come, anche in questo caso, in corrispondenza di alcune soluzioni hardware, l'utilizzazione delle risorse risulta essere la medesima. In particolare, è possibile già effettuare una riflessione riguardo i risultati ottenuti dopo il processo di sintesi. Si può notare come, in corrispondenza delle soluzioni hardware con clk=9ns e clk=8ns si ha lo stessa latency, iteration latency e utilizzazione delle risorse. Ma soprattutto, facendo riferimento alla soluzione con periodo di clock pari a 7ns, la latency e l'iteration latency aumentano di un fattore pari a 11 (rispetto alle altre soluzioni dove i valori sono di gran lunga maggiori) e l'utilizzazione di FF e LUT (dal momento che il numero di BRAM e DSP utilizzati è il medesimo) risulta essere aumentata rispettivamente del 0.6% e del 1.32%. Pertanto, a fronte del processo di sintesi, la soluzione hardware avente periodo di clock pari a 7ns si potrebbe già considerare come candidata al trade-off sopra citato.

Solution	RTL	Status	Latency			Interval		
			min	avg	max	min	avg	max
clk=10ns	VHDL	Pass	43	43	44	43	43	44
clk=9ns	VHDL	Pass	54	54	55	54	54	55
clk=8ns	VHDL	Pass	54	54	55	54	54	55
clk=7ns	VHDL	Pass	65	65	66	65	65	66
clk=6ns	VHDL	Pass	87	87	88	87	87	88
clk=5ns	VHDL	Pass	87	87	88	87	87	88
clk=4ns	VHDL	Pass	130	130	131	130	130	131
clk=4ns	VHDL	Pass	130	130	131	130	130	131
clk=3ns	VHDL	Pass	130	130	131	130	130	131

Table 17: HLS Operation Chaining Solution C/RTL Cosimulation Report

Solution	SLICE	LUT	FF	DSP	BRAM	CP required	CP achieved post-synthesis	CP achieved post-implementation
clk=10ns	81	275	160	2	0	10	5.745	6.41
clk=9ns	88	276	226	2	0	9	5.745	6.059
clk=8ns	89	276	226	2	0	8	5.745	6.034
clk=7ns	75	186	222	2	0	7	5.745	5.692
clk=6ns	76	186	241	2	0	6	4.823	4.364
clk=5ns	96	186	258	2	0	5	4.823	4.282
clk=4ns	114	188	376	2	0	4	3.667	3.724
clk=3ns	107	192	377	4	0	3	3.667	3.627

Table 18: HLS Operation Chaining Solution Export RTL Report

Si può notare come la soluzione ottimale precedentemente proposta risulta essere la migliore tra quelle citate. Infatti, è possibile evidenziare che rispetto alle altre soluzioni presenta il minor numero di risorse stimate dopo aver effettuato l'Export RTL in HLS.

Pertanto, si può affermare che il tool, in corrispondenza di un periodo di clock pari a 7ns, è riuscito ad effettuare delle ottimizzazioni che hanno permesso di ottenere un minimo aumento della latenza e una diminuzione delle risorse utilizzate.

Quindi, considerando constraint di clock differenti per ogni IP importato in Vivado, si procede con la sintesi, l'implementazione, la Post-Implementation Timing Simulation e la generazione del file .saif per ogni soluzione corrispondente.

```

1 create_clock -period 10.000 -name myclk -waveform {0.000 5.000} [get_ports ap_clk]
1 create_clock -period 9.000 -name myclk -waveform {0.000 4.500} [get_ports ap_clk]
1 create_clock -period 8.000 -name myclk -waveform {0.000 4.000} [get_ports ap_clk]
1 create_clock -period 7.000 -name myclk -waveform {0.000 3.500} [get_ports ap_clk]
1 create_clock -period 6.000 -name myclk -waveform {0.000 3.000} [get_ports ap_clk]
1 create_clock -period 5.000 -name myclk -waveform {0.000 2.500} [get_ports ap_clk]
1 create_clock -period 4.000 -name myclk -waveform {0.000 2.000} [get_ports ap_clk]
1 create_clock -period 3.000 -name myclk -waveform {0.000 1.500} [get_ports ap_clk]
```

In particolare, qui di seguito verranno mostrati e analizzati soltanto i report corrispondenti alle soluzioni avente periodo di clock compreso nell'intervallo [10ns, 4ns]. Infatti, è possibile notare come, in corrispondenza del periodo di clock pari a 3ns, i constraint di timing siano violati.

 [Timing 38-282] The design failed to meet the timing requirements. Please see the timing summary report for details on the timing violations.

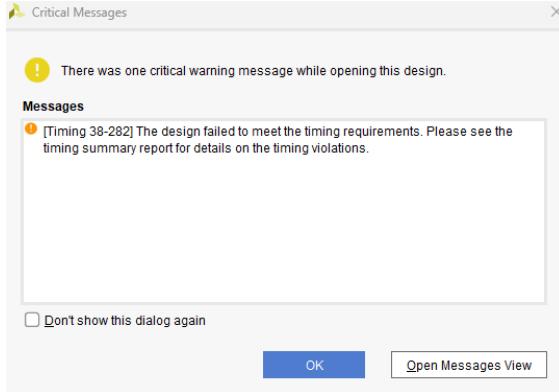


Figure 9: Vivado Operation Chaining Solution
clk=3ns Timing Constraint Violated



Figure 10: Vivado Operation Chaining Solution
clk=3ns Timing Constraint Violated

Pertanto, si procede con l'analisi dei report delle soluzioni hardware implementate in Vivado aventi periodo di clock coompresso nell'intervallo [10ns, 4ns].

Solution	LUT	LUTRAM	FF	BRAM	DSP	IO	BUFG
clk=10ns	275	32	160	0	2	71	1
clk=9ns	276	32	226	0	2	71	1
clk=8ns	276	32	226	0	2	71	1
clk=7ns	186	32	222	0	4	71	1
clk=6ns	186	32	241	0	4	71	1
clk=5ns	186	32	258	0	4	71	1
clk=4ns	188	32	376	0	4	71	1

Table 19: Vivado Operation Chaining Solution Utilization Report [#]

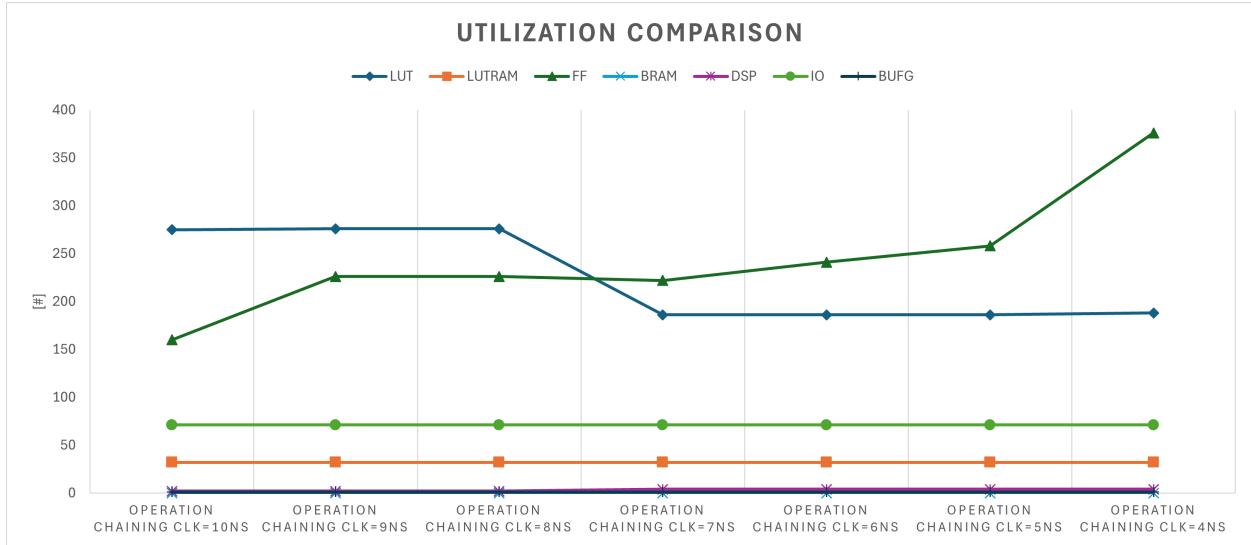


Figure 11: Vivado Operation Chaining Utilization Plot

Solution	Cycles [#]	Clock Constraint [ns]	WNS [ns]	Maximum Clock Frequency [Mhz]
clk=10ns	44	10	3.654	157.5795777
clk=9ns	55	9	3.06	168.3501684
clk=8ns	55	8	2.277	174.7335314
clk=7ns	66	7	1.33	176.366843
clk=6ns	88	6	1.573	225.8866049
clk=5ns	88	5	0.374	216.1694769
clk=4ns	131	4	0.454	282.0078962

Table 20: Vivado Operation Chaining Solution Timing Report

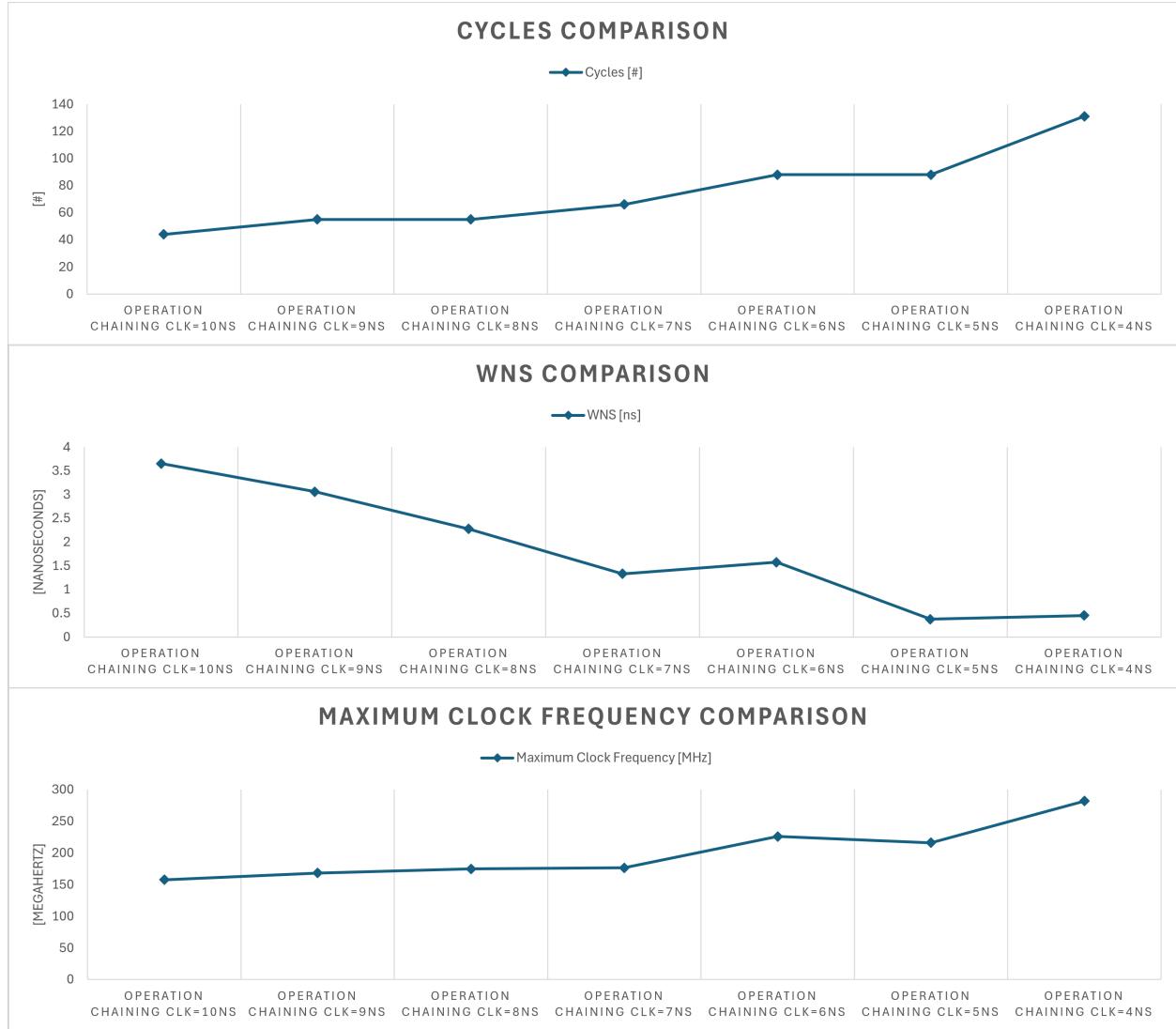


Figure 12: Vivado Operation Chaining Timing Plot

Si può notare come in corrispondenza dell'architettura hardware a cui corrisponde un periodo di clock pari a 7ns si ha l'utilizzazione delle risorse minore a meno del numero di LUT e FF che risulta essere aumentato rispettivamente di un'unità e di 66 unità. Questo aumento di risorse può essere giustificato dal momento

che, rispetto alla soluzione iniziale non ottimizzata a 10ns, il periodo del clock risulta essere diminuito e, pertanto, il tool prevede una gestione delle istruzioni in maniera più complessa comportando un'aumento delle risorse stesse. Tanto è vero che al diminuire del periodo di clock, si ha un aumento delle risorse utilizzate e dei cicli di clock considerati per ottenere un nuovo risultato. In particolare, per quanto riguarda la massima frequenza di clock associata ad ogni soluzione proposta in questa sezione, si può notare come quella a cui corrisponde un periodo di clock pari a 7ns, presenta una maximum clock frequency maggiore rispetto a quelle avente periodo maggiore. Questo aumento lo si può, inoltre, evidenziare nelle successive soluzioni aventi periodo di clock ancora minore. Tanto è vero che il WNS, al diminuire del clock constraint, anch'esso diminuisce.

Solution	Clock Enable	Clocks	DSP	Logic	Set/Reset	Data
clk=10ns	0.455035944	1.215706812	0.343570428	0.925468281	0.00349783	1.014495501
clk=9ns	0.371109403	1.601127908	0.308091694	0.871003722	0.003504661	0.845550094
clk=8ns	0.432465255	1.790320966	0.364705513	1.03614782	0.003386509	1.085355412
clk=7ns	0.4179838	1.991891069	0.305155962	0.881534419	0.002575182	0.757947506
clk=6ns	0.367850007	2.092533279	0.244985771	0.788475969	0.00273619	0.666663051
clk=5ns	0.556948828	3.174162703	0.291668839	0.943421735	0.004532358	0.906153
clk=4ns	0.57213963	3.804846667	0.253423554	0.799402245	0.002476984	0.756326073

Table 21: Vivado Operation Chaining Solution Dynamic Power Report [mW]

Solution	Dynamic Total
clk=10ns	3.957774796
clk=9ns	4.000387481
clk=8ns	4.712381475
clk=7ns	4.357087937
clk=6ns	4.163244267
clk=5ns	5.876887463
clk=4ns	6.188615153

Table 22: Vivado Operation Chaining Solution Dynamic Power Report [mW]

Solution	Energy Single Operation
clk=10ns	39.57774796
clk=9ns	36.00348733
clk=8ns	37.6990518
clk=7ns	30.49961556
clk=6ns	24.9794656
clk=5ns	29.38443731
clk=4ns	24.75446061

Table 23: Vivado Operation Chaining Solution Energy Single Operation Report [pJ]

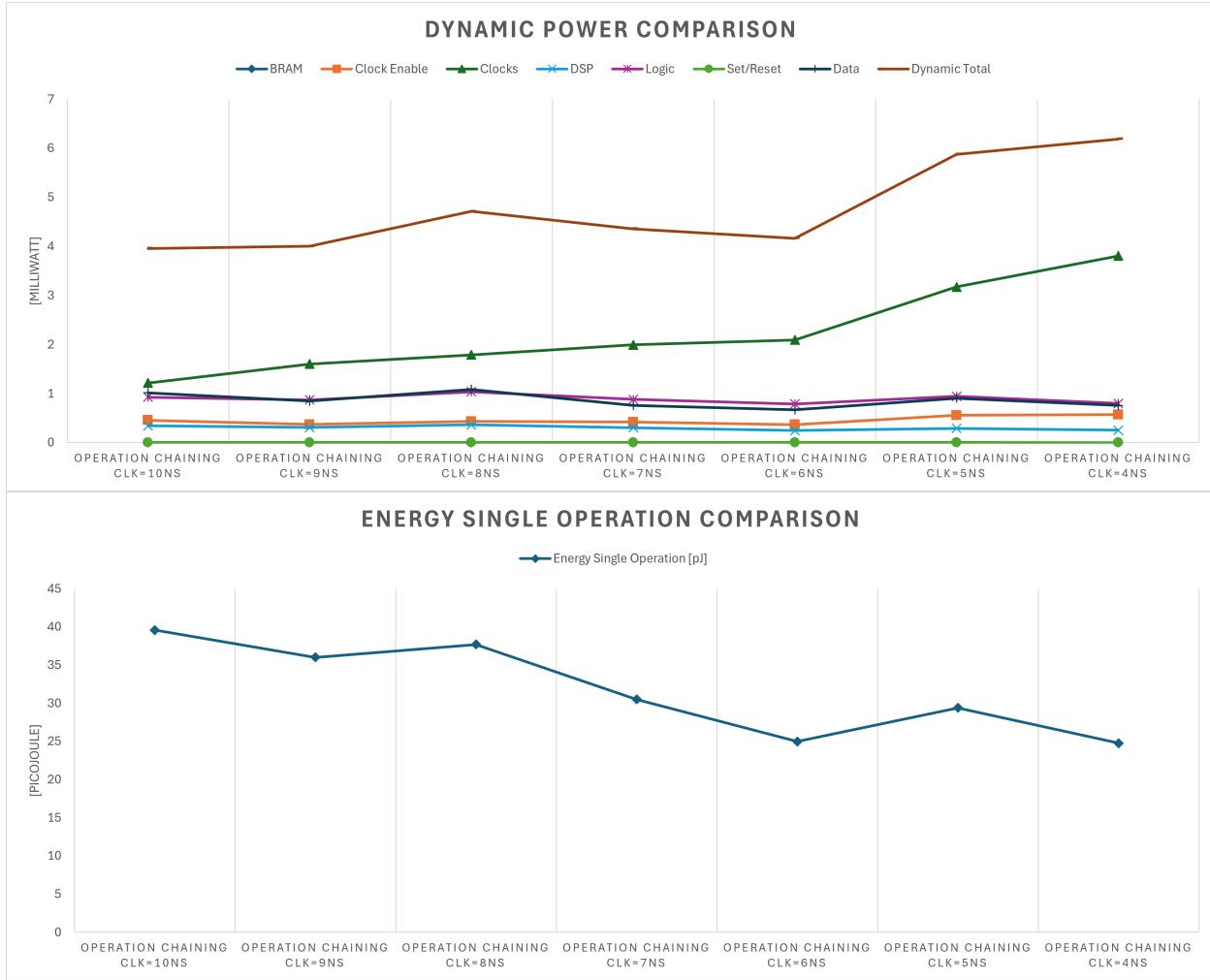


Figure 13: Vivado Operation Chaining Power Plot

Si può notare come la potenza dinamica aumenti al diminuire del periodo di clock considerato. Questo si evince dal fatto che l'utilizzazione delle risorse (in particolar modo per quello che riguarda i FF) aumenta al diminuire del constraint di clock. Tanto è vero che la potenza dinamica associata all'albero di distribuzione del clock (BUFG) e del clock enable sia di gran lunga maggiore rispetto alla soluzione iniziale dove è stato considerato un constraint pari a 10ns. In particolare, questo aumento di potenza dinamica può essere associato al maggior numero di Flip Flop utilizzato nelle soluzioni aventi periodo di clock sempre minore. Pertanto, considerando i report di utilizzazione delle risorse, di timing, di potenza dinamica ed energia per singola operazione, si può evincere che la soluzione hardware che permette un trade-off ottimale è quella ottenuta in corrispondenza del periodo di clock pari a 7ns.

5.3 Code Hoisting Solution

Considerando, ora, la stessa architettura software precedentemente utilizzata, si potrebbe mettere in luce un possibile problema: impiego di risorse aggiuntive a causa di istruzioni logiche condizionali. In particolare, questo problema è dovuto al fatto che, all'interno del ciclo for, è prevista l'istruzione condizionale *if*(*i* == 0) che comporta l'utilizzo di risorse hardware addizionali. Bisogna precisare che, istruzioni logiche quali *if*, *else* e *then* rendono i loop inefficienti limitando le performance dal punto di vista hardware. Fondamentalmente, l'istruzione condizionale *if*(*i* == 0) a livello hardware si traduce nel leggere il valore dell'indice *i* in una certa struttura dati ed effettuare il confronto, tramite un comparatore, tra tale valore appena citato e il valore zero.

```

1 #include "definitions.h"
2
3 void firConvolutionCodeHoisting(samplesType inputFilter, samplesType* outputFilter) {
4     coeffsType coefficientsFilter[SIZE] = {53, 0, -91, 0, 313, 500, 313, 0, -91, 0, 53};
5     static samplesType shiftRegister[SIZE];
6     accType accumulator;
7     accumulator = 0;
8     int i;
9     loop: for( i=SIZE-1; i>0; --i ) {
10         shiftRegister[i] = shiftRegister[i-1];
11         accumulator += shiftRegister[i] * coefficientsFilter[i];
12     }
13     accumulator += inputFilter * coefficientsFilter[0];
14     shiftRegister[0] = inputFilter;
15     *outputFilter = accumulator;
16 }
```

Pertanto, si potrebbe pensare di eliminare l'*if* in questione e spostarlo al di fuori del loop. In questo modo, anche l'istruzione condizionale *else*, che ricoprirebbe il resto delle iterazioni, verrà gestita direttamente all'interno del loop avendo quest'ultimo un'iterazione in meno dovuta ai motivi appena citati. Pertanto, quello che ci si aspetta è una diminuzione dal punto di vista della latenza totale, dal punto di vista delle risorse e per ciò che riguarda la potenza dinamica associata alla parte di logica dell'architettura hardware in questione.

Pertanto, si mostrano qui di seguito i report, generati da HLS, riguardo la sintesi, la C/RTL Cosimulation e l'Export RTL.

Clock	Target	Estimated	Uncertainty
ap_clk	10.00	8.510	1.25

Table 24: HLS Code Hoisting Solution Timing Summary (ns)

Latency		Interval	
min	max	min	max
42	42	42	42

Table 25: HLS Code Hoisting Solution Latency Summary (clock cycles)

Loop Name	Latency		Iteration Latency		Initiation	Interval	Trip Count
	min	max	min	max	achieved	target	
- loop	40	40	4	4	-	-	10

Table 26: HLS Code Hoisting Solution Latency Loops Summary

Si può notare come il numero di iterazioni del loop, in questo caso, sia pari a 10 tale che la latenza totale del ciclo risulta essere pari a 40.

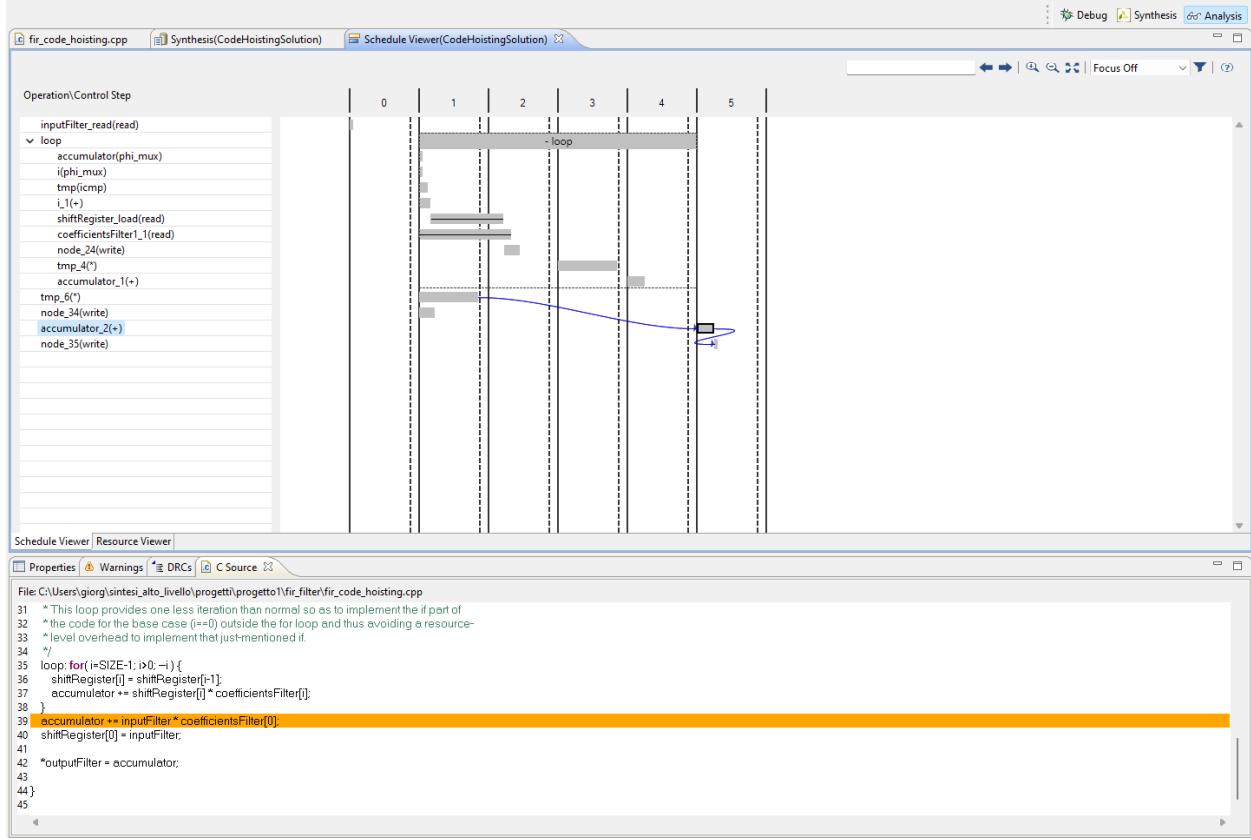


Figure 14: HLS Code Hoisting Solution Trip Count Analysis

In particolare, è opportuno fare un'attenta riflessione riguardo la latenza totale dell'architettura hardware in questione. Il report di sintesi, effettivamente, ha restituito una stima di tale valore pari a 42 dal momento bisogna considerare una latenza totale del loop pari a 40, un ciclo di latency ulteriori dovuti alla lettura iniziale dell'input (come avviene per ogni soluzione progettata) e, come ultimo, un altro ciclo finale dovuto alle due operazioni in parallelo effettuate: accumulo associato al coefficiente in posizione 0 e scrittura del dato in uscita. Precedentemente, nella soluzione iniziale non ottimizzata, tale latenza finale non era prevista dal momento che il caso $i == 0$ veniva gestito all'interno del loop mediante istruzione condizionale e la scrittura del dato in uscita veniva effettuata in parallelo all'operazione appena citata.

Name	BRAM_18K	DSP48E	FF	LUT
DSP	-	-	-	-
Expression	-	4	0	140
FIFO	-	-	-	-
Instance	-	-	-	-
Memory	0	-	74	8
Multiplexer	-	-	-	92
Register	-	-	156	-
Total	0	4	230	240
Available	280	220	106400	53200
Utilization (%)	0	1	~0	~0

Table 27: HLS Code Hoisting Solution Utilization Estimates Summary

Si può notare come il report di C/RTL Cosimulation evidenzi un ciclo di latenza in meno rispetto alla

soluzione hardware iniziale non ottimizzata. Questo è dovuto al fatto che, come precedentemente citato, il loop presenta un’iterazione in meno tale per cui il numero di cicli di clock, affinché venga processato si abbia un risultato in uscita, è pari a 43 rispetto ai 44 corrispondenti alla soluzione non ottimizzata.

RTL	Status	Latency			Interval		
		min	avg	max	min	avg	max
VHDL	Pass	42	42	43	42	42	43

Table 28: HLS Code Hoisting Solution C/RTL Cosimulation Summary

Inoltre, si può notare come il numero di risorse relative alla soluzione hardware in questione risulta essere pressocché il medesimo tranne per il numero di FF che risulta essere ridotto di circa il 18%. In aggiunta, si evidenzia come, rispetto alla soluzione iniziale, il periodo di clock raggiunto nella fase di post-implementation (stimato) risulta essere maggiore e, pertanto, avente una corrispondente incertezza minore.

Resource	VHDL
SLICE	75
LUT	270
FF	131
DSP	2
BRAM	0
SRL	0

Timing	VHDL
CP required	10.000
CP achieved post-synthesis	5.745
CP achieved post-implementation	6.847

Table 29: HLS Code Hoisting Solution Export RTL Resource Usage

Table 30: HLS Code Hoisting Solution Export RTL Final Timing

Pertanto, importando l’IP in Vivado e impostando un clock constraint pari a 10ns è possibile analizzare i seguenti report di risorse, timing, potenza dinamica ed energia per singola operazione.

```
1 create_clock -period 10.000 -name myclk -waveform {0.000 5.000} [get_ports ap_clk]
```

LUT	LUTRAM	FF	BRAM	DSP	IO	BUFG
270	32	134	0	2	71	1

Table 31: Vivado Code Hoisting Solution Utilization Report [#]

Dal momento che il WNS è diminuito, si può notare come la maximum clock frequency sia diminuita rispetto alla soluzione non ottimizzata.

Cycles [#]	Clock Constraint [ns]	WNS [ns]	Maximum Clock Frequency [Mhz]
43	10	3.074	144.3834825

Table 32: Vivado Code Hoisting Solution Timing Report

Per quanto riguarda la potenza dinamica e l’energia per singola operazione, si evidenzia un aumento rispetto alla soluzione iniziale. Si può ipotizzare che, dal punto di vista della potenza, il tool riesca meglio ad ottimizzare direttamente tutte le operazioni relative ad accumulo e shifting nello stesso ciclo, anziché averne 10 nel ciclo ed una al di fuori.

BRAM	Clock Enable	Clocks	DSP	Logic	Set/Reset [mW]	Data
0	0.370487687	1.756788697	0.41467679	0.838255044	3.35E-03	1.381990616

Table 33: Vivado Code Hoisting Solution Dynamic Power Report [mW]

Dynamic Total
4.765546238

Table 34: Vivado Code Hoisting Solution Dynamic Power Report [mW]

Energy Single Operation
47.65546238

Table 35: Vivado Code Hoisting Solution Energy Single Operation Report [pJ]

Pertanto, se da una parte si è ottenuto una piccola diminuzione delle risorse avendo gestito l'overhead di risorse dovuto alle istruzioni condizionali presenti all'interno del loop, dall'altra parte si è ottenuto un aumento di circa il 21% della potenza dinamica rispetto alla soluzione iniziale non ottimizzata.

5.4 Loop Fission Solution

Considerando la soluzione precedente, si è risolto la possibile problematica di overhead delle risorse dovuta alle istruzioni condizionali presenti all'interno del loop. Bisogna notare, però, che all'interno di tale ciclo for siano presenti istruzioni potenzialmente complesse dal punto di vista computazionale. In particolare, il tool, durante la compilazione del codice, cerca in qualche modo di apportare qualche tipologia di ottimizzazione al loop in questione. Pertanto, si potrebbe pensare di implementare due loop: uno per l'operazione di shifting ed un altro per l'operazione di accumulo. Quindi, effettuando la scissione del loop, relativo alla soluzione precedente, si potrebbe pensare che il tool possa effettuare tali ottimizzazioni in maniera più efficiente dal momento che le ottimizzazioni verrebbero effettuate su due cicli distinti. Quello che ci si aspetta, però, è un aumento significativo della latenza totale dovuto al fatto che, in questo caso, bisogna considerare la latenza di due loop.

```

1 #include "definitions.h"
2
3 void firConvolutionLoopFission(samplesType inputFilter, samplesType* outputFilter) {
4     coeffsType coefficientsFilter[SIZE] = {53, 0, -91, 0, 313, 500, 313, 0, -91, 0, 53};
5     static samplesType shiftRegister[SIZE];
6     accType accumulator;
7     accumulator = 0;
8     int i;
9     loopShifting: for( i=SIZE-1; i>0; --i ) {
10         shiftRegister[i] = shiftRegister[i-1];
11     }
12     shiftRegister[0] = inputFilter;
13     accumulator = 0;
14     loopAccumulator: for( i=SIZE-1; i>=0; --i ) {
15         accumulator += shiftRegister[i] * coefficientsFilter[i];
16     }
17     *outputFilter = accumulator;
18 }
```

Effettuando la sintesi, la C/RTL Cosimulation ed Export RTL si ottengono i seguenti report.

Clock	Target	Estimated	Uncertainty
ap_clk	10.00	8.510	1.25

Table 36: HLS Loop Fission Solution Timing Summary (ns)

Latency		Interval	
min	max	min	max
66	66	66	66

Table 37: HLS Loop Fission Solution Latency Summary (clock cycles)

Si può notare come la latenza sia aumentata rispetto alle soluzioni precedenti in cui era presente un solo loop. Questo può essere, inoltre, analizzato nella tabella "Table 38: HLS Loop Fission Solution Latency Loops Summary" sotto allegata dove è possibile osservare che i due loop vengono riconosciuti come due loop indipendenti tra loro e con una latenza associata. In particolare, "loopShifting" presenta un trip count pari a 10 dal momento che il caso $i == 0$ viene gestito secondo la tecnica del code hoisting, cioè al di fuori del ciclo stesso.

Loop Name	Latency		Iteration Latency		Initiation Interval	Trip Count
	min	max	min	max	achieved	target
- loopShifting	20	20	2	2	-	-
- loopAccumulator	44	44	4	4	-	-

Table 38: HLS Loop Fission Solution Latency Loops Summary

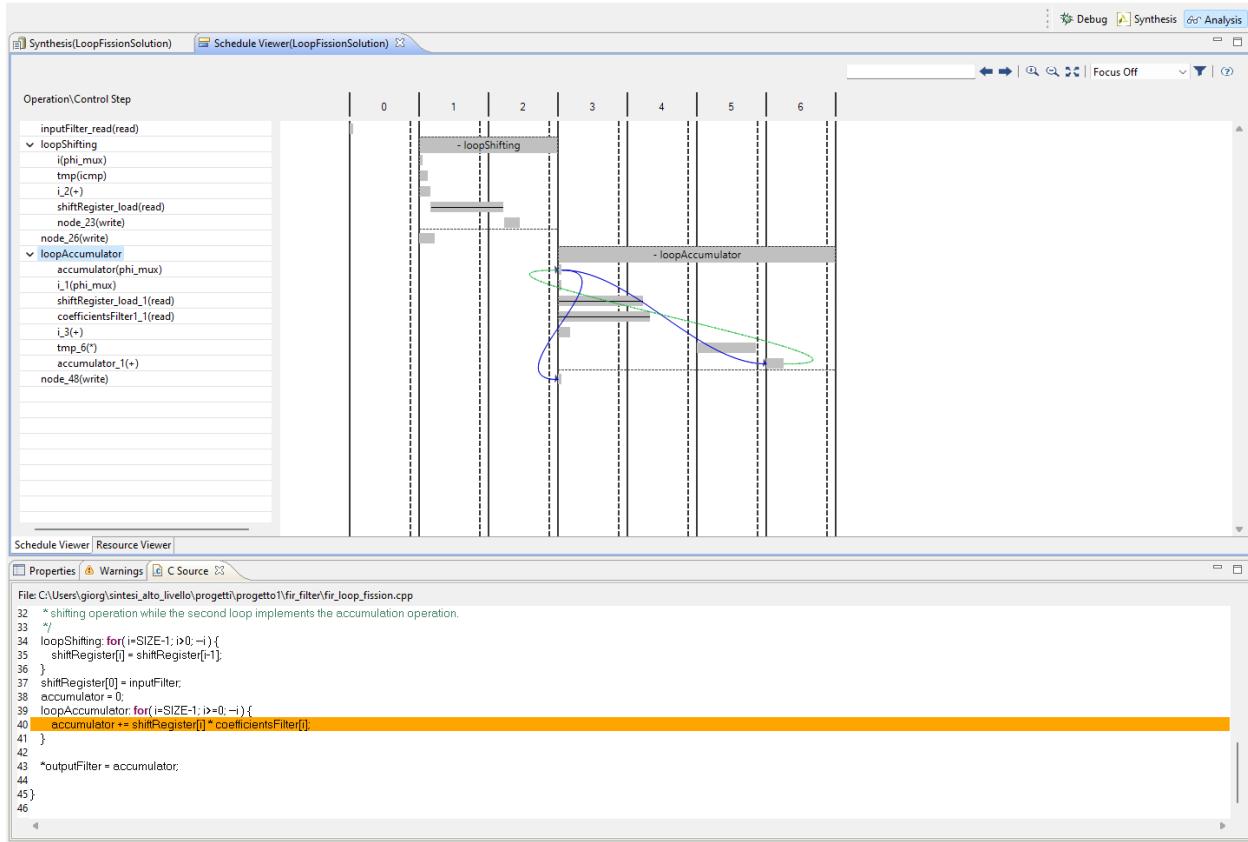


Figure 15: HLS Loop Fission Solution Analysis

La scissione del loop nei cicli "loopShifting" e "loopAccumulator" la si può notare, inoltre, nella figura sopra allegata che mostra la finestra di *Analysis*. Bisogna di nuovo precisare che, tale tecnica ha l'obiettivo di favorire il tool di ottimizzare i due loop in maniera indipendente effettuando magari uno scheduling delle operazioni differente, attraverso procedure proprietarie, rispetto alle soluzioni precedenti dove era previsto un ciclo for unico.

Name	BRAM_18K	DSP48E	FF	LUT
DSP	-	-	-	-
Expression	-	2	0	96
FIFO	-	-	-	-
Instance	-	-	-	-
Memory	0	-	74	8
Multiplexer	-	-	-	110
Register	-	-	131	-
Total	0	2	205	214
Available	280	220	106400	53200
Utilization (%)	0	~0	~0	~0

Table 39: HLS Loop Fission Solution Utilization Estimates Summary

RTL	Status	Latency			Interval		
		min	avg	max	min	avg	max
VHDL	Pass	66	66	67	66	66	67

Table 40: HLS Loop Fissiong Solution C/RTL Cosimulation Summary

Resource	VHDL
SLICE	44
LUT	157
FF	106
DSP	2
BRAM	0
SRL	0

Timing	VHDL
CP required	10.000
CP achieved post-synthesis	5.745
CP achieved post-implementation	6.363

Table 41: HLS Loop Fissiong Solution Export
RTL Resource Usage

Si può evidenziare come, dopo aver effettuato l'Export RTL, il numero di cicli di clock affinché un risultato venga processato in uscita sia aumentato a 67 rispetto a quello ottenuto nella soluzione precedente pari a 43. Di fondamentale importanza è la diminuzione dell'utilizzazione delle risorse. Infatti, si può notare, rispetto alla soluzione hardware basata sul code hoisting, come il numero di SLICE sia diminuito di circa il 41%, quello delle LUT di circa il 42% e quello dei FF di circa il 19%.

Pertanto, importando l'IP in Vivado e impostando un clock constraint pari a 10ns è possibile analizzare i seguenti report di risorse, timing, potenza dinamica ed energia per singola operazione.

```
1 create_clock -period 10.000 -name myclk -waveform {0.000 5.000} [get_ports ap_clk]
```

LUT	LUTRAM	FF	BRAM	DSP	IO	BUFG
158	32	106	0	2	71	1

Table 43: Vivado Loop Fission Solution Utilization Report [#]

Si può notare come l'utilizzazione delle risorse sia minore (–42% LUT e –21% FF) rispetto alla solution precedente e, inoltre, la maximum clock frequency sia aumentata (dato dal fatto che il WNS è maggiore). Ovviamente, a tali vantaggi corrisponde, come già anticipato precedentemente, un aumento dei cicli di clock necessari (+55.8%) affinché un risultato venga processato in uscita.

Cycles [#]	Clock Constraint [ns]	WNS [ns]	Maximum Clock Frequency [Mhz]
67	10	3.464	152.998776

Table 44: Vivado Loop Fission Solution Timing Report

BRAM	Clock Enable	Clocks	DSP	Logic	Set/Reset [mW]	Data
0	0.292603218	1.098937588	0.321194879	0.590288255	0.004160448	0.937705627

Table 45: Vivado Loop Fission Solution Dynamic Power Report [mW]

Dynamic Total
3.244890014

Table 46: Vivado Loop Fission Solution Dynamic Power Report [mW]

Energy Single Operation
32.44890014

Table 47: Vivado Loop Fission Solution Energy Single Operation Report [pJ]

Si evidenzia che, la potenza dinamica totale e l'energia per singola operazione risultano essere diminuite, rispetto alla soluzione hardware precedente, di circa il 32%. In particolare, le maggiori diminuzioni si hanno in corrispondenza dei seguenti contributi: Clocks (-38%), Logic (-30%) e Data (-32%). Infatti, bisogna ricordare che l'utilizzazione delle risorse, come già precedentemente citato, è anch'essa notevolmente diminuita rispetto alla solution precedente.

5.5 Loop Unrolling Solution

Considerando l'architettura hardware precedente, è possibile, effettivamente, implementare un approccio più veloce dal punto di vista delle operazioni, cioè una soluzione basata su parallelismo. In particolare, allocando un maggior numero di risorse, è possibile eseguire delle operazioni in parallelo riducendo di conseguenza la latenza, cioè il tempo di esecuzione totale.

Tale approccio è possibile implementarlo in due modi differenti:

- **Manuale**

Unrolling ottenuto tramite rimodulazione software

- **Automatico**

Unrolling ottenuto tramite direttiva proprietaria (pragma)

In particolare, sia l'unrolling manuale sia quello automatico, sono stati ottenuti considerando sia un fattore di parallelismo pari a 2 sia pari a 4. Inoltre, bisogna specificare che, per quanto riguarda l'unrolling automatico (tramite pragma), è stata realizzata un'ulteriore soluzione hardware, tenendo conto sia del fattore 2 sia del fattore 4, dove è stato considerato sia il pragma di unrolling sia il pragma di partitioning. Il partizionamento serve per risolvere un problema tipicamente causato dagli array. Gli array sono implementati come BRAM, solitamente progettate per un dual-port massimo. Questo può limitare il throughput di un algoritmo ad alta intensità di read/write. La larghezza di banda può essere migliorata dividendo l'array (una singola BRAM) in array più piccoli (più BRAM), aumentando di fatto il numero di porte. Gli array vengono partizionati utilizzando la direttiva ARRAY_PARTITION. Vivado HLS offre tre tipi di partizionamento degli array. I tre tipi di partizionamento sono:

- **block**

L'array originale viene suddiviso in blocchi di uguali dimensioni di elementi consecutivi dell'array originale.

- **cyclic**

L'array originale viene suddiviso in blocchi di uguali dimensioni che interlacciano gli elementi dell'array originale.

- **complete**

L'operazione predefinita consiste nel dividere l'array nei suoi singoli elementi. Ciò corrisponde alla risoluzione di una memoria in registri.

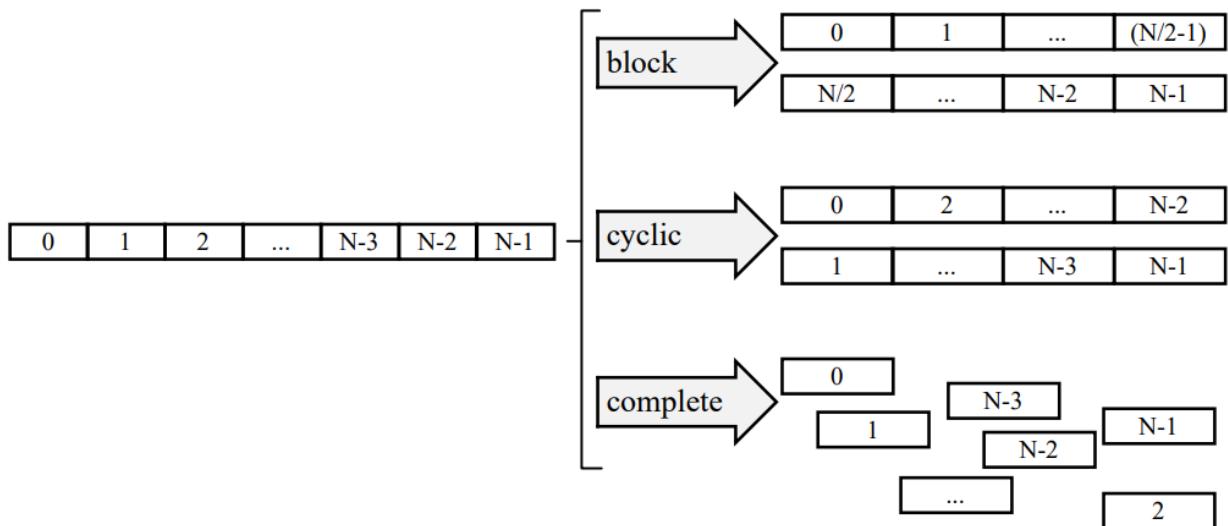


Figure 16: HLS Array Partitioning

Pertanto, è possibile elencare qui di seguito le soluzioni hardware progettate per l'approccio del parallelismo:

- Unrolling Manuale Fattore=2
- Unrolling Manuale Fattore=4
- Unrolling Automatico Fattore=2
- Unrolling Automatico Fattore=4
- Unrolling Automatico con Partitioning Fattore=2
- Unrolling Automatico con Partitioning Fattore=4

```
void top(...) {
    ...
    for_mult:for (i=3;i>0;i--) {
        a[i] = b[i] * c[i];
    }
    ...
}
```

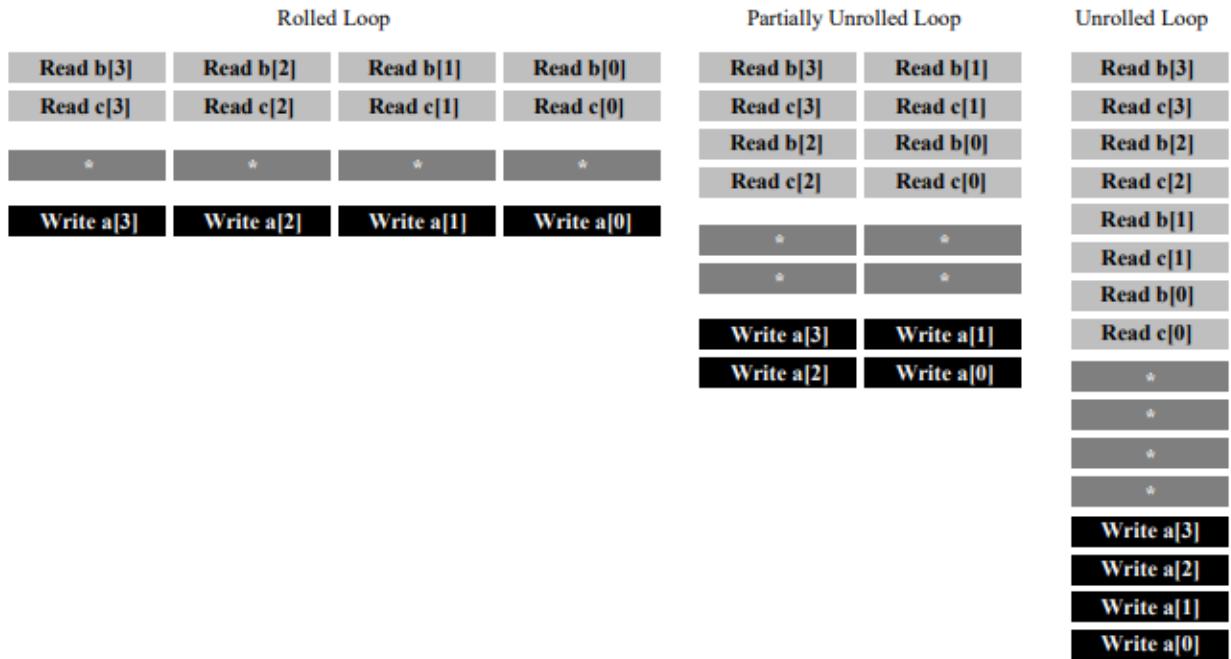


Figure 17: HLS Loop Unrolling

A questo proposito, verranno analizzati i report considerando le soluzioni hardware aventi stesso fattore così da poter effettuare confronti adeguati tra le varie implementazioni. Bisogna precisare, che in questo caso è stato effettuando un loop unrolling parziale, cioè tenendo conto della figura sopra allegata è stato considerato un fattore 2 o un fattore 4 anziché utilizzare un fattore che potesse *ricoprire* la totalità degli elementi processabili (capacità del filtro). Ovviamente, la figura sopra allegata è solo un esempio per rappresentare meglio il concetto di unrolling e partitioning.

5.5.1 Loop Unrolling Factor=2

```

1 #include "definitions.h"
2
3 void firConvolutionLoopUnrollingFactor2(samplesType inputFilter, samplesType* outputFilter)
4 {
5     coeffsType coefficientsFilter[SIZE] = {53, 0, -91, 0, 313, 500, 313, 0, -91, 0, 53};
6     static samplesType shiftRegister[SIZE];
7     accType accumulator;
8     accumulator = 0;
9     int i;
10    loopShifting: for( i=SIZE-1; i>1; i=i-2 ) {
11        shiftRegister[i] = shiftRegister[i-1];
12        shiftRegister[i-1] = shiftRegister[i-2];
13    }
14    if( i==1 ) {
15        shiftRegister[1] = shiftRegister[0];
16    }
17    shiftRegister[0] = inputFilter;
18    loopAccumulator: for( i=SIZE-1; i>=0; --i ) {
19        accumulator += shiftRegister[i] * coefficientsFilter[i];
20    }
21    *outputFilter = accumulator;
22 }

23 #include "definitions.h"
24
25 void firConvolutionLoopUnrollingFactor2Pragma(samplesType inputFilter, samplesType*
26 outputFilter) {
27     coeffsType coefficientsFilter[SIZE] = {53, 0, -91, 0, 313, 500, 313, 0, -91, 0, 53};
28     static samplesType shiftRegister[SIZE];
29     accType accumulator;
30     accumulator = 0;
31     int i;
32     loopShifting: for( i=SIZE-1; i>0; --i ) {
33         #pragma HLS unroll factor=2
34         shiftRegister[i] = shiftRegister[i-1];
35     }
36     shiftRegister[0] = inputFilter;
37     accumulator = 0;
38     loopAccumulator: for( i=SIZE-1; i>=0; --i ) {
39         accumulator += shiftRegister[i] * coefficientsFilter[i];
40     }
41     *outputFilter = accumulator;
42 }

43 #include "definitions.h"
44
45 void firConvolutionLoopUnrollingFactor2PP(samplesType inputFilter, samplesType* outputFilter
46 ) {
47     coeffsType coefficientsFilter[SIZE] = {53, 0, -91, 0, 313, 500, 313, 0, -91, 0, 53};
48     static samplesType shiftRegister[SIZE];
49     accType accumulator;
50     accumulator = 0;
51     int i;
52     loopShifting: for( i=SIZE-1; i>0; --i ) {
53         #pragma HLS array_partition variable=shiftRegister complete
54         #pragma HLS unroll factor=2
55         shiftRegister[i] = shiftRegister[i-1];
56     }
57     shiftRegister[0] = inputFilter;
58     accumulator = 0;
59     loopAccumulator: for( i=SIZE-1; i>=0; --i ) {
60         accumulator += shiftRegister[i] * coefficientsFilter[i];
61     }
62     *outputFilter = accumulator;
63 }
```

Solution	Clock	Target	Estimated	Uncertainty
Manuale	ap_clk	10.00	8.510	1.25
Automatico	ap_clk	10.00	8.510	1.25
Automatico con Partitioning	ap_clk	10.00	8.510	1.25

Table 48: HLS Loop Unrolling Factor=2 Solution Timing Summary (ns)

Si può notare come la latenza associata alla soluzione hardware basata su unrolling manuale sia la medesima di quella basata su unrolling automatico (effettuato mediante pragma). Questo è dovuto al fatto che l'implementazione è la stessa: nel primo caso è effettuato manualmente mentre nel secondo caso è effettuato mediante direttiva proprietaria del tool. Invece, nel caso della soluzione hardware basata su unrolling automatico e partitioning, si riscontra una latenza maggiore. Questo risultato è dovuto al fatto che, tramite il pragma di partizionamento, si è dovuto gestire, inoltre, letture e scritture in parallelo.

Solution	Latency		Interval	
	min	max	min	max
Manuale	56	56	56	56
Automatico	56	56	56	56
Automatico con Partitioning	61	66	61	66

Table 49: HLS Loop Unrolling Factor=2 Solution Latency Summary (clock cycles)

Si può evidenziare come, nel caso della soluzione hardware con urolling automatico e partizionamento, l'aumento della latenza si ha soltanto in corrispondenza del loop di shifting, cioè proprio dove è stato collocato il pragma di partitioning.

Solution	Loop Name	Latency		Iteration Latency		Initiation achieved	Interval target	Trip Count
		min	max	min	max			
Manuale	- loopShifting	10	10	2	2	-	-	5
	- loopAccumulator	44	44	4	4	-	-	11
Automatico	- loopShifting	10	10	2	2	-	-	5
	- loopAccumulator	44	44	4	4	-	-	11
Automatico con Partitioning	- loopShifting	15	20	3	4	-	-	5
	- loopAccumulator	44	44	4	4	-	-	11

Table 50: HLS Loop Unrolling Factor=2 Solution Latency Loops Summary

Aprendo la sezione Analysis del tool si può vedere meglio nel dettaglio che, all'interno del loop di shifting, sono presenti 2 shifting in parallelo (dal momento che è stato considerato un fattore di parallelismo pari a 2) sia peer la soluzione hardware di unrolling manuale sia per quella di unrolling automatico (mediante pragma).

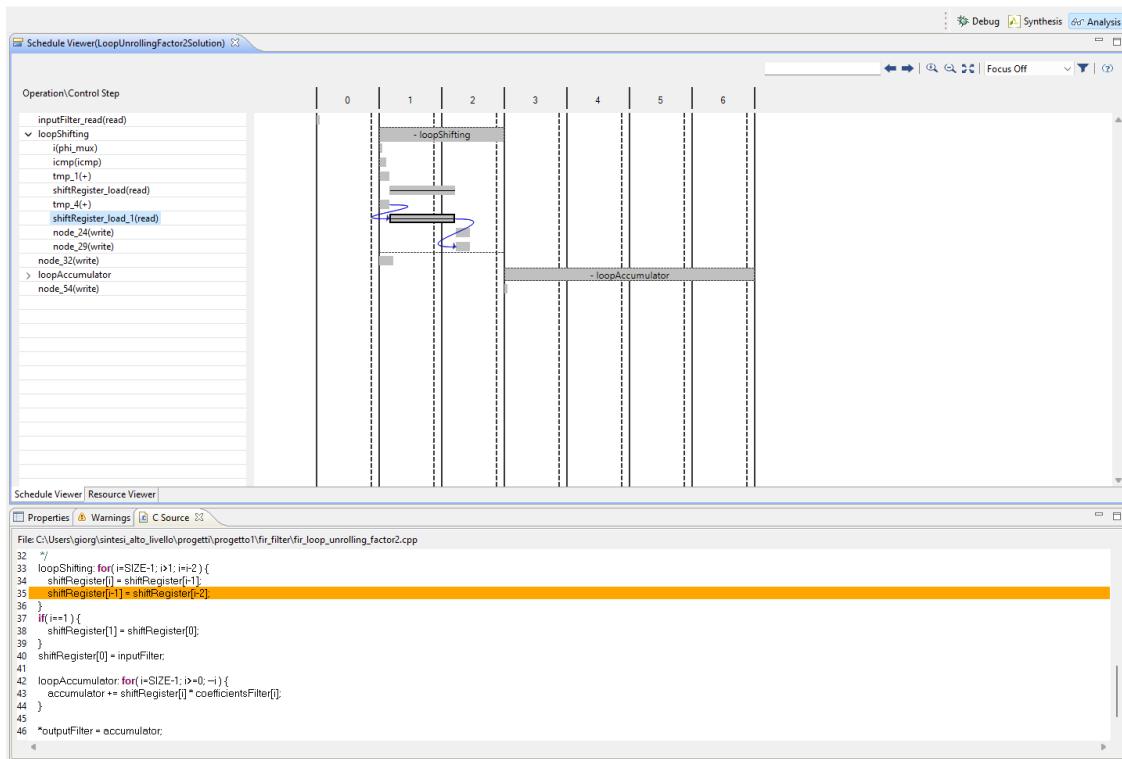


Figure 18: HLS Loop Unrolling Manual Factor=2 Analysis

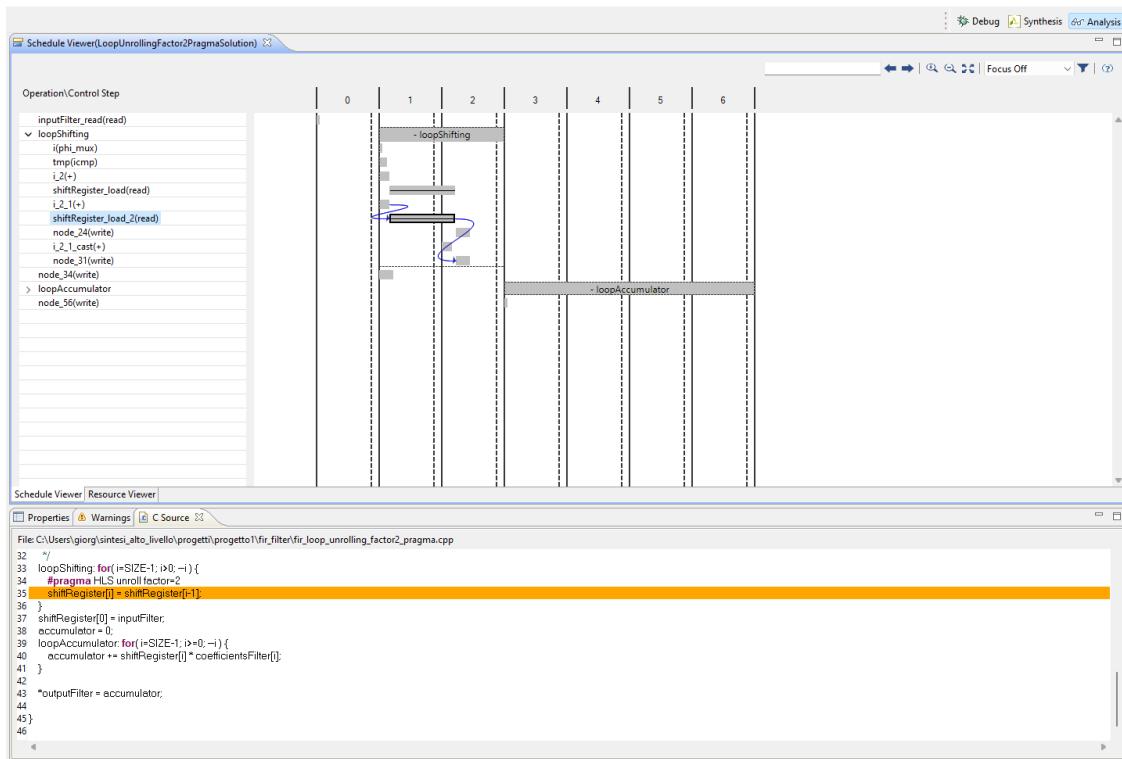


Figure 19: HLS Loop Unrolling Automatic Factor=2 Analysis

Analizzando nel dettaglio il report si può notare come il numero di DSP e LUT sia pressocché il medesimo per tutte e tre le soluzioni hardware proposte. Per quanto riguarda, invece, il numero di BRAM utilizzate, le soluzioni basate su unrolling manuale e automatico presentano un'utilizzazione pari a 2, mentre quella basata su unrolling automatico e partitioning presenta un numero di BRAM utilizzato pari a 0. Questo è dovuto al fatto che il pragma di partizionamento, al fine di garantire accessi di read/write in memoria in parallelo, impone di non utilizzare le BRAM dual-port (che limiterebbero tali accessi). In particolare, impostando la tipologia di partizionamento *complete*, la direttiva fa in modo che la soluzione hardware impieghi come memoria dei singoli registri, così che siano garantiti gli accessi in memoria in parallelo. Tanto è vero che si può notare un aumento considerevole dei FF in corrispondenza dell'architettura basata su unrolling automatico e partizionamento.

Solution	BRAM_18K	DSP48E	FF	LUT
Manuale	2	2	145	236
Automatico	2	2	141	251
Automatico con Partitioning	0	2	912	289

Table 51: HLS Loop Unrolling Factor=2 Solution Utilization Estimates [#]

Qui di seguito vengono riportati i report relativi alla C/RTL Cosimulation, dove è possibile analizzare il numero di cicli di clock che servono per ottenere un risultato in uscita, e quello relativo a Export RTL.

Solution	RTL	Status	Latency			Interval		
			min	avg	max	min	avg	max
Manuale	VHDL	Pass	56	56	57	56	56	57
Automatico	VHDL	Pass	56	56	57	56	56	57
Automatico con Partitioning	VHDL	Pass	65	65	66	65	65	66

Table 52: HLS Loop Unrolling Factor=2 Solution C/RTL Cosimulation Report

Solution	SLICE	LUT	FF	DSP	BRAM	CP required	CP achieved post-synthesis	CP achieved post-implementation
Manuale	31	97	72	2	2	10	5.745	5.692
Automatico	29	97	72	2	2	10	5.745	5.692
Automatico con Partitioning	294	413	843	2	0	10	5.745	6.188

Table 53: HLS Loop Unrolling Factor=2 Solution Export RTL Report

Pertanto, importando l'IP in Vivado e impostando un clock constraint pari a 10ns è possibile analizzare i seguenti report di risorse, timing, potenza dinamica ed energia per singola operazione.

```
1 create_clock -period 10.000 -name myclk -waveform {0.000 5.000} [get_ports ap_clk]
```

Analizzando il report di utilizzazione delle risorse generato da Vivado dopo aver effettuato il processo di implementazione, è possibile notare come il numero di risorse stimato dal tool HLS è leggermente mutato. Però, si può evidenziare come il cambiamento delle risorse tra le varie soluzioni hardware è il medesimo, cioè il numero di LUT e FF nel caso dell'implementazione con unrolling automatico e partitioning sia aumentato considerevolmente e il numero di BRAM si è ridotto a zero.

Solution	LUT	LUTRAM	FF	BRAM	DSP	IO	BUFG
Manuale	98	0	72	1	2	71	1
Automatico	98	0	72	1	2	71	1
Automatico con Partitioning	413	0	843	0	2	71	1

Table 54: Vivado Loop Unrolling Factor=2 Solution Utilization Report [#]

Si è, inoltre, analizzato l'utilizzazione delle risorse effettuando un confronto con la soluzione hardware basata sulla scissione del loop dal momento che tali architetture basate su unrolling sono basate sulla prima citata. In particolare, si può notare come il numero di risorse utilizzate sia pressoché il medesimo tra la soluzione hardware basata su loop fission e quelle basate rispettivamente sull'unrolling manuale e automatico di fattore 2. Più nello specifico, considerando le implementazioni basate sul parallelismo manuale e automatico rispetto a quella basata sul loop fission, il numero di LUT è diminuito di circa il 40%, il numero di FF è diminuito di circa il 32% e il numero di BRAM è aumentato di un'unità. Invece, effettuando un confronto tra quella basata su partitioning e quella basata sulla scissione del loop, il numero di LUT è aumentato di circa il 161% e il numero di FF di circa il 695%.

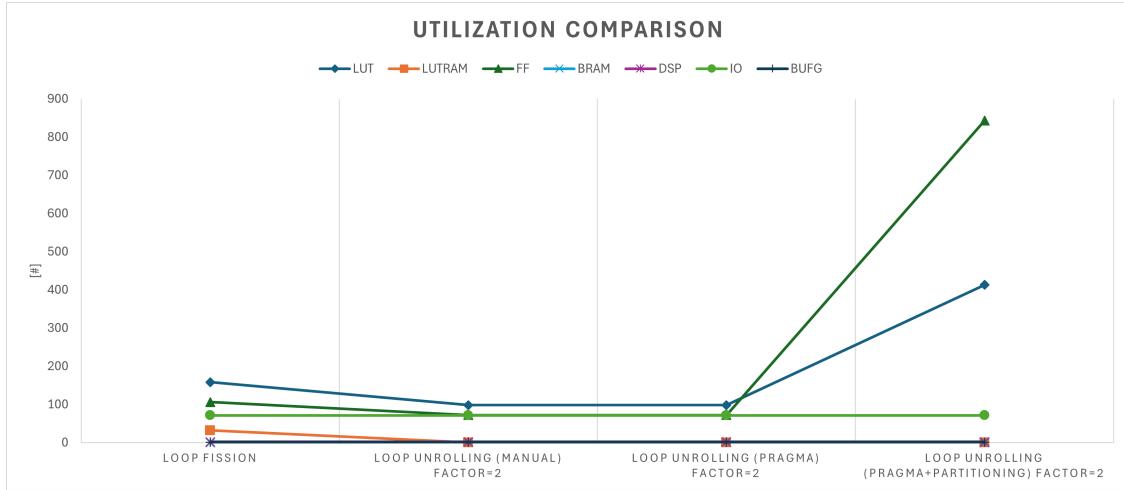


Figure 20: Vivado Loop Unrolling Factor=2 Utilization Plot

Infine, per quanto riguarda il timing, è possibile notare come il numero di cicli sia pressoché il medesimo tra la soluzione basata su unrolling manuale e quella basata su unrolling automatico. Invece, come precedentemente citato, la soluzione basata su partizionamento presenta un numero di cicli di clock, tali per garantire un risultato, maggiore.

Solution	Cycles [#]	Clock Constraint [ns]	WNS [ns]	Maximum Clock Frequency [MHz]
Manuale	57	10	4.33	176.366843
Automatico	57	10	4.33	176.366843
Automatico con Partitioning	66	10	3.469	153.1159087

Table 55: Vivado Loop Unrolling Factor=2 Solution Timing Report

Effettuando un confronto grafico e tenendo conto anche della soluzione basata sulla scissione del loop, è possibile evidenziare un numero di cicli di clock, tali per garantire un risultato, un WNS e una maximum clock frequency pressoché uguali tra l'implementazione basata su loop fission e quella basata su partizionamento. Questo accade dal momento che la soluzione hardware basata su unrolling automatico e partitioning fa

diminuire il numero di cicli di clock tramite il parallelismo e lo fa aumentare allo stesso tempo a causa del partizionamento.

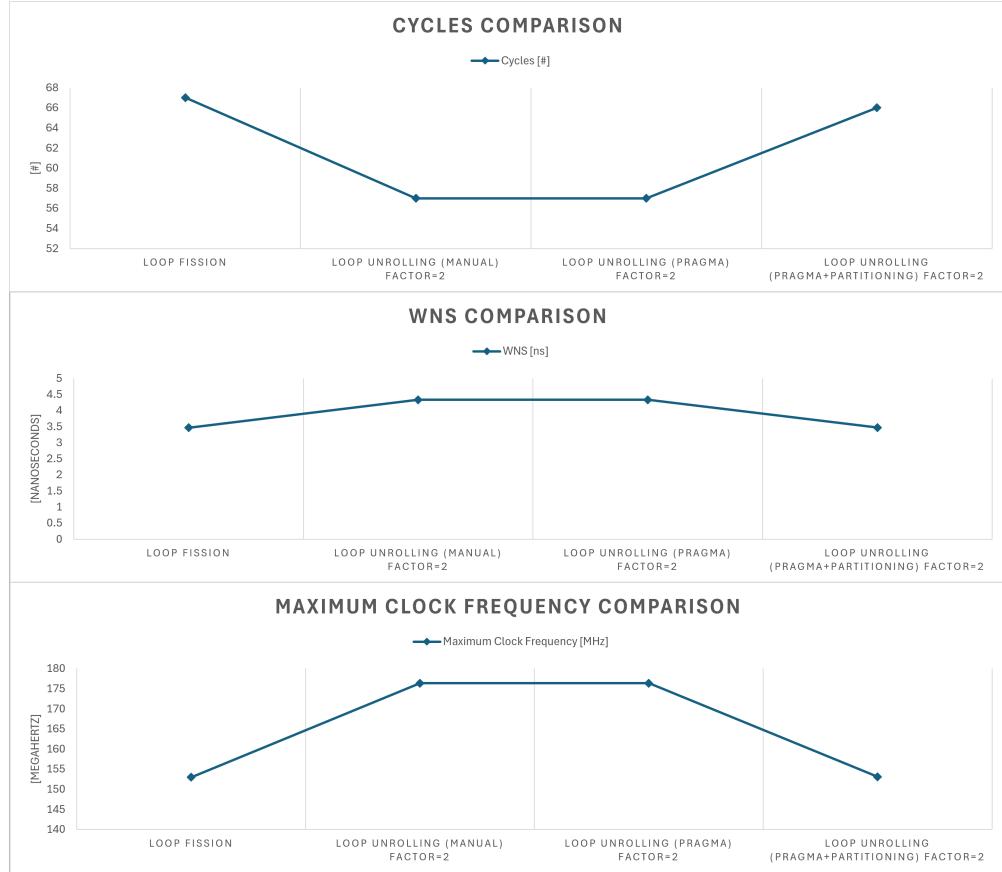


Figure 21: Vivado Loop Unrolling Factor=2 Timing Plot

Per quanto riguarda, invece, la potenza dinamica associata alle tre soluzioni hardware proposte in questa sezione, è possibile notare come i contributi associati al Clock Enable e Clocks risultano essere notevolmente più grandi in corrispondenza dell'implementazione basata su partizionamento. Questo potrebbe essere causato dal notevole aumento dell'utilizzazione dei FF per questa solution.

Solution	BRAM	Clock Enable	Clocks	DSP	Logic	Set / Reset	Data
Manuale	1.250551548	0.096387172	0.900532817	0.268251897	0.260709843	0.003146866	0.423992984
Automatico	1.240851358	0.082524632	0.960682868	0.272355421	0.266662013	0.00428147	0.425589533
Automatico con Partitioning	0	0.325270201	2.352835611	0.263715046	0.575191109	0.007010513	0.750690058

Table 56: Vivado Loop Unrolling Factor=2 Solution Dynamic Power Report [mW]

Infatti, è possibile riscontrare un aumento della potenza dinamica totale e dell'energia per singola operazione in corrispondenza dell'architettura basata su unrolling e partitioning. Invece, per quanto riguarda le altre due solution (unrolling manuale e automatico), entrambi i parametri appena citati risultano essere i medesimi.

Solution	Dynamic Total
Manuale	3.203573127
Automatico	3.252947294
Automatico con Partitioning	4.274712538

Table 57: Vivado Loop Unrolling Factor=2 Solution Dynamic Power Report [mW]

Solution	Energy Single Operation
Manuale	32.03573127
Automatico	32.52947294
Automatico con Partitioning	42.74712538

Table 58: Vivado Loop Unrolling Factor=2 Solution Energy Single Operation Report [pJ]

Analizzando graficamente queste tre soluzioni e in aggiunta anche l'implementazione basata sulla scissione del loop, è possibile notare come quella basata su loop fission presenta un medesimo valore di potenza dinamica totale ed energia per singola operazione rispetto alle due soluzioni di unrolling manuale e automatico.

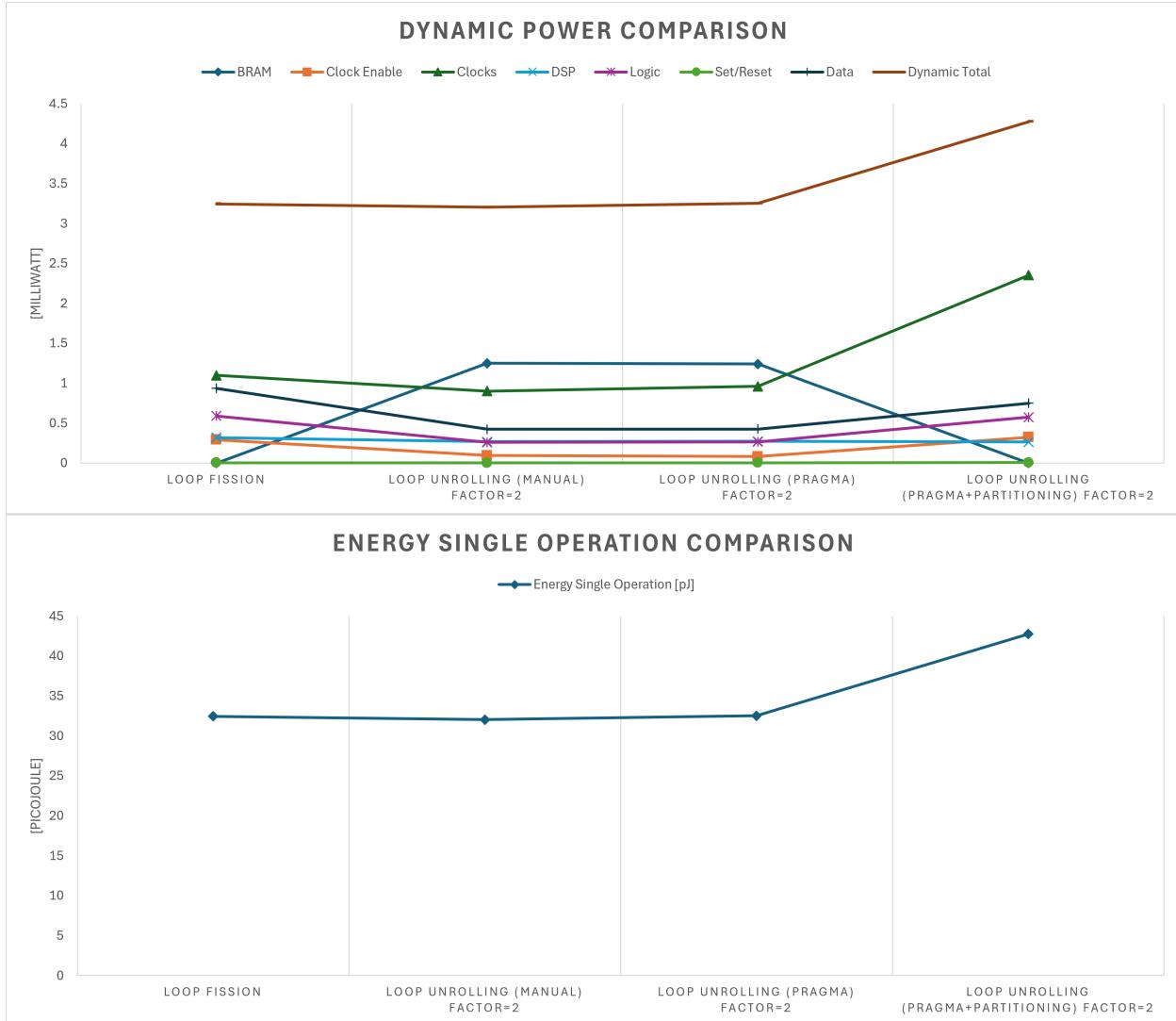


Figure 22: Vivado Loop Unrolling Factor=2 Dynamic Power Plot

5.5.2 Loop Unrolling Factor=4

```

1 #include "definitions.h"
2
3 void firConvolutionLoopUnrollingFactor4(samplesType inputFilter, samplesType* outputFilter)
4 {
5     coeffsType coefficientsFilter[SIZE] = {53, 0, -91, 0, 313, 500, 313, 0, -91, 0, 53};
6     static samplesType shiftRegister[SIZE];
7     accType accumulator;
8     accumulator = 0;
9     int i;
10    loopShifting: for( i=SIZE-1; i>1; i=i-4 ) {
11        shiftRegister[i] = shiftRegister[i-1];
12        shiftRegister[i-1] = shiftRegister[i-2];
13        shiftRegister[i-2] = shiftRegister[i-3];
14        shiftRegister[i-3] = shiftRegister[i-4];
15    }
16    if( i==1 ) {
17        shiftRegister[1] = shiftRegister[0];
18    }
19    shiftRegister[0] = inputFilter;
20    loopAccumulator: for( i=SIZE-1; i>=0; --i ) {
21        accumulator += shiftRegister[i] * coefficientsFilter[i];
22    }
23    *outputFilter = accumulator;
24 }

1 #include "definitions.h"
2
3 void firConvolutionLoopUnrollingFactor4Pragma(samplesType inputFilter, samplesType*
4     outputFilter) {
5     coeffsType coefficientsFilter[SIZE] = {53, 0, -91, 0, 313, 500, 313, 0, -91, 0, 53};
6     static samplesType shiftRegister[SIZE];
7     accType accumulator;
8     accumulator = 0;
9     int i;
10    loopShifting: for( i=SIZE-1; i>0; --i ) {
11        #pragma HLS unroll factor=4
12        shiftRegister[i] = shiftRegister[i-1];
13    }
14    shiftRegister[0] = inputFilter;
15    accumulator = 0;
16    loopAccumulator: for( i=SIZE-1; i>=0; --i ) {
17        accumulator += shiftRegister[i] * coefficientsFilter[i];
18    }
19    *outputFilter = accumulator;
20 }

1 #include "definitions.h"
2
3 void firConvolutionLoopUnrollingFactor4PP(samplesType inputFilter, samplesType* outputFilter
4 ) {
5     coeffsType coefficientsFilter[SIZE] = {53, 0, -91, 0, 313, 500, 313, 0, -91, 0, 53};
6     static samplesType shiftRegister[SIZE];
7     accType accumulator;
8     accumulator = 0;
9     int i;
10    loopShifting: for( i=SIZE-1; i>0; --i ) {
11        #pragma HLS array_partition variable=shiftRegister complete
12        #pragma HLS unroll factor=4
13        shiftRegister[i] = shiftRegister[i-1];
14    }
15    shiftRegister[0] = inputFilter;
16    accumulator = 0;
17    loopAccumulator: for( i=SIZE-1; i>=0; --i ) {
18        accumulator += shiftRegister[i] * coefficientsFilter[i];
19    }
20    *outputFilter = accumulator;
21 }
```

Solution	Clock	Target	Estimated	Uncertainty
Manuale	ap_clk	10.00	8.510	1.25
Automatico	ap_clk	10.00	8.510	1.25
Automatico con Partitioning	ap_clk	10.00	8.510	1.25

Table 59: HLS Loop Unrolling Factor=4 Solution Timing Summary (ns)

Solution	Latency		Interval	
	min	max	min	max
Manuale	58	58	58	58
Automatico	62	66	62	66
Automatico con Partitioning	62	64	62	64

Table 60: HLS Loop Unrolling Factor=4 Solution Latency Summary (clock cycles)

Solution	Loop Name	Latency		Iteration Latency		Initiation achieved	Interval target	Trip Count
		min	max	min	max			
Manuale	- loopShifting	12	12	4	4	-	-	3
	- loopAccumulator	44	44	4	4	-	-	11
Automatico	- loopShifting	15	18	5	5	-	-	3
	- loopAccumulator	44	44	4	4	-	-	11
Automatico con Partitioning	- loopShifting	15	17	5	5	-	-	3
	- loopAccumulator	44	44	4	4	-	-	11

Table 61: HLS Loop Unrolling Factor=4 Solution Latency Loops Summary

In questo caso si può notare come il trip count sia, ovviamente, il medesimo per tutte e tre le soluzioni. Ciò che cambia effettivamente, come nel caso del fattore pari a 2, è la latency associata ai loop. In particolare, analizzando più nel dettaglio, tramite l'interfaccia Analysis, si può evidenziare come il loopShifting, in cui viene implementato un parallelismo pari a 4, viene gestito in maniera differente. Infatti, nel caso del loop unrolling manuale, vengono effettuate in parallelo due read alla volta e poi successivamente vengono eseguite in parallelo due write alla volta. Invece, nel caso della soluzione hardware basata su unrolling automatico, vengono effettuate una read e una write in parallelo per ogni ciclo di latenza. Nello specifico, considerando quattro shifting per ogni iterazione, viene effettuata la read al primo ciclo di latenza e, successivamente, viene eseguita la write corrispondente in parallelo alla read relativo al secondo shifting e così via. Di conseguenza, si avranno cinque colpi di latenza per ogni iterazione come riportato dal valore di iteration latency del report di sintesi.

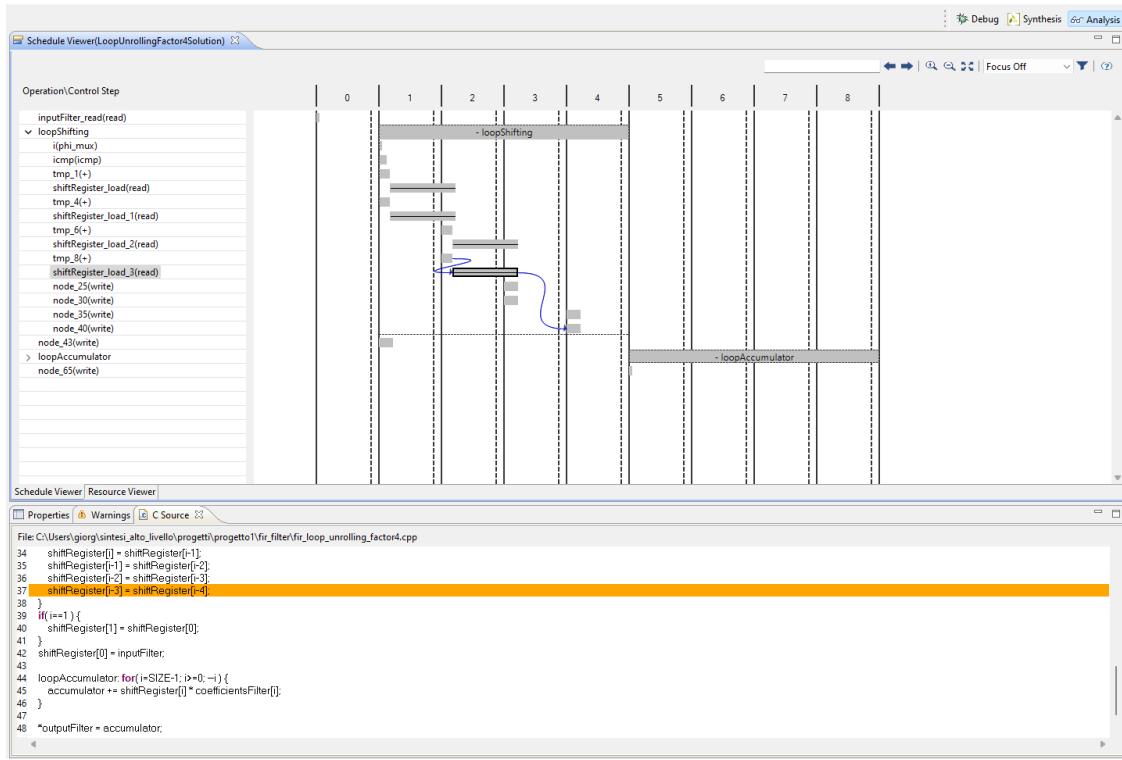


Figure 23: HLS Loop Unrolling Manual Factor=4 Analysis

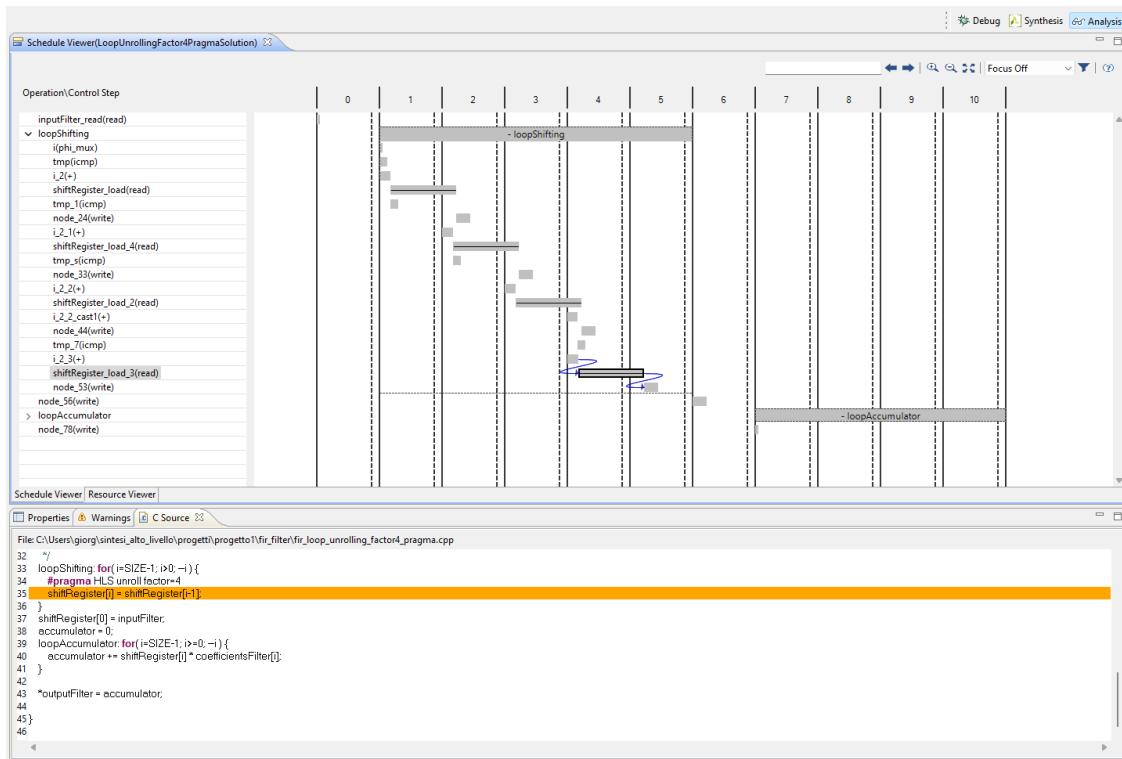


Figure 24: HLS Loop Unrolling Automatic Factor=4 Analysis

Analogamente a quanto è stato analizzato e descritto nel caso di parallelismo di fattore pari a 2, anche in questo caso si può notare come l'utilizzazione delle risorse sia quasi la medesima tra il loop unrolling manuale e automatico. Invece, per quanto riguarda il loop unrolling automatico con partitioning, si ha un notevole aumento dell'utilizzazione sia di FF sia di LUT.

Solution	BRAM_18K	DSP48E	FF	LUT
Manuale	2	2	221	318
Automatico	2	2	194	367
Automatico con Partitioning	0	2	928	748

Table 62: HLS Loop Unrolling Factor=4 Solution Utilization Estimates [#]

Qui di seguito vengono riportati i report relativi alla C/RTL Cosimulation, dove è possibile analizzare il numero di cicli di clock che servono per ottenere un risultato in uscita, e quello relativo a Export RTL.

Solution	RTL	Status	Latency			Interval		
			min	avg	max	min	avg	max
Manuale	VHDL	Pass	58	58	59	58	58	59
Automatico	VHDL	Pass	60	60	61	60	60	61
Automatico con Partitioning	VHDL	Pass	58	58	59	58	58	59

Table 63: HLS Loop Unrolling Factor=4 Solution C/RTL Cosimulation Report

Solution	SLICE	LUT	FF	DSP	BRAM	CP required	CP achieved post-synthesis	CP achieved post-implementation
Manuale	46	160	147	2	2	10	5.745	5.692
Automatico	40	145	93	2	2	10	5.745	5.692
Automatico con Partitioning	412	1145	864	2	0	10	5.745	7.109

Table 64: HLS Loop Unrolling Factor=4 Solution Export RTL Report

Pertanto, importando l'IP in Vivado e impostando un clock constraint pari a 10ns è possibile analizzare i seguenti report di risorse, timing, potenza dinamica ed energia per singola operazione.

```
1 create_clock -period 10.000 -name myclk -waveform {0.000 5.000} [get_ports ap_clk]
```

Analogamente all'unrolling di fattore pari a 2, anche in questo caso in corrispondenza della soluzione hardware basata su unrolling e partizionamento si ha un notevole aumento di utilizzazione delle risorse. Inoltre, si può notare come il numero di FF utilizzati per la solution di unrolling automatico risulta essere circa il 37% in meno rispetto a quella basata su unrolling manuale. Evidentemente il tool, tramite la direttiva proprietaria, è riuscito ad effettuare delle ottimizzazioni tali da garantire una diminuzione dei Flip Flop.

Solution	LUT	LUTRAM	FF	BRAM	DSP	IO	BUFG
Manuale	159	0	147	1	2	71	1
Automatico	145	0	93	1	2	71	1
Automatico con Partitioning	1145	0	864	0	2	71	1

Table 65: Vivado Loop Unrolling Factor=4 Solution Utilization Report [#]

Effettuando un confronto grafico e tenendo conto dei dati precedentemente ottenuti per la soluzione basata su scissione del loop, si può notare come l'utilizzazione delle risorse sia pressocché la medesima tra loop fission e loop unrolling automatico.

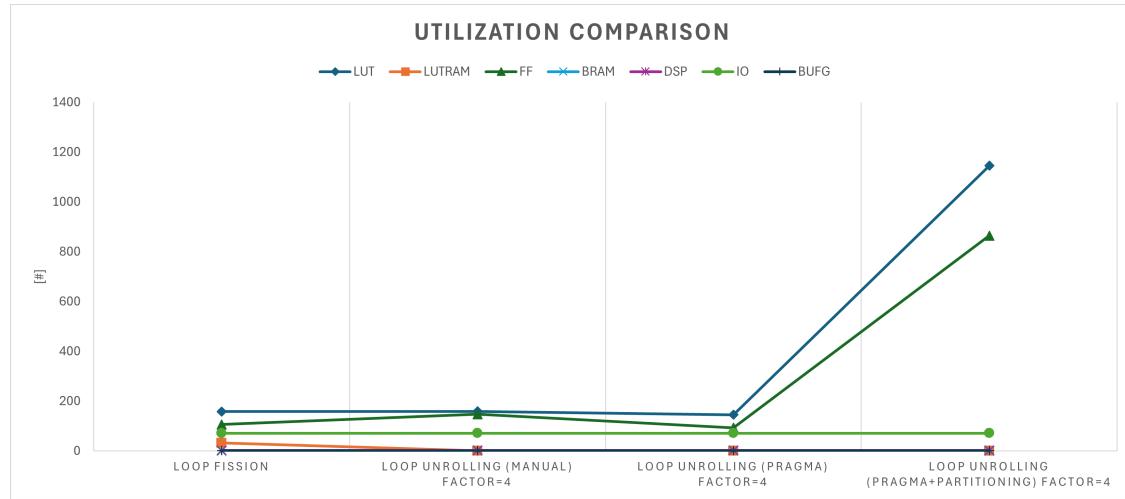


Figure 25: Vivado Loop Unrolling Factor=4 Utilization Plot

Per quanto riguarda il report di timing, è possibile notare, analogamente a quanto successo per l'unrolling di fattore pari a 2, una diminuzione della maximum clock frequency in corrispondenza della soluzione basata su loop unrolling automatico con partizionamento.

Solution	Cycles [#]	Clock Constraint [ns]	WNS [ns]	Maximum Clock Frequency [MHz]
Manuale	59	10	4.257	174.1250218
Automatico	61	10	4.33	176.366843
Automatico con Partitioning	59	10	3.097	144.8645516

Table 66: Vivado Loop Unrolling Factor=4 Solution Timing Report

Per quanto riguarda, invece, il numero di cicli di clock per garantire un risultato in uscita, in questo caso, per le tre soluzioni basate su unrolling di fattore pari a 4, si è ottenuto un valore minore rispetto alla soluzione basata su scissione del loop.



Figure 26: Vivado Loop Unrolling Factor=4 Timing Plot

Analogamente all'unrolling di fattore pari a 2, anche in questo caso in corrispondenza della soluzione hardware basata su unrolling e partizionamento si ha un notevole aumento del contributo di potenza dinamica relativo al *Clocks* comportando un aumento della potenza dinamica totale.

Solution	BRAM	Clock Enable	Clocks	DSP	Logic	Set/Reset	Data
Manuale	1.32976193	0.382625905	0.253616716	0.921033497	0.543549308	0.010887122	0.585376518
Automatico	1.23103871	0.106978332	0.972227077	0.247065967	0.258892891	0.002585625	0.41881192
Automatico con Partitioning	0	0.328624592	2.560390625	0.298048515	1.146363211	0.0030065	1.324957004

Table 67: Vivado Loop Unrolling Factor=4 Solution Dynamic Power Report [mW]

Solution	Dynamic Total
Manuale	4.026850996
Automatico	3.237600523
Automatico con Partitioning	5.661390447

Table 68: Vivado Loop Unrolling Factor=4 Solution Dynamic Power Report [mW]

Solution	Energy Single Operation
Manuale	40.26850996
Automatico	32.37600523
Automatico con Partitioning	56.61390447

Table 69: Vivado Loop Unrolling Factor=4 Solution Energy Single Operation Report [pJ]

Bisogna notare un aspetto molto interessante. La potenza dinamica totale e l'energia per singola operazione associata alla soluzione hardware basata su loop unrolling automatico di fattore pari a 4 risultano essere pressocché le medesime di quelle relative all'unrolling automatico di fattore pari a 2. Questo potrebbe essere dovuto a ottimizzazioni effettuate dal tool tramite la direttiva proprietaria dell'unrolling.

5.6 Loop Pipelining Solution

Il pipelining consente di eseguire operazioni in modo simultaneo: ogni fase di esecuzione non deve completare tutte le operazioni prima di iniziare quella successiva. Pertanto, questo approccio permette di scindere le operazioni complesse in più operazioni semplici. In questo modo si può far lavorare l'architettura con dati temporalmente differenti. Ad esempio, tenendo conto della figura sotto allegata, questa funzione senza pipeline dovrebbe eseguire tutte le iterazioni previste in maniera sequenziale e, quindi, eseguendo read, compute, write, read, compute, write e così via. Considerando, invece, un approccio mediante pipeline, quello che succede è che alla prima iterazione viene eseguita la prima micro-operazione, alla seconda iterazione viene eseguita la seconda micro-operazione e in contemporanea si può eseguire la prima micro-operazione che nel caso precedente senza pipeline era prevista per il quarto ciclo di clock. L'idea è quella di far lavorare l'architettura sfruttando la scissione della logica complessa in più micro-operazioni.

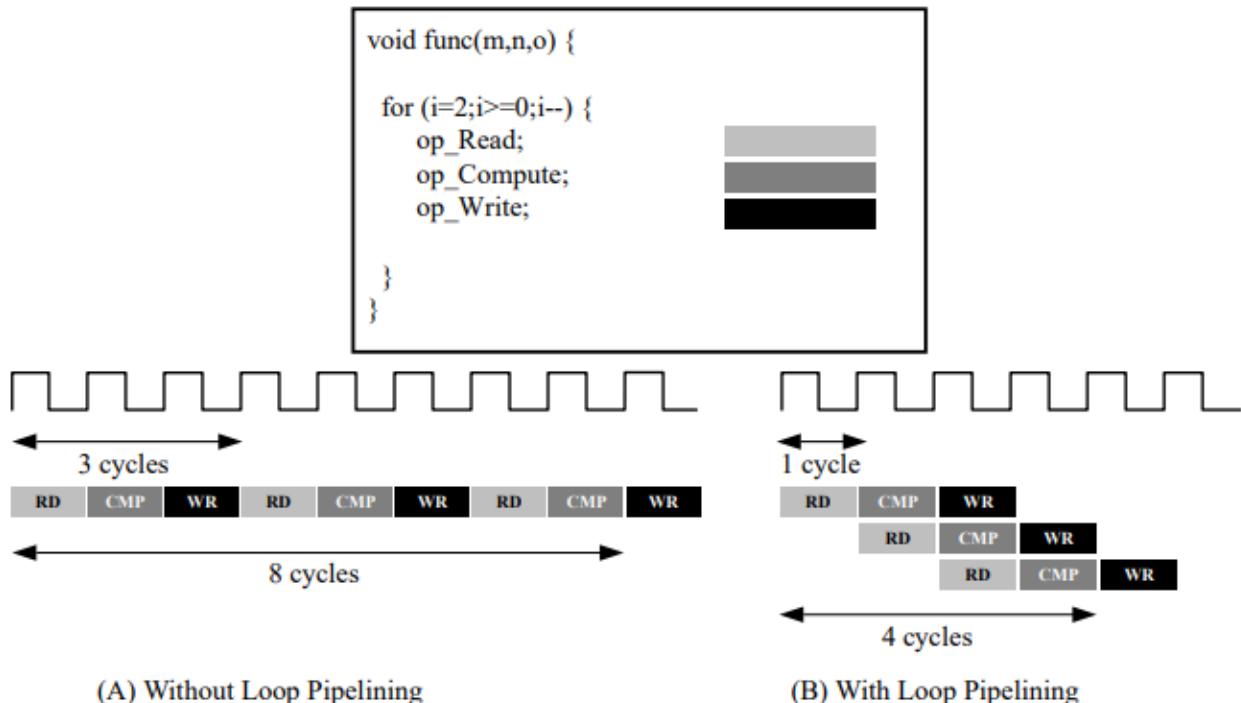


Figure 27: HLS Loop Pipelining

In particolare, è possibile evidenziare due parametri fondamentali per questa tipologia di approccio:

- **Iteration Latency (IL)**

Numero di cicli che servono per eseguire un'iterazione del corpo del ciclo.

- **Initiation Interval (II) := Throughput**

Numero di cicli di clock fino alla prossima iterazione del loop.

Idealmente, ci si aspetta un throughput pari a 1 perché corrisponderebbe al caso migliore per il pipelining. Ovviamente non sempre è possibile ottenere un throughput pari a 1 ma in questo caso l'obiettivo è quello di ottenerlo.

Il pipelining in HLS è possibile implementarlo mediante una direttiva proprietaria che permette, inoltre, di specificare l'Initiation Interval desiderato. Tanto è vero che all'interno del report di sintesi viene specificato l'II target, cioè quello specificato nella direttiva, e l'II raggiunto dal tool. Teoricamente, si dovrebbe ottenere II_target=II_achieved. Se così non fosse allora il tool non è riuscito a raggiungere l'obiettivo prefissato. Pertanto, si potrebbe modificare il codice per vedere se effettivamente il problema può essere risolto nella

seguente maniera. Se così non fosse, evidentemente non dipende dal codice ma dal tipo di applicazione che si sta cercando di implementare che non lo permette.

```

1 #include "definitions.h"
2
3 void firConvolutionLoopPipelining(samplesType inputFilter, samplesType* outputFilter) {
4     coeffsType coefficientsFilter[SIZE] = {53, 0, -91, 0, 313, 500, 313, 0, -91, 0, 53};
5     static samplesType shiftRegister[SIZE];
6     accType accumulator;
7     accumulator = 0;
8     int i;
9     loop: for( i=SIZE-1; i>=0; --i ) {
10         #pragma HLS pipeline II=1
11         #pragma HLS array_partition variable=shiftRegister complete
12         if( i==0 ) {
13             accumulator += inputFilter * coefficientsFilter[0];
14             shiftRegister[0] = inputFilter;
15         } else {
16             shiftRegister[i] = shiftRegister[i-1];
17             accumulator += shiftRegister[i] * coefficientsFilter[i];
18         }
19     }
20     *outputFilter = accumulator;
21 }
```

In questo caso, la soluzione hardware presa come riferimento è stata quella iniziale non ottimizzata. In particolare, all'interno del loop è stata definita sia la direttiva di pipeline, specificando un $II=1$, sia la direttiva di partizionamento. Nello specifico, è stato considerato il pragma di partitioning al fine di realizzabilità della soluzione in questione. Questo è dovuto al fatto che, senza considerare tale direttiva, si otterrebbe un $II_{achieved}$ differente da un II_{target} . Infatti, utilizzare lo shift register vuol dire non poter accedere ai dati intermedi ma soltanto a quelli finali. Pertanto, per far fronte a questo problema si potrebbe pensare di realizzare tali shift register mediante Flip Flop (tramite la direttiva di partizionamento). Effettuando la sintesi otteniamo i seguenti report.

Clock	Target	Estimated	Uncertainty
ap_clk	10.00	8.510	1.25

Table 70: HLS Loop Pipelining Solution Timing Summary (ns)

Latency		Interval	
min	max	min	max
15	15	15	15

Table 71: HLS Loop Pipelining Solution Latency Summary (clock cycles)

Si può notare come la latenza totale del loop sia notevolmente diminuita rispetto alla soluzione iniziale non ottimizzata. Inoltre, si può evidenziare come l'Initiation Interval raggiunto sia il medesimo di quello target, cioè pari a 1. Questo è stato possibile per i motivi precedentemente citati.

In particolare, considerando il valore di Iteration Latency, il valore del Trip Count e il valore di Initiation Interval, è possibile ricavare il valore della latenza relativa al loop:

$$IL + [(N - 1) * II] - 1 = 4 + [(11 - 1) * 1] - 1 = 4 + [10 * 1] - 1 = 4 + 10 - 1 = 13 \quad (7)$$

Loop Name	Latency		Iteration Latency		Initiation Interval		Trip Count
	min	max	min	max	achieved	target	
- loop	13	13	4	4	1	1	11

Table 72: HLS Loop Pipelining Solution Latency Loops Summary

Per quanto riguarda l'utilizzazione stimata delle risorse, si può notare come il numero di FF sia aumentato rispetto alla soluzione non ottimizzata dal momento che è stato utilizzato il pragma di partizionamento tale da partizionare shiftRegister in diverse strutture da un certo numero di FF ciascuna anziché utilizzare BRAM.

Name	BRAM_18K	DSP48E	FF	LUT
DSP	-	-	-	-
Expression	-	4	0	480
FIFO	-	-	-	-
Instance	-	-	-	-
Memory	0	-	10	2
Multiplexer	-	-	-	75
Register	-	-	700	64
Total	0	4	710	621
Available	280	220	106400	53200
Utilization (%)	0	1	~0	1

Table 73: HLS Loop Pipelining Solution Utilization Estimates Summary

In particolare, si può notare come la direttiva proprietaria di partitioning ha partizionato shiftRegister in 10 strutture da 32 FF ciascuna.

Name	FF	LUT	Bits	Const Bits
accumulator_reg_114	32	0	32	0
ap_CS_fsm	3	0	3	0
ap_enable_reg_pp0_iter0	1	0	1	0
ap_enable_reg_pp0_iter1	1	0	1	0
ap_enable_reg_pp0_iter2	1	0	1	0
ap_enable_reg_pp0_iter3	1	0	1	0
ap_phi_reg_pp0_iter1_p_pn_reg_138	32	0	32	0
ap_phi_reg_pp0_iter2_p_pn_reg_138	32	0	32	0
ap_phi_reg_pp0_iter3_p_pn_reg_138	32	0	32	0
coefficientsFilter1_1_reg_458	10	0	10	0
i_reg_127	5	0	5	0
shiftRegister_0	32	0	32	0
shiftRegister_1	32	0	32	0
shiftRegister_2	32	0	32	0
shiftRegister_3	32	0	32	0
shiftRegister_4	32	0	32	0
shiftRegister_5	32	0	32	0
shiftRegister_6	32	0	32	0
shiftRegister_7	32	0	32	0
shiftRegister_8	32	0	32	0
shiftRegister_9	32	0	32	0
shiftRegister_load_p_reg_453	32	0	32	0
tmp_1_reg_426	1	0	1	0
tmp_2_reg_417	32	0	32	0
tmp_4_reg_430	4	0	4	0
tmp_6_reg_463	32	0	32	0
tmp_reg_422	1	0	1	0
tmp_1_reg_426	64	32	1	0
tmp_reg_422	64	32	1	0
Total	700	64	574	0

Table 74: HLS Loop Pipelining Solution Register Detail Estimates Summary

Effettuando la C/RTL Cosimulation ed Export RTL si ottengono i seguenti report.

RTL	Status	Latency			Interval		
		min	avg	max	min	avg	max
VHDL	Pass	15	15	16	15	15	16

Table 75: HLS Loop Pipelining Solution C/RTL Cosimulation Summary

Resource	VHDL
SLICE	204
LUT	311
FF	455
DSP	2
BRAM	0
SRL	0

Timing	VHDL
CP required	10.000
CP achieved post-synthesis	5.745
CP achieved post-implementation	6.512

Table 77: HLS Loop Pipelining Solution Export RTL Final Timing

Table 76: HLS Loop Pipelining Solution Export RTL Resource Usage

Pertanto, importando l'IP in Vivado e impostando un clock constraint pari a 10ns è possibile analizzare i seguenti report di risorse, timing, potenza dinamica ed energia per singola operazione.

```
1 create_clock -period 10.000 -name myclk -waveform {0.000 5.000} [get_ports ap_clk]
```

Si può notare come, rispetto alla soluzione iniziale non ottimizzata, si è ottenuto un aumento di circa il 13% delle LUT e un aumento di circa il 186% dei FF.

LUT	LUTRAM	FF	BRAM	DSP	IO	BUFG
311	0	458	0	2	71	1

Table 78: Vivado Loop Pipelining Solution Utilization Report [#]

Inoltre, si può evidenziare come la maximum clock frequency sia aumentata di circa 15MHz rispetto alla soluzione non ottimizzata di riferimento. Tanto è vero che il WNS, relativo alla soluzione realizzata mediante approccio di pipelining, risulta essere leggermente maggiore.

Cycles [#]	Clock Constraint [ns]	WNS [ns]	Maximum Clock Frequency [Mhz]
16	10	4.208	172.6519337

Table 79: Vivado Loop Pipelining Solution Timing Report

In particolare, si può notare come, rispetto alla soluzione hardware iniziale non ottimizzata, i contributi di potenza dinamica sono aumentati tale da far aumentare l'energia per singola operazione. Nello specifico, si ha un aumento del contributo *Clocks* di circa il 46% dovuto principalmente all'aumento delle risorse associate alla solution in questione.

BRAM	Clock Enable	Clocks	DSP	Logic	Set/Reset [mW]	Data
0	0.196236724	1.778833685	0.814738101	1.275097136	0.012051522	1.206569896

Table 80: Vivado Loop Pipelining Solution Dynamic Power Report [mW]

Dynamic Total
5.283527065

Table 81: Vivado Loop Pipelining Solution Dynamic Power Report [mW]

Energy Single Operation
52.83527065

Table 82: Vivado Loop Pipelining Solution Energy Single Operation Report [pJ]

5.7 Bitwidth Optimization Solution

5.8 AXI Solution

6 Conclusions

Qui di seguito verranno riportati i report relativi a tutte le solution implementate sia in HLS sia in Vivado con i corrispondenti plot.

Synthesis Report										
Solution	Utilization Estimates [#]				Timing Summary [ns]			Latency Summary [#]		
	BRAM	DSP	FF	LUT	Target	Estimated	Uncertainty	Latency	Interval	Loop
Unoptimized	0	4	287	233	10	8.51	1.25	23	45	loop
Code Hoisting	0	4	230	240	10	8.51	1.25	42	42	loop
Loop Fission	0	2	205	214	10	8.51	1.25	66	66	loopShifting
Loop Unrolling (manual) factor=2	2	2	145	236	10	8.51	1.25	56	56	loopAccumulator
Loop Unrolling (pragmas) factor=2	2	2	141	251	10	8.51	1.25	56	56	loopShifting
Loop Unrolling (pragmas+partitioning) factor=2	0	2	912	289	10	8.51	1.25	61	66	loopAccumulator
Loop Unrolling (manual) factor=4	2	2	221	318	10	8.51	1.25	58	58	loopShifting
Loop Unrolling (pragmas) factor=4	2	2	194	367	10	8.51	1.25	62	66	loopAccumulator
Loop Unrolling (pragmas+partitioning) factor=4	0	2	928	748	10	8.51	1.25	62	64	loopShifting
Operation Chaining clk=10ns	0	4	287	233	10	8.51	1.25	23	45	loopAccumulator
Operation Chaining clk=9ns	0	4	619	301	9	6.912	1.12	24	57	loop
Operation Chaining clk=8ns	0	4	619	301	8	6.912	1	24	57	loop
Operation Chaining clk=7ns	0	4	623	305	7	5.745	0.88	25	69	loop
Operation Chaining clk=6ns	0	4	725	221	6	4.644	0.75	27	93	loop
Operation Chaining clk=5ns	0	4	725	221	5	4.321	0.62	27	93	loop
Operation Chaining clk=4ns	0	4	832	286	4	3.254	0.5	40	139	loop
Loop Pipelining	0	4	710	621	10	8.51	1.25	15	15	loop
Bitwidth Optimization	0	2	81	107	10	6.38	1.25	31	31	loop
Axi	0	14	767	583	10	8.742	1.25	?	?	loop 1

Figure 28: HLS Synthesis Report

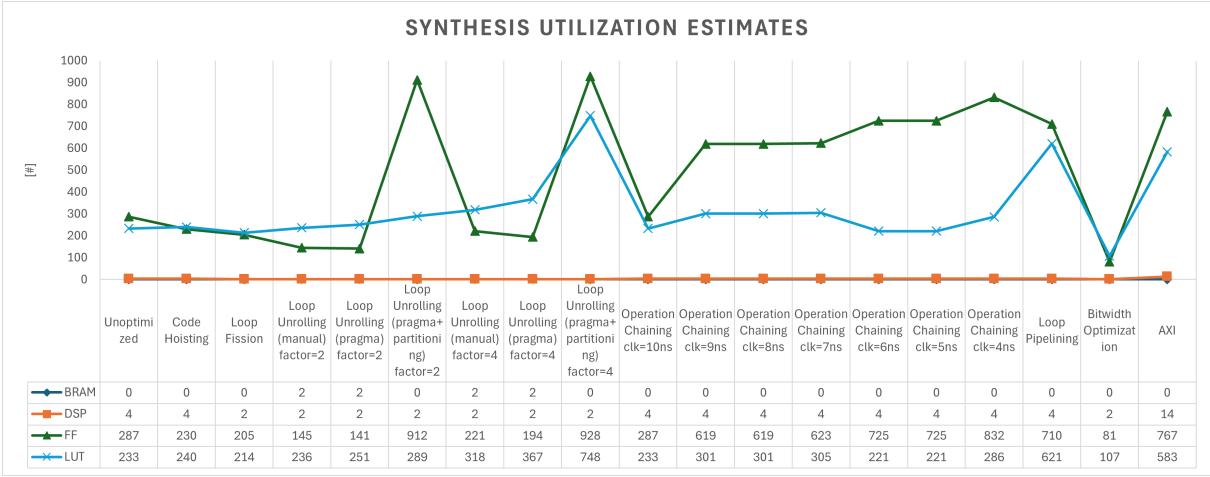


Figure 29: HLS Synthesis Utilization Estimates Report Plot

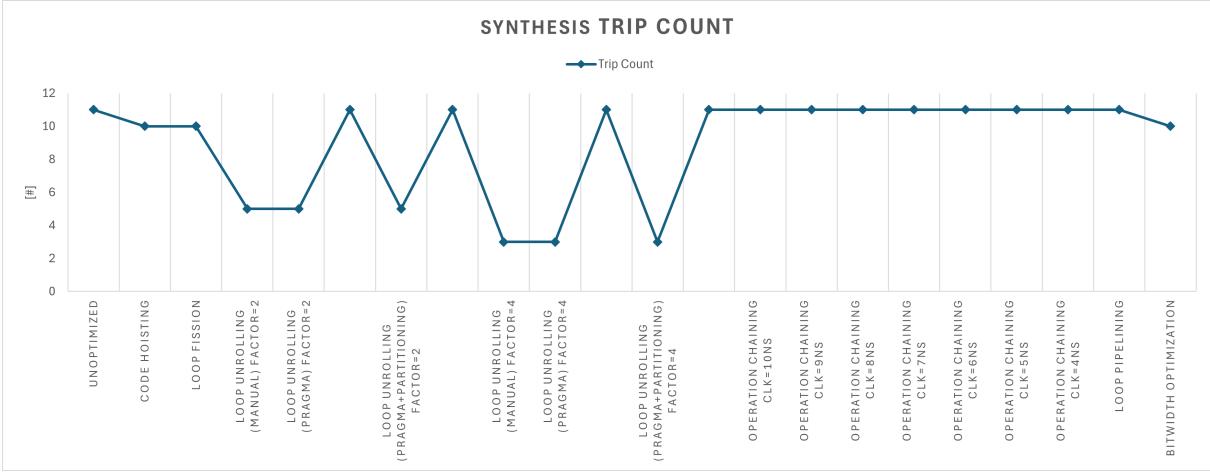


Figure 30: HLS Synthesis Trip Count Report Plot

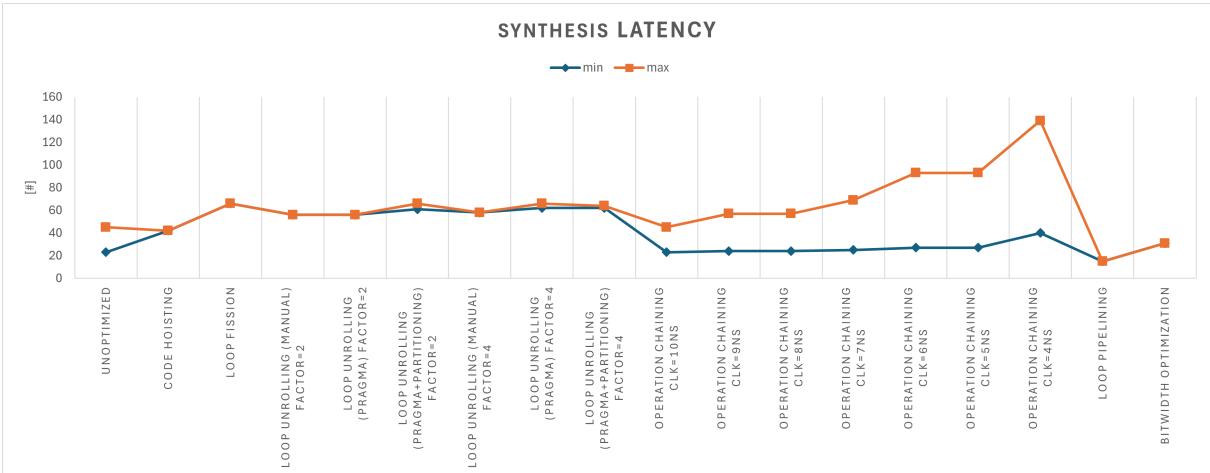


Figure 31: HLS Synthesis Latency Report Plot

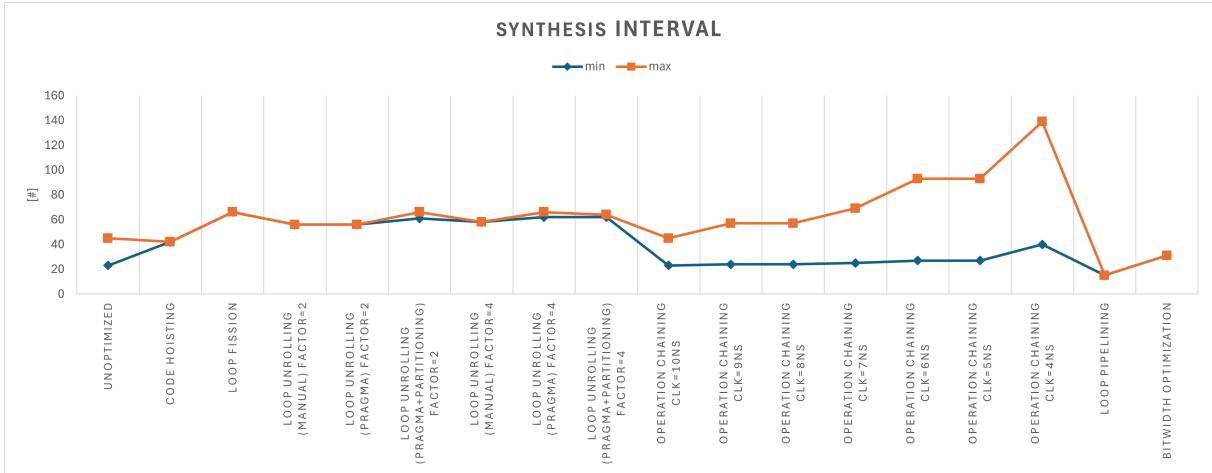


Figure 32: HLS Synthesis Interval Report Plot

Co-Simulation Report						
Solution	Latency			Interval		
	min	avg	max	min	avg	max
Unoptimized	43	43	44	43	43	44
Code Hoisting	42	42	43	42	42	43
Loop Fission	66	66	67	66	66	67
Loop Unrolling (manual) factor=2	56	56	57	56	56	57
Loop Unrolling (pragma) factor=2	56	56	57	56	56	57
Loop Unrolling (pragma+partitioning) factor=2	65	65	66	65	65	66
Loop Unrolling (manual) factor=4	58	58	59	58	58	59
Loop Unrolling (pragma+partitioning) factor=4	60	60	61	60	60	61
Loop Unrolling (manual) factor=4	58	58	59	58	58	59
Operation Chaining clk=10ns	43	43	44	43	43	44
Operation Chaining clk=9ns	54	54	55	54	54	55
Operation Chaining clk=8ns	54	54	55	54	54	55
Operation Chaining clk=7ns	65	65	66	65	65	66
Operation Chaining clk=6ns	87	87	88	87	87	88
Operation Chaining clk=5ns	87	87	88	87	87	88
Operation Chaining clk=4ns	130	130	131	130	130	131
Loop Pipelining	15	15	16	15	15	16
Bitwidth Optimization	31	31	32	31	31	32
AXI	13	13	13	NA	NA	NA

Figure 33: HLS C/RTL Cosimulation Report



Figure 34: HLS C/RTL Cosimulation Report Report Plot

Export RTL Report								
Solution	Resource Usage					Final Timing		
	SLICE	LUT	FF	DSP	BRAM	SRL	CP required	CP achieved post synthesis
Unoptimized	81	275	160	2	0	0	10	5.745
Code Hoisting	75	270	131	2	0	0	10	5.745
Loop Fission	44	157	106	2	0	0	10	5.745
Loop Unrolling (manual) factor=2	31	97	72	2	2	0	10	5.692
Loop Unrolling (pragma) factor=2	29	97	72	2	2	0	10	5.745
Loop Unrolling (pragma+partitioning) factor=2	294	413	843	2	0	0	10	5.745
Loop Unrolling (manual) factor=4	46	160	147	2	2	0	10	5.745
Loop Unrolling (pragma) factor=4	40	145	93	2	2	0	10	5.745
Loop Unrolling (pragma+partitioning) factor=4	412	1145	864	2	0	0	10	5.745
Operation Chaining clk=10ns	81	275	160	2	0	0	10	5.745
Operation Chaining clk=9ns	88	276	226	2	0	0	9	5.745
Operation Chaining clk=8ns	89	276	226	2	0	0	8	5.745
Operation Chaining clk=7ns	75	186	222	4	0	0	7	5.745
Operation Chaining clk=6ns	76	186	241	4	0	0	6	4.823
Operation Chaining clk=5ns	96	186	258	4	0	0	5	4.823
Operation Chaining clk=4ns	114	188	376	4	0	0	4	3.667
Loop Pipelining	204	311	455	2	0	0	10	5.745
Bitwidth Optimization	15	50	24	2	0	0	10	5.429
AXI	231	739	699	0	0	0	10	7.269
								8.202

Figure 35: HLS Export RTL Report

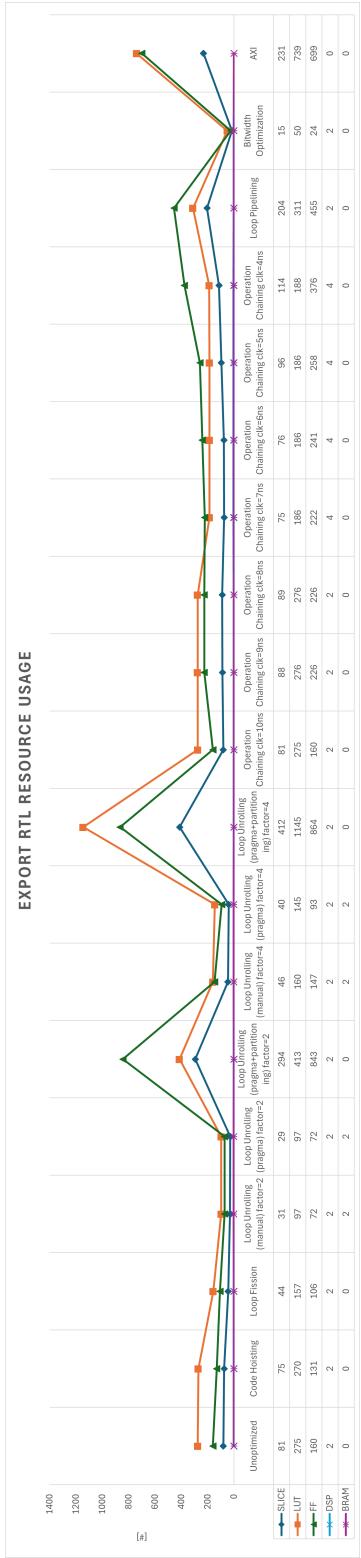


Figure 36: HLS Export RTL Report Report Plot