

# POLITECNICO DI TORINO

---

Master's Degree in Computer Engineering

Master's Degree Thesis

## Acceleration by Separate-Process Cache for Memory-Intensive Algorithms on FPGA via High-Level Synthesis



**Supervisor**  
Prof. Luciano Lavagno

**Candidate**  
Giovanni Brignone  
ID: 274148

Academic year 2020-2021



# Abstract

The end of the Moore's Law validity is making the performance advance of software run on general purpose processors more challenging than ever. Since current technology cannot scale anymore it is necessary to approach the problem from a different point of view: application-specific hardware can provide higher performance and lower power consumption, while requiring higher design efforts and higher deployment costs.

High-Level Synthesis is aimed at reducing design efforts, shrinking the gap between hardware and software design.

FPGAs allow to reduce deployment costs, making possible to implement special-purpose hardware modules on general-purpose underlying architecture.

FPGAs memory system is composed of three main kind of resources: registers, Block-RAMs and external DRAMs. Current HLS tools allow to exploit this memory hierarchy manually, in a scratchpad-like fashion: the objective of this thesis work is to automate the memory management by providing a easily integrable and fully customizable cache system for High-Level Synthesis.

# Contents

<b>1</b>	<b>Background</b>	<b>1</b>
1.1	Cache memory . . . . .	1
1.2	High-Level Synthesis . . . . .	3
1.3	Previous work . . . . .	3
<b>2</b>	<b>Thesis contribution</b>	<b>5</b>
2.1	Architecture . . . . .	5
2.2	Implementation . . . . .	5
2.3	Results . . . . .	5
<b>3</b>	<b>Conclusion</b>	<b>7</b>

# Chapter 1

## Background

### 1.1 Cache memory

Memory devices are usually the performance bottleneck in the execution of memory-bound algorithms. The ideal memory should be fast, large and cheap, but current technology forces the designer to choose a tradeoff between the metrics.

A common solution to this problem is to setup a memory hierarchy in which fast but small memories are paired with large but slow memories, which allows to get good performance on average while containing costs.

This hierarchy can be managed by two main approaches:

- *Scratchpad*: different memories belongs to different addressing spaces: the user is in charge of manually choosing what memory to access: this approach allows to optimally exploit the hierarchy at the cost of high design effort.
- *Cache*: different memories belongs to the same addressing space: the system automatically uses the whole hierarchy exploiting spatial locality (accessed data is likely physically close to previously accessed data) and temporal locality (accessed data has likely recently been accessed), which are typical of many algorithms.

### Structure

A cache memory is logically split into *sets* containing *lines* (or *ways*) which are in turn made up of *words*, as shown in Figure 1.1.

Whenever a word  $w$  is requested there are two possibilities:

- *Hit*:  $w$  is present in the cache: the request can be immediately fulfilled.
- *Miss*:  $w$  is not present in the cache: it is necessary to retrieve it from lower level memory before fulfilling the request.

During the data retrieving a cache line is filled with a block of contiguous words loaded from the lower level memory, trying to exploit spatial locality of future accesses, while

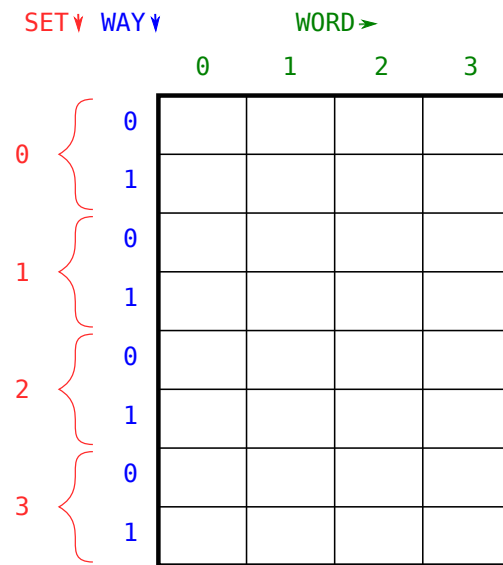


Figure 1.1: Cache logic structure.

mapping policies and replacement policies determine which cache line to overwrite, trying to exploit temporal locality.

If the cache memory is writable, data consistency is ensured by a consistency policy.

## Policies

### Mapping policy

The mapping policy is in charge of statically associating a lower level memory line to a cache set.

The *set associative* policy is the most common mapping policy: given a cache memory with  $s$  sets of  $w$  words, the word address (referred to the lower level memory) bits are split into three parts (as shown in Figure 1.2):

1.  $\log_2(w)$ : offset of the word in the line.
2.  $\log_2(s)$ : set.
3. remaining MSBs: tag identifying the specific line.

Special cases of this policy are:

- *Direct mapped* policy: each set is composed of a single line: the set bits identify a specific cache line, therefore there is no need for a replacement policy.



Figure 1.2: Set associative policy address bits meaning.

- *Fully associative* policy: there is only a single set, therefore the line is fully determined by the replacement policy.

### Replacement policy

The replacement policy is in charge of dynamically associating a lower level memory line to a cache line of a set.

Multiple solutions of this problem have been developed, trying to maximize the temporal locality exploitation. Among the most commonly used solutions there are:

- *First-In First-Out*: the line to be replaced is the first one that has been inserted to the cache.
- *Least recently used*: the line to be replaced is the one that has least recently been accessed.

### Consistency policy

The consistency policy is in charge of ensuring data consistency between memories belonging to different hierarchy levels.

The most common solutions to this problem are:

- *Write-back*: write accesses are performed to the highest level memory and lower level memories are updated when the cache line is replaced only.
- *Write-through*: each write access is propagated along the whole hierarchy.

## 1.2 High-Level Synthesis

### Multi-process modeling

### Process modeling

### Communication modeling

## 1.3 Previous work





## Chapter 2

# Thesis contribution

### 2.1 Architecture

L2 cache architecture

L1 cache architecture

Multi-port architecture

### 2.2 Implementation

### 2.3 Results



## Chapter 3

## Conclusion