# Cache for HLS

## A multi-process architecture

Brignone Giovanni

June 7, 2021

**Politecnico di Torino**

1859

Introduction
Architecture
Implementation
Results
Future work

Motivation

Introduction
Architecture
Implementation
Results
Future work

Motivation

## Motivation

- **Problem**:
  big arrays are mapped to <u>DRAM</u> $\Rightarrow$ performance <u>bottlenecks</u>
- **Proposed solution**:
  cache module which should:
  - store data to <u>BRAMs</u>
  - require <u>minimal effort</u> to be integrated into any HLS design

Introduction
**Architecture**
Implementation
Results
Future work

Inlined architecture
Multi-process architecture

**1** Introduction

**2** Architecture
Inlined architecture
Multi-process architecture

**3** Implementation

**4** Results

**5** Future work

Introduction
Architecture
Implementation
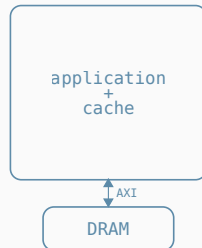Results
Future work

Inlined architecture
Multi-process architecture

# Inlined architecture (Ma Liang)

Array access $\rightarrow$ Inlined cache logic

- straightforward architecture
- cache logic mixed with application logic
- single-port only

Introduction
Architecture
Implementation
Results
Future work

Inlined architecture
**Multi-process architecture**

# Multi-process architecture

Array access $\rightarrow$ Request to separate module

- decoupling between application and cache logic (communication through FIFOs)
- may support multiple ports (work in progress)

Introduction
Architecture
**Implementation**
Results
Future work

**Tools**
Internal architecture
Problems and solutions

**1** Introduction

**2** Architecture

**3** Implementation
   Tools
   Internal architecture
   Problems and solutions

**4** Results

**5** Future work

Introduction
Architecture
**Implementation**
Results
Future work

**Tools**
Internal architecture
Problems and solutions

## Tools

C++ code compatible with *Vitis HLS* (2020.2).

Introduction
Architecture
**Implementation**
Results
Future work

**Tools**
Internal architecture
Problems and solutions

## Tools

C++ code compatible with *Vitis HLS* (2020.2).

- **Multiple processes modeling**:
  infinite loops parallelized by:
  - <u>SW simulation:</u> std::thread
  - <u>Synthesis:</u> DATAFLOW with *start propagation* disabled

Introduction
Architecture
**Implementation**
Results
Future work

**Tools**
Internal architecture
Problems and solutions

## Tools

C++ code compatible with *Vitis HLS* (2020.2).

- **Multiple processes modeling**:
  infinite loops parallelized by:
  - <u>SW simulation:</u> `std::thread`
  - <u>Synthesis:</u> `DATAFLOW` with *start propagation* disabled
- **Inter-process communication**: `hls::stream` FIFOs

Introduction
Architecture
**Implementation**
Results
Future work

**Tools**
Internal architecture
Problems and solutions

## Tools

C++ code compatible with *Vitis HLS* (2020.2).

- **Multiple processes modeling**:
  infinite loops parallelized by:
    - <u>SW simulation:</u> `std::thread`
    - <u>Synthesis:</u> `DATAFLOW` with *start propagation* disabled
- **Inter-process communication**: `hls::stream` FIFOs
- **Performance optimization**:
    - loop pipelining: new request each cycle
    - automatic port widening: access one line at a time in DRAM

Introduction
Architecture
**Implementation**
Results
Future work

**Tools**
Internal architecture
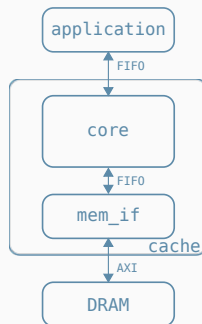Problems and solutions

# Tools

C++ code compatible with *Vitis HLS* (2020.2).

- **Multiple processes modeling**:
  infinite loops parallelized by:
    - <u>SW simulation:</u> `std::thread`
    - <u>Synthesis:</u> `DATAFLOW` with *start propagation* disabled
- **Inter-process communication**: `hls::stream` FIFOs
- **Performance optimization**:
    - loop pipelining: new request each cycle
    - automatic port widening: access one line at a time in DRAM
- **Customization:** cache characteristics set through `template`

Introduction
Architecture
**Implementation**
Results
Future work

**Tools**
Internal architecture
Problems and solutions

## Tools

C++ code compatible with *Vitis HLS* (2020.2).

- **Multiple processes modeling**:
  infinite loops parallelized by:
  - <u>SW simulation</u>: `std::thread`
  - <u>Synthesis</u>: `DATAFLOW` with *start propagation* disabled
- **Inter-process communication**: `hls::stream` FIFOs
- **Performance optimization**:
  - loop pipelining: new request each cycle
  - automatic port widening: access one line at a time in DRAM
- **Customization:** cache characteristics set through `template`
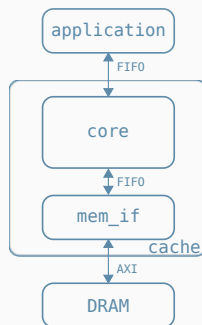- **Limitations**: cannot override "`[] operator`" for `set` due to automatic `class` disaggregation

Introduction
Architecture
**Implementation**
Results
Future work

Tools
**Internal architecture**
Problems and solutions

## Internal architecture

- **Problem**: persuade HLS to schedule response writing to FIFO in case of HIT early in the pipeline

Introduction
Architecture
**Implementation**
Results
Future work

Tools
**Internal architecture**
Problems and solutions

## Internal architecture

- **Problem**: persuade HLS to schedule response writing to FIFO in case of HIT early in the pipeline
- **Proposed solution**: split the cache into two processes:
  1. core:
     - manage requests from application
     - keep cache data structures up-to-date
  2. mem_if:
     - manage DRAM accesses

$\Rightarrow$ synthesizer job is simplified:
if HIT, request is managed in 1 cycle

Introduction
Architecture
**Implementation**
Results
Future work

Tools
**Internal architecture**
Problems and solutions

# mem_if implementation
Functionality

Manage DRAM accesses:

❶ read request from core

❷ access DRAM

❸ write response to core (if read request)

Introduction
Architecture
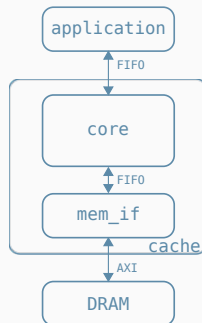**Implementation**
Results
Future work

Tools
**Internal architecture**
Problems and solutions

# core implementation

Functionality

❶ read request from application

Introduction
Architecture
**Implementation**
Results
Future work

Tools
Internal architecture
Problems and solutions

# core implementation
Functionality

**1** read request from application
**2** check if it is an HIT or a MISS

Introduction
Architecture
**Implementation**
Results
Future work

Tools
**Internal architecture**
Problems and solutions

## core implementation
Functionality

1. read request from application
2. check if it is an HIT or a MISS
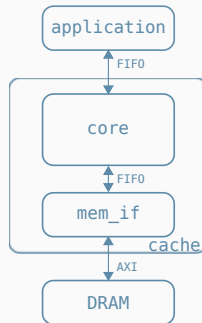3. if MISS:

```
        application
            ↕ FIFO

            core

            ↕ FIFO

          mem_if
            ↕        cache
            ↕ AXI

           DRAM
```

Introduction
Architecture
**Implementation**
Results
Future work

Tools
**Internal architecture**
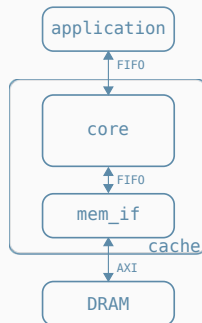Problems and solutions

## core implementation
Functionality

1. read request from `application`
2. check if it is an HIT or a MISS
3. if MISS:
   - if the cache line to be overwritten $line_{old}$ has been modified, issue a `write` request of $line_{old}$ to `mem_if`

Introduction
Architecture
**Implementation**
Results
Future work

Tools
**Internal architecture**
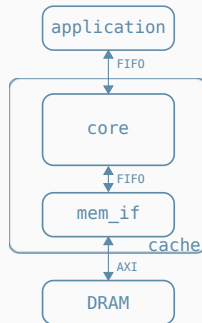Problems and solutions

# core implementation

Functionality

1. read request from application
2. check if it is an HIT or a MISS
3. if MISS:
   - if the cache line to be overwritten $line_{old}$ has been modified, issue a write request of $line_{old}$ to mem_if
   - issue a read request of the requested line to mem_if

```
        application
             ↑ FIFO
    ┌─────────────────┐
    │      core       │
    │                 │
    │        ↑ FIFO   │
    │     mem_if      │
    └──────────── cache┘
             ↑ AXI
          DRAM
```

Introduction
Architecture
**Implementation**
Results
Future work

Tools
**Internal architecture**
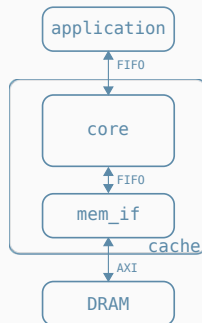Problems and solutions

## core implementation

Functionality

1. read request from application
2. check if it is an HIT or a MISS
3. if MISS:
   - if the cache line to be overwritten $line_{old}$ has been modified, issue a write request of $line_{old}$ to mem_if
   - issue a read request of the requested line to mem_if
   - read mem_if response and update BRAM

Introduction
Architecture
**Implementation**
Results
Future work

Tools
**Internal architecture**
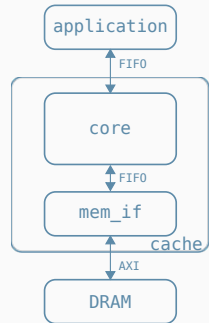Problems and solutions

# core implementation

Functionality

1. read request from application
2. check if it is an HIT or a MISS
3. if MISS:
   - if the cache line to be overwritten $line_{old}$ has been modified, issue a write request of $line_{old}$ to mem_if
   - issue a read request of the requested line to mem_if
   - read mem_if response and update BRAM
4. access BRAM

Introduction
Architecture
**Implementation**
Results
Future work

Tools
**Internal architecture**
Problems and solutions

## core implementation

Functionality

1. read request from application
2. check if it is an HIT or a MISS
3. if MISS:
   - if the cache line to be overwritten $line_{old}$ has been modified, issue a write request of $line_{old}$ to mem_if
   - issue a read request of the requested line to mem_if
   - read mem_if response and update BRAM
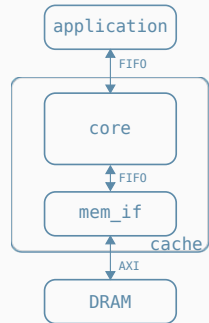4. access BRAM
5. write response to application (if read)

application

FIFO

core

FIFO

mem_if

cache

AXI

DRAM

Introduction
Architecture
**Implementation**
Results
Future work

Tools
Internal architecture
**Problems and solutions**

# BRAM dependencies - RAW `fill`

- **Problem**: reading BRAM immediately after it is written by the `fill` process causes an increase of II

Introduction
Architecture
**Implementation**
Results
Future work

Tools
Internal architecture
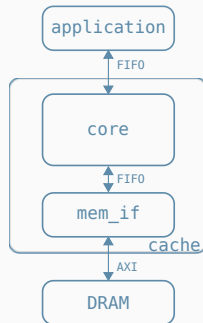**Problems and solutions**

# BRAM dependencies - RAW `fill`

- **Problem**: reading BRAM immediately after it is written by the `fill` process causes an increase of II
- **Proposed solution**: store the `mem_if` response in a buffer which can be immediately accessed and update BRAM afterwards
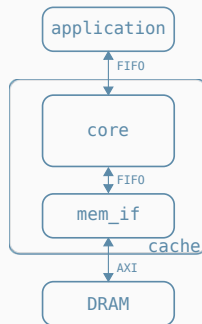
Introduction
Architecture
**Implementation**
Results
Future work

Tools
Internal architecture
**Problems and solutions**

# BRAM dependencies - RAW request

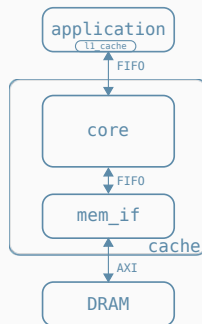- **Problem**: reading BRAM immediately after it is written by a previous write request causes an increase of II

Introduction
Architecture
**Implementation**
Results
Future work

Tools
Internal architecture
**Problems and solutions**

# BRAM dependencies - RAW `request`

- **Problem**: reading BRAM immediately after it is written by a previous `write` request causes an increase of II
- **Proposed solution**: `raw_cache` inside `cache`:
    - `write` request: store the written line to a buffer
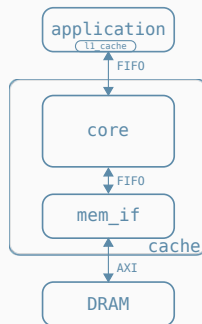    - subsequent `read` request to the same line: access the buffer instead of BRAM

Introduction
Architecture
**Implementation**
Results
Future work

Tools
Internal architecture
**Problems and solutions**

# Request optimization

- **Problem**:
  - each FIFO access costs 1 cycle
  - accesses to arrays are often sequential

Introduction
Architecture
**Implementation**
Results
Future work

Tools
Internal architecture
**Problems and solutions**

## Request optimization

- **Problem**:
    - each FIFO access costs 1 cycle
    - accesses to arrays are often sequential
- **Proposed solution**: l1_cache in the interface (application side):
    - read request: get the whole cache line, store it in a buffer and return the requested element
    - subsequent read request to the same line: access the buffer
    $\Rightarrow$ avoid passing through FIFOs

Introduction
Architecture
Implementation
**Results**
Future work

Performance

**1** Introduction

**2** Architecture

**3** Implementation

**4** Results
   Performance

**5** Future work

Introduction
Architecture
Implementation
**Results**
Future work

Performance

# Performance

- **Loop pipelining**:
  both core and mem_if loops pipelined with II=1
- **L1 HIT**: latency of 1 cycle
- **HIT**: latency of 6 cycles

Introduction
Architecture
Implementation
Results
**Future work**

Multi-port cache

**1** Introduction

**2** Architecture

**3** Implementation

**4** Results

**5** Future work
  Multi-port cache

Introduction
Architecture
Implementation
Results
Future work

Multi-port cache

## Multi-port cache

- **Problem**: cache manages one request per cycle $\Rightarrow$ unrolling a loop which uses the cache is not trivial

Introduction
Architecture
Implementation
Results
Future work

Multi-port cache

# Multi-port cache

- **Problem**: cache manages one request per cycle ⇒ unrolling a loop which uses the cache is not trivial
- **Possible workarounds**: get multiple data elements for each request and do all the unroll in `application`; two possible ways:
    - `get_line`: return a whole cache line
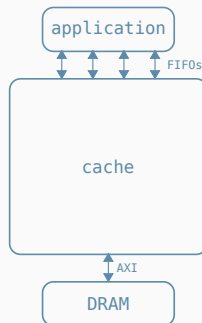    - pack multiple data elements in a single cache element (e.g. `hls::vector`)

Introduction
Architecture
Implementation
Results
Future work

Multi-port cache

# Multi-port cache

- **Problem**: cache manages one request per cycle $\Rightarrow$ unrolling a loop which uses the cache is not trivial
- **Possible workarounds**: get multiple data elements for each request and do all the unroll in application; two possible ways:
    - get_line: return a whole cache line
    - pack multiple data elements in a single cache element (e.g. hls::vector)
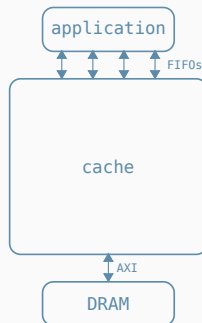- **Proposed solution**: provide $N$ ports, where $N$ is the unroll factor

Introduction
Architecture
Implementation
Results
Future work

Multi-port cache

# Multi-port cache
Implementation

- **Straightforward implementation**:
  - provide $N$ FIFOs on the interface
  - unroll core loop by a factor $N$

Introduction
Architecture
Implementation
Results
Future work
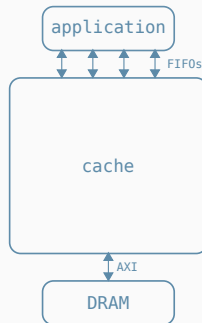
Multi-port cache

# Multi-port cache
Implementation

- **Straightforward implementation**:
  - provide $N$ FIFOs on the interface
  - unroll core loop by a factor $N$
- **Problem**: all the $N$ requests may refer to the same line $\rightarrow$ each iteration must wait completion of previous one

Introduction
Architecture
Implementation
Results
Future work

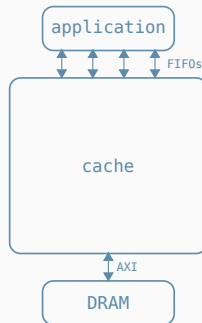Multi-port cache

# Multi-port cache

Implementation

- **Proposed solution**:
    1. read all the $N$ requests
    2. check if they are all HIT
    3. if at least one MISS issue all the requests to mem_if
    4. access BRAM
    5. write all the responses

Introduction
Architecture
Implementation
Results
Future work

Multi-port cache

# Multi-port cache
Implementation

- **Proposed solution**:
  1. read all the $N$ requests
  2. check if they are all HIT
  3. if at least one MISS issue all the requests to mem_if
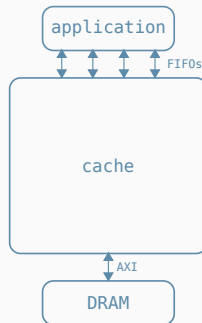  4. access BRAM
  5. write all the responses

Introduction
Architecture
Implementation
Results
Future work

Multi-port cache

# Multi-port cache

Implementation

- **Proposed solution**:
  1. read all the $N$ requests
  2. check if they are all HIT
  3. if at least one MISS issue all the requests to mem_if
  4. access BRAM
  5. write all the responses

  - cache addressing mode converted to an associative one

Introduction
Architecture
Implementation
Results
Future work

Multi-port cache

# Multi-port cache
Implementation

- **Proposed solution**:
    1. read all the $N$ requests
    2. check if they are all HIT
    3. if at least one MISS issue all the requests to mem_if
    4. access BRAM
    5. write all the responses

    - cache addressing mode converted to an associative one
    - possibly implement core at RTL to have full control on scheduling