# Hi-DMM:
# High-Performance Dynamic Memory Management in High-Level Synthesis

- **Tingyuan Liang, Jieru Zhao, Liang Feng, Sharad Sinha, and Wei Zhang**

- **Hong Kong University of Science and Technology (HKUST)**
- **Indian Institute of Technology Goa (IIT Goa)**

# Outline

- Motivation
- Overview of Hi-DMM
- Implementation of Software
- Implementation of Hardware
- Evaluation of Hi-DMM
- Open-Source Hi-DMM Platform
- Conclusion

# Outline

- Motivation
- Overview of Hi-DMM
- Implementation of Software
- Implementation of Hardware
- Evaluation of Hi-DMM
- Open-Source Hi-DMM Platform
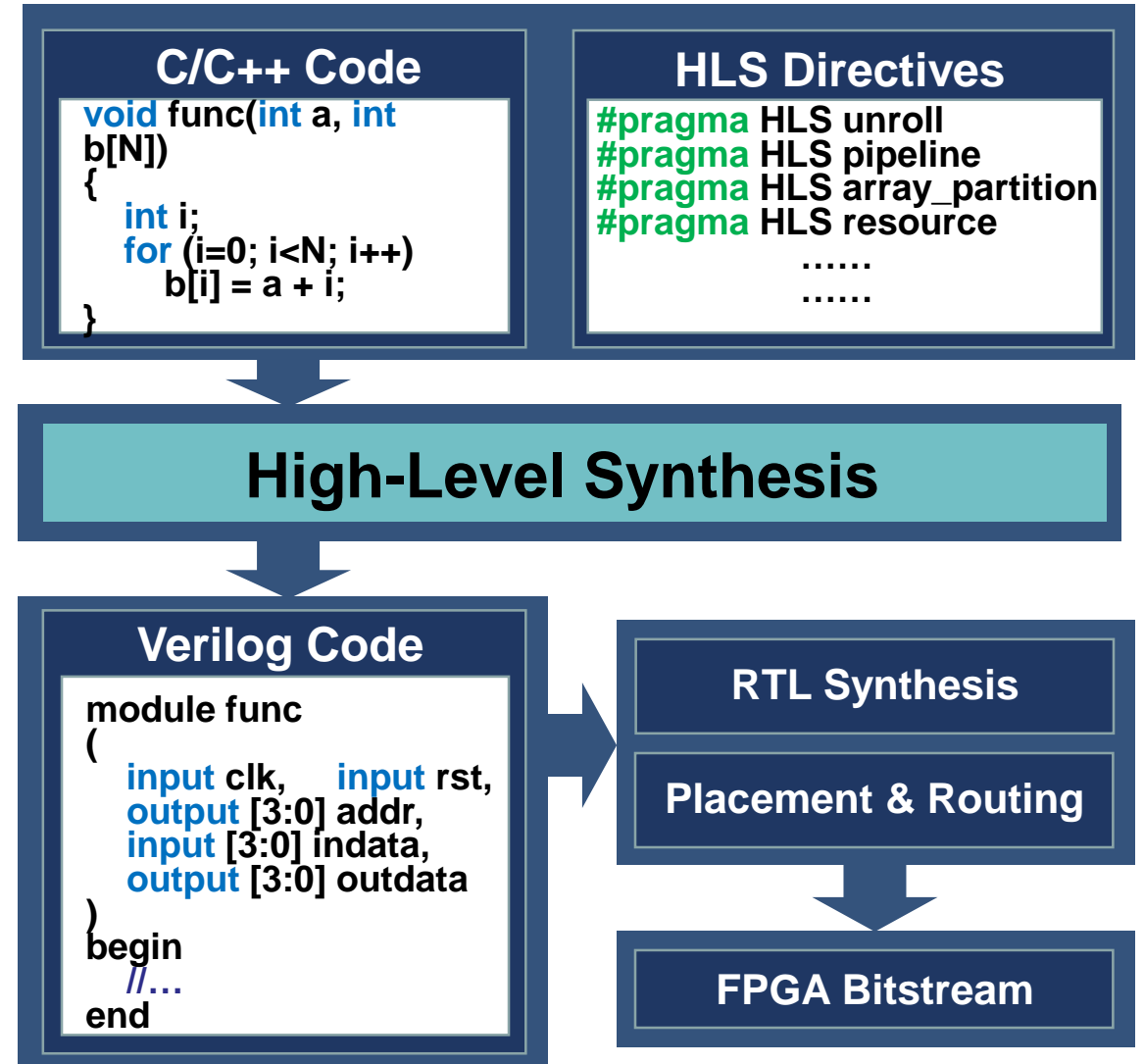- Conclusion

# Motivation: High-Level Synthesis

1. Describe Hardware at High Level

   e.g. from C/C++ to Verilog

2. Fast Development of FPGA design

3. Friendly to Complex Applications

**C/C++ Code**

```
void func(int a, int b[N])
{
    int i;
    for (i=0; i<N; i++)
        b[i] = a + i;
}
```

**HLS Directives**

```
#pragma HLS unroll
#pragma HLS pipeline
#pragma HLS array_partition
#pragma HLS resource
......
......
```

**High-Level Synthesis**

**Verilog Code**

```
module func
(
    input clk,     input rst,
    output [3:0] addr,
    input [3:0] indata,
    output [3:0] outdata
)
begin
    //...
end
```

**RTL Synthesis**

**Placement & Routing**

**FPGA Bitstream**

# Motivation: Dynamic Memory Management

## 1. Feature of High-Level Language

e.g. malloc(), free(), new, delete

## 2. Flexible and Efficient

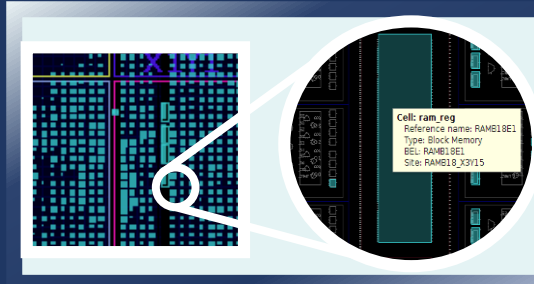make full use of memory

## 3. Unsupported by current HLS

if DMM in HLS realized,
the utilization of BRAMs will be raised



**C/C++ Code with DMM**

```
int *a = (int *) malloc(n * sizeof(int));
...
free(a);
```
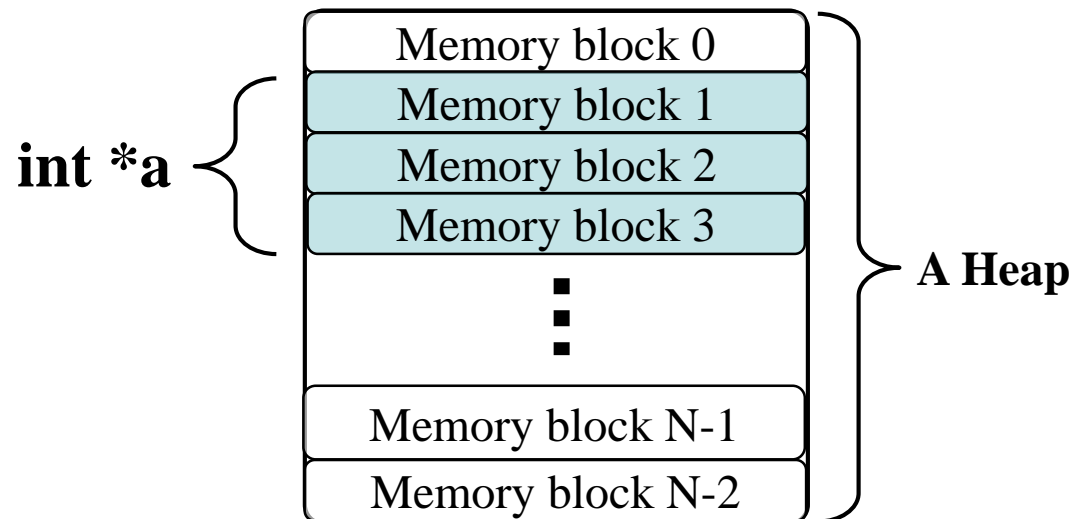
**HLS Failure**

**Block-RAM on FPGA**

Cell: ram_reg
Reference name: RAMB18E1
Type: Block Memory
BEL: RAMB18E1
Site: RAMB18_X3Y15

**High Performance**

**Resource Constraint**

# Motivation: Challenges for DMM in HLS

**int \*a    =    (int \*)    malloc    (n \* sizeof(int));**

**Where to store** dynamic data on FPGA?

**Heaps (Block-RAMs on FPGA)**

int \*a

| |
|---|
| Memory block 0 |
| Memory block 1 |
| Memory block 2 |
| Memory block 3 |
| ⋮ |
| Memory block N-1 |
| Memory block N-2 |

**A Heap**

# Motivation: Challenges for DMM in HLS

int *a   =   (int *)   malloc   (n * sizeof(int));

Who will manage the dynamic memory on FPGA?

Record which blocks are used.

Allocate memory according to required size

**Hardware Allocator**

1. Low Latency
2. High Utilization of BRAM
3. Large Management Capability with Low Overhead of Area

# Motivation: Challenges for DMM in HLS

$$\text{int *a} \quad = \quad \text{(int *)} \quad \text{malloc} \quad \text{(n * sizeof(int));}$$

**Who will manage** the dynamic memory on FPGA?

**Hardware Allocator**
1. Low Latency
2. High Utilization of BRAM
3. Large Management Capability with Low Overhead of Area

SysAlloc / DMM-HLS

Hundreds of cycles?

Fixed-size Allocation?
DOMMU

AMMU

High overhead?

# Motivation: Challenges for DMM in HLS

**int *a** = (int *) **malloc** (n * sizeof(int));

Who will manage the dynamic memory on FPGA?

Hardware Allocator

1. Low Latency
2. High Utilization of BRAM
3. Large Management Capability with Low Overhead of Area

Hi-DMM Allocators ✓

**int *a    =    (int *)    malloc    (n * sizeof(int));**

How can it **be synthesized** into a **high-performance** design with **HLS features**?

Previous Works

Automatic Transformation? ✗
Resource Mapping? ✗
HLS Directive Compatibility? ✗

**Source Code Compiler**

1. **Source Code Transformation**
2. **Resource Mapping**
3. **Adaption to HLS Directives**

# Motivation:  Challenges for DMM in HLS

int *a    =    (int *)    malloc    (n * sizeof(int));

How can it **be synthesized** into a **high-performance** design with **HLS features**?

**Source Code Compiler**

1.  Source Code Transformation
2.  Element Mapping
**Hi-DMM Compiler** ✓
3.  Adaption to HLS Directives

# Motivation: Hi-DMM

int *a    =    (int *)    malloc   (n * sizeof(int));

Who will manage the dynamic memory on FPGA?

How can it be synthesized into a high-performance design with HLS features?

**Hardware Allocator**

Hi-DMM Allocators ✓

**Source Code Compiler**

Hi-DMM Compiler ✓

Hi-DMM

# Outline

- Motivation
- **Overview of Hi-DMM**
- Implementation of Software
- Implementation of Hardware
- Evaluation of Hi-DMM
- Open-Source Hi-DMM Platform
- Conclusion

# Overview:  Workflow of Hi-DMM

**Compilation-time**

**Run-time**

**Source Code of Accelerator with DMM**

```
int *a = (int *) malloc(n * sizeof(int));
…
free(a);
```
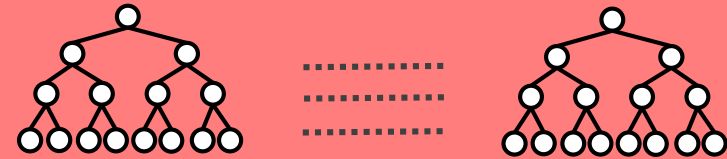
↓

**Hi-DMM Compiler**

↓

**Source Code Compatible with HLS**

```
int offset_a = Hi_malloc(n, allocator1);
…
Hi_free(offset_a);
```

↓

**High-Level Synthesis**

**Hi-DMM Allocator IP Blocks Basde on Buddy Tree**

Heap 0

Heap M

Computation Logic

Static Memory

**Accelerator IP Block integrated with Heaps**

# Overview: Workflow of Hi-DMM

## Compilation-time

**Source Code of Accelerator with DMM**

```
int *a = (int *) malloc(n * sizeof(int));
…
free(a);
```
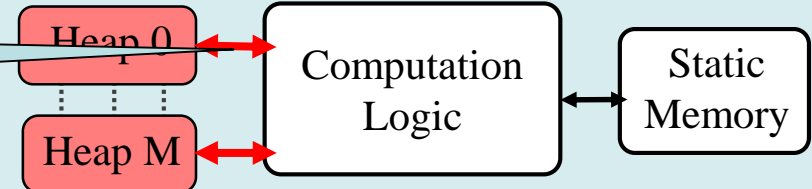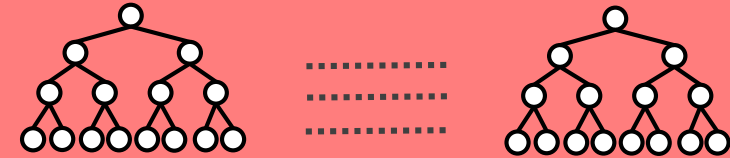
↓

**Hi-DMM Compiler**

↓

**Source Code Compatible with HLS**

```
int offset_a = Hi_malloc(n, allocator1);
…
Hi_free(offset_a);
```

↓

**High-Level Synthesis**

## Run-time

Un-synthesizable ✗

Based on Clang

Synthesizable ✓

Static Memory

**Accelerator IP Block integrated with Heaps**

# Overview: Workflow of Hi-DMM

**Compilation-time**

**Run-time**

**Source Code of Accelerator with DMM**

int *a = (int *) malloc(n * sizeof(int));
...

Both of the accelerator and the allocators are **described in C** and synthesized by Vivado HLS.

**Source Code Compatible with C**

int offset_a = Hi_malloc(n, allocator1);
...
Hi_free(offset_a);

**High-Level Synthesis**

**Hi-DMM Allocator IP Blocks Basde on Buddy Tree**

Heap 0

Heap M

Computation Logic

Static Memory

**Accelerator IP Block integrated with Heaps**

**Compilation-time**

**Run-time**

**Source Code of Accelerator with DMM**

```
int *a = (int *) malloc(n * sizeof(int));
...
free(a);
```

Accelerator can **send request** to the allocators and **get an available address**.

**Source Code Compatible with HLS**

```
int offset_a = Hi_malloc(n, allocator1);
...
Hi_free(offset_a);
```

**High-Level Synthesis**

**Hi-DMM Allocator IP Blocks Basde on Buddy Tree**

Heap 0

Heap M

Computation Logic

Static Memory

**Accelerator IP Block integrated with Heaps**

# Overview: Workflow of Hi-DMM

**Compilation-time**

**Run-time**

**Source Code of Accelerator with DMM**

```
int *a = (int *) malloc(n * sizeof(int));
…
free(a);
```

**Hi-DMM Compiler**

**Source Code Compatible with HLS**

Then the accelerator can **access the heaps** with the address directly.

**High-Level Synthesis**

**Hi-DMM Allocator IP Blocks Basde on Buddy Tree**

Heap 0

Heap M

Computation Logic

Static Memory

**Accelerator IP Block integrated with Heaps**

# Outline

- Motivation
- Overview of Hi-DMM
- **Implementation of Software**
- Implementation of Hardware
- Evaluation of Hi-DMM
- Open-Source Hi-DMM Platform
- Conclusion

# Software: Hi-DMM Compiler

Detection

↓

Analysis
&
Resource Mapping

↓

Transformation
&
Adaption

# Software: Hi-DMM Compiler

Detection

Analysis & Resource Mapping

Transformation & Adaption

---

**Definition of Pointers**

| | |
|---|---|
| int *a, *b, c;<br>ap_unint<13> *e, f, *g; | float *x, *y, z;<br>User_Struct *h, *k; |

1. Name
2. Type
3. Width

---

**Allocation Function Call**

a = (int*) malloc(n*sizeof(int));
e = (ap_unint<13> *) malloc(100*sizeof(ap_unint<13> ));
h = (User_Struct *) malloc(sizeof(User_Struct ));

4. Allocation Requester
5. Granularity

---

**Access to Pointers**

| | |
|---|---|
| a[i] = b[j] + c;<br>e = g; | h->next = k;<br>h->val = 123; |

6. Type of Access
7. Dependencies

---

**HLS Directives**

#pragma HLS array_partition variable=xxx factor=xxx
#pragma HLS unroll factor=xxx

8. Directive for Pointers
9. Directive involving Pointers

# Software: Hi-DMM Compiler

Detection

Analysis
&
Resource Mapping

Transformation
&
Adaption

Pointers

A
B
C
D
E

①Mapping?

Heaps

Heap 0

Heap 1

② Depth of heap?

Allocators

③ Allocator ?

# Software: Hi-DMM Compiler



Pointers

Heaps

Allocators

A
B
C
D
E

Heap 0

①Mapping?

Heap 1

Detection

Analysis
&
Resource Mapping

Transformation
&
Adaption

1. **Map Pointers to heaps**

Co-operations
between pointers:

**A**[i] = **B**[j] + C[i];
…
D[k] = **A**[i] + **B**[m];
…
**2 co-operations**
between A and B
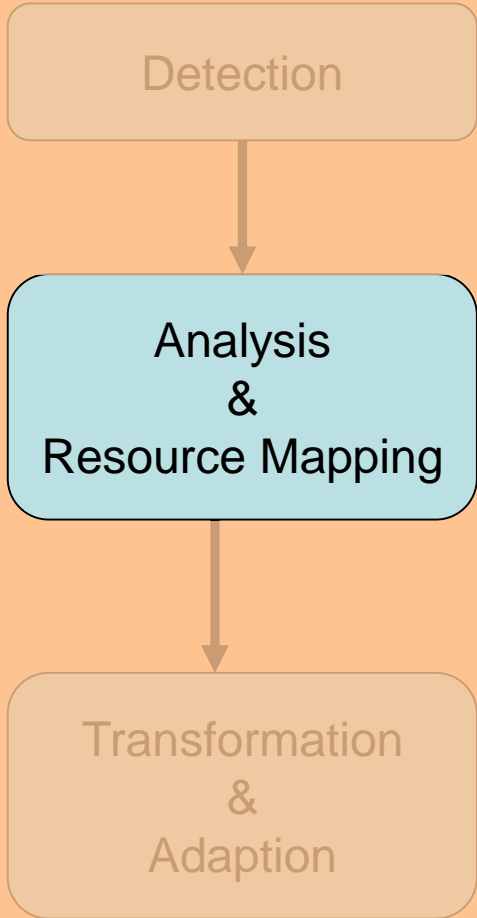
Fully-Connected
Graph

**2**

A
B
E
C
D

Karger
Algorithm

Max-Cuts

Separated
Sub-graphs

A
B
E
C
D

Max-Cut

Pointers                    Heaps                    Allocators

Detection

Analysis
&
Resource Mapping

Transformation
&
Adaption

A
B
C
D
E

①Mapping

Heap 0

Heap 1

1. **Map Pointers to heaps**

Karger
Algorithm

Max-Cuts

Separated
Sub-graphs

Max-Cut

A
B
E
C
D

Pointer
Mapping

Heap 0    Heap 1

# Software: Hi-DMM Compiler

Detection

Analysis
&
Resource Mapping

Transformation
&
Adaption

Pointers

Heaps

Allocators

A
B
C
D
E

Heap 0

Heap 1

①Mapping ✓

② Depth of heap ✓

2. Assign BRAM resource to heaps

Pointer
Mapping

Heap 0
Heap 1

Resource
Mapping

Heap 0

Heap 1

# Software: Hi-DMM Compiler

Detection

Analysis
&
Resource Mapping

Transformation
&
Adaption

Pointers

A
B
C
D
E

① **Mapping** ✓

Heaps

Heap 0

Heap 1

② **Depth of heap** ✓

Allocators

**K-way Tree**

**Hybrid Tree**

③ **Allocator** ✓

3. **Map heaps to allocators**

Depths of Heaps

Heap 0

Heap 1

Type
Width
Granularity
.........

Map Heaps to Allocators

Heap 0 — **K-way Tree**

Heap 1 — **Hybrid Tree**

# Software: Hi-DMM Compiler

## 1. Transformation

**Detection**

↓

**Analysis & Resource Mapping**

↓

**Transformation & Adaption**

### Definition of DMM Heaps

```
int Hi_DMM_Heap_0[8192];
ap_uint<7> Hi_DMM_Heap_1[2048];
#pragma HLS array_partition variable=Hi_DMM_Heap_1 cyclic  factor=4
```

### DMM Interface

```
void TOP(hidmm_alloc_port *Hi_DMM_allocator_0)
{
#pragma HLS interface ap_hs port = Hi_DMM_allocator_0
    ....
}
```

*Synthesizable* ✓

### Function Calls

```
offset_local_dis = HLS_malloc<8192>(n, Hi_DMM_allocator_1_Super_HTA8k;
```

### Accesses to "Struct" Pointers

```
//head->VAL = data[0];
head[OFFSET_LIST_VAL] = data[0];
```

### Assignment of Pointers (from one to another one)

```
// now = tail;
offset_now = offset_tail;
```

Detection

Analysis
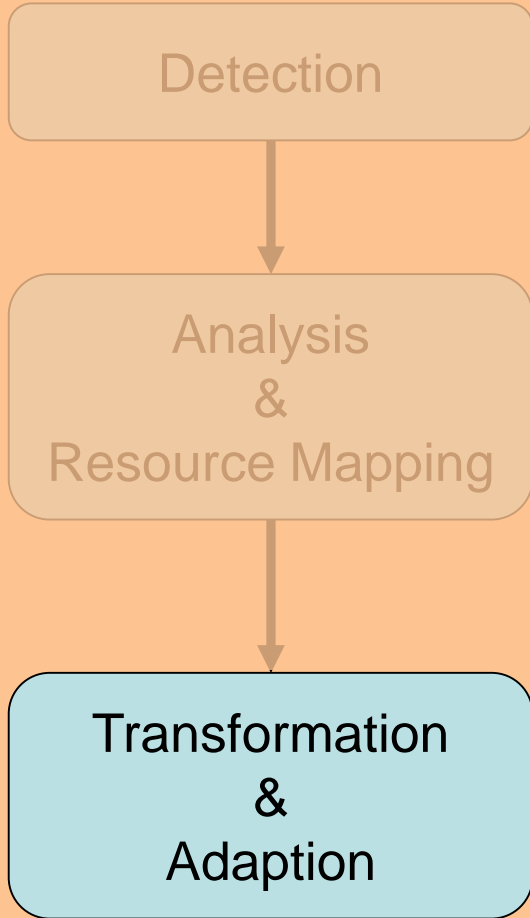&
Resource Mapping

Transformation
&
Adaption

## 2. Adaption to HLS Directives

e.g. Loop Transformation for Loop Unrolling

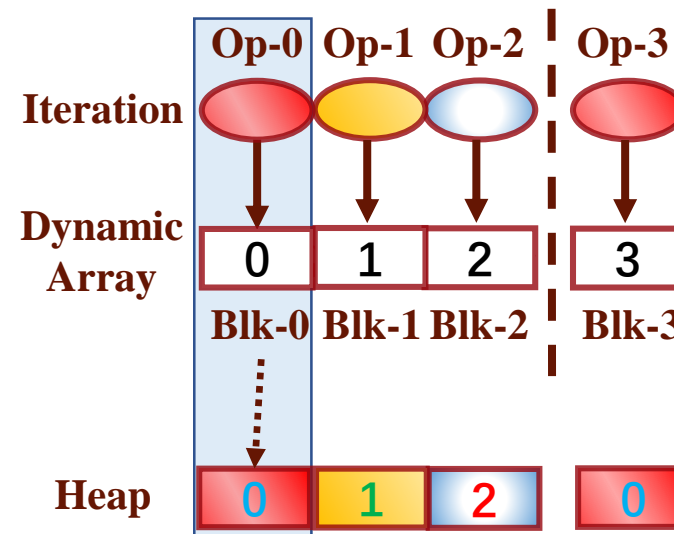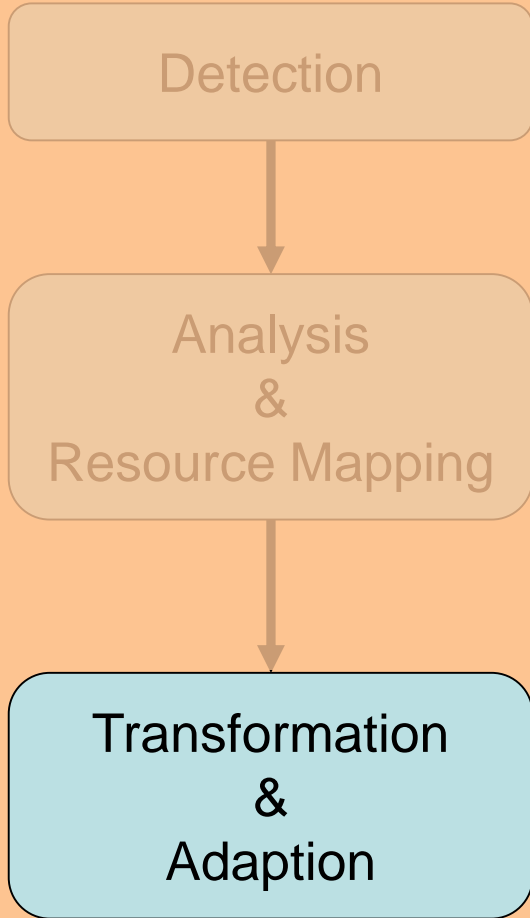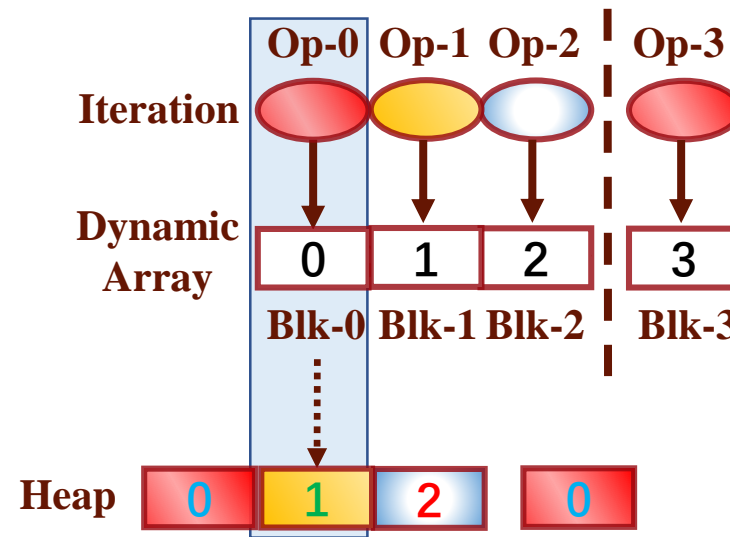**Example without DMM: Operations mapped to corresponding partitions**



Op-0  Op-1  Op-2    Op-3

Iteration

Static
Array     0    1    2     0

Blk-0 Blk-1 Blk-2  Blk-3

**Loop Unrolling** ✓

Detection

Analysis
&
Resource Mapping

Transformation
&
Adaption

## 2. Adaption to HLS Directives

**e.g. Loop Transformation for Loop Unrolling**

**Example with DMM: Operations mapped to all partitions**

| | Op-0 | Op-1 | Op-2 | Op-3 |
|---|---|---|---|---|

Iteration

Dynamic Array: | 0 | 1 | 2 | | 3 |

Blk-0  Blk-1  Blk-2   Blk-3

**Loop Unrolling** ✗

Heap: | 0 | 1 | 2 | | 0 |

**Location of dynamic array is unknown.**

29

Detection

Analysis
&
Resource Mapping

Transformation
&
Adaption

## 2. Adaption to HLS Directives

**e.g. Loop Transformation for Loop Unrolling**

**Example with DMM: Operations mapped to all partitions**



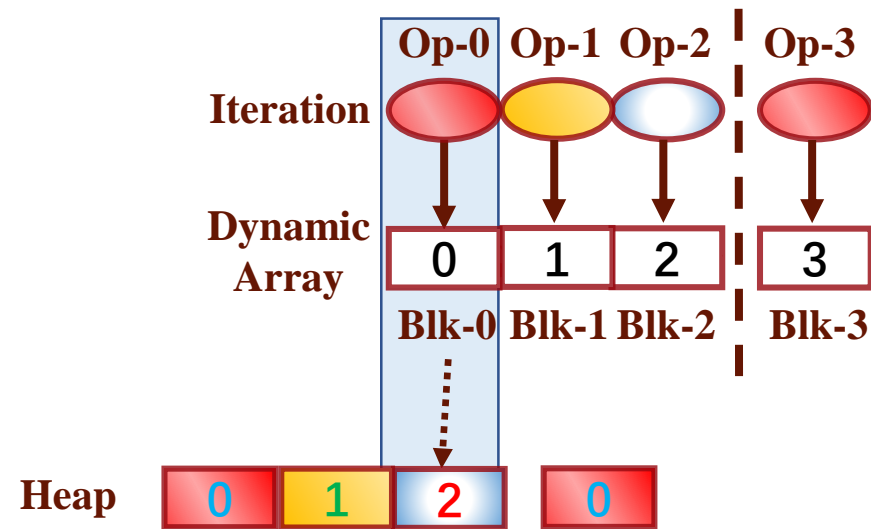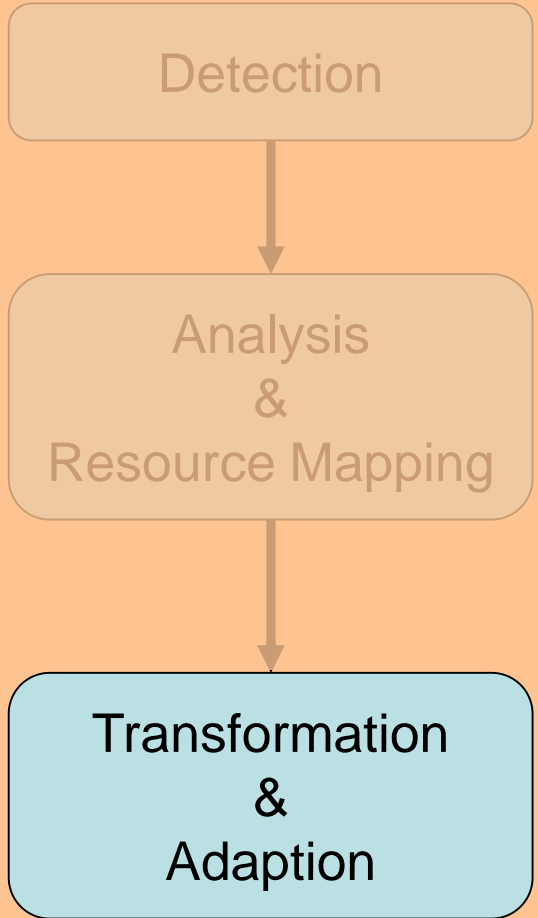**Loop Unrolling** ✗
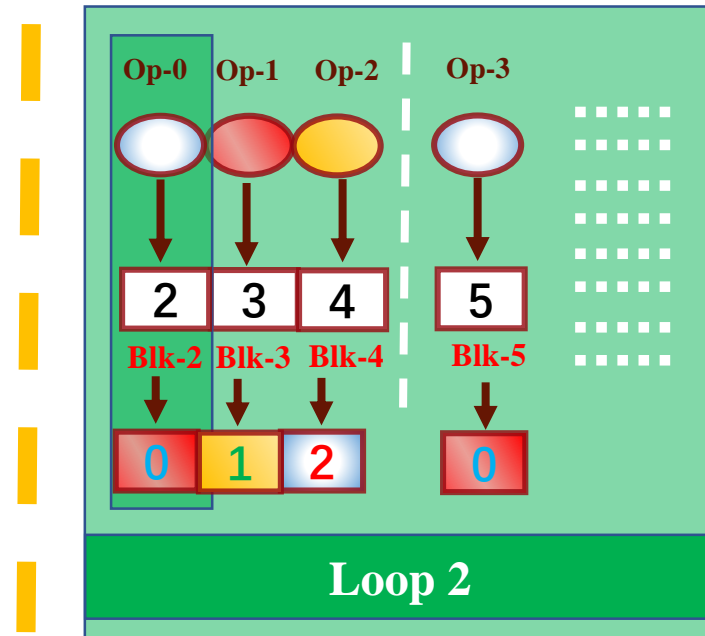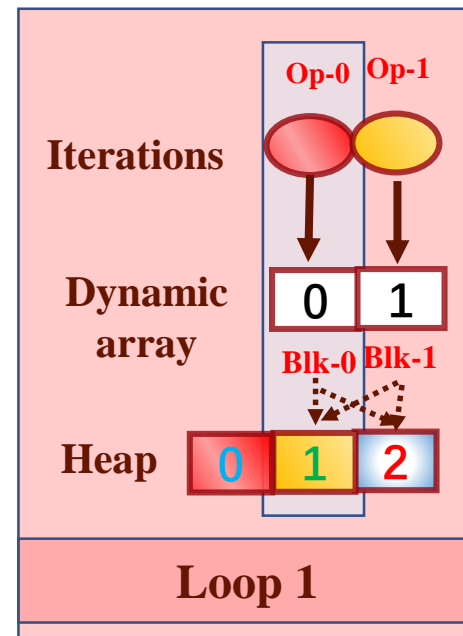
Location of dynamic array is unknown.

# Software: Hi-DMM Compiler

Detection

Analysis
&
Resource Mapping

Transformation
&
Adaption

## 2. Adaption to HLS Directives

**e.g. Loop Transformation for Loop Unrolling**

**Example with DMM: Operations mapped to all partitions**

Op-0  Op-1  Op-2     Op-3

Iteration

Dynamic
Array

| 0 | 1 | 2 | | 3 |

Blk-0 Blk-1 Blk-2    Blk-3

**Loop Unrolling** ✗

Heap

| 0 | 1 | 2 | | 0 |

**Location of dynamic array is unknown.**

# Software: Hi-DMM Compiler

Detection

Analysis
&
Resource Mapping

Transformation
&
Adaption

## 2. Adaption to HLS Directives

### e.g. Loop Transformation for Loop Unrolling

**Solution: Loop Splitting**



**Loop Unrolling**

# Software: Hi-DMM Compiler

Detection

Analysis
&
Resource Mapping

Transformation
&
Adaption

Accelerator IP
with DMM

HLS

Hi-DMM
Interconnect

## Generated Vivado Project

Allocator 1

Allocator 0

Allocator 2

Accelerator

# Outline

- Motivation
- Overview of Hi-DMM
- Implementation of Software
- **Implementation of Hardware**
- Evaluation of Hi-DMM
- Open-Source Hi-DMM Platform
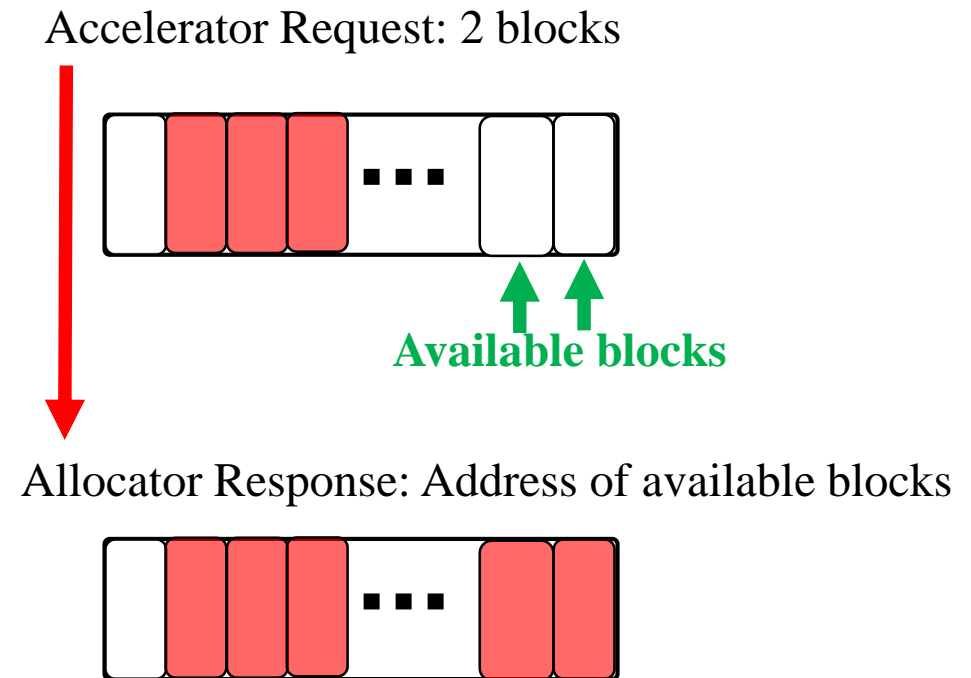- Conclusion

# Hardware: Hi-DMM Allocators

## What do allocators do?

### 1. Record which blocks in the heap are used

int *a
- Memory block 0
- Memory block 1
- Memory block 2
- Memory block 3
- ...
- Memory block N-2
- Memory block N-1

A Heap

### 2. Allocate memory according to the size

Accelerator Request: 2 blocks

**Available blocks**

Allocator Response: Address of available blocks

# Hardware: Hi-DMM Allocators

## Compare with Previous Works

| Previous Mechanism | Allocation Latency | Resource Cost | Memory Efficiency |
|---|---|---|---|
| Buddy Tree | medium | high | high |
| Fixed-Size Blocks | high | low | low |
| Free List | high | medium | high |

**Hi-DMM allocators are based on buddy tree.**

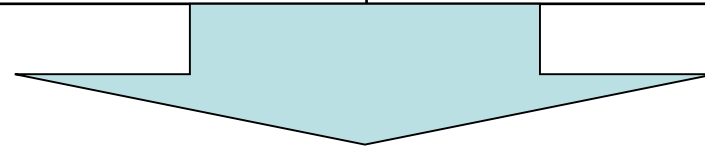| Proposed Mechanism | Allocation Latency | Resource Cost | Memory Efficiency |
|---|---|---|---|
| Hi-DMM Allocators | low | low | high |

# Hardware: Hi-DMM Allocators

## Compare among Hi-DMM Allocators

| Hi-DMM Mechanism | Allocation Latency | Resource Cost | Highlights |
|---|---|---|---|
| Fast Buddy Tree | Very low (~10 cycles) | Low | Fast allocator for small heap |
| Pre-Allocation Tree | Very low (~5 cycles) | Low | Pre-allocate a block for next request |
| Hybrid Tree | Low (~20 cycles) | Very low | Very large management capability |
| K-Way Tree | Extremely low (~1 cycles) | Extremely Low | Fixed-size allocator for user-defined struct |

**Meet the requirements of various applications**

# Hardware: Buddy Tree

- **Conventional Buddy Tree Allocation:**
  - <u>Splits</u> the entire space of heap <u>repetitively in half</u> to find an available memory block <u>best fitting the size of request</u>.
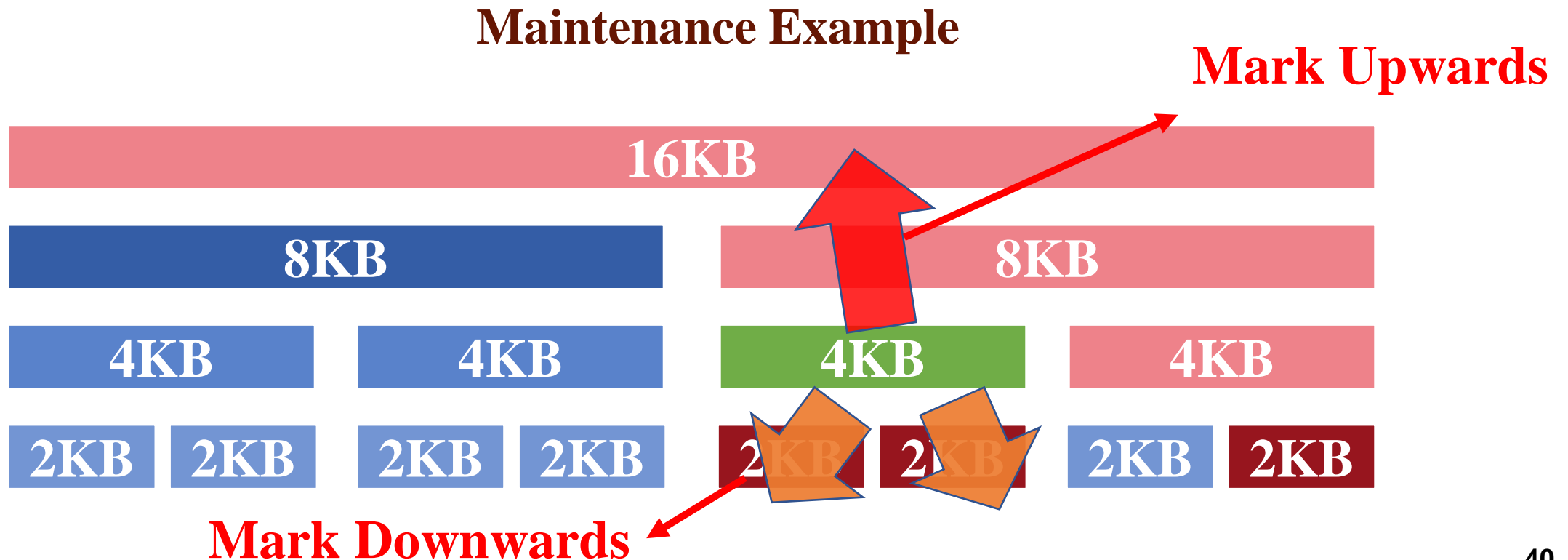
| 16KB | | | | | | | |
|---|---|---|---|---|---|---|---|
| 8KB | | | | 8KB | | | |
| 4KB | | 4KB | | 4KB | | 4KB | |
| 2KB | 2KB | 2KB | 2KB | 2KB | 2KB | 2KB | 2KB |

# Hardware: Buddy Tree

- **Conventional Buddy Tree Allocation:**

**Allocation Example**

# Hardware: Buddy Tree

- **Conventional Buddy Tree Allocation:**

**Maintenance Example**

**Mark Upwards**

| 16KB | |
| 8KB | 8KB |
| 4KB | 4KB | 4KB | 4KB |
| 2KB | 2KB | 2KB | 2KB | 2KB | 2KB | 2KB | 2KB |

**Mark Downwards**

# Hardware: Hi-DMM Allocators

- **Fast Buddy Tree Allocator (FBTA):**

  **1.** Allocation *without searching* layer by layer, based on *bit operation*

  **Bit-Vector (BV)**

  0

  1  0

  1  1  0  0

  1  1  1  1  0  0  1  0

- **Fast Buddy Tree Allocator (FBTA):**

  **1. Allocation *without searching* layer by layer, based on *bit operation***

  $$BV = \boxed{1\ 1\ 1\ 1\ 0\ 0\ \boxed{1}\ 0}$$

  **How to find the lowest set (i.e. 1) bit in BV?**

  $$-\,BV = \boxed{0\ 0\ 0\ 0\ 1\ 1\ 1\ 0} \quad \textit{(two's complement)}$$

  $$(-BV)\ \&BV = \boxed{0\ 0\ 0\ 0\ 0\ 0\ \boxed{1}\ 0} \longrightarrow \textbf{Index}$$

  **Log2 MUX**

# Hardware: Hi-DMM Allocators

- **Fast Buddy Tree Allocator (FBTA):**

  2. Maintenance *parallelized*

- **Fast Buddy Tree Allocator (FBTA):**

    2. Maintenance *parallelized*



Pipeline

16KB

8KB

8KB

4KB   4KB   4KB   4KB

2KB   2KB   2KB   2KB   2KB   2KB   2KB   2KB

# Hardware: Hi-DMM Allocators

- ## Pre-Allocation Tree Allocator (PATA):

    **Based on FBTA but can *pre-allocate before the request***

    **Locality of Allocation**

    **Temporal:** An allocation request is usually **followed closely** by another one.

    **Spatial:** Those allocation requests close to each other usually ask for **similar size.**



**Close**

**Similar Size**

- **Hybrid Tree Allocator (HTA):**

  **Based on FBTA but can manage those *wide bit-vectors with thousands of bits.***

*Resource / Latency*

**Wide BV:** 0 0 1 0 1 1 1 1 · · · · · · · · · · · · 0 0 0 0

# Hardware: Hi-DMM Allocators

- ## Hybrid Tree Allocator (HTA):

    **Based on FBTA but can manage those *wide bit-vectors with thousands of bits.***

    **Solution:** *use BV to mange a wide BV*

Group BV:  **1  1  0  1** ⋯⋯⋯ **1  1  0**

Wide BV:  **0  0  1  0**    **1  1  1  1**  — — — — —  **0  0  0  0**

# Hardware: Hi-DMM Allocators

- ## K-Way Tree Allocator (KWTA):

  **Manage fine-grained fixed-size user-defined struct variables**

  **Scenario: Dynamic Data Structure**

  **Characteristics:**

  **1. Fixed-size:**
  - *tree_node*
  - *queue_element*

  **Waste of Resource** ⟵

  **Buddy Tree?**

  **2. Frequent (De-)Allocation Operations**
  - *insert*_tree_node()
  - queue_*pop*()

  **Further Lower Latency Required** ⟵

- **K-Way Tree Allocator (KWTA):**

Manage user-defined struct variables with extremely low latency



Following requests will be handled by the cache.

An block will be allocated from the cache.

Mini-Heap: a cache of fixed-size blocks

Summary BV

8 7 6 5 4 3 2 1

Group BV

40 39 38 37 36 35 34 33

# Outline

- Motivation
- Overview of Hi-DMM
- Implementation of Software
- Implementation of Hardware
- **Evaluation of Hi-DMM**
- Open-Source Hi-DMM Platform
- Conclusion

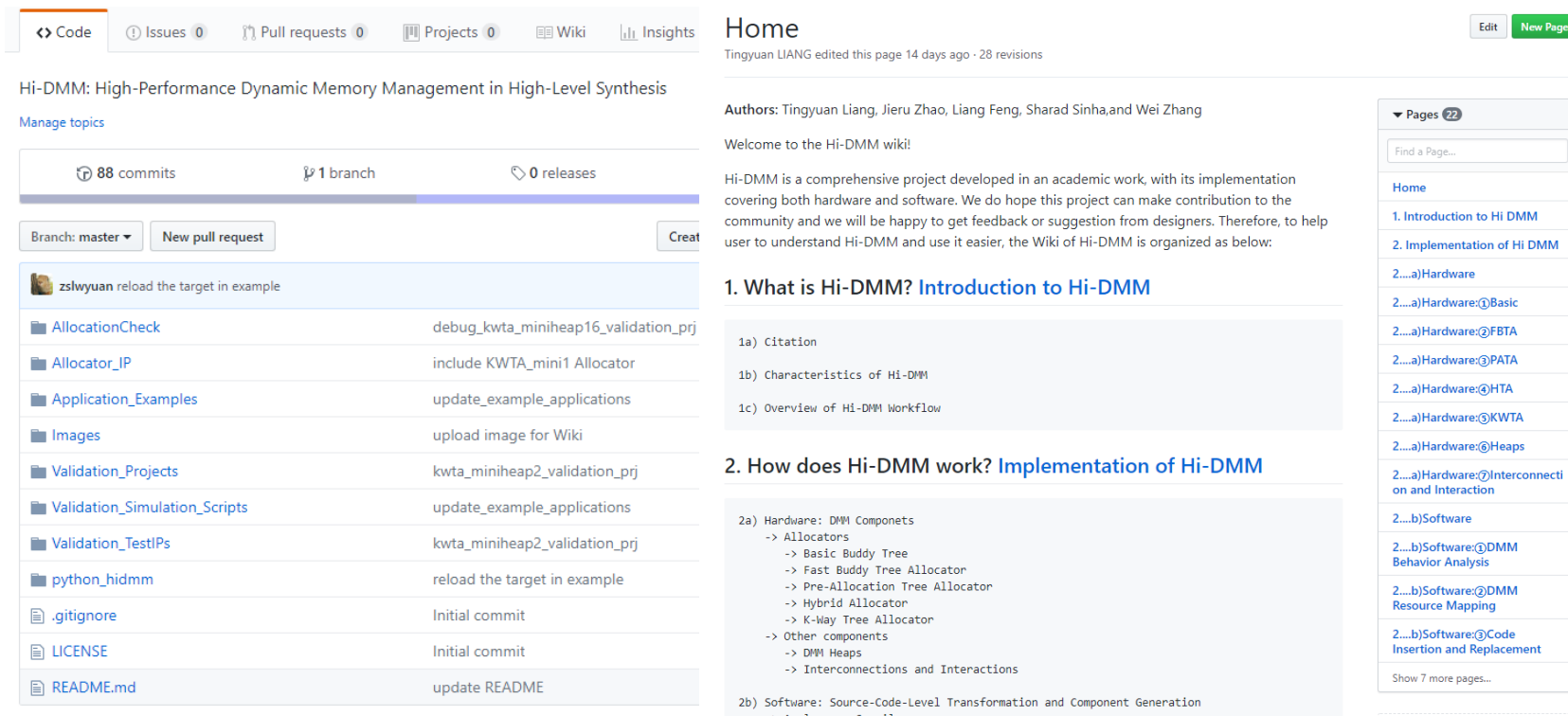# Evaluation: Allocator Performance

## @ 100MHz with Zynq-7020



**Allocation Latency**

Request handled by SysAlloc (Hundreds of cycles)

Request handled by Hi-DMM (Tens of cycles)

Legend: FBTA, PATA(min), PATA(max), KWTA(min), KWTA(max), HTA(min), HTA(max), SysAlloc

X-axis: Management Capability / MAUs — 32, 64, 128, 256, 512, 1K, 2K, 4K, 8K, 16K, 32K, 64K

Y-axis: Latency of Allocation Stage / cycles — 1, 10, 100

# Evaluation: Allocator Resource

## @ 100MHz with Zynq-7020

**LUT Usage**

**BRAM Usage**

**Resource cost by HTA is much lower than FBTA.**

# Evaluation: Memory Efficiency of KWTA

@ 100MHz with Zynq-7020



Reuse

Utilization

# Evaluation: Source Code Optimization

## Pointer Mapping

*Calculation with Multiple Matrices:   ABCD + EF*

**Automatically distributes pointers to 2 heaps with Hi-DMM** : **125033 cycles**

**← 6.0%**

**Assign all pointers to 1 heap:** **133033 cycles**

## Loop Transformation

*Reduction Operation based on Dynamic Arrays*
*(unroll_factor = 4)*

With loop transformation:  412~812 cycles

51.1%

Without loop transformation:  842 cycles

# Evaluation: Source Code Optimization

## Allocation Selection

*Shortest Path Faster Algorithm (SPFA) with a queue*

0.84%

KWTA: 15890 cycles

HTA with pre-allocation: 16563 cycles → 4.72%

HTA: 16840 cycles → 6.34%

# Outline

- Motivation
- Overview of Hi-DMM
- Implementation of Software
- Implementation of Hardware
- Evaluation of Hi-DMM
- Open-Source Hi-DMM Platform
- Conclusion

# Hi-DMM is open to the community

https://github.com/zslwyuan/Hi-DMM

# Outline

- Motivation
- Overview of Hi-DMM
- Implementation of Software
- Implementation of Hardware
- Evaluation of Hi-DMM
- Open-Source Hi-DMM Platform
- Conclusion

# Conclusion

- Software: Hi-DMM Compiler
  - Automatic Transformation
  - HLS-Friendly
  - Couple with commercial tools (Vivado / Vivado HLS)

- Hardware: Hi-DMM Allocator
  - High Performance
  - HLS-Friendly
  - Adaptive to various applications

- Future Works
  - Consider more DMM characteristics
  - Further improve the performance of allocators

# Thanks!