

# HeteroRefactor: Refactoring for Heterogeneous Computing with FPGA

Jason Lau\*, Aishwarya Sivaraman\*, Qian Zhang\*,  
Muhammad Ali Gulzar, Jason Cong, and Miryung Kim  
University of California, Los Angeles

{lau, dcssiva, zhangqian, gulzar, cong, miryung}@cs.ucla.edu

\*Equal co-first authors in alphabetical order

## ABSTRACT

Heterogeneous computing with field-programmable gate-arrays (FPGAs) has demonstrated orders of magnitude improvement in computing efficiency for many applications. However, the use of such platforms so far is limited to a small subset of programmers with specialized hardware knowledge. High-level synthesis (HLS) tools made significant progress in raising the level of programming abstraction from hardware programming languages to C/C++, but they usually cannot compile and generate accelerators for kernel programs with pointers, memory management, and recursion, and require manual refactoring to make them HLS-compatible. Besides, experts also need to provide heavily handcrafted optimizations to improve resource efficiency, which affects the maximum operating frequency, parallelization, and power efficiency.

We propose a new dynamic invariant analysis and automated refactoring technique, called HETEROREFACTOR. First, HETEROREFACTOR monitors FPGA-specific dynamic invariants—the required bit-width of integer and floating-point variables, and the size of recursive data structures and stacks. Second, using this knowledge of dynamic invariants, it refactors the kernel to make traditionally HLS-incompatible programs synthesizable and to optimize the accelerator’s resource usage and frequency further. Third, to guarantee correctness, it selectively offloads the computation from CPU to FPGA, only if an input falls within the dynamic invariant. On average, for a recursive program of size 175 LOC, an expert FPGA programmer would need to write 185 more LOC to implement an HLS compatible version, while HETEROREFACTOR automates such transformation. Our results on Xilinx FPGA show that HETEROREFACTOR minimizes BRAM by 83% and increases frequency by 42% for recursive programs; reduces BRAM by 41% through integer bitwidth reduction; and reduces DSP by 50% through floating-point precision tuning.

## KEYWORDS

heterogeneous computing, automated refactoring, FPGA, high-level synthesis, dynamic analysis

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

ICSE '20, May 23–29, 2020, Seoul, Republic of Korea

© 2020 Copyright held by the owner/author(s). Publication rights licensed to Association for Computing Machinery.

ACM ISBN 978-1-4503-7121-6/20/05... \$15.00

<https://doi.org/10.1145/3377811.3380340>

## ACM Reference format:

Jason Lau\*, Aishwarya Sivaraman\*, Qian Zhang\*, Muhammad Ali Gulzar, Jason Cong, Miryung Kim. 2020. HeteroRefactor: Refactoring for Heterogeneous Computing with FPGA. In *Proceedings of 42nd International Conference on Software Engineering, Seoul, Republic of Korea, May 23–29, 2020 (ICSE '20)*, 13 pages.

<https://doi.org/10.1145/3377811.3380340>

## 1 INTRODUCTION

In recent years, there has been a growing interest in architectures that incorporate heterogeneity and specialization to improve performance, e.g., [12, 14, 16, 22]. FPGAs are reprogrammable hardware that often exceeds the performance of general-purpose CPUs by several orders of magnitude [8, 33, 57] and offer lower cost across a wide variety of domains [7, 9, 17]. To support the development of such architectures, hardware vendors support CPU+FPGA multi-chip packages (e.g., Intel Xeon [35, 58]) and cloud providers support virtual machines with FPGA accelerators and application development frameworks (e.g., Amazon F1 [3]).

Although FPGAs provide substantial benefits and are commercially available to a broad user base, they are associated with a high development cost [64]. Programming an FPGA is a difficult task; hence, it is limited to a small subset of programmers with specialized knowledge on FPGA architecture details. To address this issue, there has been work on high-level synthesis (HLS) for FPGAs [21]. HLS tools take a kernel written in C/C++ as input and automatically generates an FPGA accelerator. However, to meet the HLS synthesizability requirement, significant code rewriting is needed. For example, developers must manually remove the use of pointers, memory management, and recursion, since such code is not compilable with HLS. To achieve high efficiency, the users must heavily restructure the kernel to supply optimization information manually at the synthesis time. Carefully handcrafted HLS optimizations are non-trivial and out of reach for software engineers who usually program with CPUs [13, 15].

Our observation is that software kernels are often over-engineered in the sense that a program is generalized to handle more inputs than what is necessary for common-case inputs. While this approach has no or little impact on the program efficiency on a CPU, in an FPGA accelerator, the design efficiency could be impacted considerably by the compiled size that depends on actual ranges of values held by program variables, the actual size of recursive data structures observed at runtime, etc. For example, a programmer may choose a 32-bit integer data type to represent a human age, whose values range from 0 to 120 in most cases. Consider another example, where in 99% of executions, the size of a linked list is

bounded by 2k; however, the programmer may manually flatten it to an array with an overly-conservative size of 16k.

We propose a novel combination of dynamic invariant analysis, automated refactoring, and selective offloading approach, called HETEROREFACTOR to guide FPGA accelerator synthesis. This approach guarantees correctness—*behavior preservation*, as it selectively offloads the computation from CPU to FPGA, only if the invariant is met, but otherwise keeps the computation on CPU. It also does not require having a representative data set for identifying dynamic invariants, as its benefit is to aggressively improve FPGA accelerator efficiency for a common case input without sacrificing correctness. In this approach, a programmer first implements her kernel code in a high-level language like C/C++. Then she executes the kernel code on existing tests or a subset of input data to identify FPGA-specific dynamic invariants. HETEROREFACTOR automatically refactors the kernel with pointers into a pointerless, non-recursive program to make it HLS-compatible and to reduce resource usage by lowering bitwidth for integers and floating-points, which in turn reduces resource usages and increases the frequency at the FPGA level.

We evaluate HETEROREFACTOR on ten programs, including five handwritten recursive programs, three integer-intensive programs from Rosetta benchmark [84], and two floating-point-intensive programs from OpenCV [6]. We generate kernels targeting to a Xilinx Virtex UltraScale+ XCVU9P FPGA on a VCU1525 Reconfigurable Acceleration Platform [80] and achieve the following results:

- (1) For recursive programs that are traditionally unsynthesizable, HETEROREFACTOR refactors pointers and recursion with the accesses to a flattened, finite-size array, making them HLS-compatible. On average, for a recursive program of size 175 LOC, an expert FPGA programmer would need to write 185 *more* LOC to implement an HLS-compatible version, while HETEROREFACTOR requires no code change. Using a tight bound for a recursive data structure depth, the resulting accelerator is also resource-efficient—an accelerator with a common-case bound of 2k size can achieve 83% decrease in BRAM and 42% increase in frequency compared to the baseline accelerator with an overly conservative size of 16k.
- (2) For integers, HETEROREFACTOR performs transparent optimization and reduces the number of bits by 76%, which leads to 25% reduction in flip-flops (FF), 21% reduction in look-up tables, 41% reduction in BRAM, and 52% decrease in DSP.
- (3) For floating-points, HETEROREFACTOR automatically reduces the bitwidth while providing a probabilistic guarantee for a user-specific quality loss and confidence level. The optimized accelerator can achieve up to 61% reduction in FF, 39% reduction in LUT, and 50% decrease in DSP when an acceptable precision loss is specified as  $10^{-4}$  at 95% confidence level.

In summary, this work makes the following contributions:

- Traditionally, automated refactoring has been used to improve software maintainability. We adapt and expand automated refactoring to lower the barriers of creating customized circuits using HLS and to improve the efficiency of the generated FPGA accelerator.
- While both dynamic invariant analysis and automated refactoring have a rich literature in software engineering, we design a novel combination of dynamic invariant analysis,

automated kernel refactoring, and selective offloading, for transparent FPGA synthesis and optimization with correctness guarantee, which is unique to the best of our knowledge.

- We demonstrate the benefits of FPGA-specific dynamic invariant and refactoring in three aspects: (1) conversion of recursive data structures, (2) integer optimization, and (3) floating-point tuning with a probabilistic guarantee.

HETEROREFACTOR’s source code and experimental artifacts are publicly available at <https://github.com/heterorefactor/heterorefactor>.

## 2 BACKGROUND

This section overviews a developer workflow when using a high-level synthesis (HLS) tool for FPGA and describes the types of manual refactoring a developer must perform to make their kernel synthesizable and efficient on FPGA.

### 2.1 Overview of FPGA Programming with HLS

Modern FPGAs include millions of look-up tables (LUTs), thousands of embedded block memories (BRAMs), thousands of digital-signal processing blocks (DSPs), and millions of flip-flop registers (FFs) [78]. Each k-input LUT can implement any Boolean function up to k inputs. An FPGA must be programmed with a specific binary *bitstream* to specify all the LUT, BRAM, DSP, and programmable switch configurations to achieve the desired behavior. Fortunately, HLS has been developed in recent years to aid the translation of algorithmic descriptions (e.g., kernel code in C/C++) to application-specific bitstreams [21, 28, 50]. Specifically, HLS raises the abstraction of hardware development by automatically generating RTL (Register-Transfer Level) descriptions from algorithms. Generation of FPGA-specific bitstream consists of a *frontend* responsible for C simulation and a *backend* responsible for hardware synthesis. In the frontend, after analysis of C/C++ code, HLS *schedules* each operation from the source code to certain time slots (clock cycles). Next, it allocates *resources*, i.e., the number and type of hardware units used for implementing functionality, like LUTs, FFs, BRAMs, DSPs, etc. Finally, the *binding* stage maps all operations to the allocated hardware units. This frontend process generates an RTL, which is sent to a backend to perform logic synthesis, placement, and routing to generate FPGA bitstreams. Software simulation is fast; however, hardware synthesis can take anywhere from *a few hours to a couple of days*, depending on the complexity of the algorithm. For example, even for tens of lines of code, hardware synthesis can take hours for our subjects in Section 4.

Therefore, such long hardware synthesis time justifies the cost of manual rewriting of kernels for optimized resource allocation, frequency, and power utilization. In other words, this motivates HETEROREFACTOR to invest time in a-priori dynamic analysis as opposed to just-in-time compilation to optimize FPGA, as frequent iterations of hardware synthesis are prohibitively expensive.

### 2.2 Refactoring for High-Level Synthesis

HLS tools aim to narrow the gap between the software program and its hardware implementation. While HLS tools take kernel code in C or C++, a developer must perform a substantial amount of manual refactoring to make it synthesizable and efficient on an FPGA chip. Such refactoring is error-prone and time-consuming

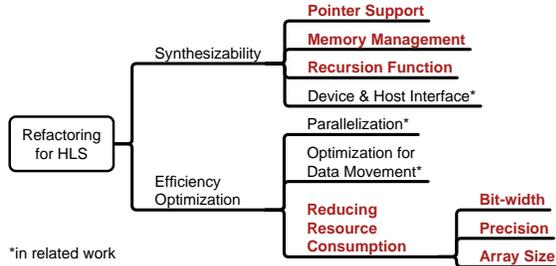


Figure 1: Overview of refactoring for high-level synthesis.

since certain language constructs for readability and expressiveness in C/C++ are not allowed in HLS [25]. A developer must have interdisciplinary expert knowledge in both hardware and software and know obscure platform-dependent details [15]. Below, we categorize manual refactorings for HLS into two kinds: (1) synthesizability and (2) efficiency optimization. In this paper, we focus on improving the Vivado HLS tool from Xilinx [21, 79], which is the most widely used FPGA HLS in the community, although our techniques can be easily generalized to other HLS tools, such as Intel HLS Compiler, Catapult HLS from Mentor, and CyberWorkBench from NEC.

### 2.2.1 Synthesizability.

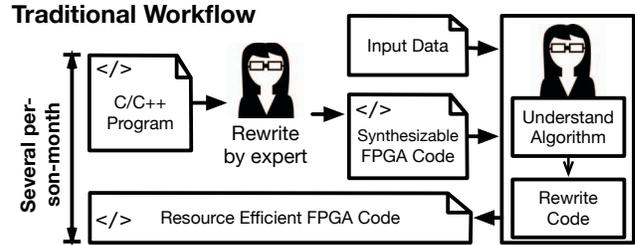
**Pointer support.** To transform kernel code into its equivalent HLS synthesizable version, a developer must manually eliminate pointer declarations and usages; there are only two types of pointers that are natively supported in HLS—pointers to hardware interfaces such as device memory or pointers to variables. Pointer reinterpretation is limited to primitive data types. Arrays of pointers or recursive data structures are strictly forbidden in Vivado HLS.

**Memory management and recursion.** Because Vivado HLS has no capability of *memory management*, function calls to memory allocation such as `malloc` cannot be synthesized. Thus, developers must create an overly conservative, large-sized static array in advance and manage data elements manually. Similarly, Vivado HLS cannot synthesize recursions. Thus, developers must manually convert recursions into iterations or create a large stack to store program states and manage function calls manually.

**Device and host interface.** Vivado HLS requires a strict description of parameters of the top-level function that acts as the *device and host interface*. The function is called from the host and is offloaded into FPGA. A function parameter can be either a scalar or pointer to the device memory with a data size in the power of 2 bytes, and a developer must write specific pragmas—e.g., `#pragma HLS interface m_axi port=input` to use AXI4 interconnect interface for passing the parameter named `input` to the FPGA design.

### 2.2.2 Efficiency Optimization.

**Parallelization.** Reprogrammable hardware provides an inherent potential to implement *parallelization*. Such parallelization can be done through pipelining of different computation stages and by duplicating processing elements or data paths to achieve an effect similar to multi-threading. To guide such parallelization, a developer must manually write HLS pragmas such as `#pragma HLS pipeline` and `#pragma HLS unroll` for suitable loops or must expose parallelization opportunities through polyhedral model-based loop transformations [5, 23, 56, 85].



### HETEROREFACTOR Workflow

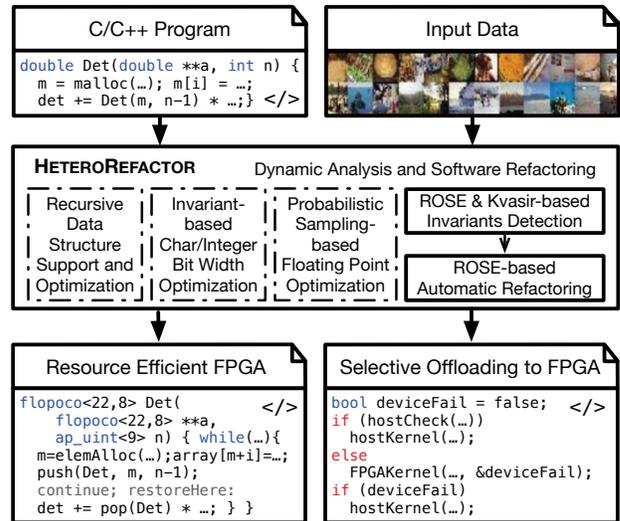
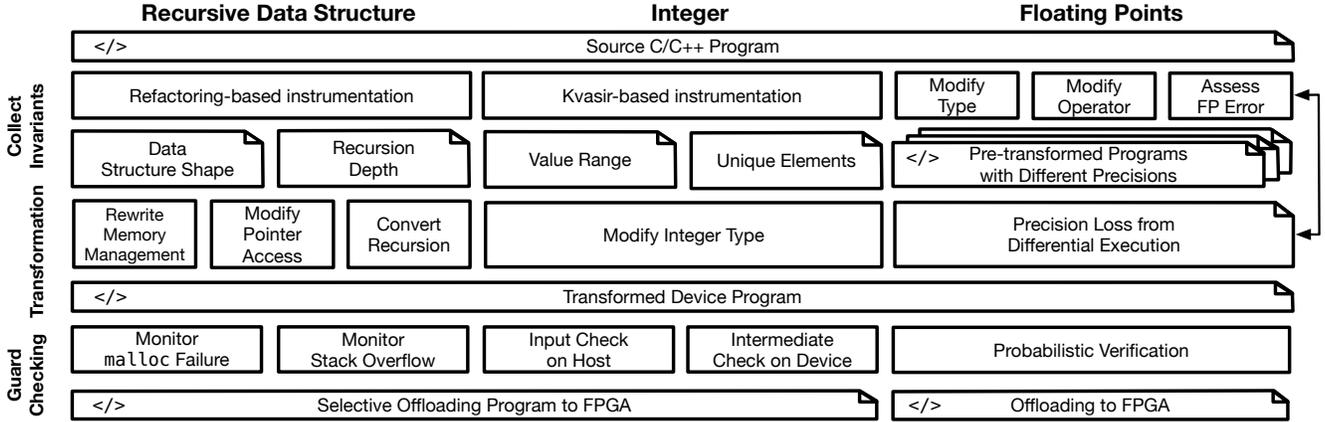


Figure 2: Approach overview of HETEROREFACTOR.

**Optimization of data movement.** Accessing of the device memory can be more efficient by packing bits into the width of DRAM access of 512 bits. To overlap communication with computation, a developer could explicitly implement a double buffering technique [15]. To cache data, developers need to explicitly store them on chip through data tiling, batching or reusing [10, 56, 65].

**Reducing resource consumption.** Provisioning more processing elements or a larger cache will require using more on-chip *resources*, limiting the potential of parallelization and data movement optimizations by duplicating processing elements or adding cache. A higher resource utilization ratio can lower the maximum operating frequency and consume more power; thus, it degrades the performance and efficiency. Besides, a resource-efficient design is economical as it can be implemented on a smaller FPGA chip. Traditionally, developers allocate integers and floating-point variables with a fixed size bitwidth large enough for all possible input values, or create a static array for the largest-possible size. Such a practice may cause wasting on-chip resources. In particular, in modern applications such as big data analytics and ML applications where on-chip resource usage is input-dependent, FPGA resource optimization becomes increasingly difficult.

Figure 1 illustrates our new contributions, highlighted with **bold** and **red**, relative to the prior HLS literature. There exists many automated approaches for generating device and host interfaces [20, 61, 83], exploring parallelization opportunities [24, 34, 46, 83],



**Figure 3: HETEROREFACTOR incorporates three techniques—dynamic invariant detection, kernel refactoring, and selective offloading with guard checking. Its profiling concerns three aspects: (1) the length of recursive data structures, (2) required integer bitwidth, and (3) required floating-point bitwidth to meet a specified precision loss.**

and optimizing data movement [10, 20, 24, 46, 55, 56]. But general methods for reducing resource consumption, pointer support, memory management and recursion support remain as *open research questions* and no automated kernel refactoring exists yet. HETEROREFACTOR addresses three important scopes of such refactoring transformations: (1) converting a program with pointers and recursion to a pointerless and non-recursive program by rewriting memory management and function calls, (2) reducing on-chip resource consumption of integer bitwidth, and (3) reducing on-chip resource consumption by tuning floating-point precision.

### 3 APPROACH

HETEROREFACTOR, as shown in Figure 2, is a novel end-to-end solution that combines dynamic invariant analysis, automated refactoring, and selective offloading for FPGA. It addresses three kinds of HLS refactorings: rewriting a recursive data structure to an array of finite size (Section 3.1); reducing integer bitwidth (Section 3.2); and tuning variable-width floating-point operations (Section 3.3). All three refactorings are based on the insight that a-priori dynamic analysis improves FPGA synthesizability and resource efficiency and that dynamic, input-dependent offloading can guarantee correctness. Figure 3 details the three components that work in concert: (A) instrumentation for FPGA-specific dynamic invariant analysis, (B) source-to-source transformation using dynamic invariants, and (C) selective offloading that checks the guard condition when offloading from CPU to FPGA. The first two kinds of refactorings follow similar implementation for selective offloading using a guard condition check, described in Section 3.4. For floating-point operations, our dynamic analysis provides a probabilistic guarantee that the precision loss is within a given bound.

#### 3.1 Recursive Data Structure

Many applications use recursive data structures built on `malloc`, `free`, and recursive function calls. As mentioned in Section 2.2, HLS tools have strict restrictions on the types of pointers allowed and do not support memory allocation and recursion. For example,

Valdo HLS throws the following error for Figure 4a: *an unsynthesizable type '[10 x %struct.Node.0.1.2]\**. This severely limits the type of programs that can be automatically ported for heterogeneous computing. Expert FPGA developers manually rewrite the recursive data structure into a flattened array to be HLS-compliant; however, as they may not know the common maximum size required for the application, they often over-provision and declare an unnecessarily large size. They also have to manually convert recursion into loop iterations and over-provision the stack required for keeping track of program state involved in recursive calls.

HETEROREFACTOR uses a source to source compiler framework, ROSE [59] to instrument code for identifying the size of recursive data structures and the corresponding stack depth and performs source-to-source transformation based on the size.

**3.1.1 Refactoring-based Instrumentation.** HETEROREFACTOR instruments memory allocation and de-allocation function calls (e.g., allocation of a linked list node), and adds tracing points at the entry and exit of recursive functions to monitor a stack depth. HETEROREFACTOR then determines the number of elements allocated for each data structure based on the collected sizes. In Figure 4a, HETEROREFACTOR sets a tracing point at line 3 to record the number of allocated nodes and another tracing point at line 16 to record the released count. To monitor the recursion depth, HETEROREFACTOR inserts tracing points `call` at the function entry point and `ret` at the function exit point of each recursive function. In Figure 4a, `call` is inserted before line 6, and `ret` is inserted at line 6 and after line 9. HETEROREFACTOR then maintains a variable `stack_size` for each function, which is incremented every time the program reaches `call` and decremented when it reaches `ret`. The highest value attained by `stack_size` during execution is reported and used as the bound for a flattened array and the corresponding stack.

**3.1.2 Refactoring.** HETEROREFACTOR is implemented based on ROSE [59] to rewrite recursive data structures. It takes C/C++ kernel code and the array sizes and recursion depths found via dynamic analysis, and outputs an HLS-compatible version with on-chip memory allocation, removes all pointers except for those with native

```

1 struct Node { Node *left, *right; int val; };
2 void init(Node **root) {
3   *root = (Node *)malloc(sizeof(Node)); }
4 void insert(Node **root, int n, int *arr);
5 void traverse(Node *curr) {
6   if (curr == NULL) return;
7   visit(curr->val);
8   traverse(curr->left);
9   traverse(curr->right); }
10 void top(int n, int *output_if) {
11 #pragma HLS interface m_axi port=output_if
12 Node *root; init(&root); // ...
13 int values[3] = {5, 4, 3};
14 insert(&root, 3, values);
15 int *curr = output_if; traverse(root); // ...
16 free(root); }

```

(a) Original kernel code using pointers and memory allocation.

```

1 bool guard_error = false;
2 struct Node { Node_ptr left, right; int val; };
3 struct Node Node_arr[NODE_SIZE];
4 typedef unsigned int Node_ptr;
5 Node_ptr Node_malloc(size_t size);
6 void Node_free(Node_ptr); // buddy allocation
7 void init(Node_ptr *root) {
8   *root = (Node_ptr)Node_malloc(sizeof(Node));
9   if (!root) guard_error = true; }
10 void insert(Node_ptr *root, int n, int *arr);
11 void traverse(Node_ptr curr) {
12   stack<context> s(TRVERSE_STACK_SIZE, {curr:curr,loc:0});
13   while (!s.empty()) { context c = s.pop(); goto L(c.loc);
14   L0:if (c.curr == NULL) continue;
15       visit(Node_arr[c.curr-1]._data.val);
16       if (s.full()) { guard_error=true; return; }
17       c.loc = 1; s.push(c); s.push({curr: Node_arr[
18         c.curr-1]._data.left, loc: 0}); continue; // traverse(left)
19   L1:// traverse(right) ...
20   L2;; } }
21 void top(int n, int *output_if, bool *fail) {
22 #pragma HLS interface m_axi port=output_if
23 Node_ptr root; init(&root); // ...
24 int values[3] = {5, 4, 3};
25 insert(&root, 3, values);
26 int *curr=output_if; traverse(root); // ...
27 Node_free(root); *fail = guard_error; }

```

(b) Refactored kernel code (schematic).

Figure 4: Example of recursive data structures: binary tree.

HLS support (to be explained further under Rule 2), and rewrites recursive functions. The transformation is semantics-preserving and consists of the following transformation rules:

**RULE 1: REWRITE MEMORY MANAGEMENT.** To replace calls to `malloc` and `free`, for each data type, we pre-allocate an array whose size is guided by instrumentation (line 3 in Figure 4b). The per-type allocation strategy with an array is based on two reasons—HLS only supports pointer reinterpretation on primitive data types, and it can optimize array accesses if the size of one element is known. For each node allocation and de-allocation, we implement a buddy memory system [54] and allocate from the array. The buddy memory system requires less overhead and has little external fragmentation [77], making it suitable for FPGA design. We identify all calls to `malloc` and `free`, the requested types and element counts, and transform them into calls to our library function `Node_malloc` (line 8 in Figure 4b) which returns an available index from the array. Section 4.1 details performance benefits in terms of increased frequency and reduced resource utilization using an array size guided by dynamic analysis rather than declaring an overly conservative size.

**RULE 2: MODIFY POINTER ACCESS TO ARRAY ACCESS.** There are only two types of pointers *natively supported* in HLS, and we do not need to convert them into array access. One is a pointer of interfaces,

which we can identify by looking up pragmas in the code (line 22 in Figure 4b). Second is a pointer to variables, which can be detected by finding all address-of operators or array references in the code. Before modifying pointer access to array access, we identify these natively supported pointers using a breadth-first search on the data flow graph and exclude them from our transformation.

We transform the pointers to an unsigned integer type that takes value less than the size of the pre-allocated array from dynamic analysis. This integer represents the offset of the pointed element in the pre-allocated array. There are three locations where this type of transformation is applied: (1) variable declarations (line 23 Figure 4b), typecasting (line 8 Figure 4b), and function parameters (line 10 Figure 4b) and the return value in both declarations and the definition. We perform a breadth-first search on the data flow graph to propagate the type changes. Since we use an array offset to reference allocated elements, we need to change all pointer dereferences into array accesses with the relative index. We transform indirection operator (`*ptr`) and structure dereference operators (`ptr->`, `ptr->*`) into array accesses with pointer integer as the array index. Similarly, the subscript operators (`ptr[]`) are transformed into array accesses with the pointer integer added with the given offset as the array index. For example, we modify pointer access (line 7 in Figure 4a) to array access (line 15 in Figure 4b).

**RULE 3: CONVERT RECURSION TO ITERATION.** To transform recursive functions into non-recursive ones, we create a stack (line 12 Figure 4b) for each function with all local variables. The depth of the allocated stack is determined through the dynamic analysis step. All references to local variables are transformed into references to elements on the top of the stack (line 14 Figure 4b). To simulate the saved context of the program counter and return value in a CPU call stack, we reserve two member variables in our stack to store the location indicating which line of code we need to restore to, and the return value of the called function.

With a stack, we can implement function calls like in CPU. Entering a function pushes the current context and new parameters to the top of the stack (line 17 Figure 4b), then continue to the first line of the function (line 18 Figure 4b). A function return writes the return value to the stack, pops the top item from the stack, and returns to the saved context (lines 13 Figure 4b).

## 3.2 Integer

**3.2.1 Kvasir-based Instrumentation.** Daikon is a dynamic invariant detection tool [27] that reports likely program invariants during a program’s execution. It consists of two parts: (a) a language-specific front-end and (b) a language-independent inference engine. A front end instruments the program and extracts the program state information by running the program. FPGA kernels are programmed in C/C++ for HLS; hence, we use Kvasir [27], a C/C++ front-end for Daikon, to instrument the target program’s binary.

**3.2.2 FPGA-Specific Invariants.** Daikon is often used for general program comprehension and testing, and therefore it outputs invariants such as an array size or binary comparison, e.g.,  $i > 0$ ,  $i < 0$ ,  $\text{size}(\text{array}) = 0$ ,  $\text{size}(\text{array}) > 0$ . However, such general invariants must be adapted for the purpose of FPGA synthesis. For example, reducing a variable bitwidth leads to resource reduction in FPGA directly [45].

```

1 int weakClassifier(int stddev, int coord[12], int haarC, int w_id);
2 int cascadeClassifier(int SUM1_data[IMG_HEIGHT][IMG_WIDTH],
3 int SQSUM1_data[IMG_HEIGHT][IMG_WIDTH], MyPoint pt) { // ...
4 int stddev = int_sqrt(stddev); // ...
5 }

```

(a) Original kernel code using int.

```

1 bool guard_error = false;
2 void guard_check(ap_int<65> value, int size, int sign) {
3 #pragma HLS inline off
4 if (sign==1) { if (value<0) {
5 if (value < -(1LL<<(size-1)) guard_error = true;
6 } else { /*...*/ } } else { /*...*/ } }
7 int weakClassifier(ap_uint<9> stddev, ap_uint<23> coord[12], ap_uint<7>
8 haarC, ap_uint<8> w_id);
9 int cascadeClassifier(ap_uint<18> SUM1_data[HEIGHT][WIDTH],
10 ap_uint<18> SQSUM1_data[HEIGHT][WIDTH], MyPoint pt) { //...
11 ap_uint<18> stddev = int_sqrt(stddev);
12 guard_check(ap_int<65>(int_sqrt(stddev)),18,0); // ...
13 }

```

(b) Refactored kernel code using ap\_(u)int.

Figure 5: Example of integers: face detection.

Therefore, to optimize FPGA synthesis, we design three types of FPGA-specific invariants: (1) the minimum and maximum value of a variable based on a range analysis, (2) the number and type of unique elements in an array, and (3) the size of an array. For example, first consider Figure 5a. A programmer may over-engineer and use a 32-bit integer by default, which is a higher bitwidth than what is actually necessary. While the instruction set architecture (ISA) for CPU defines integer arithmetics at 32 bits by default, in FPGA, individual bitwidths could be programmed.

**3.2.3 Refactoring.** RULE MODIFY VARIABLE TYPE. To convert an integer to an arbitrary precision integer, we leverage ap\_uint<k> or ap\_int<k> provided by Vivado HLS, which defines an arbitrary precision integer of k bits. As an example, the input haar\_counter to method weakClassifier in Figure 5a is declared as a 32-bit integer by the programmer. However, suppose that HETEROREFACTOR finds that it has a min value of 0 and a max value of 83—it then only needs 7 bits instead of 32 bits. It parses the program’s AST using ROSE [59], identifies the variable declaration node for stddev, coord, haarC, and w\_id, and then modifies the corresponding type as shown in Figure 5b.

### 3.3 Floating Point

Unlike the reduction of integer bitwidth in Section 3.2, reducing the bitwidth for floating-point (FP) variables can lead to FP precision loss. Estimating the error caused by lowering a FP bitwidth can be done reliably only through differential execution, because existing static analysis tends to over-approximate FP errors. Therefore, we design a new probabilistic, differential execution-based FP tuning approach, which consists of four steps: (1) source-to-source transformation for generating program variants with different biwidths, (2) estimation of the required number of input data samples based on Hoeffding’s inequality [37], (3) test generation and differential execution, and (4) probabilistic verification for FP errors.

Prior work on reducing FP precision in CPU [62, 63] used dynamic analysis; however, since they use a golden test set, they do not provide any guarantee on running the reduced precision program on unseen inputs. The key insight behind HETEROREFACTOR’s

```

1 float l2norm(float query[], float data[], int dim) {
2 float dist = 0.0;
3 for (int j = 0; j < dim; j++)
4 dist += ((query[j] - data[j]) * (query[j] - data[j]));
5 return sqrt(dist); }

```

(a) Original kernel code using float.

```

1 using namespace this; typedef policy_flopoco<16,5>::value_t LOWBIT;
2 float low_l2norm(float query[], float data[], int dim) {
3 LOWBIT dist = 0.0;
4 for (int j = 0; j < dim; j++) {
5 LOWBIT fp_query_j = to<LOWBIT, policy>(query[j]);
6 LOWBIT fp_data_j = to<LOWBIT, policy>(data[j]);
7 LOWBIT fp_neg_1 = neg(fp_data);
8 dist += (fp_query + fp_neg_1) * (fp_query + fp_neg_1); }
9 return sqrt(to<float>(dist)); }
10 int main() {
11 for (...) { // ...
12 float highValue = l2norm(args[]);
13 float lowValue = low_l2norm(args[]);
14 float error = highValue - lowValue;
15 if (fabs(error) > acceptableError) Failed++; else Passed++; }
16 if (double(Passed) / Samples > requiredProbability) {
17 /* Passed verification */ } else { /* Failed verification */ } }

```

(b) Refactored kernel code that performs differential execution and probabilistic verification (schematic).

Figure 6: Example of FP numbers: l2norm from KNN.

probabilistic verification approach is that we can draw program input samples to empirically assess whether the relative error between a low precision program and a high precision program is within a given acceptable precision loss  $e$  with probability  $p$ . Given a program with high-precision FP operations,  $hp$ , we construct a lower precision copy of the program,  $lp$ , by changing the corresponding type of all FP variables, constants and operations. For each input  $i \in I$ , we compute the actual error between the high and the low bitwidth variants,  $hp(i) - lp(i)$ . We then check whether this FP error is within the acceptable precision loss  $e$ , indicated by a predicate  $c_i = (hp(i) - lp(i) < e)$  which forms a distribution  $B$ . When the empirical measurement  $\bar{c}_i$  of the given input samples is higher than the target probability  $p$ , the verification is passed.

HETEROREFACTOR takes as inputs: (1) a program, (2) a set of sampled inputs  $I$  or a statistical distribution, (3) an acceptable loss (error)  $e$ , (4) a required probability  $p$ , (5) a required confidence level  $(1 - \alpha)$  and (6) deviation  $\epsilon$ . We use Hoeffding’s inequality [37] to compute the minimum number of samples required to satisfy the given confidence level  $(1 - \alpha)$  and deviation  $\epsilon$ . Equation 1 shows the probability that the empirical measurement  $\bar{c}_i$  of the distribution  $B$  deviates from its actual expectation  $E[\bar{c}_i]$  by  $\epsilon$ , which should be less than  $\alpha$  to achieve our target. Similar to Sampson et al.’s probabilistic assertion [66], we use Hoeffding’s inequality since it provides a conservative, general bound for expectations of *any arbitrary distribution* and relies only on probability and deviation. Therefore, it is suitable for our situation where we have no prior knowledge about the FP loss distribution, incurred by reducing the bitwidth. Equation 2 calculates the minimum number of samples required to verify whether the error is within the acceptable loss.

$$P[|\bar{c}_i - E[\bar{c}_i]| \geq \epsilon] \leq 2e^{-2n\epsilon^2} \quad (1)$$

$$n \geq \ln(2/\alpha)/(2\epsilon^2) \quad (2)$$

For example, when a user wants the actual FP error between the original version (Figure 6a) and the low precision variant (Figure 6b) to be less than  $10^{-4}$  with 95% probability, 95% confidence level and 0.03 deviation, the minimum number of samples required is 2049.

During differential testing with respect to input  $I$ , if the proportion  $\bar{c}_I$  of passing samples to  $|I|$  is greater than  $p$ , we probabilistically guarantee that it is safe to lower the FP precision to the given lower bitwidth. The following transformation rules are applied to identify a lower precision configuration for FP variables.

**RULE 1: DUPLICATE METHOD AND MODIFY TYPE.** To create multiple copies of method `l2norm` in Figure 6a, `HETEROREFACTOR` traverses its AST and redefines the type of variable `query`, `data`, and `dim` originally declared as `float` using `thls::fp_flopoco<E, F>`, whose library is based on Thomas’ work on *templated soft floating-point type for HLS* [75].  $E$  is the number of exponent bits (excluding 1 implicit bit). For example, `thls::fp_flopoco<8, 23>` is 32 bit float type, and `thls::fp_flopoco<5, 16>` uses 22 bits in total (5 for exponent, 16 for fraction and 1 for implicit bit) instead.

**RULE 2: MODIFY ARITHMETIC OPERATORS.** While addition, multiplication, and division operators are implemented by `thls::fp_flopoco<E, F>`, subtraction is not supported [75]. Hence, we convert subtraction in `l2norm` (line 4 in Figure 6a) to corresponding `neg` and `add`, i.e., `subtract(a, b) = add(a, neg(b))` using a variable `fp_neg_1` to store the intermediate result (lines 7-8 in Figure 6b).

**RULE 3: ASSESS FP ERROR FOR DIFFERENTIAL EXECUTION.** We define a skeleton method that computes the relative error and probabilistically verifies if the error is within the user given acceptable loss (lines 11-17 Figure 6b). This involves adding code to invoke the original and generated low precision variants of the function.

### 3.4 Selective Offloading with Guard Check

To selectively offload the computation that fits the reduced size, we insert guard conditions in the host (function sending data from CPU to FPGA) and the kernel (algorithm) to be mapped to FPGA.

For recursive programs, as illustrated at line 9 and line 16 in Figure 4b, we insert a guard condition at `Node_malloc`. The condition sets a global variable `guard_error` to true, if the array is full and more allocation is required. Similarly, the global variable is set to true, if the stack size grows beyond the reduced size. For integer-intensive programs, as shown at line 11 in Figure 5b, we add a guard condition in the kernel and host program. We guard the use of each input, output, and intermediate value in the kernel to proactively prevent overflow (lines 4-6 in Figure 5b). For this, we first identify all expressions containing the reduced bitwidth variables, and if the expression contains binary operations, we insert a guard.

## 4 EVALUATION

Our evaluation seeks to answer the following research questions:

- RQ1** Does `HETEROREFACTOR` effectively enlarge the scope of HLS synthesizability for recursive data structures?
- RQ2** How much manual effort can `HETEROREFACTOR` save by automatically creating an HLS-compatible program?
- RQ3** How much resource reduction does `HETEROREFACTOR` provide for recursive data structures, integer optimization, and floating-point optimization?

**Benchmarks.** We choose ten programs, listed in Table 1 as benchmarks for our main evaluation. For recursive data structures, we use the following five kernels: (R1) Aho-Corasick [2] is a string pattern searching algorithm that uses breadth-first search with a dynamic queue, a recursive Trie tree [26] and a finite state machine. (R2) DFS is depth-first search implemented with recursion. (R3) `Linked List` is insertion, removal, and sorting on a linked list. (R4) `Merge Sort` is performed on a linked list. (R5) Strassen’s [40] is a recursive algorithm for matrix multiplication. For integer optimization, we use face detection and 3D rendering from Rosetta [69, 84] (I6 and I7). We also write (I8) `bubble sort`. For FP bitwidth reduction, we modify two programs—(F9) `KNN-l2norm` and (F10) `RGB2YUV` from OpenCV examples [6].

These subject programs demonstrate `HETEROREFACTOR`’s capability on improving *synthesizability* and *resource efficiency*. For recursive data structures, the original programs are not synthesizable and cannot run on FPGA prior to our work. Thus, we compare our results against manually ported code in terms of human effort and resource utilization. The hand-optimized programs are written by experienced graduate students from an FPGA research group at UCLA. Original programs for integer and floating-point can already run on FPGA. For integers, we compare resource utilization to both original (unoptimized) and manually optimized programs, which are directly from Rosetta [69, 84]. For floating-point, there is no comparison with hand-optimized versions, because a manual optimization attempt will be similar to the verification procedure of `HETEROREFACTOR`.

Though the code size of subject programs looks small, these programs are rather sizable compared against well-known FPGA HLS benchmarks [36, 60]. Similar to creating a new instruction type in the CPU instruction set architecture, the role of FPGA is to create high performance, custom operators at the hardware circuit level. In fact, in a usual FPGA development workflow, developers instrument software on CPU, find out its hotspot corresponding to tens of lines of code, and extract it as a separate kernel for FPGA synthesis. Therefore, our work cannot be judged under the same scalability standard used for pure software refactoring (e.g., handling GitHub projects with millions of lines of code).

**Experimental Environment.** All experiments are conducted on a machine with Intel(R) Core(TM) i7-8750H 2.20GHz CPU and 16 GB of RAM running Ubuntu 16.04. The dynamic invariant analysis is based on instrumentation using Daikon version 5.7.2 with Kvasir as front-end. The automated refactoring is implemented based on ROSE compiler’s version 0.9.11.0. The refactored programs are synthesized to RTL to estimate the resource utilization by Vivado Design Suite 2018.03. The generated kernels are targeted to a Xilinx Virtex UltraScale+ XCVU9P FPGA on a VCU1525 Reconfigurable Acceleration Platform.

### 4.1 Recursive Data Structure

To answer RQ1, we assess how many recursive data structure programs are now synthesizable using `HETEROREFACTOR` that fail compilation with Vivado HLS. For RQ2, we measure manual porting effort as LOC and characters in the code. For RQ3, we assess reduction in resource utilization and increase in frequency of the resulting FPGA

**Table 1: Resource utilization for HETEROREFACTOR**

ID/Program		#LUT	#FF	BRAM	DSP
R1/ Aho- Corasick	Orig	Not Synthesizable			
	Manual	3287	4666	1939	7
	HR-8K	5492	5085	678	10
	HR-2K	5234	5006	206	10
R2/ DFS	Orig	Not Synthesizable			
	Manual	1471	1961	221	0
	HR-8K	2634	2901	254	0
	HR-2K	2563	2881	69	0
R3/ Linked List	Orig	Not Synthesizable			
	Manual	2993	3732	534	0
	HR-8K	3771	4044	318	0
	HR-2K	3655	3936	83	0
R4/ Merge Sort	Orig	Not Synthesizable			
	Manual	2755	2878	519	0
	HR-8K	2751	2958	367	0
	HR-2K	2603	2951	105	0
R5/ Strassen's	Orig	Not Synthesizable			
	Manual	21631	13722	919	12
	HR-8K	20303	14899	223	12
	HR-2K	19591	14654	68	12
I6/ Face Detection	Orig	11325	5784	49	39
	Manual	10158	4800	49	37
	HR	10298	4770	47	28
I7/ 3D Rendering	Orig	3828	2033	123	36
	Manual	2239	1357	67	12
	HR	1907	878	39	9
I8/ Bubble Sort	Orig	313	125	2	0
	Manual	306	125	1	0
	HR	302	125	1	0
F9/ KNN- l2norm	Orig	88843	18591	30	32
	$p$	$e = 10^{-2}$			
	0.95	80163	15257	30	16
	0.99	82228	15626	30	16
	0.999	82228	15626	30	16
	$p$	$e = 10^{-4}$			
	0.95	88952	17102	30	32
	0.99	88952	17102	30	32
	0.999	88952	17855	30	32
	$p$	$e = 10^{-6}$			
	0.95	88843	18591	30	32
	0.99	88843	18591	30	32
	0.999	88843	18591	30	32
F10/ RGB2YUV	Orig	398444	73437	30	288
	$p$	$e = 10^{-4}$			
	0.95	243516	28379	30	144
	0.99	250044	28827	30	144
	0.999	250044	28827	30	144
	$p$	$e = 10^{-5}$			
	0.95	304956	49468	30	144
	0.99	304956	49468	30	144
	0.999	311532	49964	30	144
	$p$	$e = 10^{-6}$			
	0.95	372236	66381	30	288
	0.99	398444	73437	30	288
	0.999	398444	73437	30	288

**Table 2: Recursive data structure kernels, no extra code with HETEROREFACTOR v.s. effort for manual refactoring**

ID/Program	Orig.	Manual	$\Delta$	Orig.	Manual	$\Delta$
	LOC	LOC	LOC	Chars	Chars	Chars
R1/A.-C.	190	291	33%	5673	8776	35%
R2/DFS	86	198	57%	2236	5699	61%
R3/L. List	131	235	44%	3061	6686	54%
R4/M. Sort	128	342	63%	3267	9124	64%
R5/Strassen's	342	735	53%	10026	40971	76%
Geomean			49%			56%

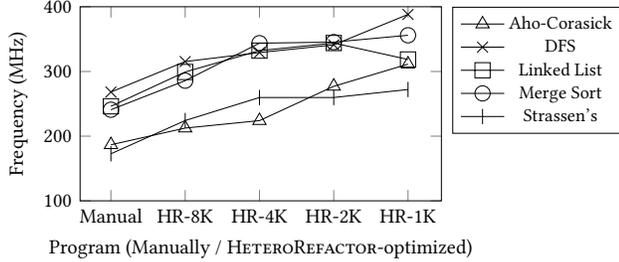
design code, compared to the FPGA design based on a manually written kernel with a conservative size.

Table 2 shows how many lines of code (Manual LOC) and characters (Manual Chars) we need to write in total, if we manually refactor a synthesizable version in Vivado HLS. These manual versions have only a naïve allocator that returns the first unallocated element in the array. If we add a buddy memory system to the manually refactored code to achieve the same functionality as in HETEROREFACTOR, about 100 additional lines of code are required, and thus manual refactoring effort would be even greater.

To evaluate reduction in resource utilization, we instrument the programs using randomly generated input data with typical size of 1k, 2k, 4k or 8k. The profiled information is then passed to HETEROREFACTOR, which automatically generates Vivado HLS-compileable variants of the original program. As mentioned in Section 3.1, FPGA programmers manually transform pointer to non-pointer programs with an overly conservative estimate for the size of the data structure. To compare traditional code rewriting to HETEROREFACTOR, we manually convert and optimize the programs in Table 2 for a conservative data structure size of 16k.

Rows R1-R5 in Table 1 summarize reduction in pre-allocated array size and resource utilization for each of these variants. Manual shows resource usage numbers for the hand-optimized program with a conservative size of 16k, and HR-8k and HR-2k show resource usage of HETEROREFACTOR with 8k and 2k typical data size. If the typical input data size is 2k, there is 83% reduction on average in BRAM, compared to the manually refactored program with a conservative array size. This decrease is significant because Vivado HLS stores most of the large array in BRAM. On the other hand, there is an increase of 302 units in LUT and 494 units in FF on average compared to the hand-optimized version. This small overhead is caused by the fixed-sized buddy memory system which does not increase, as the user design scales.

We implement FPGA accelerator on-board with a target frequency of 300 MHz. Figure 7 reports the maximum operating frequency after placement and routing by Xilinx Vivado for each typical input data size. The frequency is calculated statically by using the worst negative slack (WNS) in the report file:  $F_{max} = 1/(1/300\text{MHz} + \text{WNS})$ . If the input recursive data structure size is 2k, on average, there is 42% increase in frequency compared to the hand-written code with a conservative size of 16k. The frequency improvement comes from reducing communication time among distributed storage resources. When the array is large, the required storage is more distributed, and thus the routing paths are longer, which harms timing.



**Figure 7: Operating frequency of a hand-optimized version with a conservative size, and HETEROREFACTOR-optimized versions with different sizes.**

If an expert FPGA programmer were to use an overly conservative size of 32k, two kernels Merge Sort and Strassen’s will even fail to generate any bitstream, as they require too many resources, justifying the needs of dynamic analysis in creating a custom circuit.

#### Summary 1

By identifying an empirical bound for the recursive data structure size, HETEROREFACTOR makes programs HLS-synthesizable. The accelerators optimized for common-case inputs are 83% more memory-efficient with 42% higher frequency than hand-written code with a conservative size.

## 4.2 Integer

We assess the hypothesis that reducing bitwidth based on dynamic invariants leads to reduction in resource utilization for integers. We measure resource utilization for each program using Vivado HLS 2018.03 targeting a Xilinx XCVU9P FPGA.

For integer bitwidth reduction, HETEROREFACTOR takes as input (1) the kernel under analysis and (2) input data. We use dynamic invariants (Table 3) to create a bitwidth optimized program (e.g., Figure 5b). Table 3 reports the *FPGA-specific dynamic invariants* for integer variables. In terms of input data, we either generate synthetic data of a fixed size or use an existing test set. For Face Detection we use hex images generated from [69] and resize them to the dimension of 16x16. The program uses pre-trained weights declared as an integer array. HETEROREFACTOR identifies that one of the weight arrays requires only unsigned 14 bits based on the max and min value and has only two unique values. For 3D Rendering, we use the test input available in the benchmark [84] and split it into subsets of 100 for each instrumented run. HETEROREFACTOR identifies that the input model has a range of (38, 150) and size of 100. For

**Table 3: FPGA’s specific invariants for integer optimization**

Program	Variable	FPGA-specific Invariants			
		Min	Max	Unique	Size
I6/F. D.	Weights Array 1	8192	12288	2	2913
I6/F. D.	Stddev Variable	305	369	N/A	N/A
I6/F. D.	Coord	0	6746969	21	12
I7/3D R.	Triangle 3D (x0)	38	255	49	100
I8/B.Sort	Input Array	0	10	11	400

Bubble Sort, we generate 400 integers based on Chi Square distribution [47]. The invariants identified by HETEROREFACTOR reconcile with the distribution parameters and fixed size of the input set.

Rows I6-I8 in Table 1 summarize the bitwidth reduction and resource utilization. For each resource type, we report the numbers for (1) an original, unoptimized program in row Orig, (2) a manually optimized program in row Manual, and (3) a HETEROREFACTOR optimized version in row HR. On average, HETEROREFACTOR leads to 25% reduction in FF, 21% reduction in LUT, 41% reduction in BRAM, and 52% decrease in DSP compared to an unoptimized program. Compared to carefully hand-crafted programs by experts, it leads to 12% reduction in FF, 5% reduction in LUTs, 15% reduction in BRAM, and 16% decrease in DSP. Due to the area reduction, more processing elements can be synthesized in one single chip.

We then implement these FPGA accelerators on-board with a target frequency of 300 MHz. All of the refactored programs can meet this target; however, the original version of 3D Rendering fails the timing constraints and can only work with a final frequency of 240.6 MHz. This validates that the frequency improvement can be achieved by HETEROREFACTOR.

#### Summary 2

HETEROREFACTOR reduces the manual refactoring effort by automatically finding the optimized bit width for integers. It reduces 25% FF, 21% LUTs, 41% BRAM, and 52% DSP in resource utilization, which are better than hand-optimized kernels written by experts.

## 4.3 Floating Point

We evaluate the effectiveness of HETEROREFACTOR in providing a probabilistic guarantee while lowering a bitwidth, and reducing resource utilization compared to original programs.

We begin with the given float (32-bit) precision and generate program variants with a reduced operand bitwidth. Reducing mantissa bits leads to precision loss, whereas reducing exponent leads to a smaller dynamic range. Hence, in our experiments, we incrementally reduce mantissa and verify if the loss is within a user given loss,  $\epsilon$ , with probability  $p$ . As described in Section 3.3, we use Hoeffdings inequality to determine the number of input data samples for given  $(1 - \alpha)$  and  $\epsilon$ . In our experiments, we fix  $\epsilon$  to be 0.03, vary  $p$  to be 0.95, 0.99, and 0.999, and keep the confidence level the same as  $p$ , i.e.,  $(1 - \alpha) = p$ , which require at least 2049, 2943 and 4222 samples, respectively. We also consider the input features of FPGA kernels to determine the final number of samples. For example, RGB2YUV requires that the inputs must be multiples of 16, so the final number of samples is 2064 rather than 2049 when  $p$  is 0.95. In our evaluation, we draw random test inputs within 0 to 255. The bitwidth reduction is verified by HETEROREFACTOR with an acceptable loss ( $\epsilon$ ) ( $10^{-2}$ ,  $10^{-4}$ , or  $10^{-6}$ ) for KNN-L2norm and an acceptable loss ( $\epsilon$ ) ( $10^{-4}$ ,  $10^{-5}$ , or  $10^{-6}$ ) for RGB2YUV.

Table 4 summarizes the probabilistic verification results for different  $\epsilon$  and  $p$  configurations. For each configuration, we report verification results for 8 and 16 bits floating-point, where N indicates a verification failure, and the column HR reports the smallest verified bitwidth. As expected, a higher precision and confidence requirement leads to a higher FP bitwidth. The results show that a

**Table 4: Probabilistic floating-point verification**

Program	$p$	$e = 10^{-2}$			$10^{-4}$			$10^{-6}$		
		8	16	HR	8	16	HR	8	16	HR
F9/KNN	0.95	N	N	24	N	N	29	N	N	32
	0.99	N	N	25	N	N	29	N	N	32
	0.999	N	N	25	N	N	30	N	N	32
	$p$	$e = 10^{-4}$			$10^{-5}$			$10^{-6}$		
F10/R2Y	0.95	N	N	21	N	N	25	N	N	30
	0.99	N	N	22	N	N	25	N	N	32
	0.999	N	N	22	N	N	26	N	N	32

32 bit floating-point variable could be reduced to using 21 bits with an acceptable loss of  $10^{-4}$  at 95% confidence level.

We then synthesize the refactored program using Vivado HLS 2018.03 targeting a Xilinx XCVU9P FPGA. Rows F9-F10 in Table 1 summarize the resource utilization for each subject program. The *Orig* row indicates the original program with 32-bit float type and  $p$  represents the probability and the confidence level ( $1 - \alpha$ ). Then we report the resource utilization of FF, LUT, and DSP for each combination of  $p$  and accuracy loss  $e$ . HETEROREFACTOR can achieve up to 61% reduction in FF, 39% reduction in LUT, and 50% decrease in DSP. As existing HLS flow does not support arbitrary floating-point type, so we could not find any hand-optimized kernels, and thus we can only compare against the default high bitwidth version.

**Summary 3**

HETEROREFACTOR reduces the floating-point bitwidth while providing a probabilistic guarantee for a user-specified quality loss, probability and confidence level. It can achieve up to 61% reduction in FF, 39% in LUT, and 50% in DSP.

**4.4 Overhead and Performance**

Table 5 summarizes the instrumentation overhead and refactoring overhead for R1-I8 compared against the synthesis time of the original programs. For recursive data structures, both the instrumentation and refactoring overhead are less than 1%. For integers, HETEROREFACTOR induces less than 1% refactoring overhead, and its instrumentation overhead comes from Kvasir. For floating-point programs F9-F10, Table 6 summarizes the differential execution overhead compared against the synthesis time of the original programs with a specific quality loss  $e$  and probability  $p$ , because there is no instrumentation required. HETEROREFACTOR induces less than 2% overhead on floating-point bitwidth tuning.

We compare the execution performance of the refactored kernel against running the original program on CPU. For floating-point programs, our experiment shows a significant speedup up to 7 $\times$  and 19 $\times$  in KNN-12norm and RGB2YUV. This is because these FP programs can benefit from inherent parallel computation. For recursive programs, our refactored kernels are slower than CPU, because HETEROREFACTOR uses a sequential memory allocation, these kernels are memory-bound, and the frequency of FPGA is lower than that of CPU. For integer intensive programs, the end-to-end performance depends on whether data parallelism could be easily utilized for integer-type data processing. The kernels we selected are slightly

**Table 5: Runtime overhead for recursions and integers**

Program	Instrumentation		Refactoring	
	time (min)	ratio	time (sec)	ratio
R1/Aho-Corasick	0.10	0.26%	5.1	0.26%
R2/DFS	0.06	0.26%	4.7	0.34%
R3/Linked List	0.12	0.49%	4.5	0.31%
R4/Merge Sort	0.05	0.20%	4.5	0.29%
R5/Strassen's	0.09	0.20%	10	0.38%
I6/Face Detection	0.15	0.62%	10	0.69%
I7/3D Rendering	13.66	64.76%	10	0.79%
I8/Bubble Sort	$10^{-3}$	$\sim 0$	$10^{-3}$	$\sim 0$

**Table 6: Differential execution overhead for FP (sec / %)**

Program	$p$	$e = 10^{-2}$		$10^{-4}$		$10^{-6}$	
		sec	%	sec	%	sec	%
F9/KNN	0.95	60.8	0.3%	29.4	0.2%	11.7	0.1%
	0.99	58.2	0.3%	30.4	0.2%	11.7	0.1%
	0.999	60.8	0.3%	25.9	0.1%	12.8	0.1%
	$p$	$e = 10^{-4}$		$10^{-5}$		$10^{-6}$	
F10/R2Y	0.95	83.5	1.8%	59.3	1.3%	26.8	0.6%
	0.99	81.5	1.7%	58.5	1.2%	12.9	0.3%
	0.999	82.3	1.7%	54.5	1.2%	13.7	0.3%

slower than running on CPU because I6 and I7 in Rosetta are designed to achieve higher energy efficiency but not higher processing throughput compared to CPU [84]. HETEROREFACTOR aims to reduce resource usage, while prior work [19, 24] achieves higher performance than CPU by leveraging more on-chip resources to achieve parallelism. HETEROREFACTOR could be used jointly with other tools to produce fast and resource-efficient FPGA accelerators.

**5 RELATED WORK**

**Automated Refactoring.** Since pioneering work on automated refactoring in the early 90s [32, 49, 53], recent studies find that real-world refactorings are generally not semantics-preserving [43, 44], are done manually [76], are error-prone [42, 51], and are beyond the scope and capability of existing refactoring engines. A recent study with professional developers finds that almost 12% of refactorings are initiated by developers' motivation to improve performance [44]. HETEROREFACTOR builds on this foundation [49] but repurposes it to improve performance in the new era of heterogeneous computing with re-programmable circuits. While HETEROREFACTOR's refactoring is not semantics-preserving, it guarantees semantics-preservation by leveraging selective offloading from CPU to FPGA in tandem.

**Dynamic Invariant.** Determining program invariants has been explored widely using both static and dynamic techniques. HETEROREFACTOR is inspired by Daikon [27], which generates invariants of 22 kinds for C/C++/Java programs. Kataoka et al. [41] detect the symptoms of a narrow interface by observing dynamic invariants and refactors the corresponding API. Different from these, HETEROREFACTOR does not require having representative data a-priori, as it leverages selective offloading to guarantee correctness. Therefore, a developer may use systematic test generation tools [29, 30, 67] or test minimization [38, 73] to infer FPGA-specific invariants, as representative data is not required for correctness.

**HLS Optimization.** Klimovic et al. [45] optimize FPGA accelerators for common-case inputs by reducing bitwidths using both bit-mask analysis and program profiling [31]. When inputs exceed the common-case range, a software fallback function is automatically triggered. Their simulation results estimate that an accelerator’s area may be reduced by 28% on average. While their approach is similar to HETEROREFACTOR, its scope is limited to monitoring integer values, and they do not implement a systematic approach to monitor bitwidth invariants and the size of a recursive data structure at the kernel level, nor automatically assess the impact of tuning variable-width floating-point precision with a given error bound. While we present real hardware results on Xilinx Virtex UltraScale+ XCVU9P FPGA, they only present estimated software simulation results. To our knowledge, HETEROREFACTOR is the *first tool* for heterogeneous computing with FPGA that incorporates dynamic invariant analysis, automated kernel refactoring, selective offloading, and synthesized FPGA.

Several approaches provide HLS libraries for implementing variable-width floating-point computation units, but leave it to the programmer to specify which parameters to use and to rewrite their kernel code manually. For example, Thomas [75] presents an HLS backend for generating a customized floating-point accelerator using C++ template-based, parameterized types. This approach requires the user to manually specify the bitwidths for an exponent and fraction, which is automated in HETEROREFACTOR.

HETEROREFACTOR differs from static analysis methods which results in over-approximation. For example, Bitwise [71] propagates bitwidth constraints to variables based on the flow graph of bits. MiniBit [48] minimizes integer and fixed-point data signals with a static method based on affine arithmetic. Cong et al. [18] uses affine arithmetic, general interval arithmetic and symbolic arithmetic methods to optimize for fixed-point data. In contrast to JIT compilation techniques [4], HETEROREFACTOR uses an ahead-of-time profiling phase, due to a long FPGA synthesis time.

**Recursion in Heterogeneous Computing.** Enabling recursive data structures in FPGA has been a long challenge because the address space for each array is separate in HLS/FPGA unlike traditional CPU architectures. SynADT [82] is an HLS library for representing linked lists, binary trees, hash tables, and vectors from pointers, and it internally uses arrays and a shared system-wide memory allocator [81]. However, SynADT supports only a limited set of data structures and requires developers to manually refactor. In contrast, HETEROREFACTOR automatically monitors an appropriate size of a recursive data structure and performs fully automated kernel transformation to convert pointer usage to operations on a finite-sized array and implements a guard-condition based offloading. Thomas et al. [74] use C++ templates to create a domain-specific language to support recursion in HLS. However, it requires extensive rewriting of control statements using lambdas.

Similar limitation existed on GPU with CUDA [52] and OpenCL [72]. For example, dynamic memory management on device global memory using `malloc` was not supported until CUDA 3.2 [52], and there is no implementation of `malloc` on shared on-chip memory.

However, one may write their own universal allocator for arbitrary types as a replacement for `malloc` on any memory [1, 39, 70] because GPU allows a single address space with regular access widths, similar to CPU, while FPGA does not. Such approaches [1, 39, 52, 70] still require manually specifying a heap size. HETEROREFACTOR automatically detects the required size of FPGA on-chip memory for recursive data structures using dynamic invariant detection, and fallbacks to CPU computation when the size invariants are violated.

**Tuning Floating-point Precision.** FPTuner [11] uses static analysis for automatic precision-tuning [68] of real valued expressions. It supports a single, double, or quadruple precision rather than an arbitrary-width FP type. Precimonious [63] is a floating-point precision tuning tool that uses dynamic analysis and delta-debugging to identify lower precision instruction that satisfies the user-specified acceptable precision loss constraint. HETEROREFACTOR’s FP tuning is inspired by the success of Precimonious. However, HETEROREFACTOR extends this idea by adding a probabilistic verification logic to provide statistical guarantee on precision loss. While Precimonious is a software-only analysis tool for FP tuning, HETEROREFACTOR is an end-to-end approach that integrates dynamic invariant analysis, automated refactoring, and FPGA synthesis.

## 6 CONCLUSION

Traditionally, automated refactoring has been used to improve software maintainability. To meet the increasing demand for developing new hardware accelerators and to enable software engineers to leverage heterogeneous computing environments, we adapt and expand the scope of automated refactoring. HETEROREFACTOR provides a novel, end-to-end solution that combines (1) dynamic analysis for identifying common-case sizes, (2) kernel refactoring to enhance HLS synthesizability and to reduce on-chip resource usage on FPGA, and (3) selective offloading with guard checking to guarantee correctness. For the transformed recursive programs, HETEROREFACTOR reduces BRAM by 83% and increases frequency by 42%. For integer optimization, it reduces the number of bits for integers by 76%, leading to 41% decrease in BRAM. For floating-point optimization, it reduces DSP usage by 50%, while guaranteeing a user-specified precision loss of 0.01 with 99.9% confidence.

## 7 ACKNOWLEDGEMENT

We would like to thank Guy Van den Broeck, Brett Chalabian, Todd Millstein, Peng Wei, Cody Hao Yu and anonymous reviewers for their valuable feedback. We thank Janice Wheeler for proofreading the draft. This work is in part supported by NSF grants CCF-1764077, CCF-1527923, CCF-1723773, ONR grant N00014-18-1-2037, Intel CAPA grant, and Samsung grant. This work is also partially supported by CRISP, one of six centers in JUMP, a Semiconductor Research Corporation (SRC) program and the contributions from the member companies under the Center for Domain-Specific Computing (CDSC) Industrial Partnership Program, including Xilinx and VMWare. Jason Lau is supported by UCLA Computer Science Departmental Fellowship and Muhammad Gulzar is supported by Google PhD Fellowship.

## REFERENCES

- [1] Andrew V Adinetz and Dirk Pleiter. 2014. Halloc: a high-throughput dynamic memory allocator for GPGPU architectures. In *GPU Technology Conference (GTC)*, Vol. 152.
- [2] Alfred V Aho and Margaret J Corasick. 1975. Efficient string matching: an aid to bibliographic search. *Commun. ACM* 18, 6 (1975), 333–340.
- [3] Amazon.com. 2019. Amazon EC2 F1 Instances: Run Custom FPGAs in the AWS Cloud. <https://aws.amazon.com/ec2/instance-types/f1>. (2019).
- [4] Matthew Arnold, Stephen J Fink, David Grove, Michael Hind, and Peter F Sweeney. 2005. A survey of adaptive optimization in virtual machines. *Proc. IEEE* 93, 2 (2005), 449–466.
- [5] Uday Bondhugula, Muthu Baskaran, Sriram Krishnamoorthy, Jagannathan Ramanujam, Atanas Rountev, and Ponnuswamy Sadayappan. 2008. Automatic transformations for communication-minimized parallelization and locality optimization in the polyhedral model. In *International Conference on Compiler Construction*. Springer, 132–146.
- [6] Gary Bradski and Adrian Kaehler. 2008. *Learning OpenCV: Computer vision with the OpenCV library*. O'Reilly Media, Inc.
- [7] Jared Casper and Kunle Olukotun. 2014. Hardware acceleration of database operations. In *Proceedings of the 2014 ACM/SIGDA international symposium on Field-programmable gate arrays*. ACM, 151–160.
- [8] Adrian M Caulfield, Eric S Chung, Andrew Putnam, Hari Angepat, Jeremy Fowers, Michael Haselman, Stephen Heil, Matt Humphrey, Puneet Kaur, Joo-Young Kim, et al. 2016. A cloud-scale acceleration architecture. In *The 49th Annual IEEE/ACM International Symposium on Microarchitecture*. IEEE Press, 7.
- [9] Zhe Chen, Hugh T Blair, and Jason Cong. 2019. LANMC: LSTM-Assisted Non-Rigid Motion Correction on FPGA for Calcium Image Stabilization. In *Proceedings of the 2019 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. 104–109.
- [10] Yuze Chi, Jason Cong, Peng Wei, and Peipei Zhou. 2018. SODA: stencil with optimized dataflow architecture. In *2018 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*. IEEE, 1–8.
- [11] Wei-Fan Chiang, Mark Baranowski, Ian Briggs, Alexey Solovyev, Ganesh Gopalakrishnan, and Zvonimir Rakamarundefined. 2017. Rigorous Floating-Point Mixed-Precision Tuning. In *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages (POPL 2017)*. Association for Computing Machinery, New York, NY, USA, 300–315.
- [12] Andrew A Chien, Allan Snaveley, and Mark Gahagan. 2011. 10x10: A general-purpose architectural approach to heterogeneity and energy efficiency. *Procedia Computer Science* 4 (2011), 1987–1996.
- [13] Young-kyu Choi and Jason Cong. 2017. HLScope: High-Level performance debugging for FPGA designs. In *2017 IEEE 25th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*. IEEE, 125–128.
- [14] Eric S Chung, Peter A Milder, James C Hoe, and Ken Mai. 2010. Single-chip heterogeneous computing: Does the future include custom logic, FPGAs, and GPGPUs?. In *2010 43rd annual IEEE/ACM international symposium on microarchitecture*. IEEE, 225–236.
- [15] Jason Cong, Zhenman Fang, Yuchen Hao, Peng Wei, Cody Hao Yu, Chen Zhang, and Peipei Zhou. 2018. Best-Effort FPGA Programming: A Few Steps Can Go a Long Way. *arXiv preprint arXiv:1807.01340* (2018).
- [16] Jason Cong, Mohammad Ali Ghodrati, Michael Gill, Beayna Grigorian, Karthik Gururaj, and Glenn Reinman. 2014. Accelerator-rich architectures: Opportunities and progresses. In *Proceedings of the 51st Annual Design Automation Conference*. ACM, 1–6.
- [17] Jason Cong, Licheng Guo, Po-Tsang Huang, Peng Wei, and Tianhe Yu. 2018. SMEM++: A Pipelined and Time-Multiplexed SMEM Seeding Accelerator for Genome Sequencing. In *2018 28th International Conference on Field Programmable Logic and Applications (FPL)*. 210–2104.
- [18] Jason Cong, Karthik Gururaj, Bin Liu, Chunyue Liu, Zhiru Zhang, Sheng Zhou, and Yi Zou. 2009. Evaluation of static analysis techniques for fixed-point precision optimization. In *2009 17th IEEE Symposium on Field Programmable Custom Computing Machines*. IEEE, 231–234.
- [19] Jason Cong, Muhuan Huang, Peichen Pan, Yuxin Wang, and Peng Zhang. 2016. Source-to-source optimization for HLS. In *FPGAs for Software Programmers*. Springer, 137–163.
- [20] Jason Cong, Muhuan Huang, Peichen Pan, Di Wu, and Peng Zhang. 2016. Software infrastructure for enabling FPGA-based accelerations in data centers. In *Proceedings of the 2016 International Symposium on Low Power Electronics and Design*. ACM, 154–155.
- [21] Jason Cong, Bin Liu, Stephen Neuendorffer, Juanjo Noguera, Kees Vissers, and Zhiru Zhang. 2011. High-level synthesis for FPGAs: From prototyping to deployment. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 30, 4 (2011), 473–491.
- [22] Jason Cong, Vivek Sarkar, Glenn Reinman, and Alex Bui. 2010. Customizable domain-specific computing. *IEEE Design & Test of Computers* 28, 2 (2010), 6–15.
- [23] Jason Cong and Jie Wang. 2018. PolySA: polyhedral-based systolic array auto-compilation. In *2018 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*. IEEE, 1–8.
- [24] Jason Cong, Peng Wei, Cody Hao Yu, and Peng Zhang. 2018. Automated accelerator generation and optimization with composable, parallel and pipeline architecture. In *Proceedings of the 55th Annual Design Automation Conference (DAC)*. IEEE, 1–6.
- [25] Jefferson Santiago da Silva, François-Raymond Boyer, and JM Langlois. 2019. Module-per-Object: a Human-Driven Methodology for C++-based High-Level Synthesis Design. *arXiv preprint arXiv:1903.06693* (2019).
- [26] Rene De La Briandais. 1959. File searching using variable length keys. In *Papers presented at the the March 3-5, 1959, Western Joint Computer Conference*. ACM, 295–298.
- [27] Michael D Ernst, Jeff H Perkins, Philip J Guo, Stephen McCamant, Carlos Pacheco, Matthew S Tschantz, and Chen Xiao. 2007. The Daikon system for dynamic detection of likely invariants. *Science of computer programming* 69, 1-3 (2007), 35–45.
- [28] Daniel D Gajski, Nikil D Dutt, Allen CH Wu, and Steve YL Lin. 2012. *High-Level Synthesis: Introduction to Chip and System Design*. Springer Science & Business Media.
- [29] Patrice Godefroid, Nils Klarlund, and Koushik Sen. 2005. DART: Directed Automated Random Testing. In *Proceedings of the 2005 ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '05)*. ACM, New York, NY, USA, 213–223. <https://doi.org/10.1145/1065010.1065036>
- [30] Patrice Godefroid, Michael Y. Levin, and David A Molnar. 2008. Automated White-box Fuzz Testing. In *Network Distributed Security Symposium (NDSS)*. Internet Society. <http://www.truststc.org/pubs/499.html>
- [31] Marcel Gort and Jason H Anderson. 2013. Range and bitmask analysis for hardware optimization in high-level synthesis. In *2013 18th Asia and South Pacific Design Automation Conference (ASP-DAC)*. IEEE, 773–779.
- [32] William G. Griswold. 1991. *Program Restructuring as an Aid to Software Maintenance*. Ph.D. Dissertation. University of Washington.
- [33] Licheng Guo, Jason Lau, Zhenyuan Ruan, Peng Wei, and Jason Cong. 2019. Hardware acceleration of long read pairwise overlapping in genome sequencing: A race between fpga and gpu. In *2019 IEEE 27th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*. IEEE, 127–135.
- [34] Licheng Guo, Jason Lau, Jie Wang, Cody Hao Yu, Yuze Chi, Zhe Chen, Zhiru Zhang, and Jason Cong. 2020. Analysis and Optimization of the Implicit Broadcasts in FPGA HLS to Improve Maximum Frequency. In *Proceedings of the 28th ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA)*. ACM.
- [35] Prabhat Gupta. 2019. Xeon+FPGA Platform for the Data Center. <https://www.archive.ece.cmu.edu/~calcm/carl/lib/exe/fetch.php?media=carl15-gupta.pdf>. (2019).
- [36] Yuko Hara, Hiroyuki Tomiyama, Shinya Honda, Hiroaki Takada, and Katsuya Ishii. 2008. Chstone: A benchmark program suite for practical c-based high-level synthesis. In *2008 IEEE International Symposium on Circuits and Systems*. IEEE, 1192–1195.
- [37] Wassily Hoeffding. 1994. Probability inequalities for sums of bounded random variables. In *The Collected Works of Wassily Hoeffding*. Springer, 409–426.
- [38] Hwa-You Hsu and Alessandro Orso. 2009. MINTS: A General Framework and Tool for Supporting Test-suite Minimization. In *Proceedings of the 31st International Conference on Software Engineering (ICSE '09)*. IEEE Computer Society, Washington, DC, USA, 419–429. <https://doi.org/10.1109/ICSE.2009.5070541>
- [39] Xiaohuang Huang, Christopher I Rodrigues, Stephen Jones, Ian Buck, and Wenmei Hwu. 2010. Xmalloc: A scalable lock-free dynamic memory allocator for many-core machines. In *2010 10th IEEE International Conference on Computer and Information Technology*. IEEE, 1134–1139.
- [40] Steven Huss-Lederman, Elaine M Jacobson, Jeremy R Johnson, Anna Tsao, and Thomas Turnbull. 1996. Implementation of Strassen's algorithm for matrix multiplication. In *Supercomputing'96: Proceedings of the 1996 ACM/IEEE Conference on Supercomputing*. IEEE, 32–32.
- [41] Yoshio Kataoka, Michael D. Ernst, William G. Griswold, and David Notkin. 2001. Automated support for program refactoring using invariants. In *ICSM 2001, Proceedings of the International Conference on Software Maintenance*. Florence, Italy, 736–743.
- [42] Miryung Kim, Dongxiang Cai, and Sunghun Kim. 2011. An Empirical Investigation into the Role of Refactorings during Software Evolution. In *ICSE' 11: Proceedings of the 2011 ACM and IEEE 33rd International Conference on Software Engineering*.
- [43] Miryung Kim, Thomas Zimmermann, and Nachiappan Nagappan. 2012. A field study of refactoring challenges and benefits. In *Proceedings of the ACM SIGSOFT 20th International Symposium on the Foundations of Software Engineering (FSE '12)*. ACM, New York, NY, USA, Article 50, 11 pages. <https://doi.org/10.1145/2393596.2393655>
- [44] Miryung Kim, Thomas Zimmermann, and Nachiappan Nagappan. 2014. An Empirical Study of Refactoring Challenges and Benefits at Microsoft. *IEEE Transactions on Software Engineering* 40, 7 (2014), 1–1. <https://doi.org/10.1109/TSE.2014.2318734>

- [45] Ana Klimovic and Jason H Anderson. 2013. Bitwidth-optimized hardware accelerators with software fallback. In *2013 International Conference on Field-Programmable Technology (FPT)*. IEEE, 136–143.
- [46] David Koepfinger, Raghu Prabhakar, Yaqi Zhang, Christina Delimitrou, Christos Kozyrakis, and Kunle Olukotun. 2016. Automatic generation of efficient accelerators for reconfigurable hardware. In *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*. Ieee, 115–127.
- [47] Henry Oliver Lancaster and Eugene Seneta. 2005. Chi-square distribution. *Encyclopedia of biostatistics 2* (2005).
- [48] Dong-U Lee, Altaf Abdul Gaffar, Oskar Mencer, and Wayne Luk. 2005. MiniBit: bit-width optimization via affine arithmetic. In *Proceedings of the 42nd annual Design Automation Conference*. ACM, 837–840.
- [49] Tom Mens and Tom Tourwe. 2004. A Survey of Software Refactoring. *IEEE Transactions on Software Engineering* 30, 2 (2004), 126–139. <https://doi.org/10.1109/TSE.2004.1265817>
- [50] Giovanni De Micheli. 1994. *Synthesis and optimization of digital circuits*. McGraw-Hill Higher Education.
- [51] Emerson Murphy-Hill, Chris Parnin, and Andrew P. Black. 2009. How we refactor, and how we know it. In *ICSE '09: Proceedings of the 31st International Conference on Software Engineering*. IEEE Computer Society, Washington, DC, USA, 287–297. <https://doi.org/10.1109/ICSE.2009.5070529>
- [52] Nvidia. 2011. Nvidia CUDA C programming guide. *Nvidia Corporation* 120, 18 (2011), 8.
- [53] William F. Opydyke. 1992. *Refactoring Object-Oriented Frameworks*. Ph.D. Dissertation. University of Illinois, Urbana-Champaign, IL, USA. [citeseer.ist.psu.edu/opydyke92refactoring.html](http://citeseer.ist.psu.edu/opydyke92refactoring.html)
- [54] James L Peterson and Theodore A Norman. 1977. Buddy systems. *Commun. ACM* 20, 6 (1977), 421–431.
- [55] Nam Khanh Pham, Amit Kumar Singh, Akash Kumar, and Mi Mi Aung Khin. 2015. Exploiting loop-array dependencies to accelerate the design space exploration with high level synthesis. In *2015 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 157–162.
- [56] Louis-Noel Pouchet, Peng Zhang, Ponnuswamy Sadayappan, and Jason Cong. 2013. Polyhedral-based data reuse optimization for configurable computing. In *Proceedings of the ACM/SIGDA international symposium on Field programmable gate arrays*. ACM, 29–38.
- [57] Andrew Putnam, Adrian M Caulfield, Eric S Chung, Derek Chiou, Kypros Constantinides, John Demme, Hadi Esmaeilzadeh, Jeremy Fowers, Gopi Prashanth Gopal, Jan Gray, et al. 2014. A reconfigurable fabric for accelerating large-scale datacenter services. *ACM SIGARCH Computer Architecture News* 42, 3 (2014), 13–24.
- [58] Weikang Qiao, Jieqiong Du, Zhenman Fang, Michael Lo, Mau-Chung Frank Chang, and Jason Cong. 2018. High-throughput lossless compression on tightly coupled CPU-FPGA platforms. In *2018 IEEE 26th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*. IEEE, 37–44.
- [59] Dan Quinlan and Chunhua Liao. 2011. The ROSE source-to-source compiler infrastructure. In *Cetus users and compiler infrastructure workshop, in conjunction with PACT*, Vol. 2011. Citeseer, 1.
- [60] Brandon Reagen, Robert Adolf, Yakun Sophia Shao, Gu-Yeon Wei, and David Brooks. 2014. Machsuite: Benchmarks for accelerator design and customized architectures. In *2014 IEEE International Symposium on Workload Characterization (IISWC)*. IEEE, 110–119.
- [61] Zhenyuan Ruan, Tong He, Bojie Li, Peipei Zhou, and Jason Cong. 2018. ST-Accel: A high-level programming platform for streaming applications on FPGA. In *2018 IEEE 26th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*. IEEE, 9–16.
- [62] Cindy Rubio-González, Cuong Nguyen, Benjamin Mehne, Koushik Sen, James Demmel, William Kahan, Costin Iancu, Wim Lavrijsen, David H Bailey, and David Hough. 2016. Floating-point precision tuning using blame analysis. In *Proceedings of the 38th International Conference on Software Engineering*. ACM, 1074–1085.
- [63] Cindy Rubio-González, Cuong Nguyen, Hong Diep Nguyen, James Demmel, William Kahan, Koushik Sen, David H Bailey, Costin Iancu, and David Hough. 2013. Precimonious: Tuning assistant for floating-point precision. In *SC'13: Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*. IEEE, 1–12.
- [64] Kyle Rupnow, Yun Liang, Yinan Li, and Deming Chen. 2011. A study of high-level synthesis: Promises and challenges. In *2011 9th IEEE International Conference on ASIC*. IEEE, 1102–1105.
- [65] Giacinto Paolo Saggese, Antonino Mazzeo, Nicola Mazzocca, and Antonio GM Strollo. 2003. An FPGA-based performance analysis of the unrolling, tiling, and pipelining of the AES algorithm. In *International Conference on Field Programmable Logic and Applications*. Springer, 292–302.
- [66] Adrian Sampson, Pavel Panchevka, Todd Mytkowicz, Kathryn S. McKinley, Dan Grossman, and Luis Ceze. 2014. Expressing and Verifying Probabilistic Assertions. In *PLDI*.
- [67] Koushik Sen, Darko Marinov, and Gul Agha. 2005. CUTE: A Concolic Unit Testing Engine for C. In *Proceedings of the 10th European Software Engineering Conference Held Jointly with 13th ACM SIGSOFT International Symposium on Foundations of Software Engineering (ESEC/FSE-13)*. ACM, New York, NY, USA, 263–272. <https://doi.org/10.1145/1081706.1081750>
- [68] Alexey Solovyev, Marek S Baranowski, Ian Briggs, Charles Jacobsen, Zvonimir Rakamarić, and Ganesh Gopalakrishnan. 2018. Rigorous estimation of floating-point round-off errors with symbolic taylor expansions. *ACM Transactions on Programming Languages and Systems (TOPLAS)* 41, 1 (2018), 1–39.
- [69] Nitish Srivastava, Steve Dai, Rajit Manohar, and Zhiru Zhang. 2017. Accelerating Face Detection on Programmable SoC Using C-Based Synthesis. In *25th ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*.
- [70] Markus Steinberger, Michael Kenzel, Bernhard Kainz, and Dieter Schmalstieg. 2012. ScatterAlloc: Massively parallel dynamic memory allocation for the GPU. In *2012 Innovative Parallel Computing (InPar)*. IEEE, 1–10.
- [71] Mark Stephenson, Jonathan Babb, and Saman Amarasinghe. 2000. Bidwidth Analysis with Application to Silicon Compilation. In *Proceedings of the ACM SIGPLAN 2000 Conference on Programming Language Design and Implementation (PLDI '00)*. ACM, New York, NY, USA, 108–120. <https://doi.org/10.1145/349299.349317>
- [72] John E Stone, David Gohara, and Guochun Shi. 2010. OpenCL: A parallel programming standard for heterogeneous computing systems. *Computing in science & engineering* 12, 3 (2010), 66–73.
- [73] Sriraman Tallam and Neelam Gupta. 2005. A Concept Analysis Inspired Greedy Algorithm for Test Suite Minimization. In *Proceedings of the 6th ACM SIGPLAN-SIGSOFT Workshop on Program Analysis for Software Tools and Engineering (PASTE '05)*. ACM, New York, NY, USA, 35–42. <https://doi.org/10.1145/1108792.1108802>
- [74] David B Thomas. 2016. Synthesizable recursion for C++ HLS tools. In *2016 IEEE 27th International Conference on Application-specific Systems, Architectures and Processors (ASAP)*. IEEE, 91–98.
- [75] David B Thomas. 2019. Templatised Soft Floating-Point for High-Level Synthesis. In *2019 IEEE 27th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*. IEEE.
- [76] Mohsen Vakilian, Nicholas Chen, Stas Negara, Balaji Ambresh Rajkumar, Brian P. Bailey, and Ralph E. Johnson. 2012. Use, disuse, and misuse of automated refactorings. In *Software Engineering (ICSE), 2012 34th International Conference on*. 233–243. <https://doi.org/10.1109/ICSE.2012.6227190>
- [77] Paul R Wilson, Mark S Johnstone, Michael Neely, and David Boles. 1995. Dynamic storage allocation: A survey and critical review. In *International Workshop on Memory Management*. Springer, 1–116.
- [78] Xilinx. 2019. UltraScale Architecture and Product Data Sheet: Overview. [https://www.xilinx.com/support/documentation/data\\_sheets/ds890-ultrascale-overview.pdf](https://www.xilinx.com/support/documentation/data_sheets/ds890-ultrascale-overview.pdf). (2019).
- [79] Xilinx. 2019. Vivado High-Level Synthesis. <https://www.xilinx.com/products/design-tools/vivado/integration/esl-design.html>. (2019).
- [80] Xilinx. 2019. Xilinx Virtex UltraScale+ FPGA VCU1525. <https://www.xilinx.com/products/boards-and-kits/vcu1525-a.html>. (2019).
- [81] Zeping Xue and David B Thomas. 2015. SysAlloc: A hardware manager for dynamic memory allocation in heterogeneous systems. In *2015 25th International Conference on Field Programmable Logic and Applications (FPL)*. IEEE, 1–7.
- [82] Zeping Xue and David B Thomas. 2016. SynADT: Dynamic Data Structures in High Level Synthesis. In *2016 IEEE 24th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*. IEEE, 64–71.
- [83] Cody Hao Yu, Peng Wei, Max Grossman, Peng Zhang, Vivek Sarker, and Jason Cong. 2018. S2FA: an accelerator automation framework for heterogeneous computing in datacenters. In *Proceedings of the 55th Annual Design Automation Conference (DAC)*. ACM, 153.
- [84] Yuan Zhou, Udit Gupta, Steve Dai, Ritchie Zhao, Nitish Srivastava, Hanchen Jin, Joseph Featherston, Yi-Hsiang Lai, Gai Liu, Gustavo Angarita Velasquez, Wenping Wang, and Zhiru Zhang. 2018. Rosetta: A Realistic High-Level Synthesis Benchmark Suite for Software-Programmable FPGAs. *Int'l Symp. on Field-Programmable Gate Arrays (FPGA)* (Feb 2018).
- [85] Wei Zuo, Peng Li, Deming Chen, Louis-Noël Pouchet, Shunan Zhong, and Jason Cong. 2013. Improving polyhedral code generation for high-level synthesis. In *2013 International Conference on Hardware/Software Codesign and System Synthesis (CODES+ ISSS)*. IEEE, 1–10.