

# GPU-Accelerated High-Level Synthesis for Bitwidth Optimization of FPGA Datapaths

Nachiket Kapre  
School of Computer Engineering  
Nanyang Technological University  
Singapore 639798  
nachiket@ieee.org

Deheng Ye  
School of Computer Engineering  
Nanyang Technological University  
Singapore 639798  
ye0014ng@e.ntu.edu.sg

## ABSTRACT

Bitwidth optimization of FPGA datapaths can save hardware resources by choosing the fewest number of bits required for each datapath variable to achieve a desired quality of result. However, it is an NP-hard problem that requires unacceptably long runtimes when using sequential CPU-based heuristics. We show how to parallelize the key steps of bitwidth optimization on the GPU by performing a fast brute-force search over a carefully constrained search space. We develop a high-level synthesis methodology suitable for rapid prototyping of bitwidth-annotated RTL code generation using gcc's GIMPLE backend. For range analysis, we perform parallel evaluation of sub-intervals to provide tighter bounds compared to ordinary interval arithmetic. For bitwidth allocation, we enumerate the different bitwidth combinations in parallel by assigning each combination to a GPU thread. We demonstrate up to 10–1000× speedups for range analysis and 50–200× speedups for bitwidth allocation when comparing NVIDIA K20 GPU implementation to an Intel Core i5-4570 CPU while maintaining identical solution quality across various benchmarks. This allows us to generate tailor-made RTL with minimum bitwidths in hundreds of milliseconds instead of hundreds of minutes when starting from high-level C descriptions of dataflow computations.

## 1. INTRODUCTION

FPGAs have long empowered the circuit designer to build hardware tailored completely to the particular performance, cost and accuracy requirements of the application. A key aspect of this flexibility is the ability to select the data representation for signals in the circuit. As different applications have varying accuracy requirements, the exact number of bits necessary may change. This is in stark contrast to ISA-based processors where a small set of types are supported (char, short, int, long, float, double). While the low-level tuning of bits may seem time-consuming and excessive, it is a well-known approach [8, 15, 22, 23, 2, 1, 24] for saving resource costs by as much as 2–4× while improving cir-

cuit performance by as much as 20–30% without compromising quality of the output signals. In particular, with the promise of approximate computing, we must ask how many bits are actually needed to represent data (or signals) to deliver an acceptable computational outcome for our designs. The objective of the bitwidth optimization problem is to determine the cheapest data representation (fewest bits) for all variables in the computation subject to a user-supplied accuracy constraint. This can be achieved through static compile-time analysis of variable bounds and error.

While the benefits are clear, unfortunately, bitwidth allocation is an NP-hard problem [8]. Broadly speaking, we can classify existing approaches for solving this problem based on the focus of their optimization – (1) **tighter analysis of error** [2, 1, 16, 17] that focuses on fast arithmetic techniques for producing tight error bounds, and (2) **intelligent search algorithms** [15, 22, 23, 12] that use specially-formulated automated heuristics. To perform bitwidth allocation on CPUs, all heuristics typically examine a limited set of bitwidth combinations during the search process to keep runtimes low. More importantly, these approaches solve the problem in a sequential manner by refining the bitwidth combination in each sequential iteration while learning from the previous trials.

Commercial tools such as Matlab HDL Coder<sup>1</sup> provide a floating-point to fixed-point conversion toolflow but exposes the accuracy analysis to a simulation-driven workflow that requires developer involvement. LegUp [7] provides no precision analysis engine while Vivado HLS merely allows expression of templated types for fixed-point arithmetic without providing the necessary automation to analyze error and select bitwidths. A key limitation that prevents integration of automated analysis engines in these tools is the large computational cost of the automation which translates into minutes to hours of runtime for even simple dataflow blocks of code. In fact, pure CPU-based fixed-point simulations to determine bitwidths of a simple FIR filter with 10<sup>5</sup>-element-long input test vector can take 40 days [4] of runtime. In this paper, we provide an automated precision analysis engine for gcc based on GIMPLE backend. We also show how to speedup the optimization problem using GPUs with a carefully constrained brute-force approach backed by an intelligent pruning of the search space. The use of GPUs in the FPGA CAD process needs to expand and grow to help tackle the series of slow, NP-hard heuristics that have traditionally constrained the design development process. The bitwidth optimization problem is one such slow heuristic that needs

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.

FPGA'16, February 21–23, 2016, Monterey, CA, USA.

Copyright is held by the owner/author(s). Publication rights licensed to ACM.

ACM 978-1-4503-3856-1/16/02...\$15.00.

DOI:http://dx.doi.org/10.1145/2847263.2847266.

<sup>1</sup><http://www.mathworks.com/products/hdl-coder/>

Table 1: Different Phases of the Precision Analysis Flow for  $y = a * x^2 + b * x + c$ . Assume  $a, b, c$  are integer constants,  $W_i$  represents the number of bits required to represent variable  $i$

Dataflow Graph (DFG)	DFG Node	Range Analysis	Error Propagation	Resource Modeling
	$x$	$[0, 1]$	$2^{(-W_x-1)}$	$W_x$
	$a$	$[a, a]$	0	$W_a$
	$b$	$[b, b]$	0	$W_b$
	$c$	$[c, c]$	0	$W_c$
	$a * x$	$[0, a]$	$\max(a * 2^{(-W_x-1)}, 2^{-W_{m1}})$	$W_a * W_x$
	$b * x$	$[0, b]$	$\max(b * 2^{(-W_x-1)}, 2^{-W_{m2}})$	$W_b * W_x$
	$a * x^2$	$[0, a]$	$\max(a * 2^{-W_x}, 2^{-W_{m3}})$	$W_a * W_x^2$
	$a * x^2 + b * x$	$[0, a + b]$	$\max((a + b/2) * 2^{-W_x}, 2^{-W_{a1}})$	$\max(W_a * W_x^2, W_b * W_x)$
	$a * x^2 + b * x + c$	$[c, a + b + c]$	$\max((a + b/2) * 2^{-W_x}, 2^{-W_{a2}})$	$\max(W_a * W_x^2, W_b * W_x, W_c)$

faster evaluation for broader adoption and seamless integration with high-level synthesis (HLS). The key enabling idea here is rapid evaluation of dynamic range, error and resource costs estimates for multiple bitwidth combinations in data-parallel fashion. By choosing *where* to explore and by evaluating multiple candidate solutions in parallel, we can accelerate the time to optimized solution. Unlike multi-core CPUs, GPUs are well-suited for this kind of brute-force exploration as they support hundreds of thousands of concurrent threads. We formulate our search space in a manner that fits entirely within the GPU memory capacity while maximizing the likelihood of finding the optimal solution fast. We develop heuristics that prune the search space to bound variable bit width ranges based on error propagation. For large problem sizes, we use simulated annealing-based heuristics running on the CPU to constrain the search space before transferring the computation to the GPU for acceleration. This proposal is in the same spirit as end-case optimal layout algorithms for ASICs [5], where the final branch of the design space is solved optimally while the parent heuristic is still retained.

In this paper, we make the following key contributions:

- Development of a high-level synthesis flow suitable for rapid prototyping of bitwidth optimization transforms that compile C descriptions of dataflow computations into RTL by exploiting gcc’s GIMPLE backend.
- Development of optimized GPU-based kernels for (1) accelerated sub-interval analysis to derive higher dynamic range bounds, and (2) parallel error analysis and FPGA resource model kernels for bitwidth optimization.
- Engineering of CPU-based pruning heuristics that intelligently constrain the search to make it feasible to perform brute-force exploration on the GPU.
- Quantification of speed and quality of bitwidth optimization when comparing the NVIDIA K20 GPU to state-of-the-art simulated annealing-based heuristics on an Intel Core i5-4570 CPU across a variety of benchmarks.

## 2. BACKGROUND

Precision analysis is typically performed before detailed RTL design and pipelining. We must determine the exact number of bits required for various datapath variables and encode these into the RTL descriptions of our computation. This static analysis procedure can be automated and potentially integrated with high-level synthesis to generate tailor-made RTL for a given accuracy constraint. We now describe the basic concepts underlying precision analysis that are important to explain our parallelization approach. The precision analysis flow for bitwidth optimization can be broken down into a sequence of three steps (1) range analysis, (2) error propagation, and (3) resource modeling. We will use the example expression  $a \cdot x^2 + b \cdot x + c$  to illustrate this flow as shown in Table 1.

**Range Analysis:** We need to first identify the dynamic range of all variables in our computation which requires calculating the smallest and the largest values they may attain. This is essential as the implemented circuit only needs to be correct over this range which is typically much smaller than the range of real numbers (or floating-point for canonical implementations). The underlying algorithm to achieve this simply propagates the interval of the input iteratively through the feed-forward dataflow computation while evaluating the range at the output of each arithmetic operator. For example, in the polynomial expression of Table 1 (column Range Analysis), if  $x \in [0, 1]$ , then  $y \in [c, a + b + c]$  can be computed through a top-down range propagation using Interval Arithmetic (IA). However, purely relying on IA often produces bounds that are over-estimated (loose), because they ignore input correlations. Loose bounds affect our final solution quality by requiring more bits than strictly necessary to meet the error constraints. One option is to use Affine Arithmetic (AA) [9] which allows us to produce accurate estimates of intervals by identifying correlations in the inputs at the expense of tracking correlation terms that are proportional to the number of inputs to the expression (thereby trickier for parallelism). However, Affine Arithmetic (AA) formulations are trickier to parallelize on GPUs due to the growth in correlation terms. In this paper, we choose to use the well-known sub-interval analysis approach to address this limitation of traditional IA. This is achieved by splitting each input interval into multiple sub-

intervals, performing range analysis on all input sub-interval combinations, and composing the final range by merging the results of each sub-interval evaluation. This allows us to trivially exploit parallel processing capacity on the GPU across sub-intervals even while an individual sub-interval is sequentially evaluated. Our formulation using sub-interval arithmetic helps us generate tight bounds while exposing GPU-friendly parallel behavior.

**Error Propagation:** Once we know the range of each variable, we can estimate rounding and truncation errors due to each arithmetic operation. These errors depend on the number of fraction bits used in the variable representation. We assume sufficient number of integer bits are allocated for the largest fixed-point number to eliminate overflow errors. While many error models are available such as mean squared error (MSE) in signal processing systems, we use the worst-case quantization error models [21] for fixed-point implementations that are evaluated over operating intervals of the variable. For our implementation, we exact and verify the error model equations as expressed in the error analysis tool Gappa [3]. The errors are propagated from the input to output of the arithmetic expression based on the number of bits for each arithmetic operation as shown in Table 1 (column Error Propagation). Here, we observe that each arithmetic operation introduces an error term that is then carried forward downstream in the expression tree. Structurally this is similar to range propagation, but here we need to repeat this analysis for each proposed bitwidth combination. For a given combination, this is a sequential operation that processes error values from the inputs and calculates output error of that operation. CPU-based analysis tools are typically driven by heuristics such as simulated annealing that iteratively generate bitwidth combinations for evaluation and terminate once a sufficiently low-cost solution is found. Our insight here is that we can generate and evaluate multiple precision combination proposals for the variables in the arithmetic expression and process them in parallel. Again, this form of large-scale exploration is inappropriate on multi-core platforms due to the limited number of threads that are possible and the large scope of the combinations that must be considered. However, GPU-based systems with thousands of lightweight data-parallel threads are ideally suited to exploit this parallel pattern.

**Resource Modeling:** In this paper, our optimization problem is geared towards minimization of resource utilization of the implemented circuit. The goal is flexible, and we can substitute resource minimization for speed, latency or power if required (or some combined figure of merit). To support resource minimization, we develop highly-accurate analytical resource models of arithmetic operators by performing a complete FPGA implementation flow (synthesis, place and route). This accounts for internal optimizations and interactions between the various FPGA CAD stages. Based on these compilations, we build analytical resource models using regression-fit for LUT and DSP count usage as a function of the number of bits in the arithmetic operation using Weka [10], a popular data mining toolkit. This approach is in stark contrast to the simplistic models in [8, 15, 22, 23, 2, 1] and helps improve our optimization accuracy. For a given precision combination, this computation is a simple accumulation of cost estimates for each operation as shown in Table 1 (column Resource Modeling). For all precision combinations processed in parallel, we filter those

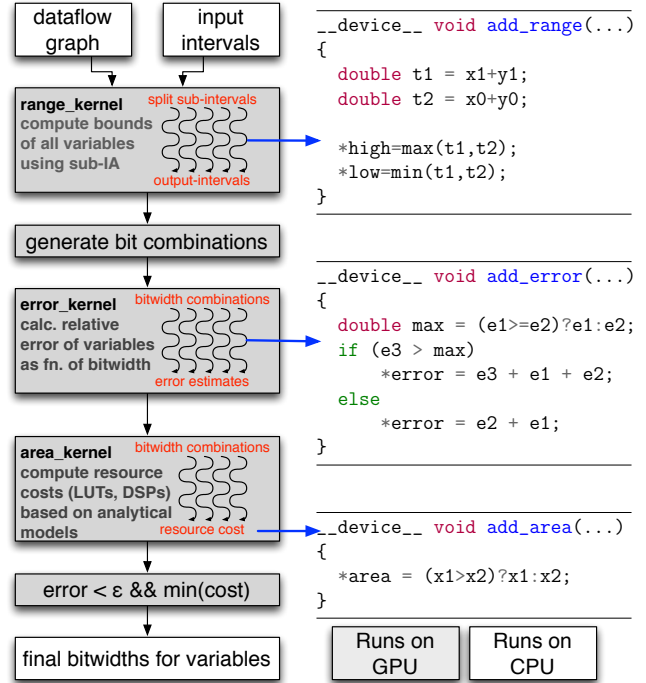


Figure 1: Compute Flow for GPU-Accelerated Precision Analysis (with sample GPU code for add operator)

that satisfy a user-supplied error constraint, and calculate the cheapest cost implementation using a GPU-optimized reduce operation.

### 3. GPU ACCELERATION

Bitwidth optimization on sequential CPUs is typically performed using heuristics such as simulated annealing (SA) that help drive the optimization towards a final solution while avoiding getting trapped in local minima through hill climbing. The state-of-the-art precision analysis tools [15, 22, 24] use this approach for formulating and tuning the optimization. They use the Adaptive Simulated Annealing (ASA) [11] package which is a highly-tuned optimization software that uses SA with intelligent selection of solution candidates to approach the optimal solutions. We use this state-of-the-art tool as reference for comparing time and solution quality.

In this section, we discuss the design of the key kernels used in our GPU approach that implements this search in a different manner. Unlike sequential CPU implementations of ASA, where a single combination of bitwidths is considered in a given iteration, on the GPU, we exploit the opportunity to consider multiple combinations in parallel.

#### 3.1 GPU Kernel Design

The parallel bitwidth allocation problem on GPUs is implemented across three kernels (1) sub-interval range analysis, (2) error propagation, and (3) resource modeling calculations. We show a high-level view of the flow in Figure 1 for the add operator. Each of these three kernels is data-parallel. In function `add_range`, `x0` and `y0` are the lower bounds of the two input operands respectively, while `x1` and `y1` are the upper bounds. In function `add_error`, `e1` and `e2`

are the inherent errors of the two operands respectively, and  $e_3$  is the propagation error. In function `add_area`,  $x_1$  and  $x_2$  are the area costs of the two operands respectively, derived from regression-fitted analytical models. Sub-interval Analysis is parallelized by trivially distributing the sub-interval combinations across the GPU threads while the Error propagation and Resource Modeling phases are parallelized across multiple bitwidth combinations that are considered. These approaches are *brute-force* approaches and are ultimately limited by the number of parallel threads that can run on your GPU and the amount of DRAM space available to store all intermediate results. With the availability of GPUs that can process thousands of threads in parallel and store GBs of data in the DRAM (e.g. NVIDIA K20 can handle  $\approx 26K$  threads and 5–6 GB of RAM), we can start tackling brute-force approaches to optimization problems with relative ease.

For range analysis, each GPU thread is programmed to evaluate a small sub-interval slice of the complete input range. For multiple inputs, this means we must consider all possible combinations of input sub-intervals across all threads. Each thread uses a unique thread index to automatically identify the sub-interval it is supposed to explore. An individual GPU thread propagates the sub-intervals through the graph to calculate the sub-ranges of the intermediate and output variables. Once this is done, we fuse the results of the multiple sub-intervals to form complete tight bounds on all variables in the program. While this procedure is highly-parallel, we are ultimately limited by the memory capacity required to hold all intermediate sub-ranges. For error analysis, we allocate each GPU thread to evaluate one combination of bitwidths. Each thread locally determines the combination of bitwidths it must explore based on the limit of precision choices possible for all variables and the unique thread index. We discuss how the precision limits for the variables are identified shortly in Section 3.2. Given a bitwidth combination, a thread can process the error equations and resource model expressions to determine the cumulative error of the output variable and resource costs. Finally, across all threads, we filter out those threads that do not satisfy the user-supplied error constraint. From this set of threads (precision combinations), we pick the thread with the cheapest implementation cost as the solution.

When developing our GPU code for arithmetic operators we derive our error models from [21] and [3], except we consider rounding to the nearest while they use truncation directly. Furthermore we adapt the models to better account for constant integer inputs. When using fixed-point to represent integer constants, we can have no rounding or truncation errors by choosing the bitwidth properly. As such, the number of bits required for such integers can be pre-determined. For area modeling of the arithmetic operators, instead of using analytical models based on theoretical expectations as used in [15], and [23], we use regression-fitted models derived from the real-world results of a complete FPGA CAD flow to capture the effect of CAD tools on implementation costs (LUTs and DSPs). In this regard, our approach requires a first-pass evaluation and mapping of all basic arithmetic building blocks through the FPGA CAD flow to help expose the DSP-LUT partitioning and allocation of resources that may differ from a simplistic theoretical models used in [15], and [23]. We quantify the impact of this accurate assessment of resource costs in Fig-

ure 4 later to show the resulting benefits. In these experiments, we observed a 30% improvement in prediction accuracy when considering the post-synthesis mapping results to select bitwidths.

In Table 2, we show the preliminary performance results for the range, error and area evaluation kernels on the GPU across various arithmetic operations. We observe 100–400× faster processing on the GPU when compared to OpenMP-optimized, multi-threaded (8-thread) implementation on the CPU. This initial experiment suggests high parallel potential for GPU evaluation of the constituent kernels. As we scale to larger GPUs, we expect this gap over multi-cores to grow as multi-core CPU trends do not align with the sheer degree of data-parallel thread capacity of modern GPU architectures and their rising DRAM capacity and associated bandwidths.

Table 2: NVIDIA K20 GPU Speedups over Intel Core i5-4570 CPU for Kernels used in Precision Analysis

Kernel	Speedup				
	add	mult	div	exp	log
Range analysis	312×	213×	119×	298×	254×
Error propagation	246×	80×	103×	261×	297×
Resource estimation	272×	266×	251×	414×	407×

### 3.2 Search Space Pruning

Instead of using sequential ASA-like heuristics directly on the GPU, we develop a GPU-optimized solution that uses a brute-force approach to explore multiple combinations in parallel. However, a naive brute force approach will quickly exhaust available GPU memory capacity, limiting the usefulness of this technique to toy problems. To allow the search space to fit within the memory limits of the GPU, we develop a pruning heuristic that restricts the combinations considered during the search. The pruning is also necessary for smooth execution of ASA running on the CPU. For our benchmark set introduced later in Section 5.1, ASA without pruning simply fails to converge and runs for hours before aborting. On the GPU, what pruning allows us to do is to optimally solve the subspace extracted from the pruning, unlike the multi-core CPUs which still explore this space sub-optimally based on heuristics. We show a high-level sketch of our pruning heuristic in Algorithm 1. Bitwidth combinations that result in error larger than the user-supplied error constraint are invalid. Our pruning heuristic must maximize the coverage of valid bitwidth combinations in our search. At the start, we choose a uniform fixed point bitwidth (*target\_fb* in line 1 in Algorithm 1) for all variables in the code and keep increasing precision until we satisfy the user-supplied error constraint (line 2–4 in Algorithm 1). The uniform bitwidth is designed to satisfy the worst case error and forces all variables to use the corresponding worst case precision. Thus, the uniform bitwidth *uniform\_bit* will be over-provisioned for some variables that could be implemented with fewer bits. Hence, we now decrease the bitwidth one variable at a time (line 6–12 in Algorithm 1) while keeping the other bitwidths equal to *uniform\_bit*. By doing this across all variables, we get the lowest possible precision (*lowest*) for each variable independently. We now replace each variable’s precision with their respective *lowest* but in this configuration we will most likely violate the error constraint. As the last step, we increase the precisions of

all variables simultaneously one bit at a time (line 14–15 in Algorithm 1) until it meets the required error criteria once again. In some instances, this pruning is excessively aggressive, and we relax the constraints by adding extra padding bits (*guard\_bit* in line 17 in Algorithm 1) to help cover potentially better solutions. Overall, our heuristic is able to compress the number of potential precision choices of each variable into ranges that make them feasible for brute-force exploration within the GPU memory space.

---

**Algorithm 1:** Search Space Pruning Heuristic

---

```

Data: The number of variables N; Targeted
         Fixed-point Precision
Result: Bounded search space
1 bit_width(0:N-1) ← target_fb;
2 while current_error > error_constraint do
3   | bit_width(0:N-1) ++;
4 end
5 uniform_bit = bit_width[0];
6 foreach n=0:N-1 do
7   | while current_error ≤ error_constraint do
8     | | bit_width(n)--;
9     | end
10    | lowest(n)← bit_width(n);
11    | bit_width(n)← uniform_bit;
12 end
13 bit_width(0:N-1) ← lowest(0:N-1);
14 while current_error ≤ error_constraint do
15   | bit_width(0:N-1)++;
16 end
17 highest(0:N-1) ← bit_width(0:N-1) + guard_bit;

```

---

For very large problems, this pruning is still insufficient to make the design feasible for brute-force exploration on the GPU. In these cases, we rely on preliminary cost-reducing moves in sequential ASA to assist in the pruning process. In this arrangement, we split the computation between CPU and GPU by allowing the first 10–100 iterations to run sequentially on the CPU and then switch to the GPU once the search space becomes feasible for GPU exploration. We keep the lower bound of the possible bitwidth identical as before and use the best observed precision combinations as the new upper bound for search. We achieved identical solution quality as the CPU-only solution after  $\approx 25$  iterations.

### 3.3 Overcoming GPU Limits

While the GPU is a great platform for rapid parallel evaluation of data-parallel problems, we consider several architectural optimizations that enhance performance of the FPGA CAD computations.

**GPU Global Memory Capacity:** A key constraint to consider is the size of the search space. For instance, on the NVIDIA K20 GPU with a 5GB DRAM, we are able to explore search spaces with precision combinations as large as  $2^{24}$  with a single GPU call. Similarly we were only able to scale from 1K–8K sub-intervals before we run out of memory resources to store intermediate sub-ranges.

**Kernel Fusion:** The GPU kernels for error propagation and resource modeling are logically separate and can be invoked independently. However, we achieve significant reductions in kernel invocation and synchronization time if we fuse them together into a single call.

**CPU-GPU Offload:** While it is tempting (and easy) to offload all data-parallel computations to the GPU, in our

case, we saw significant degradation in performance of the error kernel when evaluating truncation error of the form  $2^{-t}$  using the `pow` CUDA function. Since the possible values of precision we consider are limited ( $t < 128$  bits), we precompute these values offline on the CPU only once and simply pass them to the kernels as a lookup table.

**Memory Bandwidth:** To avoid needless data transfers between the CPU and GPU, we allocate sub-intervals, error and resource arrays directly in the GPU main memory. We only need to access summarized data, such as the unified final interval and the minimum cost implementation result from these large arrays. Furthermore, we identify thread-local state for storing per-variable structures into fast memory (registers and shared memory) instead of allocating them on the GPU global memory space. When we stored intermediate results in GPU shared memory space we were able to minimize needless memory traffic to the off-chip DRAM and obtain additional  $\approx 2\times$  speedup across our kernels.

## 4. HIGH-LEVEL SYNTHESIS TOOLFLOW

Contemporary attention in FPGA high-level synthesis is focussed on the robust and popular LLVM framework [14] (e.g. LegUp, VAST HLS). While LegUp [7] is a popular choice for development of HLS compiler transformations in the academic community, we choose an alternative methodology using `gcc`. We are able to rapidly prototype our optimization algorithms based on the `GIMPLE` backend in the `gcc` compiler, while retaining all the existing optimization benefits of the compiler flow. Transformations described on the three-register `GIMPLE` syntax can be developed independently of the compiler and need not be constrained to the language or development quirks of the compiler framework itself. Ultimately, once the technology is demonstrated and interest generated, the toolflow can and should be integrated with more robust tools such as LegUp, Vivado HLS or OpenCL compilers for broader distribution.

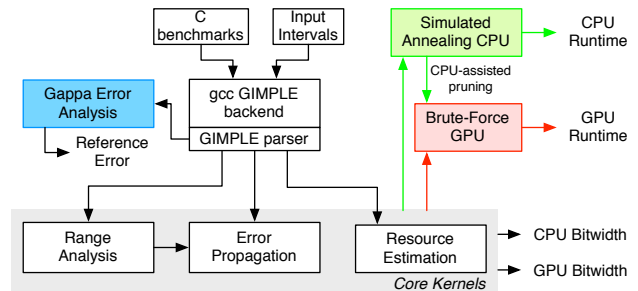


Figure 2: Compilation Flow for Precision Analysis

In Figure 2, we show a high-level block diagram of our compiler flow. We support simple feed-forward computations written in C through the `gcc`'s `GIMPLE` [6] backend. It is an Intermediate Representation (IR) provided to support development of plugins and optimization passes. For precision analysis, we process the C input files along with user-annotated range information to generate intermediate `GIMPLE` IR for post-processing. We then translate the `GIMPLE` IR into a suitable assembly-like, dataflow format for GPU

processing that captures the dependencies and range information in a compact data structure. We develop a generic interpreter on the GPU that iterates over this suitably encoded dataflow graph. This not only allows us to avoid needless recompilations on the GPU for each benchmark input, but enables optimized GPU kernel design that works across all benchmarks.

We verify the correctness of our GPU-calculated error bounds by generating Gappa scripts for the given benchmark. Once again, the simplicity of the GIMPLE IR makes it feasible to rapidly assemble a translator for Gappa syntax. Gappa [16] is a static analysis tool that proves numerical properties of programs through formal techniques. Our CPU-based ASA annealer also operates on the exact same GIMPLE IR and the exact same error and resource models as the GPU implementation for a fair comparison.

Table 3: Auto-generated Assembly and Gappa Code

(a) ASM Code	(b) Gappa Script
LD, 0, 0, 1	x.fx = <b>fx1</b> (x); a.fx = <b>fx2</b> (a);
LD, 1, 1, 1;	b.fx = <b>fx3</b> (b); c.fx = <b>fx4</b> (c);
LD, 2, 1, 1;	d.fx = <b>fx5</b> (a.fx * x.fx);
LD, 3, 1, 1;	e.fx = <b>fx6</b> (b.fx * x.fx);
MUL, 4, 0, 1;	f.fx = <b>fx7</b> (d.fx * x.fx);
MUL, 5, 0, 2;	g.fx = <b>fx8</b> (e.fx + f.fx);
MUL, 6, 0, 4;	y.fx = <b>fx9</b> (g.fx + c.fx);
ADD, 7, 5, 6;	y = a*x*x + b*x + c;
ADD, 8, 3, 7;	{ a in [1, 1] /\ b in [1, 1] /\
ST, 9, 8, -1;	c in [1, 1] /\ x in [0, 1]
	-> y.fx in ? /\ (y.fx - y) in ? }

In Table 3, we show the auto-generated GPU dataflow code and Gappa code for an example circuit  $y = a * x^2 + b * x + c$ . Here, we assume  $a, b, c$  are constants all with numerical value 1,  $x$  is an input variable with range  $[0, 1]$ . Each entry in the GPU dataflow code contains 4 fields for – instruction code, destination register, first source register, and second source register. The opcodes cover common arithmetic operations that we wish to analyze such as addition, multiplication, division, square-root, exponential, logarithm and conditionals. In addition, we support load and store instructions (shown as LD and ST respectively, in Table 3) that serve as a convenient place to pass in the various input intervals to the GPU threads. For CPU-based Gappa analysis, we generate a script that follows Gappa’s custom syntax. Gappa analysis requires custom type specifications for each variable precision as indicated by the **fx** type-casts for each operation. We also generate VHDL (not shown) from the same input C to synthesize pipelined hardware. The full FPGA CAD flow is only invoked once after the optimized bitwidth combination is determined.

#### 4.1 Compatibility with ASICs and need for FPGA-specific Flow

Like many tools developed for FPGAs, such as LegUp [7], our GPU-accelerated tool can also be made to work with ASICs while remaining a valuable and useful tool within the FPGA community. FPGA-based designs are likely to benefit to a greater extent from a complete per-variable customization of computation unlike ASIC-based designs that must factor in safety margins to handle changes to accuracy requirements after fabrication. Rapid prototyping in a high-level synthesis environment may be of stronger ap-

Table 4: Benchmark Problem Characteristics

Benchmark	#Vars	Arithmetic Operations	Inputs	Search Space
Level1 <sub>linear</sub> [24]	10	8	3	1K
Poly[15]	12	8	1	2K
Diode[24]	8	4	2	2K
Bellido[19]	13	9	3	24K
Approx1[24]	12	9	3	35K
Poly6[15]	19	12	1	52K
Level1 <sub>saturn</sub> [24]	14	10	3	2M
Caprasse[19]	16	10	4	8M
Poly8[15]	22	21	1	160M
Approx2[24]	19	17	4	15G

peal to a time-conscious FPGA developer than an ASIC developer who may also be slightly more tolerant of longer development cycles. An unconstrained bitwidth optimization problem already takes minutes to hours for modest-sized problems which is comparable to FPGA CAD times. Our approach reduces this to seconds or minutes to help make this particularly attractive as an optimization for an FPGA developer hunting for resource wins. To adapt our flow for ASICs, we would need to engineer an appropriate ASIC backend toolchain and construct resource and cost models for the particular technology.

## 5. EXPERIMENT SETUP

### 5.1 Benchmarks

We evaluate our framework using a variety of typical benchmarks that have been used frequently in previous range and precision analysis work [15], [24], [23]. These are taken from the Alias-COPRIN benchmark set, the Minibit circuits, Mix FX-SCORE examples and others. Our benchmark set contains a mixture of problems with multiple inputs and outputs, complex non-linear operations (e.g. exponential and logarithm) and varying search space sizes. We tabulate their characteristics in Table 4. While the variable count may seem small, these generate enormously large search spaces ( $\approx 15G$  data points) that are tricky to cover effectively. The key to enabling precision analysis as a routine optimization on larger problems (multiple basic blocks) in high-level synthesis flow is to make this exploration tractable – using GPUs as shown in this paper.

### 5.2 Tools and Hardware

We compare the sequential annealing (ASA) implementation on an Intel Core i5-4570 CPU @ 3.2GHz against our approach mapped to an NVIDIA K20 GPU. We use CUDA 6.5 along with Thrust library (version 1.7.0) for simplified transfer of data between CPU and GPU and also for efficient GPU reduction routines. For precision analysis on the CPU, we use two tools that we have mentioned previously (1) simulated annealing package ASA v30.15, and (2) numerical analysis package Gappa v1.1.1. ASA uses the same CPU versions of our GPU thread code for range analysis, error propagation and resource estimation. We modify Gappa to add support for **exp** and **log** operations. For FPGA compilations, we use Vivado Design Suite v2013.4 targeting the Kintex-7 XC7KLX160 FPGA to generate regression models for the various fundamental arithmetic operators. We use

off-the-shelf `gcc-4.8` compiler that ships with Ubuntu 14.04 but expect our toolflow to be compatible with even older versions of `gcc` (2007 onwards) but we have not tested this.

## 6. EVALUATION

In this section, we evaluate the experimental results of our GPU-assisted range analysis and bitwidth allocation. We show the results of FPGA resource utilization as a function of desired accuracy as well as the sensitivity of the FPGA mapping to fidelity of the resource models. For speedup calculations we compared optimized CPU and GPU implementations and include the pruning time on the CPU when calculating total GPU time for speedup calculations. Furthermore, our GPU timing calculations also include CPU-GPU memory transfer times.

### 6.1 Resource Utilization

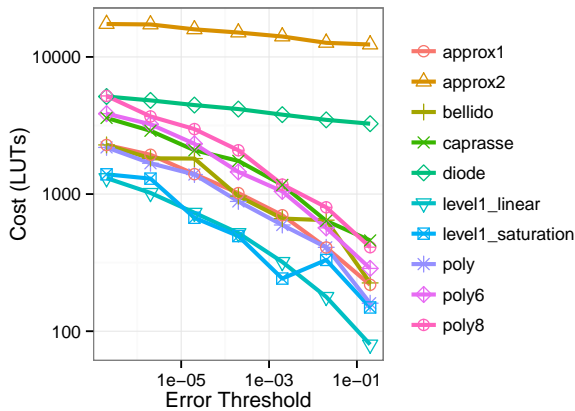


Figure 3: Impact of Error Threshold on FPGA Resource Utilization (LUTs). As we relax the error threshold, we observe that smaller designs are possible

When the resulting computation accuracy is flexible, we can achieve significant reductions in FPGA resource utilization. In our compilation framework, we are able to provide a desired level of accuracy when driving the search for datapath variable precision. We show the sensitivity of the solution quality as a function of the reference error threshold for valid solutions in Figure 3. As is evident, we can generate 2–10× reduction in LUT count when considering various desired error thresholds. We observe that most benchmarks show reduced implementation cost at varying error thresholds indicating significant resource savings are possible for real-world accuracy requirements. Certain benchmarks with `exp` and `log` operators do not show as impressive reductions in cost as seen from the `approx2` and the `diode` benchmark. In these instances, the larger resource utilization of the `exp` and `log` operator dominates the system. Thus, for applications where approximate computations are allowed, our compiler can deliver a suitably optimized solution that aims to deliver exactly the accuracy that is desired and no more. Overall, our bitwidth optimization saves FPGA resource costs by as much as 4× (mean 2.5×) when compared to a baseline double-precision implementation. These qualitative results are better than the ones reported in [15, 22,

24], as we use a superior pruning heuristic and cover every single combination within the pruned search space.

### 6.2 Resource Model Accuracy

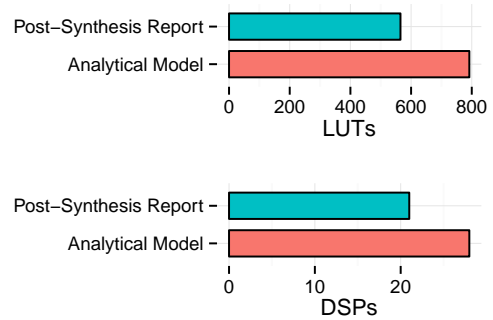


Figure 4: Comparing Quality of Results

As mentioned earlier in Section 3, due to the multi-stage nature of the FPGA CAD compilation flow, we need to build our resource model with particular care. We investigate the efficacy of using different area models on the final observed cost of the mixed precision results for the `level1-saturation` model. We compare: (1) approximate analytical area model, and (2) accurate post logic synthesis report. In each case, we run the full GPU-accelerated optimization to evaluate bitwidths and use the calculated resource numbers to filter out and select the best bitwidth. However, we generate final resource utilization after a complete place-and-route at the optimized bitwidths to compare the predicted model with actual real results. The optimizations considered at the different stages of the FPGA CAD process affect the accuracy of the predictive capabilities of the resource models in each of the CAD process. Ideally, we should use post-place-and-route result as the optimization goal of ASA, which requires invoking place-and-route calls inside ASA for every annealing iteration. This would lead to astronomically large runtimes. Therefore, we restrict our experiments until the post-logic synthesis stage. In Figure 4, we observe the reductions in LUTs and DSPs as we improve the fidelity of our resource model for `level1-saturation` by using downstream CAD tool results. If we perform bitwidth optimization solely relying on approximate analytical models, we need to spend 800 LUTs (post place-and-route). We can reduce this to 550 LUTs along with a DSP count reduction by 8 when we re-run the bitwidth optimization using the post-synthesis models. Logic synthesis tools are better at exploiting freedom of choice between LUTs vs. DSPs. Our post-synthesis models are better capable of exploiting this knowledge when driving the optimization.

### 6.3 Range Analysis

We first perform GPU-accelerated sub-interval analysis to improve the quality of the variable bounds. We obtain tighter interval bounds between 1–4× that of vanilla interval arithmetic across our benchmarks. In Figure 5, we show the performance impact of GPU acceleration on subinterval analysis as we vary the number of sub-intervals. Here, we compare optimized GPU runtime with sequential GPU Gappa runtimes when using Gappa’s *dichotomy search* feature when keeping the number of sub-intervals identical in

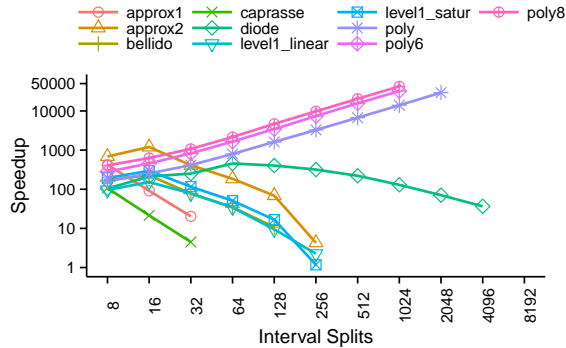


Figure 5: GPU Speedup for Sub-Interval Analysis (16-thread CPU runtimes can be minutes to hours)

both cases. While we observe peak speedups of  $3000\times$ , we achieved best pruning results below 128 sub-intervals with associated speedups of  $10\text{--}1000\times$ . Fewer sub-intervals generate limited parallelism for GPU acceleration while larger sub-intervals exceed the GPU storage limits and become infeasible. The maximum number of permissible sub-intervals for GPU-compatible acceleration varies with the number of inputs to the benchmark (listed in the **Inputs** column of Table 4). Benchmarks with few inputs 1–2 such as `poly8`, `poly6`, and `poly` deliver continued speedups upto 8192 sub-intervals while the rest saturate between 8–64 sub-intervals.

## 6.4 Bitwidth Allocation

We now show results for the bitwidth exploration phase of the problem. The CPU-based method typically ran for dozens of minutes on average and was reduced to dozens of seconds when using GPU acceleration. Certain instances of ASA did not terminate even after running for hours, and we inserted an early exit condition to restrict runtimes. These runtimes are for the small dataflow kernels listed in Table 4 which roughly correspond to the size of basic blocks in modern compilers. Larger programs in HLS must handle hundreds of such basic blocks resulting in a large overall runtime for the complete program when using CPUs alone. In Figure 6, we show the overall speedups for GPU accelerated brute-force exploration compared to sequential CPU implementation of ASA. We make the following key observations on analyzing the nature of speedups:

- For smallest benchmarks with very small search space sizes such as the `poly`, `diode`, and `level1_linear`, we observe that the single-shot GPU evaluation delivers a performance improvement greater than  $10\times$  without needing pruning assistance from ASA running on the CPU. This is to be expected and is unsurprising.
- For medium-sized benchmarks such as the `poly6`, `approx1`, `level1_saturation`, and `bellido`, the GPU-only approach delivers larger but still modest speedups. We record substantial speedup improvements when using the first 10–25 ASA iterations as pruning assistance. Switching to the GPU after  $\approx 25$  iterations generally delivers the best balance between spending sequential pruning time on the CPU and fast evaluation of smaller search space on GPU.

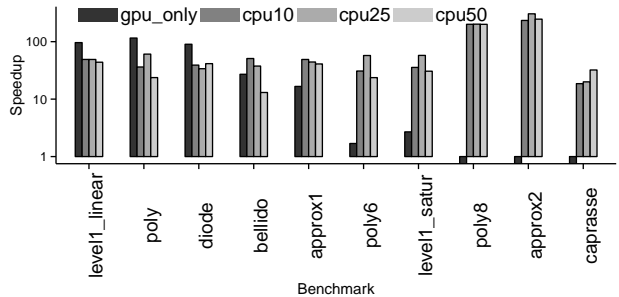


Figure 6: Bitwidth Allocation Speedups. Here `cpuX` means X iterations of ASA on CPU run before GPU invoked), CPU runtimes again run into dozens of minutes and sometimes ASA must be terminated early as it runs for hours without termination.

- For the largest benchmarks such as the `approx2`, `poly8` and `caprasse`, if we purely rely on the pruning heuristic in Algorithm 1, the search space is too large to fit in the GPU memory capacity. In this setting, the GPU-only approach is simply infeasible (shown as 0 speedup in Figure 6). For these cases, we use the ASA-assisted pruning to move the search space into a feasible region. We observe significant speedups due the sheer size of the search space to be explored on the CPU.

## 6.5 Quality-Time Tradeoffs

Finally, we show the quality-time tradeoff trends in Figure 7. Here, we report the best-observed resource cost ( $\text{cost} = 1/\text{quality}$ ) of the sequential ASA search after each trial as a function of CPU runtime across various benchmarks. It is useful to see the convergence towards the final results (normalized costs converge close to 0) often shows significant gains in the first 10% of the overall runtime. This enables us to utilize the fast convergence property of ASA to perform search space pruning, as shown in Section 6.4. While the runtime of ASA running on the CPU is tunable, the thousands of threads running on the GPU allow highly data-parallel search. For the GPU-accelerated evaluation, our one-shot GPU approach gives us the best answer as a single co-ordinate in this space at the same solution quality.

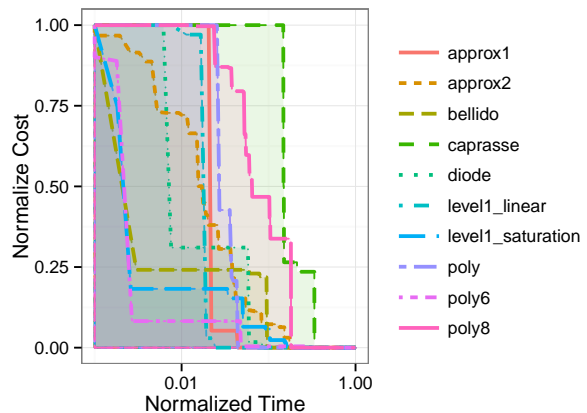


Figure 7: Quality-Time Tradeoffs for ASA (CPU-only)



## 7. DISCUSSION AND CONTEXT

### 7.1 Review of Previous Work

Research on range and precision analysis has evolved from simulation-based methods [13] to analytical approaches [15], [23], due to their speed in error modeling and range calculation. Interval Arithmetic (IA) [20] is the most well-known analytical method to calculate true bounds for general numerical algorithms. Some previous studies develop improved arithmetic techniques for producing more accurate error estimations. Lopez et al. [17] present quantized affine arithmetic (AA) [9] which only considers inputs' uncertainty into affine expansions, thereby reducing the complexity of AA representation. Kinsman et al. [12] use SAT-Modulo Theory (SMT) to break input ranges into sub-ranges for allocating optimal bit-width targeting a given precision specification, but is limited to range analysis only. Boland et al. [1] invent a polynomial algebra-based analytical approach to find provable and tight error bounds. This line of work mainly handles the accuracy issue of bit-width error modeling, but does not help with the optimization speed directly. Another set of work focuses on how to speed up the process of bit-width optimization using custom heuristics. In [15] and [24], the authors utilize simulated annealing for picking global optimized bits that minimize area cost. In [22], Osborne et al. perform partitioning followed by simulated annealing to improve runtime of [15]. [18] uses greedy algorithm for a coarse-grained search to find initial solutions, in tandem with Tabu search for refining the results' quality. Vakili et al. [23] propose two semi-analytical heuristics: progressive selection algorithm and tree-based search algorithm, dealing with fastness and optimality of precision analysis respectively. A GPU-based approach in the context of fixed-point analysis of VLSI circuits is presented in [4]. However, the speedups are compared against single-core CPUs on a single toy FIR-filter benchmark and the speedup is lower than the ones reported in this paper. They are only  $\approx 58\times$  which is at the low end of our speedups and our speedups are over a faster baseline 16-threaded implementation.

### 7.2 Case for GPU-based Brute-Force Acceleration

More broadly, the use of GPU-based acceleration provides a model for limited brute-force exploration in FPGA CAD. At multiple heuristics stages of the CAD flow (*i.e.* placement, routing, synthesis), we are often forced to make sub-optimal decisions due to our inability to fully explore potential solutions. A GPU-accelerated approach can open the door to using single-shot brute-force techniques at leaf stages of the search tree inherent in the CAD algorithms. In the context of ASIC cell layout, the idea of fast, optimal end-case placement exploration has been explored earlier in [5]. This approach still uses branch-and-bound at the leaf but optimally solves the small end-case solution when the problem size becomes small enough. In our approach, we explore all-possible combinations in the end-case using parallel GPU-based threads. In the long run, exploring the entire search space will stay intractable even with advances in GPU technology (as FPGA problem sizes will keep growing), we expect this technique to be used at stages of the search trees when the search spaces are sufficiently small for brute-force exploration and optimal solutions to the sub-problem are desirable.

## 8. CONCLUSIONS

In this paper, we show how to accelerate precision analysis for bitwidth optimization by 10–1000 $\times$  for range analysis and 50–200 $\times$  for bitwidth allocation and when comparing an NVIDIA K20 GPU to an Intel i5-4570 CPU. We demonstrate the design and engineering a high-level synthesis approach suitable for rapid prototyping of bitwidth optimization algorithms. Using our HLS framework with tuneable accuracy targets, we are able to reduce FPGA resource utilization by as much as 10 $\times$  if suitable approximate results are acceptable to the computation. To exploit GPU potential, we parallelize sub-interval analysis to improve interval analysis bounds and carefully prune the search space to enable single-shot exhaustive exploration of the search space. For range analysis, speedups correlated with the number of inputs and bitwidth allocation speedups tracked the size of the search space. As part of future work, we intend to extend this work to support Monte-Carlo sampling-based methods and integrate Affine Arithmetic models.

Our framework is available for download from <http://yedeheng.github.io/bitgpu>. The LLVM-based [14] LegUp license specifically restricts the use of the open-source compiler to non-commercial, not-for-profit development<sup>2</sup>. In contrast, our own tool is released under the GPLv3 copyleft license<sup>3</sup> and includes customized ASA code as part of the distribution. The ASA package is originally available under a permissive BSD 3-Clause license.

The Tesla K20 used for this research was donated by the NVIDIA Corporation.

## 9. REFERENCES

- [1] D. Boland and G. A. Constantinides. Automated precision analysis: A polynomial algebraic approach. In *FCCM*, pages 157–164, 2010.
- [2] D. Boland and G. A. Constantinides. A scalable approach for automated precision analysis. In *FPGA*, pages 185–194, 2012.
- [3] S. Boldo, J.-C. Filliâtre, and G. Melquiond. Combining coq and gappa for certifying floating-point programs. *Intelligent Computer Mathematics*, pages 59–74, 2009.
- [4] G. Caffarena and D. Menard. Many-core parallelization of fixed-point optimization of vlsi circuits through gpu devices. In *Design and Architectures for Signal and Image Processing (DASIP), 2012 Conference on*, pages 1–8, Oct 2012.
- [5] A. E. Caldwell, A. B. Kahng, and I. L. Markov. Optimal partitioners and end-case placers for standard-cell layout. *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on*, 19(11):1304–1313, 2000.
- [6] S. Callanan, D. J. Dean, and E. Zadok. Extending gcc with modular gimple optimizations. In *Proceedings of the 2007 GCC Developers' Summit*, pages 31–37, 2007.
- [7] A. Canis, J. Choi, M. Aldham, V. Zhang, A. Kammoona, J. H. Anderson, S. Brown, and T. Czajkowski. Legup: high-level synthesis for fpga-based processor/accelerator systems. In *Proceedings of the 19th ACM/SIGDA international*

<sup>2</sup><http://legup.eecg.utoronto.ca/license.php>

<sup>3</sup><https://www.gnu.org/licenses/licenses.html>

- symposium on Field programmable gate arrays*, pages 33–36. ACM, 2011.
- [8] G. Constantinides and G. Woeginger. The complexity of multiple wordlength assignment. *Applied Mathematics Letters*, 15(2):137–140, 2002.
- [9] L. H. De Figueiredo and J. Stolfi. Affine arithmetic: concepts and applications. *Numerical Algorithms*, 37(1-4):147–158, 2004.
- [10] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten. The weka data mining software: an update. *ACM SIGKDD explorations newsletter*, 11(1):10–18, 2009.
- [11] L. Ingber. Very fast simulated re-annealing. *Mathematical and computer modelling*, 12(8):967–973, 1989.
- [12] A. Kinsman and N. Nicolici. Finite precision bit-width allocation using sat-modulo theory. In *Design, Automation Test in Europe Conference Exhibition, 2009. DATE '09.*, pages 1106–1111, April 2009.
- [13] K.-I. Kum and W. Sung. Combined word-length optimization and high-level synthesis of digital signal processing systems. *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on*, 20(8):921–930, Aug 2001.
- [14] C. Lattner and V. Adve. Llvm: A compilation framework for lifelong program analysis & transformation. In *Code Generation and Optimization, 2004. CGO 2004. International Symposium on*, pages 75–86. IEEE, 2004.
- [15] D.-U. Lee, A. A. Gaffar, R. C. C. Cheung, O. Mencer, W. Luk, and G. A. Constantinides. Accuracy-guaranteed bit-width optimization. *IEEE Trans. on CAD of Integrated Circuits and Systems*, 25(10):1990–2000, 2006.
- [16] M. D. Linderman, M. Ho, D. L. Dill, T. H. Meng, and G. P. Nolan. Towards program optimization through automated analysis of numerical precision. In *Proc. IEEE/ACM Int. Symp. on Code Generation and Optimization*, CGO '10, pages 230–237, New York, NY, USA, 2010. ACM.
- [17] J. Lopez, C. Carreras, and O. Nieto-Taladriz. Improved interval-based characterization of fixed-point lti systems with feedback loops. *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on*, 26(11):1923–1933, Nov 2007.
- [18] D. Menard, N. Herve, O. Sentieys, and H.-N. Nguyen. High-level synthesis under fixed-point accuracy constraint. *JECE*, 2012, 2012.
- [19] J.-P. Merlet. The COPRIN benchmarks. <http://www-sop.inria.fr/coprin/logiciels/ALIAS/Benches/benches.html>.
- [20] R. E. Moore, R. B. Kearfott, and M. J. Cloud. *Introduction to Interval Analysis*. SIAM, 2009.
- [21] A. Nayak, M. Haldar, A. Choudhary, and P. Banerjee. Precision and error analysis of MATLAB applications during automated hardware synthesis for FPGAs. In *Proc. of Conf. on Design, Automation and Test in Europe*, DATE '01, pages 722–728, 2001.
- [22] W. G. Osborne, R. C. C. Cheung, J. Coutinho, W. Luk, and O. Mencer. Automatic accuracy-guaranteed bit-width optimization for fixed and floating-point systems. In *Int. Conf. on Field Programmable Logic and Applications (FPL)*, 2007.
- [23] S. Vakili, J. Langlois, and G. Bois. Enhanced precision analysis for accuracy-aware bit-width optimization using affine arithmetic. *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on*, 32(12):1853–1865, 2013.
- [24] D. Ye and N. Kapre. MixFX-SCORE: Heterogeneous fixed-point compilation of dataflow computations. In *Proceedings of the 2014 IEEE 22nd International Symposium on Field-Programmable Custom Computing Machines*, FCCM '14, 2014.