# Invited: Accelerator Design with Decoupled Hardware Customizations: Benefits and Challenges

Debjit Pal[1*], Yi-Hsiang Lai[4], Shaojie Xiang[1], Niansong Zhang[1], Hongzheng Chen[1], Jeremy Casas[2], Pasquale Cocchini[2], Zhenkun Yang[2], Jin Yang[2], Louis-Noël Pouchet[3], Zhiru Zhang[1*]

[1]Cornell University, [2]Strategic CAD Labs, Intel Corporation,[3]Colorado State University,[4]Amazon Web Services, Inc.

*{debjit.pal,zhiruz}@cornell.edu

## ABSTRACT

The past decade has witnessed increasing adoption of high-level synthesis (HLS) to implement specialized hardware accelerators targeting either FPGAs or ASICs. However, current HLS programming models entangle algorithm specifications with hardware customization techniques, which lowers both the productivity and portability of the accelerator design. To tackle this problem, recent efforts such as HeteroCL propose to decouple algorithm definition from essential hardware customization techniques in compute, data type, and memory, increasing productivity, portability, and performance.

While the decoupling of the algorithm and customizations provides benefits to the compilation/synthesis process, they also create new hurdles for the programmers to productively debug and validate the correctness of the optimized design. In this work, using HeteroCL and realistic machine learning applications as case studies, we first explain the key advantages of the decoupled programming model brought to a programmer to rapidly develop high-performance accelerators. Using the same case studies, we will further show how seemingly benign usage of the customization primitives can lead to new challenges to verification. We will then outline the research opportunities and discuss some of our recent efforts as the first step to enable a robust and viable verification solution in the future.

## 1 INTRODUCTION

Targeted specialization of functionality in hardware has become unarguably the best means to achieve improved compute performance and energy efficiency for a plethora of emerging applications. Unfortunately, it is a very unproductive practice to design and implement special-purpose accelerators using the conventional RTL methodology. For this reason, both academia and industry are seeing increasing use of HLS to automatically generate hardware accelerators from software programs [6, 8]. However, more widespread adoption of HLS is currently held back by its deficiencies in the quality of results (QoR) and ease of programming.

Programming high-performance hardware accelerators with HLS tools requires a deep understanding of hardware details and is a significant departure from traditional software programming. In

particular, current HLS programming models entangle algorithm specifications with hardware customization techniques. This approach has several limitations: (1) To achieve good QoRs, HLS programmers need to considerably restructure the source program to guide the HLS tool to realize specialized architectures such as data reuse buffers and systolic arrays; (2) Programmers are furthermore required to use various vendor-specific data types and pragmas [25], reducing the portability of accelerators across target hardware; (3) Existing HLS programming models fail to capture the interdependence among different hardware customization techniques, thus weakening the support of user-guided or automatic design space exploration (DSE).

One promising direction in modern language designs for heterogeneous computing is to apply the principle of separation of concerns. This principle refers to the decoupling of algorithm and performance optimization for enhanced modularity, composability, productivity, and performance. Halide is the first domain-specific language (DSL) to propose such decoupling for image processing applications [16]. TVM builds on Halide to support decoupled customization for machine learning applications [4]. Inspired by Halide and TVM, HeteroCL [9], T2S-Tensor [18], and SuSy [10] also separate algorithm definition from hardware customizations, aiming to make accelerator designs much more productive, performant, and portable. In the rest of this paper, we discuss how HeteroCL, an open-source Python-based programming abstraction for accelerator-rich computing that fully decouples the algorithm from hardware customization, increases the productivity of accelerator designs using two realistic case studies. We further outline several research challenges and opportunities that are worth pursuing to enable productive and pervasive hardware specialization.

## 2 DECOUPLING ALGORITHM FROM HARDWARE CUSTOMIZATIONS

In this section, we first introduce HeteroCL's decoupled programming model. Then, we outline how the programming model increases the productivity and performance of hardware accelerators using two case studies.

### 2.1 HeteroCL Accelerator Programming Model

HeteroCL is a multi-paradigm programming framework targeting accelerator-rich heterogeneous architectures. It is composed of a Python-based DSL and an automated compilation flow that maps the input algorithm into efficient accelerators. Similar to Halide [16] and TVM [4], HeteroCL separates an algorithm specification from a temporal compute schedule. Unlike the previous approaches that

mainly focus on CPUs/GPUs, HeteroCL further decouples the algorithm from memory architectures and data quantization schemes, which are both essential for efficient hardware customization, as shown in Figure 1. Another major advantage of HeteroCL is the mixed-paradigm support of both declarative tensor-based code and imperative programming constructs such as if-then-else and for loops. This offers two significant benefits: (1) The tensorized compute graph derived from the declarative code (extended from TVM) exposes many high-level domain-specific optimization opportunities to the compiler (*e.g.*, mapping stencil or systolic code to spatial architecture templates); (2) the flexibility of imperative code allows the description of more general-purpose algorithms with less-regular and fine-grained parallelism.
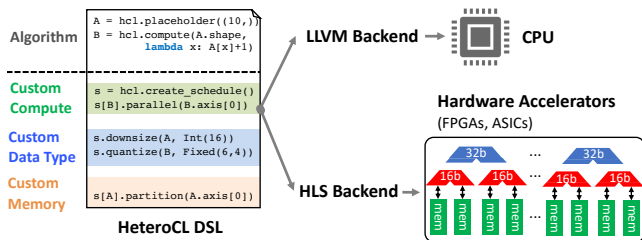


**Figure 1: Decoupled Hardware Customizations in HeteroCL.**

The current HeteroCL compilation framework has multiple backend supports, including CPU execution and HLS-based flows targeting Intel or AMD Xilinx FPGAs and ASIC accelerators using SystemC. It can produce highly efficient hardware implementations for a variety of popular image processing and DNN workloads by targeting spatial architecture templates such as systolic arrays and stencils with dataflow architectures.

**Table 1: Sample customization primitives in HeteroCL [9, 20]** – `s` is the temporal schedule of a HeteroCL program**.**

**(a) Compute Customization**

`s[kernel].unroll(axis, factor)`: unroll loop with target `factor`
`s[kernel].tile(i, factors)`: tile loop with `factors`

**(b) Memory Customization**

`s.reuse_at(tensor, kernel)`: create reuse buffer for tensor in `kernel`

**(c) Data Type Customization**

`s.quantize(tensor, dtype)`: quantize tensor to fixed-point type

**(d) Data Placement Customization**

`s.to(data, destination, mode)`: place data to destination using mode

Table 1(a)-(c) show a subset of HeteroCL customization primitives. In HeteroCL, a kernel contains a loop nest or function for computations. The decoupled customization primitives are applied to either the computation of a kernel (*i.e.*, compute customization) or the memory and data used by a kernel (*e.g.*, memory and data type customization).

Recently, HeteroFlow has extended HeteroCL to decouple data placement specification from the algorithm description and other hardware customizations [20]. Data placement refers to orchestrating the placement and movement of data between memory buffers inside the accelerator and between the accelerator and the CPU host. HeteroFlow provides a unified programming interface named `.to()` for customizing (1) host-accelerator data placement between

```
1  def ultranet(image, weight1, weight2, ...):
2    out1 = layer1_conv2d_im2col(image, weight1)
3    out2 = layer2_conv2d_im2col(out1, weight2)
4    out3 = layer3_conv2d_im2col(out2, weight3)
5    ...
6  s.to(out2, layer3_conv2d_im2col)
7  yo, yi, xo, xi = s[out3].tile(axis=[0,1], factor=[4,4])
8  s[out3].reorder(yo, xo, yi, xi)
9  PEs = s[out3].unroll(axis=[yi, xi])
10 for i in range(4):
11   s.to(out2[i][:].X, PEs[i,0]).to(PEs[i,1]).to(PEs[i,2])...
12   s.to(out2[:][i].W, PEs[0,i]).to(PEs[1,i]).to(PEs[2,i])...
13 for PE in PEs:
14   s[PE].vectorize(axis=PE.j, factor=32)
15   s.quantize(PE.X, PE,W], hf.Int(4))
```

| | # LUT/FF | # BRAM/DSP | Freq(MHz) | RT(ms) | LoC |
|---|---|---|---|---|---|
| Original HLS | 60.2K/39.6K | 377/508 | 231 | 2.97 (1.0x) | 2872 |
| HeteroCL+SA | 69.8K/39.4K | 375/594 | 233.8 | 2.27 (1.3x) | 204 |

**Figure 2: Evaluation on UltraNet in HeteroFlow** – RT refers to the total run-time, and LoC refers to the lines of code.

the CPU host memory and accelerator device memory; (2) interkernel data placement, to specify data streaming between different compute kernels within an accelerator; (3) intra-kernel data placement to specify fine-grained dataflow patterns commonly used in spatial architectures such as systolic arrays. The `.to()` primitive also enables programmers to co-optimize data placement schemes with other hardware customization techniques such as tiling and data quantization and integrate non-systolic kernels with optimized systolic arrays. Table 1(d) shows the usage of the `.to()` primitive.

## 2.2 Case Studies

To demonstrate HeteroCL's effectiveness in designing hardware accelerators, we discuss an application-specific accelerator and a domain-specific accelerator as case studies.

**Object Detection Accelerator**: UltraNet [21] is an object detection neural network implemented on FPGAs, and it is the winner of the DAC System Design Contest in 2020. Figure 2 shows its implementation in HeteroCL. UltraNet has 9 convolution layers (L1-4). In the decoupled customization, we map the third layer to a 4x4 output-stationary systolic array (L6-9), and connect it with neighborhood layers using FIFOs (L6). Each PE in the systolic array is vectorized and quantized to achieve better throughput (L11-15). The evaluation shows that the HeteorCL design achieves **1.3x** speedup of overall latency while only using **14.1x** fewer lines of code compared with the manually optimized HLS design. It is worth noting that the HeteroCL design employs decoupled customization, which makes it much easier for programmers to explore the trade-offs of different hardware architectures without altering the algorithm.

**Versatile Tensor Accelerator**: This case study showcases HeteroCL's capability of modeling a domain-specific accelerator design with a simple instruction set. We modeled an open-source deep-learning accelerator called Versatile Tensor Accelerator (VTA) [13] in HeteroCL. VTA is composed of four modules – fetch, load, compute, and store. There are three benefits of using HeteroCL to implement VTA. First, The HeteroCL programming model helped to quickly develop an executable specification for early software and compiler development. Especially, the declarative programming model allowed the programmer to quickly integrate tensor-based
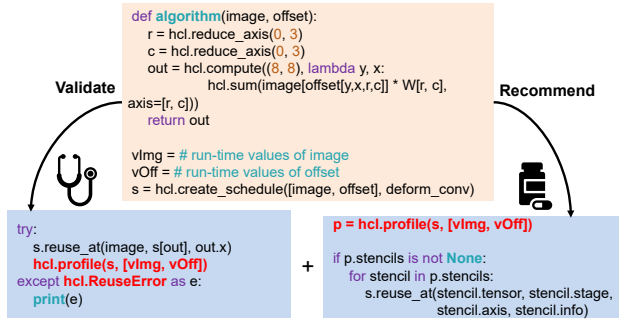
**Figure 3: Example of trace-based profiling.**

```
1 Conv2D = lambda Image, Filter: hcl.compute((N, M), lambda x, y:
2      hcl.sum(Image[x+r, y+c] * Filter[r, c], axis=[r, c]))
3 A = Conv2D(I, FA); B = Conv2D(A, FB); s = hcl.create_schedule()
4 s.to(A, s[B]).reuse_at(B.axis[0]) # streaming and reuse buffer
5 s[B].reorder(B.axis[1], B.axis[0]) # reorder outermost loops
```

**(a) HeteroCL code snippet**



**(b) Correct read-write access**          **(c) Wrong read-write access**

**Figure 4: A buggy HeteroCL example** – the RAW data dependency between two Conv2D kernels is violated after applying data streaming and loop reordering customization.

operations without worrying about implementation details of complex tensor operations. Second, the same HeteroCL executable specification was repurposed to derive a high-performance hardware implementation using decoupled customizations. Finally, the declarative programming, imperative programming, and the decoupled customizations together generated HLS code whose performance is similar to the original VTA with significantly fewer lines of code.

## 3 CHALLENGES AND OPPORTUNITIES

It is clear that with decoupled customizations, we can improve productivity, performance, and portability. However, there remain multiple challenges and opportunities with respect to the correctness verification of decoupled customizations and integration. In the following sections, we discuss two challenges and opportunities.

### 3.1 Automated Generation and Validation of Customization Primitives

The current HeteroCL compiler relies on user-specified customization primitives for optimizing the input program. To achieve high QoR, programmers need to have in-depth knowledge of the target accelerators, which creates a high threshold for most software programmers. Take memory customization as an example. Since programmers are more used to implicit memory orchestration such as caches on CPUs, it is non-trivial for them to explicitly design and manage custom memory hierarchy on FPGAs. Although HeteroCL provides customization primitives such as .reuse_at() to simplify the optimization process, it is still difficult for programmers to tell how and where to apply such a primitive. Even worse, the misapplication of primitives may worsen the performance or end up with incorrect results. Therefore, there is an urgent need for techniques that resolve the above challenges by providing programmers recommendations and validations.
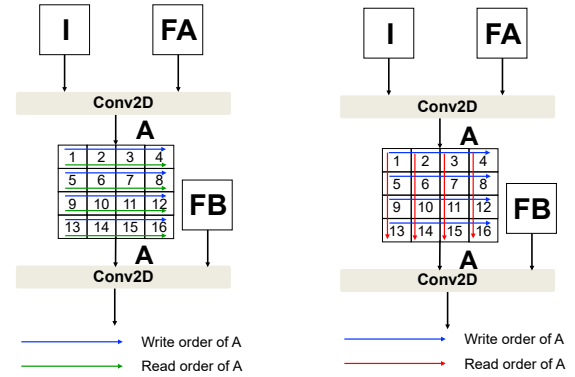
AutoTVM [5] and FlexTensor [24] leverage template-based methods to search the optimal parameters for the schedule for TVM. Ansor [23] and Halide's AutoScheduler [1] employ a template-free approach to directly generate schedule sequences from different primitive combinations. However, the above approaches only work for programs with static behaviors that are known at compile-time targeting CPUs or GPUs. There are also considerable number of works to automate DSE for HLS. Most of the existing effort develop different kinds of heuristics to automatically insert HLS pragmas into C/C++ programs [19]. For example, by leveraging graph learning models, a recent HLS DSE framework can make accurate QoR

prediction and achieve better performance [17]. While these works have considerable compute pattern optimizations, there are limited automation for data reuse and date access optimization.

We further propose a run-time trace-based profiling technique that provides automated validation and recommendation for application-specific hardware customization on FPGAs. Figure 3 shows an example of using the trace-based profiling techniques for validating and recommending data reuse by introducing a new primitive .profile() to HeteroCL. With decoupled customizations, programmers do not need to modify the algorithm specification (*i.e.*, the orange box). In addition, they can validate the specified data reuse primitives in combination (*i.e.*, the left blue box) and/or improve the QoR by applying the recommended primitives generated from the profiling results (*i.e.*, the right blue box). Moreover, with the trace-based technique, we can handle both regular and irregular data access patterns for memory customization.

### 3.2 Formal Verification of Decoupled Customizations

Although the decoupled hardware customizations in HeteroCL provide a multitude of benefits to the compilation/synthesis process they also create new hurdles for programmers to productively debug and validate the correctness of the transformed design. Figure 4a shows a HeteroCL design where the output from one convolution kernel is streamed to another (typically implemented by a FIFO in hardware). Such inter-kernel data streaming imposes a RAW dependency on FIFO for the correctness of the final result. We use .to() (L4) to specify data streaming and .reuse_at() (L4) to create a reuse buffer for the second convolution at the receiving end of the stream. One can easily break the design by (only) applying a loop reordering primitive (L6) to the second convolution of Figure 4a. As shown in Figure 4b, the write (→) and read (→) accesses of the streaming FIFO are in-order before applying the reorder primitive. However, after applying .reorder() to second convolution, the

write ($\rightarrow$) and read ($\rightarrow$) access orders of the streaming FIFO disagree as shown in Figure 4b, thereby violating RAW dependency. A seemingly benign usage of `.reorder()` breaks the correctness of the kernel functionality. While software simulation may show erroneous outputs, it is nontrivial to pinpoint the bug as the loop nests are implicit in declarative programming. Additionally, if the compiler incorrectly infers the size of the reuse buffer, the error may not manifest until the time-consuming hardware emulation is invoked. This would be particularly hard to debug by a programmer at the source level. Hence, we argue that there is an urgent need to develop new techniques and tools to automatically verify the correctness of the decoupled customizations specified in a modern DSL like HeteroCL. Specifically, the tool shall verify (1) the validity of a given sequence of decoupled customization primitives, and (2) the semantics equivalence between the original and transformed code after the compiler implements the specified customizations.

Proving the correctness of an optimizing compiler is an extremely hard task, as illustrated with the CompCert project [12]. Translation validation [22] is also difficult to deploy. A decoupled approach, where the validity of primitives is checked quickly by static analysis, independently of a later and more costly phase of validation of the generated code, can enable more efficient verification. There is a large class of programs typically candidates for acceleration that can be accurately represented using polyhedral compilation [15]. A large class of numerical computations with regular control flow such as dense linear algebra including tensor computations, image processing algorithms, n-dimensional convolutions, etc. can be exactly represented and analyzed in the polyhedral model (the affine dialect [11]), or over-approximated to it otherwise [3]. This enables the use of powerful and exact static analyses for data and control-flow dependences [7], and complex static and dynamic analyses for advanced program equivalence [2]. Many customizations can be expressed as loop transformations in the imperative form. It is therefore possible to formulate validity conditions for HeteroCL customizations directly in the polyhedral representation enabling real-time user feedback on the legality of customization during development. It also paves the way for automated techniques to generate legal customization sequences, inspired from results on legal transformation sets for loop-based programs [14].

## 4 CONCLUSION

In this paper, we have discussed the shortcomings of the HLS programming model, which seriously hinder its widespread adoption in designing high-performance hardware accelerators. We also discuss how emerging programming model like HeteroCL addresses those shortcomings and increases the productivity and performance of hardware accelerators using two realistic case studies. We outline various research challenges and opportunities that these new programming models present. Finally, we conclude with some of our recent efforts as initial steps to address those research challenges.

## ACKNOWLEDGEMENTS

## REFERENCES

[1] A. Adams, K. Ma, L. Anderson, R. Baghdadi, T.-M. Li, M. Gharbi, B. Steiner, and S. e. a. Johnson. Learning to Optimize Halide with Tree Search and Random Programs. *ACM Trans. on Graphics (TOG)*, 2019.

[2] W. Bao, S. Krishnamoorthy, and L.-N. e. a. Pouchet. PolyCheck: Dynamic Verification of Iteration Space Transformations on Affine Programs. *ACM SIGPLAN Symp. on Principles of Programming Languages (POPL)*, 2016.

[3] M.-W. Benabderrahmane, L.-N. Pouchet, A. Cohen, and C. Bastoul. The Polyhedral Model is More Widely Applicable than You Think. *Int'l Conf. on Compiler Construction (CC)*, 2010.

[4] T. Chen, T. Moreau, Z. Jiang, L. Zheng, E. Yan, M. Cowan, H. Shen, L. Wang, Y. Hu, L. Ceze, C. Guestrin, and A. Krishnamurthy. TVM: An Automated End-to-End Optimizing Compiler for Deep Learning. *USENIX Conf. on Operating Systems Design and Implementation (OSDI)*, 2018.

[5] T. Chen, L. Zheng, E. Yan, Z. Jiang, T. Moreau, L. Ceze, C. Guestrin, and A. Krishnamurthy. Learning to Optimize Tensor Programs. *Int'l Conference on Neural Information Processing Systems (NeurIPS)*, 2018.

[6] J. Cong, B. Liu, S. Neuendorffer, J. Noguera, K. Vissers, and Z. Zhang. High-Level Synthesis for FPGAs: From Prototyping to Deployment. *IEEE Trans. on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*, 2011.

[7] P. Feautrier. Dataflow Analysis of Array and Scalar References. *Int'l Journal of Parallel Program (JPP)*, 1991.

[8] Y. Lai, E. Ustun, S. Xiang, Z. Fang, H. Rong, and Z. Zhang. Programming and Synthesis for Software-defined FPGA Acceleration: Status and Future Prospects. *ACM Trans. on Reconfigurable Technology and Systems (TRETS)*, 2021.

[9] Y.-H. Lai, Y. Chi, Y. Hu, J. Wang, C. H. Yu, Y. Zhou, J. Cong, and Z. Zhang. HeteroCL: A Multi-paradigm Programming Infrastructure for Software-defined Reconfigurable Computing. *Int'l Symp. on Field-Programmable Gate Arrays (FPGA)*, 2019.

[10] Y.-H. Lai, H. Rong, S. Zheng, W. Zhang, X. Cui, Y. Jia, J. Wang, B. Sullivan, Z. Zhang, Y. Liang, et al. SuSy: A Programming Model for Productive Construction of High-Performance Systolic Arrays on FPGAs. *Int'l Conf. on Computer-Aided Design (ICCAD)*, 2020.

[11] C. Lattner, J. A. Pienaar, M. Amini, U. Bondhugula, R. Riddle, A. Cohen, T. Shpeisman, A. Davis, N. Vasilache, and O. Zinenko. MLIR: A Compiler Infrastructure for the End of Moore's Law. *arXiv*, 2020.

[12] X. Leroy. Formal Verification of a Realistic Compiler. *Commun. ACM*, 2009.

[13] T. Moreau, T. Chen, Z. Jiang, L. Ceze, C. Guestrin, and A. Krishnamurthy. VTA: An Open Hardware-Software Stack for Deep Learning. *CoRR*, 2018.

[14] L.-N. Pouchet, U. Bondhugula, C. Bastoul, A. Cohen, J. Ramanujam, P. Sadayappan, and N. Vasilache. Loop Transformations: Convexity, Pruning and Optimization. *ACM SIGPLAN Symp. on Principles of Programming Languages (POPL)*, 2011.

[15] L.-N. Pouchet, P. Zhang, P. Sadayappan, and J. Cong. Polyhedral-Based Data Reuse Optimization for Configurable Computing. *Int'l Symp. on Field-Programmable Gate Arrays (FPGA)*, 2013.

[16] J. Ragan-Kelley, C. Barnes, A. Adams, S. Paris, F. Durand, and S. Amarasinghe. Halide: A Language and Compiler for Optimizing Parallelism, Locality, and Recomputation in Image Processing Pipelines. *ACM SIGPLAN Notices*, 2013.

[17] A. Sohrabizadeh, Y. Bai, Y. Sun, and J. Cong. Automated Accelerator Optimization Aided by Graph Neural Networks. *Design Automation Conf. (DAC)*, 2022.

[18] N. Srivastava, H. Rong, P. Barua, G. Feng, H. Cao, Z. Zhang, D. Albonesi, V. Sarkar, W. Chen, P. Petersen, et al. T2S-Tensor: Productively Generating High-performance Spatial Hardware for Dense Tensor Computations. *IEEE Symp. on Field Programmable Custom Computing Machines (FCCM)*, 2019.

[19] Q. Sun, T. Chen, S. Liu, J. Miao, J. Chen, H. Yu, and B. Yu. Correlated Multi-objective Multi-fidelity Optimization for HLS Directives Design. *Design, Automation, and Test in Europe (DATE)*, 2021.

[20] S. Xiang, Y. Lai, Y. Zhou, H. Chen, N. Zhang, D. Pal, and Z. Zhang. HeteroFlow: An Accelerator Programming Model with Decoupled Data Placement for Software-Defined FPGAs. *Int'l Symp. on Field-Programmable Gate Arrays (FPGA)*, 2022.

[21] K. Zhan, J. Guo, B. Song, W. Zhang, and Z. Bao. UltraNet: An FPGA-based Object Detection for the DAC-SDC 2020. https://github.com/heheda365/ultra_net. Accessed: June 23, 2022.

[22] J. Zhao, S. Nagarakatte, M. M. Martin, and S. Zdancewic. Formalizing the LLVM Intermediate Representation for Verified Program Transformations. *ACM SIGPLAN Symp. on Principles of Programming Languages (POPL)*, 2012.

[23] L. Zheng, C. Jia, M. Sun, Z. Wu, C. H. Yu, A. Haj-Ali, Y. Wang, J. Yang, D. Zhuo, K. Sen, J. E. Gonzalez, and I. Stoica. Ansor: Generating High-Performance Tensor Programs for Deep Learning. *OSDI*, 2020.

[24] S. Zheng, Y. Liang, S. Wang, R. Chen, and K. Sheng. FlexTensor: An Automatic Schedule Exploration and Optimization Framework for Tensor Computation on Heterogeneous System. *Int'l Conf. on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, 2020.

[25] Y. Zhou, U. Gupta, S. Dai, R. Zhao, N. Srivastava, H. Jin, J. Featherston, Y.-H. Lai, G. Liu, G. A. Velasquez, W. Wang, and Z. Zhang. Rosetta: A Realistic High-Level Synthesis Benchmark Suite for Software Programmable FPGAs. *Int'l Symp. on Field-Programmable Gate Arrays (FPGA)*, 2018.