

# BraggHLS

Maksim Levental  
University of Chicago  
Email: test@test.tes

Ryan Chard  
Ecole Supérieure  
Nantes, France  
Email: second@second.fr

Kyle Chard  
and Ian Foster  
Star Academy  
San Francisco, California 99999-9999  
Telephone: (800) 555-5555  
Fax: (888) 555-5555

**Abstract**—In many experiment-driven scientific domains, such as high-energy physics, material science, and cosmology, very high data rate experiments impose hard constraints on the corresponding data acquisition systems: collected data must either be indiscriminately stored for post-processing and analysis, thereby necessitating large storage capacity, or accurately filtered in real-time, thereby necessitating low latency execution. Deep neural networks, effective in many other filtering tasks, have not been widely employed in such data acquisition systems, due to design and deployment difficulties. This paper presents an open source, lightweight, compiler framework BraggHLS, based on high-level synthesis techniques, for translating high-level representations of deep neural networks to low-level representations, suitable for deployment to near-sensor devices such as field-programmable gate arrays. We evaluate BraggHLS on various workloads and present a case-study implementation of a deep neural network for Bragg peak detection in the context of high-energy diffraction microscopy. We show BraggHLS is able to produce an implementation of the network with a throughput 4.8  $\mu$ s/sample, which is approximately a 4 $\times$  improvement over the existing implementation.

## CONTENTS

<b>I</b>	<b>Introduction</b>	<b>1</b>
<b>II</b>	<b>Background</b>	<b>2</b>
II-A	Compilers: the path from high to low . . .	2
II-A1	PyTorch and TorchScript . . . . .	3
II-A2	MLIR . . . . .	3
II-B	High-level synthesis and FPGA design . . .	3
II-B1	High-level synthesis . . . . .	3
II-C	FPGA design . . . . .	4
<b>III</b>	<b>BraggHLS compiler and HLS framework</b>	<b>5</b>
III-A	Symbolic interpretation for fun and profit .	5
III-B	AST transformations and behavioral verification . . . . .	6
III-C	Scheduling . . . . .	7
<b>IV</b>	<b>Evaluation</b>	<b>8</b>
IV-A	DNN layers . . . . .	8
IV-B	BraggNN case study . . . . .	9
<b>V</b>	<b>Conclusion</b>	<b>12</b>
	<b>References</b>	<b>12</b>
	<b>Appendix</b>	<b>13</b>

## I. INTRODUCTION

Very high data rates are observed and, consequently, large datasets are generated across a broad range of experiments in scientific domains, such as high-energy physics, material science, and cosmology. For example, in high-energy physics, the LHCb detector, at the CERN Large Hadron Collider, is tasked with observing the trajectories of particles produced in proton-proton collisions at a rate of 40 million per second (i.e., 40 MHz) [1]. With a packet size of approximately 50 kB (per collision), this implies a data rate of approximately 2 TB/s. Ultimately, in combination with other detectors, the LHC processes approximately 100 EB of data a year. In materials science, high-energy diffraction microscopy (HEDM) techniques, which provide non-destructive characterization of structure and its evolution in a broad class of single-crystal and polycrystalline materials, can have collection rates approaching 1 MHz [2], with a corresponding packet size of 80 kB. In cosmology, the Square Kilometer Array, a radio telescope projected to be completed in 2024 and to be operational by 2027 [3], will sustain data rates in excess of 10 TB/s [4].

Naturally, for high data rate experiments, directly storing and distributing such large quantities of data to the associated research communities for further analysis is cost prohibitive. Thus, either compression (in the case of storage and transmission) or outright filtering is necessary, i.e., only a small fraction of the most “interesting” data is selected at time of collection, with the remainder being permanently discarded. In this work we focus on the filtering approach. Note, the tradeoff made in employing filtering should be clear: reduced storage at the expense of more stringent latency constraints (on the filtering mechanisms). In addition, the risk of discarding meaningful data introduces accuracy (of the filtering mechanisms) as a critical new dimension of the data acquisition systems. Typically, these filtering mechanisms consist either of physics based models [5] or machine learning models [6]; in either case maximally efficient and effective use of the target hardware platform is tantamount to accuracy. Irrespective of the type of technique employed, almost universally, for the ultra-low latency use cases (e.g., sub-microsecond latency constraints), the implementation is deployed to either field-programmable gate arrays (FPGAs) or application-specific integrated circuits (ASICs) [7]. Here we focus primarily on FPGAs.

Deep neural networks (DNNs), a particular type of machine learning model, have been shown to be effective in many scientific and commercial domains due to their “representational capacity”, i.e., they demonstrate a capacity to (approximately) represent diverse sets of mappings [8]. DNNs “learn” to represent a mapping over the course of “training”, wherein they are iteratively evaluated on sample data while a “learning rule” periodically updates the parameters (*weights*) that parameterize the DNN. In recent years they have been investigated for near real-time scientific use cases [9], [10], [11] but their use for the lowest latency use cases has been very limited [7]. The reasons for this are threefold:

- 1) Graphics Processing Units (GPUs), the conventional hardware target for DNNs, until very recently, have not been performant enough for these very high data rate, very low latency, use cases (due to their low clock speeds and low peripheral bandwidth [12]);
- 2) DNNs, by virtue of their depth, are resource intensive, in terms of both memory (for the weights) and compute (floating-point arithmetic), thereby preventing their deployment to FPGAs, which, in particular, have limited static RAM available;
- 3) DNNs are (typically) defined, trained, and distributed using high-level frameworks (such as PyTorch [13], TensorFlow [14], MXNet [15]), which abstract all implementation details from the user, thereby making portability of existing model architectures (to e.g., FPGA) nigh impossible.

These three barriers demand of a solution that can simultaneously translate a high-level representation of a DNN to a low-level representation, suitable for deployment to FPGA, while optimizing resource usage and minimizing latency. In general, the task of *lowering* high-level representations of programs to lower-level representations is the domain of a compiler. Similarly, the task of *synthesizing* a *register-transfer level* (RTL) *design*, rendered in a *hardware description language* (HDL), from a program, is the domain of high-level synthesis (HLS) [16]. While several such HLS tools exist [17], [18], [19] and despite, often, bundling robust optimizing compilers, they struggle to effectively perform the necessary optimizations in reasonable amounts of time (see Section II-B1).

Recently, deep learning compilers (such as TVM [20], MLIR [21], and Glow [22]) have demonstrated the ability to dramatically reduce inference latencies [23], training times [24], and memory usage [25] of DNNs. These compilers function by extracting intermediate-level representations (IRs) of the DNNs, from the representations produced by the frameworks, and performing various optimizations on those IRs (such as kernel fusion [26], vectorization [27], and memory planning [25]). The highly optimized IR is then used to generate code for various target hardware platforms. Given the successes of these compilers, it’s natural to wonder whether they can be adapted to the task of sufficiently optimizing a DNN such that it might be synthesized to RTL, for deployment to FPGA.

In this paper, we present BraggHLS, an open source, lightweight, compiler and HLS framework which can lower DNNs defined as PyTorch models to FPGA compatible implementations. BraggHLS uses a combination of compiler and HLS techniques to compile the entire DNN into fully scheduled RTL, thereby eliminating all synchronization overheads and achieving ultra-low latency. BraggHLS is general and supports a wide range of DNN layer types, and thus a wide range of DNNs, but we particularly focus on optimizations relevant to a DNN designed for identifying Bragg diffraction peaks. In summary our specific contributions include:

- 1) We describe and implement a compiler framework, BraggHLS, which can effectively transform unoptimized, hardware-agnostic PyTorch models into ultra-low latency RTL designs suitable for deployment to Xilinx FPGAs. BraggHLS is thoroughly tested, open source, and available at <https://github.com/makslevental/bragghls/>;
- 2) We show that designs generated by BraggHLS achieve lower latency than Xilinx’s state-of-the-art commercial HLS tool (Vitis HLS) for a variety of DNN layer types. In particular we show that BraggHLS can produce synthesizable designs that meet placement, routing, and timing constraints for BraggNN, a DNN designed for identifying Bragg diffraction peaks;
- 3) We discuss some of the challenges faced even after successful synthesis of RTL from a high-level representation of a DNN, namely during the place and route phases of implementation.

The rest of this paper is organized as follows: Section II reviews key concepts from compilers, high-level synthesis, and RTL design for FPGA. Section III describes the BraggHLS compiler and HLS framework in detail. Section IV evaluates BraggHLS’s performance, scalability, and competitiveness with designs generated by Vitis HLS. Section IV-B describes our case study, i.e., BraggHLS applied to BraggNN, a Bragg peak detection DNN with a target latency of 1  $\mu$ s/sample. Finally, Section V concludes with a summary, and related and future work.

## II. BACKGROUND

### A. Compilers: the path from high to low

The path from a high-level, abstract, representation of a DNN to a register-transfer level representation can be neatly formulated as a series of progressive lowerings between adjacent levels of abstraction. Each level of abstraction is rendered as a programming language, IR, or HDL, and thus we describe each lowering in terms the representations and tools BraggHLS employ in manipulating those representations:

- 1) An imperative, *define-by-run*, Python representation, in PyTorch;
- 2) High-level data-flow graph representation, in TorchScript;
- 3) Low-level data and control flow graph representation, in MLIR.

1) *PyTorch and TorchScript*: Typically DNN models are represented in terms of high-level frameworks, themselves implemented within general purpose programming languages. Such frameworks are widely used because of their ease of use and large library of example implementations of various DNN model architectures. BraggHLS targets the PyTorch framework, thus we focus on relevant aspects of PyTorch. DNNs developed within PyTorch are *defined-by-run*: the author imperatively describes the DNN in terms of high-level operations, using python, which when executed materializes the (partial) high-level data-flow graph (DFG) corresponding to the DNN (e.g., for the purposes of reverse-mode automatic differentiation). From the perspective of the user, define-by-run enables fast iteration at development time, possibly at the cost of some runtime performance.

On the other hand, from the perspective of compilation, define-by-run precludes efficient extraction of the high-level DFG; since the DFG is materialized only at runtime, it cannot be inferred from the textual representation (i.e., the python source) of the DNN. Furthermore, a priori, the runtime-materialized DFG is only partially materialized<sup>1</sup>, and only as an in-memory data structure. Thus, framework support is necessary for efficiently extracting the full DFG. Indeed, PyTorch supports a Single Static Assignment (SSA) IR, called TorchScript (TS) IR and accompanying tracing mechanism (the TS JIT), which produces TS IR from conventionally defined PyTorch models. Lowering from PyTorch to TS IR enables various useful analyses and transformations on a DNN at the level of the high-level DFG but targeting FPGAs requires a broader collection of transformations. To this end, we turn to a recent addition to the compiler ecosystem.

2) *MLIR*: Multi-level Intermediate Representation [21] presents a new approach to building reusable and extensible compiler infrastructure. MLIR is composed of a set of *dialect* IRs, subsets of which are mutually compatible, either outright or by way of translation/legalization. The various dialects aim to capture and formalize the semantics of compute intensive programs at varying levels of abstraction, as well as namespace related sets of IR transformations. The entrypoint into this compiler framework, from PyTorch, is the `torch` dialect [28], a high-fidelity mapping from TS IR to MLIR native IR, which, in addition to performing the translation to MLIR, fully refines all shapes of intermediate tensors in the DNN (i.e., computes concrete values for all dimensions of each tensor); this is necessary for downstream optimizations and eliminating inconsistencies in the DNN [29].

While the `torch` dialect is necessary for lowering to MLIR and shape refinement, it is a representation of a DNN at the same level of abstraction as TS IR: it does not capture the precise data flow and control flow necessary for de novo implementations of DNN operations (e.g., for FPGA). Fortunately, MLIR supports lower-level dialects, such as the `linalg`, `affine` and `scf` (structured control flow) dialects. The

`scf` dialect describes standard control flow primitives, such as conditionals and loops, and is mutually compatible with the `arith` (arithmetic operations) and `memref` (memory buffers) dialects. The `affine` dialect, on the other hand, provides a formalization of semantics that lend themselves to polyhedral compilation techniques [30], i.e., techniques that enable loop dependence analysis and loop transformations. Such loop transformations, particularly loop unrolling, are crucial for achieving lowest possible latencies [31] because they directly inform the concurrency and parallelism of the final RTL design.

## B. High-level synthesis and FPGA design

1) *High-level synthesis*: High-level synthesis tools produce RTL descriptions of designs from high-level representations, such as C or C++ [17], [19]. In particular, Xilinx’s Vitis HLS, based on the Autopilot project [18], is a state-of-the-art HLS tool. Given a high-level, procedural, representation, HLS carries out three fundamental tasks, in order to produce a corresponding RTL design:

- 1) HLS schedules operations (such as `mulf`, `addf`, `load`, `store`) in order to determine which operations should occur during each clock cycle; such a schedule depends on three characteristics of the high-level representation:
  - a) The topological ordering of the DFG/CFG of the procedural representation (i.e., the dependencies of operations on results of other operations and resources);
  - b) The completion time for each operation;
  - c) The user’s desired clock rate/frequency;
- 2) HLS associates (called *binding*) operations to particular RTL instantiations of intellectual property (IP) for those operations; for example whether to associate an addition operation followed by a multiply operation to IPs for each, or whether to associate them both with a single IP, designed to perform a “fused” multiply-accumulate (MAC);
  - a) In the case of floating-point arithmetic operations, HLS also (with user guidance) determines the precision of the floating-point representation;
- 3) HLS builds a finite-state machine (FSM) that implements the schedule of operations as control logic, i.e., logic that initiates operations during the appropriate stages of the schedule.

In addition to fulfilling these three fundamental tasks, high-level synthesis aims to optimize the program. In particular, HLS attempts to maximize concurrency and parallelism (number of concurrent operations scheduled during a clock-cycle) in order maximize the throughput and minimize the latency of the final implementation. Maximizing concurrency entails pipelining operations: operations are executed such that they overlap in time, subject to available resources. Maximizing parallelism entails partitioning the DNN into subsets of operation that can be computed independently and simultaneously and whose results are combined upon completion.

<sup>1</sup>“...instead, every intermediate result records only the subset of the computation graph that was relevant to their computation.” [13]

```

def conv2d(
    input: MemRef(b, cin, h, w),
    output: MemRef(b, cout, h, w),
    weight: MemRef(cout, cin, k, k)
):
    for i1 in range(0, b):
        for i2 in range(0, cout):
            for i3 in range(0, h):
                for i4 in range(0, w):
                    for i5 in range(0, cin):
                        for i6 in range(0, k):
                            for i7 in range(0, k):
                                _3 = i3 + i6
                                _4 = i4 + i7
                                _5 = input[i1, i5, _3, _4]
                                _6 = weight[i2, i5, i6, i7]
                                _7 = output[i1, i2, i3, i4]
                                _8 = _5 * _6
                                _9 = _7 + _8
                                output[i1, i2, i3, i4] = _9

```

Listing 1: Python representation of a padding  $\lfloor k/2 \rfloor$ , stride 1,  $c_{out}$  filter convolution with  $k \times k$  kernel applied to  $(b, c_{in}, h, w)$ -dimensional input tensor, where  $b$  is the batch size,  $c_{in}$  is the number of channels, and  $(h, w)$  are the height and width, respectively.

While HLS aims to optimize various characteristics of a design, there are challenges associated with this kind of automated optimization. In particular, maximum concurrency and parallelism necessitates data-flow analysis in order to identify data dependencies amongst operations, both for scheduling and identifying potential data hazards. Such data-flow analysis is expensive and grows (in runtime) as better performance is pursued. This can be understood in terms of loop-nest representations of DNN operations; for example consider a convolution as in Listing 1. A schedule that parallelizes (some of) the arithmetic operations for this loop nest can be computed by first unrolling the loops up to some “trip-count” and then computing the topological sort of the operations. Using this scheduling algorithm (known as *list scheduling*), the degree to which the loops are unrolled determines how many arithmetic operations can be scheduled in parallel. The issue is that the stores and loads on the output array prevent reconstruction of explicit relationships between the inputs and outputs of the arithmetic operations across loop iterations. The conventional resolution to this loss of information is to perform *store-load forwarding*: pairs of store and load operations on the same memory address are eliminated, with the operand of the store forwarded to the uses of the load (see Listing 2). In order for this transformation to be correct (i.e., preserve program semantics), for each pair of candidate store and load operations, it must be verified that there are no intervening memory operations on the same memory address. These verifications are non-trivial since the iteration spaces of the loops need not be regular; in general it might involve solving a small constraint satisfaction program [32]. Furthermore, the number of such verifications grows

```

1 def conv2d(
2     input: MemRef(b, cin, h, w),
3     output: MemRef(b, cout, h, w),
4     weight: MemRef(cout, cin, k, k)
5 ):
6     for i1 in range(0, b):
7         for i2 in range(0, cout):
8             for i3 in range(0, h):
9                 for i4 in range(0, w):
10                    ...
11                    # e.g., i5, i6, i7 = 2, 3, 4
12                    _31 = i3 + i6
13                    _41 = i4 + i7
14                    _51 = input[i1, i5, _31, _41]
15                    _61 = weight[i2, i5, i6, i7]
16                    _71 = output[i1, i2, i3, i4]
17                    _81 = _51 * _61
18                    _91 = _71 + _81
19                    output[i1, i2, i3, i4] = _91
20                    # i5, i6, i7 = 2, 3, 5
21                    _32 = i3 + i6
22                    _42 = i4 + i7
23                    _52 = input[i1, i5, _32, _42]
24                    _62 = weight[i2, i5, i6, i7]
25                    _72 = output[i1, i2, i3, i4]
26                    _82 = _52 * _62
27                    _92 = _72 + _82
28                    output[i1, i2, i3, i4] = _92
29                    ...

```

Listing 2: Store-load forwarding across successive iterations (e.g.,  $i7 = 4, 5$ ) of the inner loop in Listing 1, after unrolling. The forwarding opportunity is from the store on line 19 to the load on line 25; both can be eliminated and  $\_91$  can replace uses of  $\_72$ , such as in the computation of  $\_92$  (and potentially many others).

polynomially in the parameters of the convolution since the loop nest unrolls into  $b \times c_{out} \times h \times w \times c_{in} \times k^2$  store-load pairs on the output array.

Finally, note, though greedy solutions to the scheduling problem solved by HLS are possible, in principle scheduling is an integer linear programming problem (ILP), instances of which are NP-hard. In summary, HLS tools solve computationally intensive problems in order to produce a RTL description of a high-level representation of a DNN. These phases of the HLS process incur “development time” costs (i.e., runtime of the tools) and impose practical limitations on the amount of design space exploration (for the purpose of achieving latency goals) which can be performed. BraggHLS addresses these issues by enabling the user to employ heuristics (during both the parallelization and scheduling phases) which, while not guaranteed to be correct, can be *behaviorally verified* (see Section III-A).

### C. FPGA design

Broadly, at the register-transfer level of abstraction, there remain two more steps prior to being able to actually deploy a design to an FPGA; one of them being a final lowering, so-called logic synthesis, and the other being place and route



(P&R). The entire process is carried out, for example, by Xilinx’s Vivado tool.

Logic synthesis is the process of mapping RTL to actual hardware primitives on the FPGA (so-called *technology mapping*), such as lookup tables (LUTs), block RAMs (BRAMs), flip-flops (FFs), and digital signal processors (DSPs). Logic synthesis produces a network list (*netlist*) describing the logical connectivity of various parts of the design. Logic synthesis effectively determines the implementation of floating-point operations in terms of DSPs; depending on user parameters and other design features, DSP resource consumption for floating-point multiplication and addition can differ greatly. The number of LUTs and DSPs that a high-level representation of a DNN corresponds to is relevant to both the performance and feasibility of that DNN when deployed to FPGA.

After the netlist has been produced, the entire design undergoes P&R. The goal of P&R is to determine which configurable logic block within an FPGA should implement each of the units of logic required by the digital design. P&R algorithms need to minimize distances between related units of functionality (in order to minimize wire delay), balance wire density across the entire fabric of the FPGA (in order to reduce route congestion), and maximize the clock speed of the design (a function of both wire delay, logic complexity, and route congestion). The final, routed design, can then be deployed to the FPGA by producing a proprietary *bitstream*, which is written to the FPGA.

### III. BRAGGHLS COMPILER AND HLS FRAMEWORK

BraggHLS is a compiler and HLS framework which employs MLIR for extracting loop-nest representations of DNNs. It is implemented in python for ease of use and extensibility. Critically, and distinctly, it handles the DNN transformations as well as scheduling, binding, and FSM extraction; there is no dependence on any commercial HLS tools. Figure 1 shows the architecture of BraggHLS. BraggHLS proceeds by first lowering DNNs from PyTorch to MLIR through TorchScript and the `torch` dialect (see Section II-A2). They are then further lowered from the `torch` dialect to the `scf` dialect (through the `linalg` dialect). Such a representation lends itself to a straightforward translation to python (compare Listing 1 to Listing 1) and indeed BraggHLS performs this translation. The benefits of translating `scf` dialect to python are manifold and discussed in the following (see Section III-A). Ultimately, BraggHLS produces a representation of the DNN that is then fully scheduled using the scheduling infrastructure in CIRCT [33] (an MLIR adjacent project). After scheduling, BraggHLS emits corresponding RTL (as Verilog).

BraggHLS delegates to the FloPoCo [34] IP generator the task of generating pipelined implementations of the standard floating-point arithmetic operations (`mul`, `div`, `add`, `sub`, `sqrt`) at various precisions. In addition, we implement a few generic (parameterized by bit width) operators in order to support a broad range of DNN operations: two-operand maximum (`max`); negation (`neg`); rectified linear

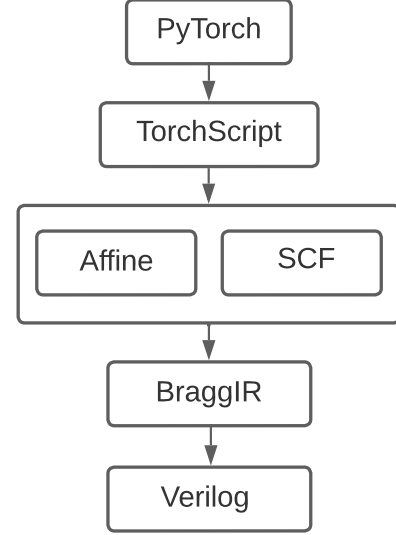


Fig. 1. BraggHLS framework overview (placeholder).

units (`relu`). Transcendental functions, such as `exp`, are implemented using a Taylor series expansion to  $k$ -th order (where  $k$  is determined on a case-by-case basis). Note, FloPoCo’s floating-point representation differs slightly from IEEE754, foregoing subnormals and differently encoding zeroes, infinities and NaNs (for the benefit of reduced complexity) and our implementations `max`, `neg`, `relu` are adjusted appropriately.

We now discuss some aspects of BraggHLS in greater detail.

#### A. Symbolic interpretation for fun and profit

As discussed in Section II-B1, maximizing concurrency and parallelism for a design entails unrolling loop nests and analyzing the data-flow of encompassed operations. As also discussed in Section II-B1, the formally correct approach to unrolling a loop nest is prohibitively expensive in terms of runtime. Indeed, for example, in the case of BraggNN, repeatedly performing this unrolling was a hindrance to effectively searching the design space for a RTL representation achieving the target latency, since it often took an enormous amount of time (see Figure 2). Translating `scf` dialect to python enables BraggHLS to overcome this barrier by enabling us to use the python interpreter as a *symbolic interpreter*. Interpreting the resulting python loop-nests (i.e., running the python program) while treating the arithmetic and memory operations on SSA values as operations on symbols (i.e., python classes with overloaded methods) enables us to:

- 1) Partially evaluate functions of iteration variables, such as `%3 = arith.addi %i3, %i6`, which enables concretely determining array index operands of all stores and loads, such as `memref.load %input[%i1,%i5,%i3,%3,%4]`, and thereupon performing memory dependence checks, thus

```

@conv2d(
  %input: memref<b × cin × h × w>,
  %weight: memref<b × cout × h × w>,
  %output: memref<cout × cin × k × k>
) {
  scf.for %i1 = %c0 to b step %c1 {
    scf.for %i2 = %c0 to cout step %c1 {
      scf.for %i3 = %c0 to h step %c1 {
        scf.for %i4 = %c0 to w step %c1 {
          scf.for %i5 = %c0 to cin step %c1 {
            scf.for %i6 = %c0 to k step %c1 {
              scf.for %i7 = %c0 to k step %c1 {
                %3 = arith.addi %i3, %i6
                %4 = arith.addi %i4, %i7
                %5 = memref.load %input[
                  %i1, %i5, %i3, %4]
                %6 = memref.load %weight[
                  %i2, %i5, %i6, %i7]
                %7 = memref.load %output[
                  %i1, %i2, %i3, %i4]
                %8 = arith.mulf %5, %6
                %9 = arith.addf %7, %8
                memref.store %9, %output[
                  %i1, %i2, %i3, %i4]
              }
            }
          }
        }
      }
    }
  }
  return %2
}

```

Listing 3: scf dialect loop representation of the convolution in Listing 1.

transforming the problem of statically verifying memory dependence into a matter of checking assertions at runtime;

- 2) Unroll loops by recording each floating-point arithmetic operation executed while enforcing SSA; e.g., for a loop whose body has repeated assignments to the same SSA value (ostensibly violating SSA), we execute the loop and instantiate new, uniquely identified, symbols for the result of each operation;
- 3) Reconstruct all data flow through arithmetic operations and memory operations by interpreting memrefs as *geometric symbol tables* (i.e., symbol tables indexed by array indices rather than identifiers/names) and stores and loads as assignments and reads on those symbol tables;
- 4) Easily swap evaluation rules in order to support various functionality modes, e.g., evaluating the floating-point arithmetic operations using (python) bindings to FloPoCo’s C++ functional models thereby enabling behavioral verification.

See Table I for the translation rules from MLIR dialects to python.

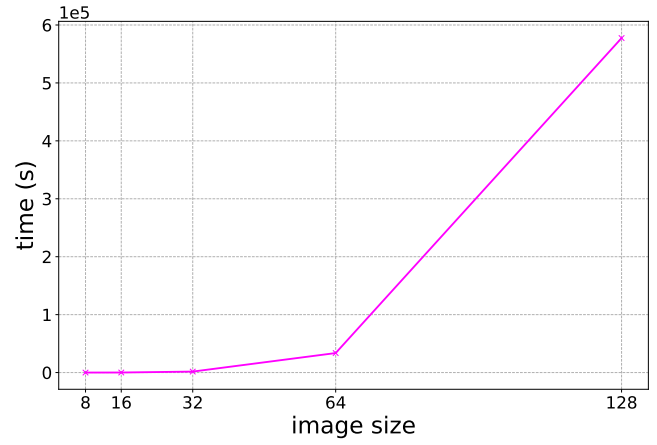


Fig. 2.  $3 \times 3$ -kernel convolution (cf. Listing 3) full unrolling time vs. input (square) image size, with store-load forwarding using MLIR’s -affine-scalrep pass. Note that the longest time is  $577419 \div 3600 \approx 160$  hours (for a loop nest with a trip count of  $128 \times 128 \times 3 \times 3 = 147,456$ ).

### B. AST transformations and behavioral verification

Prior to interpretation, BraggHLS performs some simple AST transformations on the python generated from scf dialect:

- 1) **Hoist globals**: all DNN tensors which are fixed (i.e., weights) are moved out of the body of the python<sup>2</sup> and into the parameter list, for the purpose of ultimately exposing them at the RTL module interface;
- 2) **Remove if expressions**: DNN relu operations are lowered to the scf dialect as a decomposition of arith.cmpfugt and arith.select; this transformation recomposes them into a relu;
- 3) **Remove MACs**: sequences of load-multiply-add-store are very common in DNN implementations, thus we schedule such sequences jointly (this transformation coalesces such sequences into a single FMAC);
- 4) **Reduce fors**: this transformation implements the reduction tree structure for non-parallelizable loop-nests mentioned in Section III-A.

These transformations on the python AST are simple (implemented with procedural pattern matching), extensible, and efficient (marginal runtime cost) because they are unverified: no effort is made to verify their formal correctness. Thus, BraggHLS trades formal correctness for development time performance. This tradeoff enables quick design space iteration, which for example, enabled us to achieve very low latency implementations for BraggNN (see Section IV-B). As a substitute for formal verification, BraggHLS supports behavioral verification. Specifically, BraggHLS can generate testbenches for all synthesized RTL. The test vectors for these testbenches are generated by evaluating the generated python representation of the DNN on randomly generated

<sup>2</sup>BraggHLS translates the MLIR module corresponding to the DNN into a single python function in order to simplify analysis and interpretation.

TABLE I  
SCF, ARITH, AND MEMREF DIALECT TO PYTHON TRANSLATION.

[[MLIR]]	python
[[%5]]	v5 = Val("%5")
[[memref<b × c <sub>in</sub> × h × w>]]	MemRef(b, c <sub>in</sub> , h, w)
[[%5 = memref.load %input[%i1, %i5, %3, %4]]]	[[%5]] = [[%input]].__getitem__(([[%i1]], [[%i5]], [[%3]], [[%4]]))
[[ memref.store %9, %output[%i1, %i5, %3, %4]]]	[[%output]].__getitem__(([[%i1]], [[%i5]], [[%3]], [[%4]]), [[%9]])
[[ scf.for %i1 = %c0 to b step %c1]]	for [[%i1]] in range([[%c0]], b, [[%c1]])
[[%3 = arith.addi %i3, %i6]]	[[%3]] = [[%i3]] + [[%i6]]
[[%8 = arith.mulf %5, %6]]	[[%8]] = [[%5]].__mul__(([[%6]]))
[[%9 = arith.addf %7, %8]]	[[%9]] = [[%7]].__add__(([[%8]]))
[[%63 = arith.cmpfugt %10, %c0]]	[[%63 = arith.cmpfugt %10, %c0]]
[[%64].relu([[%10]])]	[[%64]].relu([[%10]])
[[%8 = arith.mulf %5, %6]]	[[%8 = arith.addf %7, %8]]
[[%9] = fma([[%5]], [[%6]], [[%7]])]	

inputs but with floating-point operations now evaluated using functional models of the corresponding FloPoCo operators. The testbenches can then be run using any IEEE 1364 compliant simulator. For example, we run a battery of such testbenches (corresponding to various DNN operation types), using cocotb [35] and iverilog [36], as a part of our continuous integration process<sup>3</sup>.

### C. Scheduling

Recall that one of the critical functions which HLS fulfills is the scheduling of operations during each clock cycle, in such a way that they preserve the data-flow graph of a DNN; that schedule then informs the construction of a corresponding FSM. As already mentioned, scheduling arbitrary DNNs involves formulating and solving an ILP. In the resource-unconstrained case, due to the precedence relations induced by data-flow, the constraint matrix of the associated ILP is *totally unimodular matrix* and the feasible region of the problem is an integral polyhedron. Thus, in such cases, the scheduling problem can be solved optimally in polynomial time with a LP solver [37]. In the resource-constrained case it is possible to transform resource constraints into precedence constraints as well, by picking a particular (possibly heuristic) linear ordering on the particularly resource-constrained operations. This transformation partitions resource constrained operations into distinct clock cycles, thereby guaranteeing sufficient resources are available for all operations scheduled within the same clock cycle [38].

BraggHLS uses the explicit parallelism of the scf.parallel loop-nest representation to inform such

a linear ordering on resource-constrained operations. By assumption, for loop-nests which can be represented as scf.parallel loop-nests (see Listing 4), each instance of a floating-point arithmetic operation in the body corresponding to unique values of the iteration variables (e.g., %i1, %i2, %i3, %i4 for Listing 4) is independent of all other such instances<sup>4</sup>. This exactly determines total resource usage per loop-nest; for example, the convolution in Listing 4, would bind to  $2K_i$  DSPs (assuming mulf, addf bind to one DSP each), where

$$K_i := |\{\%iv1 = \%c0 + \%c1 \times \mathbb{N} \wedge \%iv1 < b\}| \times |\{\%iv2 = \%c0 + \%c1 \times \mathbb{N} \wedge \%iv2 < c_{out}\}| \times |\{\%iv3 = \%c0 + \%c1 \times \mathbb{N} \wedge \%iv3 < h\}| \times |\{\%iv4 = \%c0 + \%c1 \times \mathbb{N} \wedge \%iv4 < w\}|$$

where  $\%c1 \times \mathbb{N}$  represents all multiples of  $\%c1$ . That is to say,  $K_i$  is the cartesian product of the iteration spaces of the parallel iteration variables. Taking the maximum over such  $K := \max_i K_i$  across all scf.parallel loop-nests, we can infer peak usage of any resource. Then, after indexing available hardware resources  $j = 1, \dots, K$ , we can bind the operations of any particular loop-nest. This leads to a linear ordering on resource-constrained operations such that operations bound to the same hardware resource index  $j$  must be ordered according to their execution order during interpretation<sup>5</sup>. Note this ordering coincides with the higher-level

<sup>4</sup>Data-flow within a loop body must still be respected.

<sup>5</sup>BraggHLS only needs to construct a partial precedence ordering  $op_a < op_b$  for operations  $op_a, op_b$  which CIRCT then combines with the delays of the operations to construct constraints such as  $start_{op_a} + delay_a \leq start_{op_b}$ .

<sup>3</sup><https://github.com/makslevental/braggHLS/actions>

```

@conv2d(
  %input: memref<b × cin × h × w>,
  %weight: memref<b × cout × h × w>,
  %output: memref<cout × cin × k × k>
) {
  scf.parallel (%i1, %i2, %i3, %i4) =
    (%c0, %c0, %c0, %c0) to
      (b, cout, h, w) step
      (%c1, %c1, %c1, %c1) {
    scf.for %i5 = %c0 to cin step %c1 {
      scf.for %i6 = %c0 to k step %c1 {
        scf.for %i7 = %c0 to k step %c1 {
          %3 = arith.addi %i3, %i6
          %4 = arith.addi %i4, %i7
          %5 = memref.load %input[%i1, %i5, %i3, %3, %4]
          %6 = memref.load %weight[%i2, %i5, %i6, %i7]
          %7 = memref.load %output[%i1, %i2, %i3, %i4]
          %8 = arith.mulf %5, %6
          %9 = arith.addf %7, %8
          memref.store %9, %output[%i1, %i2, %i3, %i4]
        }
      }
    }
  }
  return %2
}

```

Listing 4: Parallel loop representation of the convolution in Listing 1.

structure of the DNN, since ordering of `scf.parallel` loop nests (and thus interpretation order during execution of the python program) is determined by the higher-level structure of the DNN.

For DNN operations that do not lower `scf.parallel` loop-nests but lower to sequential loop nests (e.g., `sum`, `max`, or `prod`), we fully unroll the loops and transform the resulting, sequential, operations to a reduction tree; we use As-Late-As-Possible scheduling [39] amongst the subtrees of such reduction trees.

#### IV. EVALUATION

We evaluate BraggHLS both on individual DNN layers, and end-to-end, on our use-case BraggNN. We compare BraggHLS to Xilinx’s Vitis HLS by comparing the latencies and resource usages of the final designs generated by each. We also compare the runtimes of the tools themselves. Both BraggHLS and Vitis HLS produce Verilog RTL, on which we run a synthesis pass using Xilinx’s Vivado. The particular FPGA target is Xilinx Alveo U280. We measure LUT, DSP, BRAM, and FF usage. For the DNN layer evaluations, we use FloPoCo (5, 11)-floating point representations (5 bits for the exponent and 11 bits for the mantissa) corresponding to Vitis HLS’s IEEE half precision IPs. We synthesize all designs for a 10 ns target clock period and report end-to-end latency as a product of the total schedule interval count of the design and achieved clock period ( $10 - WNS$ , where  $WNS$  is the worst negative slack reported). Note, in the case of Vitis HLS, which potentially explicitly pipelines the design and therefore implements with an initiation interval strictly less than the total schedule interval count, we report the best possible interval count (LatencyBest from the Vitis HLS

reports). All other measurements are collected from Vivado synthesis reports. Note, since Vitis HLS operates on C++ representations, we generate such a representation for our test cases by first lowering each to the affine dialect and then using the `scalehls-translate` tool of the ScaleHLS project [31] to emit C++. Note further, we do not make any use of `scalehls-opt` tool (of the same project).

Since our ultimate goal is low latency inference, and since the strategy that BraggHLS employs in the pursuit of this goal is loop unrolling, in order to produce a like for like comparison, we similarly unroll the representation that’s passed to Vitis HLS. Thus, all Vitis HLS measurements are reported in terms of *unroll factor*: an unroll factor of  $k$  corresponds to a  $k$ -fold increase in the number of statements in the body of a loop and commensurate  $k$ -fold decrease in the trip count of the loop. For loop nests we unroll inside out: if  $k$  is greater than the trip count  $t$  of the innermost loop, we unroll the innermost loop completely and then unroll the enclosing loop by a factor of  $k - t$ . We do not perform any store-load forwarding during this preprocessing but we annotate all arrays with the directive `array_partition complete dim=1` in order that Vitis HLS can effectively pipeline.

##### A. DNN layers

See Figure 3; we evaluate BraggHLS against Xilinx’s Vitis HLS by comparing the latency of the final design on the following, individual, DNN layer types:

- `addmm(a, b, c)` : Matrix multiplication of the matrices  $a$  and  $b$ . The matrix  $c$  is added to the final result;
- `batch_norm_2d(num_features)` : Batch Normalization over a 4D input as described in [40];
- `conv_2d(cin, cout, k)` : 2D convolution with bias, with  $k \times k$  kernel, over a  $b \times c_{in} \times h \times w$  input, producing  $b \times c_{out} \times h' \times w'$  output;
- `max_pool_2d(k, stride)` : 2D max pooling, with  $k \times k$  kernel, and striding;
- `soft_max` : The Softmax function.

The layers were chosen to cover a range of arithmetic operations and data access patterns, i.e., iteration, accumulation, and reduction. The concrete parameter values and dimensions of inputs used during evaluation are summarized in Table II.

We observe that Vitis HLS indeed approaches lower and lower end-to-end latencies as a function of unroll factor but that for none of the unroll factors is it able to match the end-to-end latency achieved by BraggHLS. Even at unroll factor equal to 1024 (which corresponds to fully unrolled for all of the loop nests comprising these layer types), Vitis HLS is only within  $10\times$  of BraggHLS. We hypothesize that in general this is due to Vitis HLS’s inability to effectively pipeline due to inability to eliminate memory dependencies, either



TABLE II  
DNN LAYERS USED FOR EVALUATION OF BRAGGHLS.

Layer	Parameter values	Input dimensions
addmm	N/A	$\dim(a) = \dim(b) = \dim(c) = (16, 16)$
batch_norm_2d	$\text{num\_features} = 2$	$\dim(\text{input}) = (10, 2, 3, 3)$
conv_2d	$c_{in} = 1, c_{out} = 3, k = 3$	$\dim(\text{input}) = (1, 1, 16, 16)$
max_pool_2d	$k = 3, \text{stride} = 2$	$\dim(\text{input}) = (1, 3, 16, 16)$
soft_max	N/A	$\dim(\text{input}) = (1, 3, 16, 16)$

through store-load forwarding or further array partitioning<sup>6</sup>. Conversely, BraggHLS’s ability to effectively perform store-load forwarding is evident in the complete lack of BRAM usage: all weights are kept on FFs or LUTs. While this is infeasible for larger designs (which would be constrained by the number of available FFs), for our particular use case, this unconstrained usage of FFs is acceptable. Note, the increasing latency (as a function of unroll factor) in the case of max\_pool\_2d is due to Vitis HLS’s failure to meet timing, i.e., while the interval count is decreased in the clock period is increasing (see Appendix Table III for exact values).

### B. BraggNN case study

High-energy diffraction microscopy techniques can provide non-destructive characterization for a broad class of single-crystal and polycrystalline materials. The critical steps in a typical HEDM experiment involve an analysis to determine precise Bragg diffraction peak characteristics (followed by reconstruction of material grain information from the peak characteristics). Peak characteristics are typically computed by fitting the peaks to a probability distribution, e.g., Gaussian, Lorentzian, Voigt, or Pseudo-Voigt. As already mentioned (in Section I) these experiments can have collection rates in excess of 80 GB/s. These data rates, though more modest than those observed at the LHC, combined with the runtime of the fitting procedure, merit exploring the same low latency approach (as explored employed by LHC experiments) in order to enable experiment modalities that depend on measurement-based feedback (i.e., experiment steering).

BraggNN [41], a DNN aimed at efficiently characterizing Bragg diffraction peaks, achieves a batched latency of approximately 22  $\mu\text{s}/\text{sample}$ . This is a large speedup over the classical pseudo-Voigt peak fitting methods, but still falls far short of the 1  $\mu\text{s}/\text{sample}$  target for handling the 1 MHz sampling rates. In addition, the current implementation of BraggNN, deployed to either a data-center class GPU such as a NVIDIA V100, or even a workstation class GPU such as a NVIDIA RTX 2080Ti, has no practicable means to being deployed at the

```

BraggNN(s) (
  (cnn_layers_1): Conv2d(s × 16, kernel=3, stride=1)
  (nlb): NLB(
    (theta_layer): Conv2d(s × 16, s × 8, kernel=1, stride=1)
    (phi_layer): Conv2d(s × 16, s × 8, kernel=1, stride=1)
    (g_layer): Conv2d(s × 16, s × 8, kernel=1, stride=1)
    (out_cnn): Conv2d(s × 8, s × 16, kernel=1, stride=1)
    (soft): Softmax()
  )
  (cnn_layers_2): Sequential(
    (0): ReLU()
    (1): Conv2d(s × 16, s × 8, kernel=3, stride=1)
    (2): ReLU()
    (3): Conv2d(s × 8, s × 2, kernel=3, stride=1)
    (4): ReLU()
  )
  (dense_layers): Sequential(
    (0): Linear(in_features=s × 50, out_features=s × 16)
    (1): ReLU()
    (2): Linear(in_features=s × 16, out_features=s × 8)
    (3): ReLU()
    (4): Linear(in_features=s × 8, out_features=s × 4)
    (5): ReLU()
    (6): Linear(in_features=s × 4, out_features=2)
    (7): ReLU()
  )
)

```

Listing 5: BraggNN for scale  $s = 1, 2$ .

edge, i.e., adjacent or proximal to the high energy microscopy equipment. Towards this end, we applied BraggHLS to the PyTorch representation of BraggNN( $s = 1$ ) (see Listing 5) and achieved a RTL implementation which synthesizes to a 1,238 interval design that places, routes, and meets timing closure for a clock period of 10 ns (for a Xilinx Alveo U280). The design consists of a three stage pipeline with the longest stage measuring 480 intervals. Thus, the throughput of the implementation is 4.8  $\mu\text{s}/\text{sample}$ . See Table IV for a summary of the resource usage of the implementation and Figure 5 for a comparison with designs generated by Vitis HLS.

The most challenging aspect of implementing BraggNN was minimizing latency while satisfying compute resource constraints (LUTs, DSPs, BRAMs) and meeting routing “closure”, i.e., not exceeding available routing resources and avoiding congestion. Two design choices were made for the purposes of reducing resource consumption. The first was reducing the precision used for the floating-point operations. We reduced the precision from half precision to FloPoCo (5, 4)-precision (5 bits for the exponent and 4 bits for the mantissa).

<sup>6</sup>Indeed Vitis HLS can be seen to produce warnings such as  
WARNING: [HLS 200-885] Unable to schedule 'load'  
operation ... due to limited memory ports (II = 1).  
Please consider using a memory core with more ports  
or partitioning the array ...

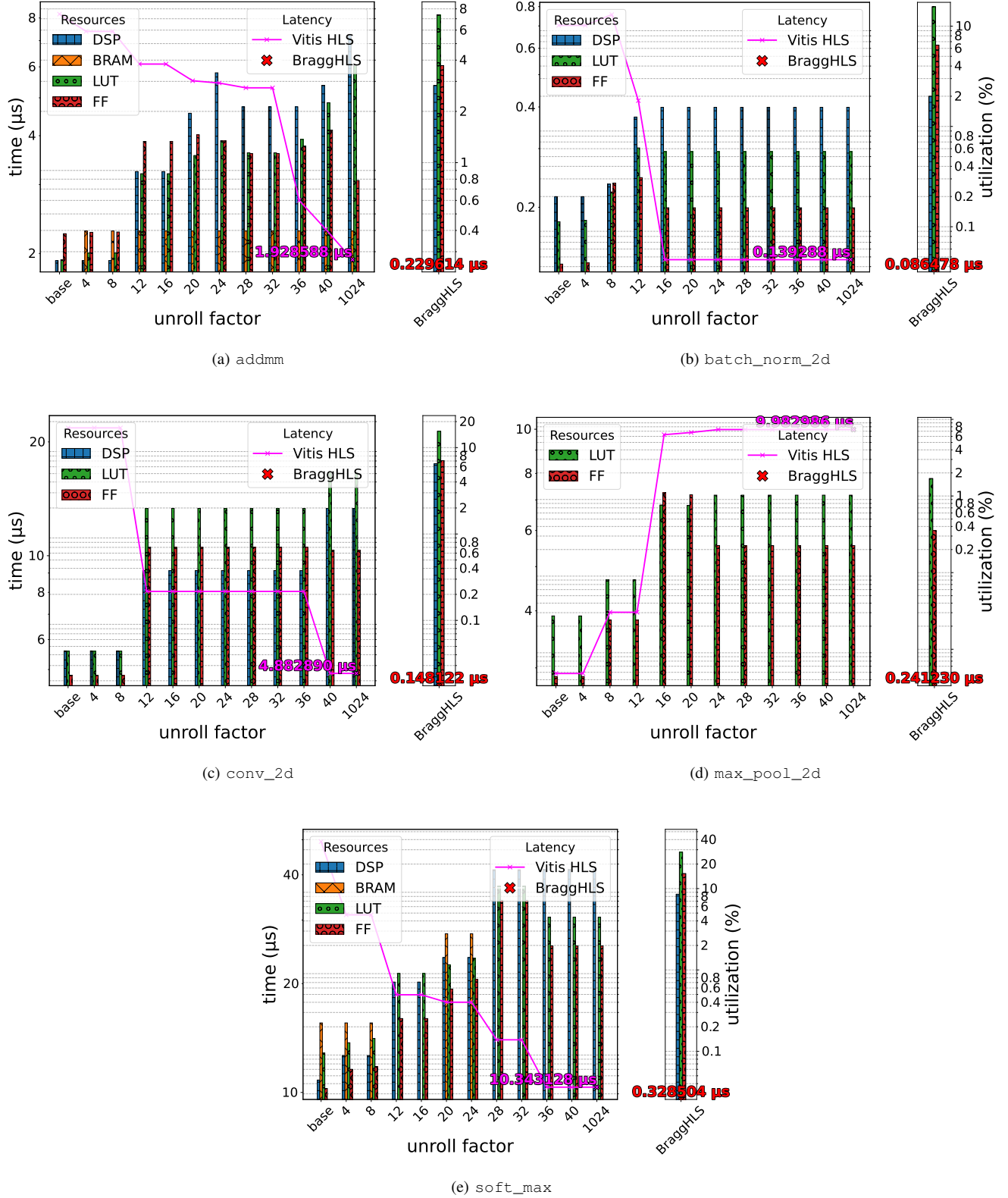


Fig. 3. Vitis HLS vs. BraggHLS resource usage and latency vs. unroll factor of various DNN modules. All  $y$ -scales are log.

This was justified by an examination of the distribution of the weights of the fully trained BraggNN (see figure 6). Reducing

the precision enabled the second design choice: we eliminate BRAMs from the design, since, at the lower precision, all

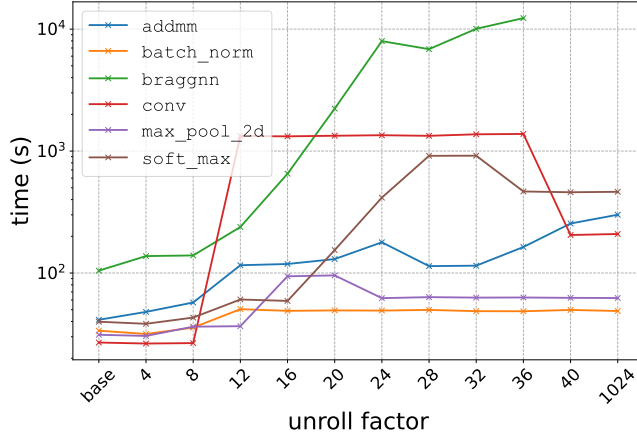


Fig. 4. Runtime of Vitis HLS vs. unroll factor.

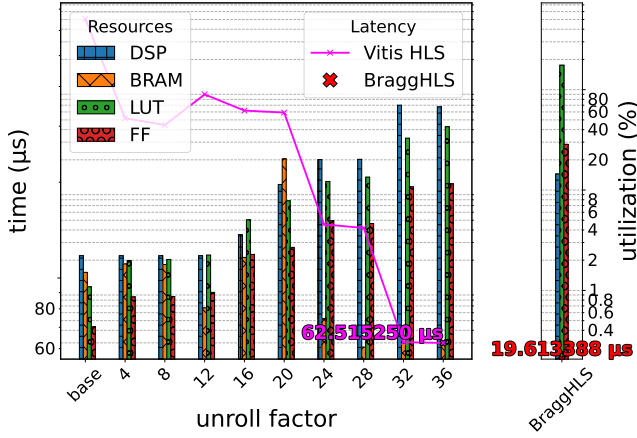


Fig. 5. BraggNN Vitis HLS vs. BraggHLS resource usage and latency vs. unroll factor.

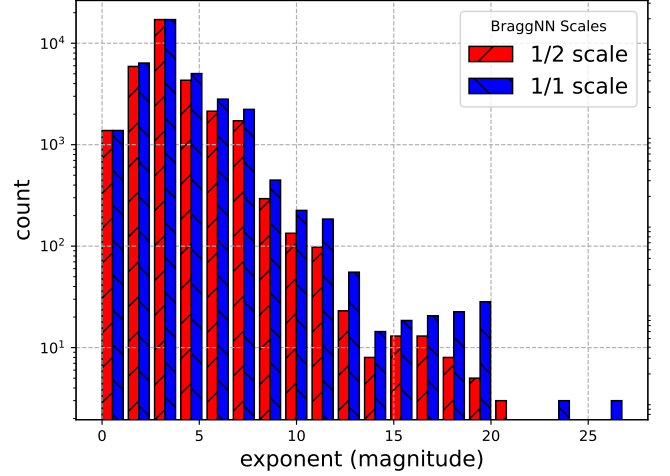


Fig. 6. BraggHLS weights exponent distribution.

weights can be represented as registered constants. The reduced precision also drives the Vivado synthesizer to infer implementations of the floating-point operations that make no use of DSPs; this was not intentional but seemingly cannot be altered. Most likely this is due to the fact that the DSP48 hardware block includes a 18-bit by 25-bit signed multiplier and a 48-bit adder [42], neither of which neatly divide the bit width of FloPoCo (5, 4)-precision cores<sup>7</sup>.

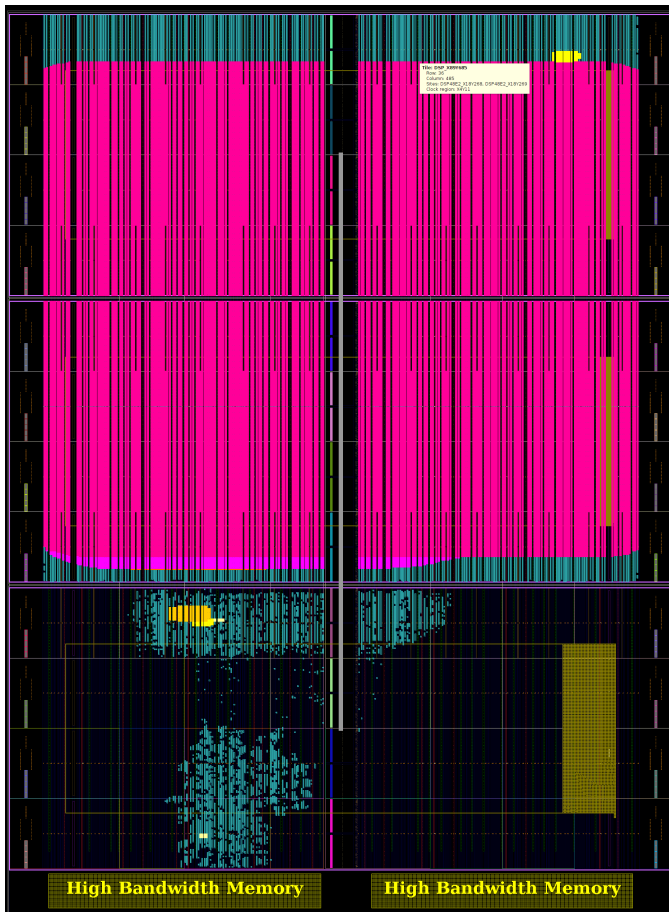
Achieving routing closure was very difficult due to the nature of the Xilinx’s UltraScale architecture, of which the Alveo U280 is an instance. The UltraScale architecture achieves its scale through “Stacked Silicon Interconnect” (SSI) technology [43], which implies multiple distinct FPGA dies, called Super Logic Regions (SLRs), on the same chip, connected by interposers. Adjacent SLRs communicate with each other using a limited set of Super Long Lines (SLLs). These SLLs

determine the maximum bus width that spans two SLRs. On the Alveo U280 there are exactly 23,040 SLLs available between adjacent SLRs and at (5, 4)-precision BraggNN needs 23,328 SLLs between SLR2 and SLR1<sup>8</sup>. Thus, we further reduced the precision to (5, 3). Finally, since multiple dies constitute independent clock domains, the SLLs that cross SLRs are sensitive hold time violations due to the higher multi-die variability [44]. This multi-die variability leads to high congestion if not addressed. Thus, routing across SLRs needs to be handled manually, using placement and routing constraints for logic in each SLR and the addition of so-called “launch” and “latch” registers in each SLR. See figure 7 for an illustration on the effect of using launch and latch registers as well as placement and routing constraints.

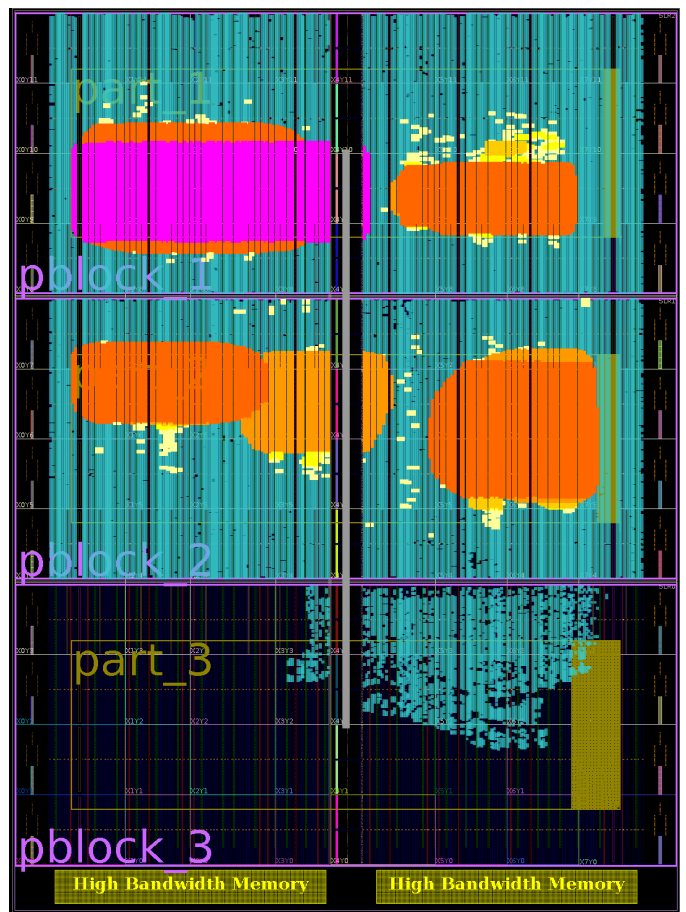
Thus, the confluence of the aforementioned design choices

<sup>7</sup>The actual bit width for FloPoCo (5, 4)-precision is 12 bits: 1 extra bit is needed for the sign and 2 bits are needed for FloPoCo’s handling of exceptional conditions.

<sup>8</sup>We route the output of `cnn_layers_1` ( $1 \times 16 \times 9 \times 9 \times 12$  wires) as well as the output of `soft(theta_layer \times phi_layer) \times g_layer` ( $1 \times 8 \times 9 \times 9 \times 12$  wires) from SLR2 to SLR1.



(a) BraggNN fails to achieve routing closure without placement and routing constraints and launch and latch registers.



(b) BraggNN achieves routing closure with use of per SLR placement and routing constraints (pblock\_1, pblock\_2, pblock\_3) and launch and latch registers (not highlighted).

Fig. 7. Congestion maps for BraggNN on a Xilinx Alveo U280. **Magenta** indicates areas of high congestion.

(in combination with compiler level optimizations performed by BraggHLS) and careful management of routing constraints enables us to lower, compile, synthesize, place and route BraggNN to Xilinx’s Alveo U280 at a throughput of 4.8  $\mu\text{s}/\text{sample}$ . This is  $\sim 5\times$  higher latency than the target 1  $\mu\text{s}/\text{sample}$  but it is a  $4\times$  improvement over the GPU PyTorch-GPU implementation. In the following we propose some future work that can further improve on this.

## V. CONCLUSION

## REFERENCES

- (in combination with compiler level optimizations performed by BraggHLS) and careful management of routing constraints enables us to lower, compile, synthesize, place and route BraggNN to Xilinx's Alveo U280 at a throughput of 4.8 us/sample. This is  $\sim 5\times$  higher latency than the target 1 us/sample but it is a  $4\times$  improvement over the GPU PyTorch-GPU implementation. In the following we propose some future work that can further improve on this.
- ## V. CONCLUSION
- ## REFERENCES
- [1] V. Gligorov, "Real-time data analysis at the lhc: present and future," in *Proceedings of the NIPS 2014 Workshop on High-energy Physics and Machine Learning*, ser. Proceedings of Machine Learning Research, G. Cowan, C. Germain, I. Guyon, B. Kegl, and D. Rousseau, Eds., vol. 42. Montreal, Canada: PMLR, 13 Dec 2015, pp. 1–18. [Online]. Available: <https://proceedings.mlr.press/v42/glig14.html>
  - [2] M. Hammer, K. Yoshii, and A. Miceli, "Strategies for on-chip digital data compression for x-ray pixel detectors," *Journal of Instrumentation*, vol. 16, no. 01, pp. P01025–P01025, jan 2021. [Online]. Available: <https://doi.org/10.1088/1748-0221/16/01/P01025>
  - [3] J. McMullin, P. Diamond, M. Caiazzo, A. Casson, T. Cheetham, P. Dewdney, R. Laing, B. Lewis, A. Schinckel, L. Stringhetti *et al.*, "The square kilometre array project update," in *Ground-based and Airborne Telescopes IX*, vol. 12182. SPIE, 2022, pp. 263–271.
  - [4] K. Grainge, B. Alachkar, S. Amy, D. Barbosa, M. Bommineni, P. Boven, R. Braddock, J. Davis, P. Diwakar, V. Francis *et al.*, "Square kilometre array: The radio telescope of the xxi century," *Astronomy reports*, vol. 61, no. 4, pp. 288–296, 2017.
  - [5] "Comparison of particle selection algorithms for the LHCb Upgrade," 2020. [Online]. Available: <https://cds.cern.ch/record/2746789>
  - [6] V. V. Gligorov and M. Williams, "Efficient, reliable and fast high-level triggering using a bonsai boosted decision tree," *Journal of Instrumentation*, vol. 8, no. 02, pp. P02013–P02013, feb 2013. [Online]. Available: <https://doi.org/10.1088/1748-0221/8/02/P02013>
  - [7] J. Duarte, S. Han, P. Harris, S. Jindariani, E. Kreinar, B. Kreis, J. Ngadiuba, M. Pierini, R. Rivera, N. Tran, and Z. Wu, "Fast inference of deep neural networks in FPGAs for particle physics," *Journal of Instrumentation*, vol. 13, no. 07, pp. P07027–P07027, jul 2018. [Online]. Available: <https://doi.org/10.1088/1748-0221/13/07/P07027>
  - [8] L. Alzubaidi, J. Zhang, A. J. Humaidi, A. Al-Dujaili, Y. Duan, O. Al-Shamma, J. Santamaría, M. A. Fadhel, M. Al-Amidie, and L. Farhan, "Review of deep learning: Concepts, cnn architectures, challenges, applications, future directions," *Journal of big Data*, vol. 8, no. 1, pp. 1–74, 2021.
  - [9] Z. Liu, T. Bicer, R. Kettimuthu, and I. Foster, "Deep learning accelerated light source experiments," in *2019 IEEE/ACM Third Workshop on Deep*



- Learning on Supercomputers (DLS)*. IEEE, 2019, pp. 20–28.
- [10] R. M. Patton, J. T. Johnston, S. R. Young, C. D. Schuman, D. D. March, T. E. Potok, D. C. Rose, S.-H. Lim, T. P. Karnowski, M. A. Ziatdinov *et al.*, “167-pflops deep learning for electron microscopy: from learning physics to atomic manipulation,” in *SC18: International Conference for High Performance Computing, Networking, Storage and Analysis*. IEEE, 2018, pp. 638–648.
  - [11] Y. Liu, K. P. Kelley, H. Funakubo, S. V. Kalinin, and M. Ziatdinov, “Exploring physics of ferroelectric domain walls in real time: deep learning enabled scanning probe microscopy,” *Advanced Science*, p. 2203957, 2022.
  - [12] R. Aaij, J. Albrecht, M. Belous, P. Billoir, T. Boettcher, A. Brea Rodríguez, D. Vom Bruch, D. Cámpora Pérez, A. Casais Vidal, D. Craik *et al.*, “Allen: A high-level trigger on gpus for lhcb,” *Computing and Software for big Science*, vol. 4, no. 1, pp. 1–11, 2020.
  - [13] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, “Automatic differentiation in PyTorch,” 2017.
  - [14] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mane, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viegas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng, “Tensorflow: Large-scale machine learning on heterogeneous distributed systems,” 2016. [Online]. Available: <https://arxiv.org/abs/1603.04467>
  - [15] T. Chen, M. Li, Y. Li, M. Lin, N. Wang, M. Wang, T. Xiao, B. Xu, C. Zhang, and Z. Zhang, “Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems,” 2015. [Online]. Available: <https://arxiv.org/abs/1512.01274>
  - [16] R. Nane, V.-M. Sima, C. Pilato, J. Choi, B. Fort, A. Canis, Y. T. Chen, H. Hsiao, S. Brown, F. Ferrandi, J. Anderson, and K. Bertels, “A survey and evaluation of fpga high-level synthesis tools,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 35, no. 10, pp. 1591–1604, 2016.
  - [17] A. Canis, J. Choi, M. Aldham, V. Zhang, A. Kammoona, T. Czajkowski, S. D. Brown, and J. H. Anderson, “Legup: An open-source high-level synthesis tool for fpga-based processor/accelerator systems,” *ACM Trans. Embed. Comput. Syst.*, vol. 13, no. 2, sep 2013. [Online]. Available: <https://doi.org/10.1145/2514740>
  - [18] Z. Zhang, Y. Fan, W. Jiang, G. Han, C. Yang, and J. Cong, *AutoPilot: A Platform-Based ESL Synthesis System*. Dordrecht: Springer Netherlands, 2008, pp. 99–112. [Online]. Available: [https://doi.org/10.1007/978-1-4020-8588-8\\_6](https://doi.org/10.1007/978-1-4020-8588-8_6)
  - [19] F. Ferrandi, V. G. Castellana, S. Curzel, P. Fezzardi, M. Fiorito, M. Lattuada, M. Minutoli, C. Pilato, and A. Tumeo, “Invited: Bambu: an open-source research framework for the high-level synthesis of complex applications,” in *2021 58th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 2021, pp. 1327–1330.
  - [20] T. Chen, T. Moreau, Z. Jiang, L. Zheng, E. Yan, H. Shen, M. Cowan, L. Wang, Y. Hu, L. Ceze *et al.*, “{TVM}: An automated {End-to-End} optimizing compiler for deep learning,” in *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*, 2018, pp. 578–594.
  - [21] C. Lattner, M. Amini, U. Bondhugula, A. Cohen, A. Davis, J. Pienaar, R. Riddle, T. Shpeisman, N. Vasilache, and O. Zinenko, “Mlir: A compiler infrastructure for the end of moore’s law,” 2020. [Online]. Available: <https://arxiv.org/abs/2002.11054>
  - [22] N. Rotem, J. Fix, S. Abdulrasool, G. Catron, S. Deng, R. Dzhabarov, N. Gibson, J. Hegeman, M. Lele, R. Levenstein, J. Montgomery, B. Maher, S. Nadathur, J. Olesen, J. Park, A. Rakhov, M. Smelyanskiy, and M. Wang, “Glow: Graph lowering compiler techniques for neural networks,” 2018. [Online]. Available: <https://arxiv.org/abs/1805.00907>
  - [23] Y. Liu, Y. Wang, R. Yu, M. Li, V. Sharma, and Y. Wang, “Optimizing cnn model inference on cpus,” 2018. [Online]. Available: <https://arxiv.org/abs/1809.02697>
  - [24] S. Zheng, R. Chen, Y. Jin, A. Wei, B. Wu, X. Li, S. Yan, and Y. Liang, “Neoflow: A flexible framework for enabling efficient compilation for high performance dnn training,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 11, pp. 3220–3232, nov 2022.
  - [25] T. Chen, B. Xu, C. Zhang, and C. Guestrin, “Training deep nets with sublinear memory cost,” 2016. [Online]. Available: <https://arxiv.org/abs/1604.06174>
  - [26] A. Ashari, S. Tatikonda, M. Boehm, B. Reinwald, K. Campbell, J. Keenleyside, and P. Sadayappan, “On optimizing machine learning workloads via kernel fusion,” vol. 50, no. 8, pp. 173–182, jan 2015. [Online]. Available: <https://doi.org/10.1145/2858788.2688521>
  - [27] S. Maleki, Y. Gao, M. J. Garzar, T. Wong, D. A. Padua *et al.*, “An evaluation of vectorizing compilers,” in *2011 International Conference on Parallel Architectures and Compilation Techniques*. IEEE, 2011, pp. 372–382.
  - [28] S. Silva and A. Elangovan, “Torch-MLIR,” <https://mlir.llvm.org/OpenMeetings/2021-10-07-The-Torch-MLIR-project.pdf>, 2021.
  - [29] M. Hattori, N. Kobayashi, and R. Sato, “Gradual tensor shape checking,” 2022. [Online]. Available: <https://arxiv.org/abs/2203.08402>
  - [30] U. Bondhugula, “Polyhedral Compilation Opportunities in MLIR,” [https://acohen.gitlabpages.inria.fr/impact/impact2020/slides/IMPACT\\_2020\\_keynote.pdf](https://acohen.gitlabpages.inria.fr/impact/impact2020/slides/IMPACT_2020_keynote.pdf), 2020.
  - [31] H. Ye, C. Hao, J. Cheng, H. Jeong, J. Huang, S. Neuendorffer, and D. Chen, “Scalehls: A new scalable high-level synthesis framework on multi-level intermediate representation,” in *2022 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, 2022.
  - [32] S. V. Rajopadhye, “Dependence analysis and parallelizing transformations,” 2002.
  - [33] J. Oppermann, M. Urbach, and J. Demme, “How to Make Hardware with Maths: An Introduction to CIRCT’s Scheduling Infrastructure,” in *2022 European LLVM Developers’ Meeting (EuroLLVM)*, 2022.
  - [34] F. de Dinechin, “Reflections on 10 years of flopeco,” in *2019 IEEE 26th Symposium on Computer Arithmetic (ARITH)*, 2019, pp. 187–189.
  - [35] B. J. Rosser, “Cocotb: a python-based digital logic verification framework,” 2018.
  - [36] S. Williams, “Icarus verilog, 1998–2020,” URL <http://iverilog.icarus.com>.
  - [37] J. Oppermann, “Advances in ilp-based modulo scheduling for high-level synthesis,” Ph.D. dissertation, Technische Universität, Darmstadt, 2019. [Online]. Available: <http://tuprints.ulb-tu-darmstadt.de/9272/>
  - [38] S. Dai, G. Liu, and Z. Zhang, “A scalable approach to exact resource-constrained scheduling based on a joint sdc and sat formulation,” in *Proceedings of the 2018 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, ser. FPGA ’18. New York, NY, USA: Association for Computing Machinery, 2018, pp. 137–146. [Online]. Available: <https://doi.org/10.1145/3174243.3174268>
  - [39] Z. Baruch, “Scheduling algorithms for high-level synthesis,” *ACM Scientific Journal*, vol. 5, no. 1-2, pp. 48–57, 1996.
  - [40] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” 2015. [Online]. Available: <https://arxiv.org/abs/1502.03167>
  - [41] Z. Liu, H. Sharma, J.-S. Park, P. Kenesei, A. Miceli, J. Almer, R. Kettimuthu, and I. Foster, “BraggNN: fast X-ray Bragg peak analysis using deep learning,” *IUCrJ*, vol. 9, no. 1, pp. 104–113, Jan 2022.
  - [42] A. S. U. Guide, “Ultrascale architecture dsp slice.”
  - [43] S. Leibson and N. Mehta, “Xilinx ultrascale: The next-generation architecture for your next-generation architecture,” *Xilinx White Paper WP435*, vol. 143, 2013.
  - [44] “Create placed and routed dcp to cross slr,” [https://www.rapidwright.io/docs/SLR\\_Crosser\\_DCP\\_Creator\\_Tutorial.html](https://www.rapidwright.io/docs/SLR_Crosser_DCP_Creator_Tutorial.html), accessed: 2022-10-15.

## APPENDIX

TABLE III: Latency, resource usage, and runtimes for all BraggHLS evaluations.

	Metric name Unroll factor	BRAM	DSP	FF	LUT	Latency	Clock Period	Runtime
addmm	0	0.00E+00	2.66E-01	3.83E-01	2.69E-01	1334	6.13E+00	4.14E+01
	4	3.97E-01	2.66E-01	3.90E-01	2.94E-01	1205	6.13E+00	4.80E+01
	8	3.97E-01	2.66E-01	3.91E-01	2.95E-01	1205	6.13E+00	5.73E+01
	12	3.97E-01	8.87E-01	1.33E+00	8.61E-01	993	6.14E+00	1.16E+02
	16	3.97E-01	8.87E-01	1.33E+00	8.61E-01	993	6.14E+00	1.19E+02
	20	3.97E-01	1.95E+00	1.46E+00	1.10E+00	900	6.14E+00	1.30E+02
	24	3.97E-01	3.37E+00	1.34E+00	1.35E+00	888	6.14E+00	1.79E+02
	28	3.97E-01	2.13E+00	1.13E+00	1.14E+00	821	6.46E+00	1.14E+02
	32	3.97E-01	2.13E+00	1.13E+00	1.14E+00	821	6.46E+00	1.15E+02
	36	3.97E-01	2.13E+00	1.25E+00	1.37E+00	425	6.43E+00	1.63E+02
	40	3.97E-01	2.84E+00	1.55E+00	2.25E+00	359	6.42E+00	2.55E+02
batch_norm	1024	3.97E-01	5.67E+00	7.86E-01	3.74E+00	314	6.14E+00	3.00E+02
	0	0.00E+00	1.99E-01	4.24E-02	1.12E-01	116	6.06E+00	3.36E+01
	4	0.00E+00	1.99E-01	4.36E-02	1.16E-01	116	6.06E+00	3.16E+01
	8	0.00E+00	2.66E-01	2.72E-01	2.22E-01	125	6.06E+00	3.57E+01
	12	0.00E+00	1.24E+00	3.09E-01	6.11E-01	69	6.06E+00	5.05E+01
	16	0.00E+00	1.55E+00	1.55E-01	5.61E-01	23	6.06E+00	4.90E+01
	20	0.00E+00	1.55E+00	1.55E-01	5.61E-01	23	6.06E+00	4.94E+01
	24	0.00E+00	1.55E+00	1.55E-01	5.61E-01	23	6.06E+00	4.92E+01
	28	0.00E+00	1.55E+00	1.55E-01	5.61E-01	23	6.06E+00	4.98E+01
	32	0.00E+00	1.55E+00	1.55E-01	5.61E-01	23	6.06E+00	4.86E+01
	36	0.00E+00	1.55E+00	1.55E-01	5.61E-01	23	6.06E+00	4.85E+01
BraggNN	40	0.00E+00	1.55E+00	1.55E-01	5.61E-01	23	6.06E+00	4.98E+01
	1024	0.00E+00	1.55E+00	1.55E-01	5.61E-01	23	6.06E+00	4.88E+01
	0	1.51E+00	2.22E+00	4.34E-01	1.09E+00	85851	7.59E+00	1.05E+02
	4	1.84E+00	2.22E+00	8.61E-01	1.99E+00	49310	6.45E+00	1.38E+02
	8	1.81E+00	2.22E+00	8.69E-01	2.03E+00	49257	6.14E+00	1.39E+02
	12	6.70E-01	2.22E+00	9.56E-01	2.25E+00	61499	6.14E+00	2.38E+02
	16	2.13E+00	3.59E+00	2.29E+00	5.04E+00	52515	6.40E+00	6.52E+02
	20	2.04E+01	1.14E+01	2.69E+00	7.82E+00	51768	6.40E+00	2.23E+03
	24	5.21E-01	2.01E+01	4.95E+00	1.22E+01	24004	6.14E+00	7.97E+03
	28	2.73E-01	2.02E+01	4.64E+00	1.35E+01	22377	6.42E+00	6.85E+03
	32	2.98E-01	7.01E+01	1.07E+01	3.29E+01	9793	6.42E+00	1.01E+04
conv	36	2.98E-01	6.73E+01	1.15E+01	4.27E+01	9730	6.42E+00	1.23E+04
	0	0.00E+00	4.43E-02	2.30E-02	4.40E-02	3551	6.13E+00	2.69E+01
	4	0.00E+00	4.43E-02	2.30E-02	4.40E-02	3551	6.13E+00	2.64E+01
	8	0.00E+00	4.43E-02	2.30E-02	4.40E-02	3551	6.13E+00	2.66E+01
	12	0.00E+00	3.77E-01	7.04E-01	1.97E+00	1310	6.14E+00	1.33E+03
	16	0.00E+00	3.77E-01	7.03E-01	1.98E+00	1310	6.14E+00	1.32E+03
	20	0.00E+00	3.77E-01	7.03E-01	1.98E+00	1310	6.14E+00	1.34E+03
	24	0.00E+00	3.77E-01	7.03E-01	1.97E+00	1310	6.14E+00	1.35E+03
	28	0.00E+00	3.77E-01	7.03E-01	1.97E+00	1310	6.14E+00	1.33E+03
	32	0.00E+00	3.77E-01	7.03E-01	1.97E+00	1310	6.14E+00	1.37E+03
	36	0.00E+00	3.77E-01	7.03E-01	1.97E+00	1310	6.14E+00	1.38E+03
max_pool_2d	40	0.00E+00	1.97E+00	6.49E-01	5.07E+00	795	6.14E+00	2.05E+02
	1024	0.00E+00	1.97E+00	6.49E-01	5.07E+00	795	6.14E+00	2.09E+02
	0	0.00E+00	0.00E+00	4.41E-03	2.72E-02	889	3.28E+00	3.11E+01
	4	0.00E+00	0.00E+00	4.41E-03	2.72E-02	889	3.28E+00	3.05E+01

Continued on next page

TABLE III: Latency, resource usage, and runtimes for all BraggHLS evaluations.

module	metric_name unroll_factor	metric_val	DSP	FF	LUT	Latency	Clock Period	Runtime
		BRAM						
soft_max	8	0.00E+00	0.00E+00	2.40E-02	8.04E-02	641	6.20E+00	3.63E+01
	12	0.00E+00	0.00E+00	2.40E-02	8.04E-02	641	6.20E+00	3.66E+01
	16	0.00E+00	0.00E+00	1.10E+00	7.63E-01	419	2.32E+01	9.40E+01
	20	0.00E+00	0.00E+00	1.04E+00	7.51E-01	412	2.39E+01	9.56E+01
	24	0.00E+00	0.00E+00	2.25E-01	1.03E+00	262	3.81E+01	6.22E+01
	28	0.00E+00	0.00E+00	2.25E-01	1.03E+00	262	3.81E+01	6.34E+01
	32	0.00E+00	0.00E+00	2.25E-01	1.03E+00	262	3.81E+01	6.29E+01
	36	0.00E+00	0.00E+00	2.25E-01	1.03E+00	262	3.81E+01	6.30E+01
	40	0.00E+00	0.00E+00	2.25E-01	1.03E+00	262	3.81E+01	6.25E+01
	1024	0.00E+00	0.00E+00	2.25E-01	1.03E+00	262	3.81E+01	6.24E+01
	0	2.23E-01	4.43E-02	3.49E-02	9.52E-02	7742	6.38E+00	3.99E+01
	4	2.23E-01	8.87E-02	6.01E-02	1.27E-01	5053	6.13E+00	3.83E+01
	8	2.23E-01	8.87E-02	6.50E-02	1.44E-01	5053	6.13E+00	4.31E+01
	12	0.00E+00	7.09E-01	2.53E-01	9.04E-01	3037	6.13E+00	6.06E+01
	16	0.00E+00	7.09E-01	2.53E-01	9.04E-01	3037	6.13E+00	5.90E+01
	20	2.78E+00	1.42E+00	5.82E-01	1.16E+00	2893	6.14E+00	1.54E+02
	24	2.78E+00	1.42E+00	7.64E-01	1.39E+00	2893	6.14E+00	4.15E+02
	28	0.00E+00	1.70E+01	7.17E+00	1.07E+01	2277	6.14E+00	9.14E+02
	32	0.00E+00	1.70E+01	7.17E+00	1.07E+01	2277	6.14E+00	9.16E+02
	36	0.00E+00	1.72E+01	1.97E+00	4.45E+00	1684	6.14E+00	4.66E+02
	40	0.00E+00	1.72E+01	1.97E+00	4.45E+00	1684	6.14E+00	4.60E+02
	1024	0.00E+00	1.72E+01	1.97E+00	4.45E+00	1684	6.14E+00	4.63E+02

TABLE IV  
RESOURCE USAGE FOR BRAGGNN WITH  $k = 1$  AND (5, 3)-PRECISION FLOPoCo

Site Type	SLR0	SLR1	SLR2	SLR0 %	SLR1 %	SLR2 %
CLB	5047	52648	53900	9.18	97.50	99.81
CLBL	2773	28613	29227	9.47	97.72	99.82
CLBM	2274	24035	24673	8.86	97.23	99.81
CLB LUTs	19797	263733	311794	4.50	61.05	72.17
LUT as Logic	19797	263733	311794	4.50	61.05	72.17
using O5 output only	277	3944	4304	0.06	0.91	1.00
using O6 output only	17176	202564	266733	3.91	46.89	61.74
using O5 and O6	2344	57225	40757	0.53	13.25	9.43
LUT as Memory	0	0	0	0.00	0.00	0.00
LUT as Distributed RAM	0	0	0	0.00	0.00	0.00
LUT as Shift Register	0	0	0	0.00	0.00	0.00
CLB Registers	12527	286226	339820	1.42	33.13	39.33
CARRY8	244	5184	5184	0.44	9.60	9.60
Block RAM Tile	0	0	0	0.00	0.00	0.00
RAMB36/FIFO	0	0	0	0.00	0.00	0.00
RAMB18	0	0	0	0.00	0.00	0.00
URAM	0	0	0	0.00	0.00	0.00
DSPs	0	0	0	0.00	0.00	0.00
Unique Control Sets	189	2641	3179	0.17	2.45	2.94

TABLE V  
SUPER LONG LINE USAGE ACROSS SUPER LOGIC REGIONS FOR BRAGGNN

	Used	Fixed	Available	Util %
SLR2 $\leftrightarrow$ SLR1	21366		23040	92.73
SLR1 $\rightarrow$ SLR2	2			< 0.01
SLR2 $\rightarrow$ SLR1	21364			92.73
SLR1 $\leftrightarrow$ SLR0	3904		23040	16.94
SLR0 $\rightarrow$ SLR1	2			< 0.01
SLR1 $\rightarrow$ SLR0	3902			16.94
Total SLLs Used	25270			