

# POSTER: Hardware Acceleration of Complex Machine Learning Models through Modern High-Level Synthesis

Serena Curzel\*  
serena.curzel@polimi.it  
Politecnico di Milano  
Milano, Italy

Nicolas  
Bohm Agostini†  
nicolas.agostini@pnnl.gov  
Northeastern University  
Boston, MA, USA

Antonino Tumeo  
antonino.tumeo@pnnl.gov  
PNNL  
Richland, WA, USA

Fabrizio Ferrandi  
fabrizio.ferrandi@polimi.it  
Politecnico di Milano  
Milano, Italy

## CCS CONCEPTS

• **Hardware** → **Electronic design automation; Reconfigurable logic and FPGAs; Emerging languages and compilers.**

## KEYWORDS

FPGA, machine learning, High-Level Synthesis, MLIR

### ACM Reference Format:

Serena Curzel, Nicolas, Bohm Agostini, Antonino Tumeo, and Fabrizio Ferrandi. 2022. POSTER: Hardware Acceleration of Complex Machine Learning Models through Modern High-Level Synthesis. In *19th ACM International Conference on Computing Frontiers (CF'22)*, May 17–19, 2022, Torino, Italy. ACM, New York, NY, USA, 2 pages. <https://doi.org/10.1145/3528416.3530870>

## INTRODUCTION

Machine learning (ML) and deep learning algorithms are well suited to process and analyze large amounts of data, as it has been repeatedly proven in applications such as image classification, natural language processing, or recommendation systems. Both ML training and inference are compute- and memory-intensive, leading to widespread adoption of heterogeneous systems containing specialized accelerators. While graphic processing units (GPUs) are the established platform of choice to accelerate training, they are often too power-hungry to run inference tasks, or cannot meet the strict latency requirements of scientific experiments. A variety of custom solutions implemented as field programmable gate arrays (FPGAs) or application-specific circuit (ASICs) have been proposed in their place, ranging from generic "neural processors" to accelerators that focus on a narrow set of models with great efficiency.

While FPGAs and ASICs promise to provide the necessary hardware specialization, experimental workflows and the related ML methods evolve quickly. Even if FPGAs can be configured after deployment, allowing to update the accelerators and support new models, they still require significant expertise in hardware design and hardware description languages (HDLs) to be effectively used.

On the other hand, domain scientists are used to leverage python-based high-level frameworks (e.g., TensorFlow, Pytorch) to quickly

develop and explore ML algorithms, and typically do not have any hardware design expertise. To bridge the two worlds, some works propose to exploit High-Level Synthesis [1] [4], a well established technique that reduces the effort required to design hardware accelerators by automatically translating a behavioral description (traditionally C or C++) into Verilog/VHDL.

Most of the HLS-based frameworks for ML use C/C++ as an intermediate representation of the input model, augmenting it with tool-specific directives that drive the synthesis to obtain an efficient design. Our proposal, instead, is a compiler-based, open source framework that exploits Multi-Level Intermediate Representation (MLIR [6]) to lower the input model, apply optimization passes at the correct level of abstraction, and exploit knowledge about the HLS process to guide them.

A limitation of the existing frameworks is that they usually focus on a narrow set of models, specifically Multi-Layer Perceptrons and Convolutional Neural Networks. While it is true that these cover a significant part of ML applications (especially in the computer vision field), there is room for exploring other types of models, for example to accelerate scientific applications that work on graphs or sparse data structures. A narrow focus limits the possibility of quickly adapting to newer algorithmic approaches, a factor which needs to be taken into higher consideration.

Finally, as the MLIR project continues to grow, our tools are going to be available to other classes of algorithms, not exclusively ML models: in fact, any domain-specific framework with a lowering to basic MLIR dialects will be able to exploit our optimizations and translate algorithms into an efficient hardware accelerator. In the context of a scientific experiment, this will allow to accelerate also other stages in the data acquisition pipeline, such as pre-processing, analysis, and simulation.

## "CLASSIC" AND "MODERN" HLS FLOWS

Current frameworks that help automating the design of ML accelerators (e.g., hls4ml [4], FINN [1]) use existing High-Level Synthesis tools (e.g. Vitis/Vivado HLS) as back-end: they parse a model exported from popular ML frameworks and replace operators with C/C++ functions taken from a library of templates, adding tool-specific directives to guide optimizations (Figure 1a). The HLS tool processes this intermediate C/C++ representation and produces an accelerator design without further manual intervention.

Our compiler-based flow SODA [7] (Figure 1b) has two fundamental differences with respect to such a "classic" approach. First, instead of translating one high-level description into another high-level language such as C/C++, we embrace the MLIR multi-level approach to work with progressive lowerings. Increased modularity

\*Also with PNNL.

†Also with PNNL.

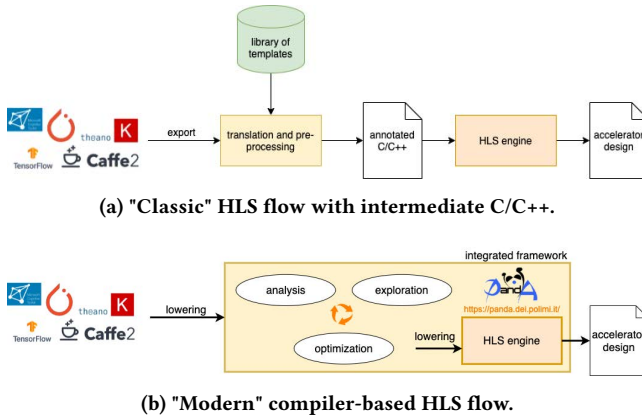
Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

CF'22, May 17–19, 2022, Torino, Italy

© 2022 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-9338-6/22/05.

<https://doi.org/10.1145/3528416.3530870>



**Figure 1: Two different ways of generating hardware accelerators through HLS.**

makes it possible to apply and explore both high-level algorithmic optimizations and low-level hardware-oriented ones. Second, we integrate the open-source HLS tool Bambu [5] in the loop, to have more control on the underlying High-Level Synthesis process.

A significant drawback of the "classic" flow is that the library of templates is necessarily tied to a single, specific HLS tool: each tool expects coding patterns, annotations, and optimization directives that are not portable to other tools; incompatible directives would be ignored, resulting in a very inefficient design. As typically there is one commercial HLS tool for each FPGA vendor, the choice of target boards is limited to the ones supported by the chosen HLS tool. This is not the case for our design flow, since Bambu is a multi-platform tool supporting multiple FPGA (Xilinx, Intel, Lattice, NanoXplore) and ASIC targets without any modification in the input code: such flexibility enables a smooth transition from high-level design, to FPGA prototyping, to ASIC manufacturing and deployment.

## MODEL DIVERSITY CHALLENGES

Machine Learning is an umbrella term that covers a broad spectrum of algorithms; research works about FPGA acceleration and HLS-based design flows are mostly focused on a subset of ML models, i.e., multi-layer perceptrons and convolutional neural networks (CNN). Sometimes their field of action can be even smaller: for example, the original implementation of hls4ml was optimized for small, fully-connected models under tight latency constraints, reflecting the needs of a high-energy physics experiment at CERN. Appropriate implementation choices for such a specific case may not be beneficial to a generic model [3].

A second challenge is supporting sparsity: large models can be compressed to reduce their computation and memory requirements, for example by employing low precision data types (quantization) or by removing operations with zero values. Quantization is well suited to hardware acceleration, with dedicated HLS libraries for custom precision operators. Sparse tensors, on the other hand, imply irregular computation, communication, and memory access patterns, which result in poor efficiency when mapped on accelerators designed for dense models.

**Table 1: Preliminary results on CNN and GCN models.**

|                | Clock (MHz) | Registers | LUTs   | Latency (s) |
|----------------|-------------|-----------|--------|-------------|
| LeNet          | 146.75      | 44171     | 44 325 | 0.689       |
| ResNet-50      | 217.82      | 20861     | 20 522 | 284.640     |
| GCN 1 (dense)  | 204.80      | 14812     | 8965   | 1302.990    |
| GCN 2 (sparse) | 210.00      | 14639     | 8685   | 6.414       |

Another class of models that could benefit from hardware acceleration and requires unique design choices is Graph Neural Networks. Existing HLS-based design flows are good at extracting data- and instruction-level parallelism (e.g. by unrolling loops), but they are not equipped to deal with the irregular task-based patterns required by graph processing. Recently, we also started to extend the SODA framework to support Spiking Neural Networks (SNN) [2]. SNNs will require a completely new MLIR abstraction to capture temporal sequences, emulate SNN execution in software, and eventually map operations onto analog hardware.

Table 1 presents some results obtained synthesizing two CNNs and a simple graph convolutional network with SODA, targeting a Xilinx Zynq-7000 FPGA. LeNet is a relatively small network, which was synthesized as a single Verilog module; for ResNet-50 instead each single layer was outlined and implemented as a single accelerator. The GCN was implemented in two slightly different versions, one for dense inputs and one for sparse inputs. These preliminary results do not include any optimization, but they show that the SODA design flow has the potential to tackle a diverse set of models; more research will allow to tune the compilation passes according to the specific needs of different ML algorithms.

## CONCLUSION

The SODA framework provides a novel approach to the design of accelerators for Machine Learning, bridging the gap between algorithmic frameworks and hardware design through MLIR and High-Level Synthesis. Its multi-level strategy allows to operate analysis and transformations as compiler passes on different intermediate representations, retrieving relevant information that might otherwise be lost in the translation to C/C++. Future developments will focus on extending the framework to support different classes of ML models.

## REFERENCES

- [1] M. Blott et al. 2018. FINN-R: An end-to-end deep-learning framework for fast exploration of quantized neural networks. *ACM Transactions on Reconfigurable Technology and Systems (TRETs)* 11, 3 (2018), 1–23.
- [2] S. Curzel et al. 2021. Automated Generation of Integrated Digital and Spiking Neuromorphic Machine Learning Accelerators. In *2021 IEEE/ACM International Conference On Computer Aided Design (ICCAD)*. 1–7.
- [3] S. Curzel et al. 2021. De-specializing an HLS library for Deep Neural Networks: improvements upon hls4ml. *arXiv:2103.13060* (2021).
- [4] J. Duarte et al. 2018. Fast inference of deep neural networks in FPGAs for particle physics. *Journal of Instrumentation* 13, 07 (2018), P07027.
- [5] F. Ferrandi et al. 2021. Bambu: an Open-Source Research Framework for the High-Level Synthesis of Complex Applications. In *DAC 2021: 58th ACM/IEEE Design Automation Conference*. 1327–1330.
- [6] C. Lattner et al. 2020. MLIR: A Compiler Infrastructure for the End of Moore's Law. *arXiv:2002.11054* (2020).
- [7] J. Zhang et al. 2021. Towards Automatic and Agile AI/ML Accelerator Design with End-to-End Synthesis. In *IEEE 32nd International Conference on Application-specific Systems, Architectures and Processors (ASAP)*. IEEE.