# From PyTorch to Calyx: An Open-Source Compiler Toolchain for ML Accelerators

Jiahan Xie
jxie84@ucsc.edu
University of California, Santa Cruz
Santa Cruz, California, USA

Evan Williams
emw236@cornell.edu
Cornell University
Ithaca, New York, USA

Adrian Sampson
asampson@cs.cornell.edu
Cornell University
Ithaca, New York, USA

## Abstract

We present an end-to-end open-source compiler toolchain that targets synthesizable SystemVerilog from ML models written in PyTorch. Our toolchain leverages the accelerator design language Allo, the hardware intermediate representation (IR) Calyx, and the CIRCT project under LLVM. We also implement a set of compiler passes for memory partitioning, enabling effective parallelism in memory-intensive ML workloads. Experimental results demonstrate that our compiler can effectively generate optimized FPGA-implementable hardware designs that perform reasonably well against closed-source industry-grade tools such as Vitis HLS.

## 1 Introduction

High-level synthesis (HLS) is a promising approach for designing accelerators for ML applications [3, 7, 9, 11]. HLS allows us to compile high-level functional specifications from software (typically written in C or C++) to synthesizable RTL suitable for hardware implementation. However, few HLS approaches exist for compiling directly from the languages and frameworks that express ML models, such as PyTorch. Moreover, the prior work relies heavily on closed-source commercial HLS toolchains, such as Xilinx Vitis HLS. Recent MLIR-based HLS systems demonstrate the potential for open-source compiler stacks for accelerator design [12].

Building such a system introduces several challenges. ML models are typically written in Python and rely heavily on floating-point arithmetic, complex tensor operators, and multi-module execution patterns, none of which map cleanly onto low-level hardware descriptions. Supporting these workloads requires a compilation flow that can translate high-level tensor programs into hardware-oriented IRs while preserving enough structure to enable parallel execution and aggressive memory optimizations. However, traditional HLS approaches are limited, as high-level software programming languages are not well-suited for direct translation to hardware semantics. Additionally, modern ML accelerators must exploit fine-grained parallelism, but parallel execution introduces hazards such as memory port contention that require careful program analysis.

Therefore, we introduce an end-to-end open-source compiler toolchain that generates synthesizable SystemVerilog from PyTorch models through a structured compilation flow. Our system uses Allo [1] to translate PyTorch programs into MLIR, leverages domain-specific MLIR dialects [10] to preserve high-level tensor structure, and relies on the CIRCT [2, 6] infrastructure to perform hardware lowering. The resulting program is expressed in Calyx, whose explicit separation of control and hardware structure allows us to encode accelerator architectures and optimizations cleanly. Finally, we compile Calyx to SystemVerilog and use standard FPGA vendor tools such as Xilinx Vivado to synthesize, place–and–route, and deploy accelerators. Our contributions are as follows:

- An end-to-end open-source compiler stack from PyTorch to synthesizable SystemVerilog using Allo, CIRCT, and Calyx.
- Memory banking and partitioning analyses in Calyx enabling safe and efficient parallel access patterns.
- FPGA evaluation showing performance gain up to 2.21× Vitis HLS.

## 2 Background

### 2.1 Allo

Allo [1] is a compiler for constructing large-scale, high-performance hardware accelerators. Its lowering pipeline translates PyTorch programs into MLIR while preserving tensor semantics, control flow, and data-layout information. The resulting structured MLIR program integrates cleanly with downstream dialects and compiler infrastructures, forming a bridge between high-level ML frameworks and hardware generation backends.

### 2.2 Calyx

Calyx [8] is an IR and compiler infrastructure designed for generating hardware accelerators from high-level programming languages. It explicitly separates control flow from structural hardware descriptions, enabling optimizations that leverage both perspectives. The Calyx compiler then lowers the program into synthesizable RTL through a series of transformation and optimization passes.

### 2.3 CIRCT

Circuit IR Compilers and Tools (CIRCT) [6] is an open-source, LLVM-based infrastructure for building hardware compilers. Built atop MLIR [5], CIRCT provides a suite of hardware-specific dialects and transformation passes to support the development of custom compilation flows, hardware synthesis pipelines, and IRs.

Calyx is integrated into CIRCT as one of its dialects and is particularly valuable because CIRCT provides an MLIR-native path down to a hardware-oriented dialect, and Calyx is the natural endpoint of that lowering because it retains both structural hardware description and software-style control, giving us a unified space to express accelerator-specific rewrites before RTL generation.

## 3 Contribution

### 3.1 Overview

The goal of this work is to run PyTorch programs on custom hardware accelerators. We use FPGAs as the prototyping platform and develop a fully open-source compilation pipeline that transforms high-level Python programs into synthesizable hardware designs.
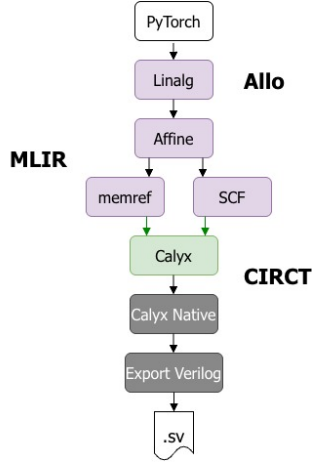
**Figure 1: Compilation pipeline from PyTorch through Allo to Calyx.**

Our toolchain is composed of several open-source components. We begin by using the Allo framework [1] to compile PyTorch models into MLIR programs. These MLIR programs are then progressively lowered via native MLIR passes to dialects supported by CIRCT. During this stage, we also apply high-level optimizations to improve performance and hardware compatibility.

Once the program reaches a form compatible with Calyx, we emit code in the Calyx IR. The Calyx compiler then performs hardware-specific transformations and generates synthesizable SystemVerilog. This hardware design is finally deployed to an FPGA, completing the software-to-hardware compilation path.

## 3.2 Lowering to Calyx

Figure 1 shows the compilation flow we developed and orchestrated to lower PyTorch to Calyx. We leverage Allo to lower to the MLIR Linalg dialect. We then lower to the MLIR Affine dialect, and then the Memref and SCF dialects. We use CIRCT to produce Calyx, and Vivado HLS to execute the design on an FPGA.

Building upon [10], which developed initial lowering support for structured control flow (SCF) constructs (e.g., `for`, `while`), we expanded this work by supporting additional SCF operations, including `if` and `parallel`. We also implemented the translation of software-level functions into Calyx components: each function becomes a hardware module with ports for scalar arguments and instantiated memories for tensor arguments, while function calls become component instantiations. We developed a full, general floating-point library for Calyx to provide floating-point support in CIRCT, including the integration of the Berkeley HardFloat [4] components. We also implemented bit-level IEEE-754 constant handling and human-readable attribute representations to enhance visibility and debuggability in CIRCT.

Collectively, these contributions form the foundation of our compiler stack, enabling high-level MLIR programs to be lowered into the Calyx IR and synthesized into RTL.

## 3.3 Memory and Parallelism Optimizations

To optimize the forward pass of the model, we leverage parallelism to increase hardware throughput. Calyx supports parallel execution as a first-class control construct, allowing us to explicitly model concurrent computation. Our task, then, is to expose and maximize parallelism - particularly in memory access patterns - while adhering to hardware constraints.

We base our approach on data parallelism: we partition memories into multiple banks so that different operations can access data concurrently. Since Calyx assumes each memory supports only a single read or write per cycle, banking allows us to duplicate storage and route accesses to different banks.

Our data-parallel strategy must avoid bank conflicts: memory banking increases the number of ports and introduces control logic to select which bank to access, often resulting in nested conditionals. If not optimized, these additional control paths increase latency and resource usage.

Our implementation supports cyclic memory partitioning, and we assume this scheme throughout. To route accesses to the correct bank, we use a `switch` statement (or nested `if-else` chains when `switch` is unavailable). A naive implementation that directly emits control branches for each bank leads to code-size blow-up. For a memory with $d$ dimensions and a partition factor $c$, the number of unique control branches scales as $c^d$. In Calyx, these branches are all instantiated as hardware, even if only one is active at runtime. This not only increases area usage but also results in deeper control finite state machines (FSMs), which hurt performance.

To illustrate the problem concretely, consider the following simple loop writing to a four-element memory with a memory banking factor of 2:

```
for (int i = 0; i < 4; ++i) {
    if (i % 2 == 0) {
        mem_bank_0[i / 2] = i;
    } else {
        mem_bank_1[i / 2] = i;
    }
}
```

To parallelize this loop, we materialize it with a factor of 2 using nested seq and par:

```
seq for (int i = 0; i < 2; ++i) {
    par for (int j = 0; j < 2; ++j) {
        int new_index = 2 * i + j;
        if (new_index % 2 == 0) {
            mem_bank_0[new_index / 2] = new_index;
        } else {
            mem_bank_1[new_index / 2] = new_index;
        }
    }
}
```

However, each par block must be unrolled into separate arms with statically known indices. In our example, this produces two parallel arms—j = 0 and j = 1—each computing a distinct new_index. Although each arm activates only one side of the banking condition at runtime, Calyx still instantiates both sides of the if−else because the predicate cannot be folded symbolically. As a result, even conflict-free access patterns may lead to multiple arms writing to the same physical bank, creating write hazards that Calyx cannot eliminate statically.

We implement two techniques to ensure safe and efficient memory parallelism: (1) we express banking by embedding the bank index into the memory's dimensional layout instead of guarding accesses with conditional logic; and (2) we rewrite loop nests whose structure would otherwise duplicate sequential controllers when executed in parallel.

For the first, rather than emitting branch logic per access, we raise the memory's dimensionality and bake the bank index into the first dimension. For instance:

```
1  seq for (int i = 0; i < 2; ++i) {
2      par for (int k = 0; k < 2; ++k) {
3          mem[k][i] = 2 * i + k;
4      }
5  }
```

Unrolling this loop yields:

```
1  seq for (int i = 0; i < 2; ++i) {
2      parallel execution {
3          execute par-arm-0 {
4              mem[0][i] = 2 * i + 0;
5          }
6          execute par-arm-1 {
7              mem[1][i] = 2 * i + 1;
8          }
9      }
10 }
```

Here, the bank index is a compile-time constant in each parallel arm, ensuring that memory accesses are disjoint and contention-free. This approach enables us to preserve parallelism without incurring the cost of unnecessary control overhead.

The second transformation operates at the level of loop structure. This example also highlights a broader point: loop transformations that are semantically equivalent in software do not necessarily yield equivalent hardware. Consider two nestings: one where seq(i) surrounds par(j), and another where par(j) surrounds seq(i). Although they are semantically equivalent to par(j) around seq(i) in software, they behave very differently in hardware. In the first form, there is a single sequential controller that iterates over i, and each iteration triggers a parallel group over j. In the second form, however, each parallel arm receives its own private sequential controller for iterating over i, effectively duplicating the entire FSM. This replication inflates the hardware area and increases control overhead. To prevent this unnecessary duplication, our compiler detects these patterns and rewrites parallel–sequential loop nests into schedules that share control logic while preserving the intended parallelism.

Through these memory and loop-aware transformations, our compiler produces parallel Calyx programs that are both correct and efficient for hardware execution.

## 4 Results and Evaluation

In this section, we evaluate our Calyx-based flow against Vitis HLS under two configurations. In Section 4.2, we compare baseline designs with no data parallelism: neither toolchain applies any banking strategy. In Section 4.3, we enable memory banking for both flows using matched partitioning factors, schemes, and dimensions, allowing a direct comparison of parallelized configurations.

Before presenting results, we note that Vitis HLS incorporates many mature and often implicit optimizations that are not currently implemented in our Calyx-based flow and cannot be disabled or inspected through pragmas. Thus, we do not expect Calyx to outperform Vitis in baseline latency or resource usage. Instead, our goal is to assess how competitive Calyx can be and whether our targeted banking transformations meaningfully narrow the gap. As shown below, while Vitis maintains an advantage in baseline configurations, Calyx becomes competitive once banking is enabled, offering an open-source and compiler-controlled foundation for further optimization.

### 4.1 Benchmark Models

We evaluate the performance and resource usage of our compiler by benchmarking three representative ML models and comparing them against a commercial HLS toolchain, Xilinx Vitis HLS. The models include a feed-forward neural network (FFNN), a convolutional neural network (CNN), and a multi-head attention (MHA) module.

The FFNN model takes an input of 64 features, followed by a fully connected layer of size $64 \times 48$, a ReLU activation, and a second fully connected layer of size $48 \times 4$.

The CNN model processes a $80 \times 60$ color image with 3 channels. The first layer performs a 2D convolution with $5 \times 5$ kernels, 3 input channels, and 8 output channels using unit strides. This is followed by a ReLU activation and a max-pooling layer with a $2 \times 3$ window. The resulting feature map is flattened and passed through a fully connected layer for binary classification.

The MHA model is derived from the Transformer architecture and uses 2 attention heads. Each head operates on a 21-dimensional subspace of a 42-dimensional embedding, with causal masking for autoregressive decoding.

### 4.2 Comparison Across Models

For comparison, we use Vitis HLS, a widely used commercial HLS tool. The Allo project provides a shared frontend that lowers PyTorch models to MLIR, which is then further compiled to either HLS C++ for Vitis or Calyx for our flow. To ensure fairness, the MLIR input is held constant across both paths. All HLS pragmas are disabled except for #pragma ARRAY_PARTITION, which is used to compare against our memory banking configuration. We ensure consistent banking factors, schemes (cyclic), and dimensions across both flows.

Wall-clock latency is computed from each design's cycle count and its maximum post–place-and-route frequency that meets timing, obtained by sweeping the timing constraint and selecting the implementation that achieves the lowest overall latency.

Figure 2 shows the wall-clock latency of each model across the two toolchains. Vitis outperforms our Calyx-based flow in all cases, particularly on the CNN model. This performance gap is largely due to limitations in Calyx's memory model: Calyx only supports single-dimensional memories, requiring us to flatten multi-dimensional tensors. As a result, access indices—often affine expressions of loop variables—are lowered to hardware as expensive arithmetic operations like multiplication and modulo, which introduce latency. Furthermore, the CNN model's convolution and pooling layers naturally lead to deeper loop nests, amplifying this cost.

Table 1 reports resource utilization across the same models. We observe that Vitis uses more BRAMs for storing weights and biases, while Calyx consumes significantly more LUTs and FFs due to
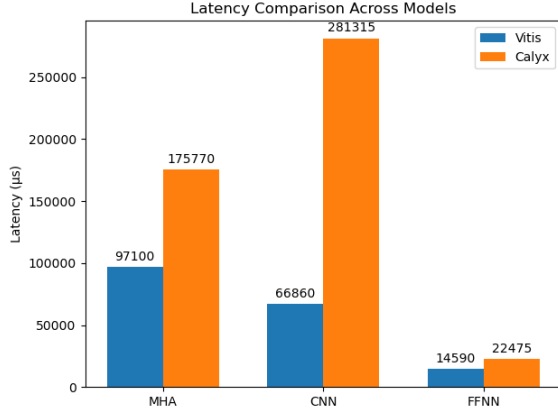
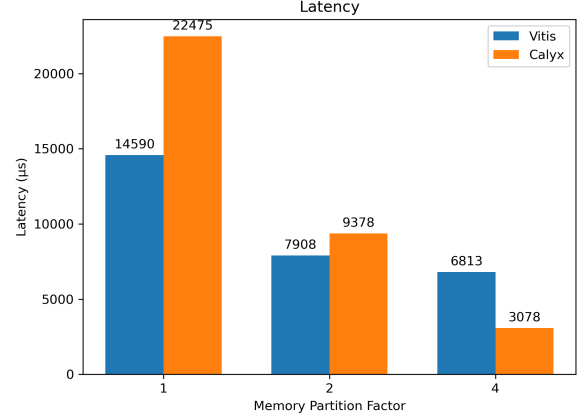Figure 2: Wall-clock latency comparison across models.

Table 1: Resource usage across models.

| Resource | MHA | | CNN | | FFNN | |
|---|---|---|---|---|---|---|
| | Vitis | Calyx | Vitis | Calyx | Vitis | Calyx |
| LUTs | 7846 | 33312 | 3136 | 4574 | 2011 | 3730 |
| FFs | 4017 | 5561 | 1815 | 1223 | 1281 | 742 |
| BRAMs | 194 | 71 | 213 | 43 | 43 | 9 |
| DSPs | 19 | 67 | 5 | 14 | 5 | 6 |

its use of explicit FSMs for control. Overall, Calyx incurs higher resource usage—particularly LUTs—because of its verbose control modeling and lack of aggressive scheduling.

## 4.3 FFNN Memory Partitioning

To further understand the impact of memory banking, we consider the FFNN model. In this study, every memory is partitioned cyclically along each dimension. In Vitis, this is done via `#pragma ARRAY_PARTITION`, while in Calyx the partitioning is implemented via our compiler pass. We compare both latency and resource usage across varying banking factors.

Figure 3 shows the latency across different partition factors. For `factor=1` and `factor=2`, Vitis performs better. However, at `factor=4`, Calyx becomes faster. Moreover, the relative speedup for Vitis is modest: increasing the banking factor from 2 to 4 yields diminishing returns, with a speedup of only $7908/6813 \approx 1.16$. Given that all matrices are two-dimensional, the theoretical maximum speedup is $2^2 = 4$ under ideal conditions. The limited gain suggests that other bottlenecks remain.

In contrast, the Calyx-based flow shows significant improvement with higher partitioning. The speedup from `factor=1` to `factor=2` is $22475/9378 \approx 2.40$, and from `factor=2` to `factor=4` is $9378/3078 \approx 3.05$. This demonstrates the effectiveness of our memory banking analysis and transformation passes, which allow Calyx to unlock more parallelism at the memory level.

Table 2 shows the resource usage corresponding to these experiments. DSP usage is comparable across both toolchains. Vitis



Figure 3: Latency vs. partition factor for FFNN.

Table 2: Resource usage vs. partition factor for the FFNN model

| Resource | Partition 1 | | Partition 2 | | Partition 4 | |
|---|---|---|---|---|---|---|
| | Vitis | Calyx | Vitis | Calyx | Vitis | Calyx |
| LUTs | 2011 | 3730 | 6021 | 13197 | 13799 | 49121 |
| FFs | 1281 | 742 | 4036 | 3145 | 15083 | 10657 |
| BRAMs | 43 | 9 | 39 | 10 | 64 | 20 |
| DSPs | 5 | 6 | 7 | 20 | 80 | 69 |

uses slightly more BRAMs and FFs, while Calyx consumes significantly more LUTs, particularly at `factor=4`, due to the additional control logic needed for managing multiple memory banks. These results illustrate the tradeoff between performance and hardware complexity introduced by fine-grained memory partitioning.

## 5 Conclusion

This work presents a complete open-source compilation toolchain that transforms PyTorch programs into synthesizable hardware designs using Allo, MLIR, CIRCT, and Calyx. We bridge the gap between software-level ML programs and FPGA-executable hardware accelerators by implementing support for structured control flow, function modeling, floating-point arithmetic, and memory layout transformations.

We focus in particular on optimizing memory access concurrency through static memory banking and control restructuring. We leave hardware pipelining, resource sharing, and compilation of the backward pass for future work. Our evaluation on representative ML models shows that while the Calyx-based flow trails behind commercial tools like Vitis in general-purpose scheduling and resource efficiency, it demonstrates promising performance gains when aggressive memory partitioning is applied. These results highlight the potential of Calyx as a research and prototyping platform for hardware accelerator compilation, especially in settings where open-source and customization are prioritized.

# References

[1] Hongzheng Chen, Niansong Zhang, Shaojie Xiang, Zhichen Zeng, Mengjia Dai, and Zhiru Zhang. 2024. Allo: A Programming Model for Composable Accelerator Design. *Proc. ACM Program. Lang.* 8, PLDI, Article 171 (June 2024), 28 pages. doi:10.1145/3656401

[2] John Demme and Aaron Landy. 2022. Using CIRCT for FPGA Physical Design. Presented at the 2nd Workshop on Languages, Tools, and Techniques for Accelerator Design (LATTE'22). https://capra.cs.cornell.edu/latte22/paper/10.pdf

[3] Javier Duarte et al. 2018. Fast inference of deep neural networks in FPGAs for particle physics. *JINST* 13, 07 (2018), P07027. arXiv:1804.06913 [physics.ins-det] doi:10.1088/1748-0221/13/07/P07027

[4] John R. Hauser. 2019. Berkeley HardFloat. https://github.com/ucb-bar/berkeley-hardfloat

[5] Chris Lattner, Mehdi Amini, Uday Bondhugula, Albert Cohen, Andy Davis, Jacques Pienaar, River Riddle, Tatiana Shpeisman, Nicolas Vasilache, and Oleksandr Zinenko. 2021. MLIR: scaling compiler infrastructure for domain specific computation. In *Proceedings of the 2021 IEEE/ACM International Symposium on Code Generation and Optimization* (Virtual Event, Republic of Korea) *(CGO '21)*. IEEE Press, 2–14. doi:10.1109/CGO51591.2021.9370308

[6] LLVM CIRCT. "CIRCT" / Circuit IR Compilers and Tools. [n. d.]. https://github.com/llvm/circt

[7] Rachit Nigam, Sachille Atapattu, Samuel Thomas, Zhijing Li, Theodore Bauer, Yuwei Ye, Apurva Koti, Adrian Sampson, and Zhiru Zhang. 2020. Predictable accelerator design with time-sensitive affine types. In *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation* (London, UK) *(PLDI 2020)*. Association for Computing Machinery, New York, NY, USA, 393–407. doi:10.1145/3385412.3385974

[8] Rachit Nigam, Samuel Thomas, Zhijing Li, and Adrian Sampson. 2021. A compiler infrastructure for accelerator generators. In *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems* (Virtual, USA) *(ASPLOS '21)*. Association for Computing Machinery, New York, NY, USA, 804 – 817. doi:10.1145/3445814.3446712

[9] Hardik Sharma, Jongse Park, Divya Mahajan, Emmanuel Amaro, Joon Kyung Kim, Chenkai Shao, Asit Mishra, and Hadi Esmaeilzadeh. 2016. From high-level deep neural models to FPGAs. In *The 49th Annual IEEE/ACM International Symposium on Microarchitecture* (Taipei, Taiwan) *(MICRO-49)*. IEEE Press, Article 17, 12 pages.

[10] Mike Urbach and Morten B. Petersen. 2022. HLS from PyTorch to SystemVerilog with MLIR and CIRCT. Presented at the 2nd Workshop on Languages, Tools, and Techniques for Accelerator Design (LATTE'22). https://capra.cs.cornell.edu/latte22/paper/2.pdf

[11] Stylianos I. Venieris and Christos-Savvas Bouganis. 2016. fpgaConvNet: A Framework for Mapping Convolutional Neural Networks on FPGAs. In *2016 IEEE 24th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*. 40–47. doi:10.1109/FCCM.2016.22

[12] Hanchen Ye, HyeGang Jun, Hyunmin Jeong, Stephen Neuendorffer, and Deming Chen. 2022. ScaleHLS: a scalable high-level synthesis framework with multi-level transformations and optimizations: invited. In *Proceedings of the 59th ACM/IEEE Design Automation Conference* (San Francisco, California) *(DAC '22)*. Association for Computing Machinery, New York, NY, USA, 1355–1358. doi:10.1145/3489517.3530631