

HPVM: Hardware-Agnostic Programming for Heterogeneous Parallel Systems

Adel Ejeh , Aaron Councilman, and Akash Kothari, *University of Illinois Urbana-Champaign, Urbana, IL, 61801, USA*

Maria Kotsifakou, *Runtime Verification, Inc., Champaign, IL, 61820, USA*

Leon Medvinsky, Abdul Rafae Noor, Hashim Sharif, Yifan Zhao, Sarita Adve, Sasa Misailovic, and Vikram Adve, *University of Illinois Urbana-Champaign, Urbana, IL, 61801, USA*

We present Heterogeneous Parallel Virtual Machine (HPVM), a compiler framework for hardware-agnostic programming on heterogeneous compute platforms. HPVM introduces a hardware-agnostic parallel intermediate representation with constructs for the hierarchical task, data, and pipeline parallelism, including dataflow parallelism, and supports multiple front-end languages. In addition, HPVM provides optimization passes that navigate performance, energy, and accuracy tradeoffs, and includes retargetable back ends for a wide range of diverse hardware targets, including central processing units, graphics processing units, domain-specific accelerators, and field-programmable gate arrays. Across diverse hardware platforms, HPVM optimizations provide significant performance and energy improvements, while preserving object-code portability. With these capabilities, HPVM facilitates developers, domain experts, and hardware vendors in programming modern heterogeneous systems.

With the slowdown of Moore's law and the end of Dennard scaling, heterogeneous architectures are increasingly dominating the systems used for modern applications. These systems have been evolving to include a plethora of computing- and energy-efficient processing elements (PEs), ranging from graphics processing units (GPUs) and field-programmable gate arrays (FPGAs) to fixed-function and programmable domain-specific accelerators. Enabling hardware-agnostic programming of these heterogeneous systems is important to facilitate their use by a broad range of software developers. By hardware-agnostic, we mean that the entire programming process, including the software itself and the iterative manual development and tuning processes, should not be specific to a particular target system. Moreover, any system-specific performance tuning should be automated by the compilation flow as much as possible, and if that is not entirely

possible, the programming process should minimize the need for changes to the application source code.

Currently, programming such heterogeneous devices present many challenges, including the need to use diverse hardware-specific programming languages, poor source-code portability, lack of object-code portability, which is essential in certain domains, such as mobile applications, and the need for system-specific performance tuning.¹ Moreover, in many emerging application domains [e.g., video analytics, augmented and virtual reality (AR/VR), mobile robotics, etc.], it is also crucial—but challenging for application developers—to increase performance and/or energy efficiency by sacrificing small amounts of accuracy or application quality.

We believe that the solution to these challenges lies in a suitable compiler and autotuning infrastructure, with the right abstractions of parallelism, that allows seamless compilation of hardware-agnostic code into multiple general-purpose and/or specialized hardware targets. In addition, the compiler must automatically optimize programs using efficient hardware-level primitives and domain-specific transformations while requiring minimal source-code tuning by the developer for performance. This can be achieved by separating the hardware-agnostic functional specification from the

hardware-specific tuning, which itself can be automated using autotuning.

We present Heterogeneous Parallel Virtual Machine (HPVM),² a compiler framework that achieves the aforementioned goals. HPVM addresses all the programming challenges by providing a *compiler IR design and code-generation framework that is retargetable* to a wide range of heterogeneous parallel architectures, *hardware-agnostic language front ends* that support ease of programming, and a *domain- and hardware-specific optimization framework* that automatically navigates performance, energy, and accuracy tradeoffs. HPVM's flexibility allows *programmers* to easily write and optimize code for heterogeneous platforms while giving *hardware vendors* the ability to easily extend the compiler infrastructure with new hardware back ends and corresponding optimizations.

No existing infrastructure we know of provides the combination of features needed for hardware-agnostic programming of a broad range of accelerator-based systems, including a flexible abstraction of parallelism, source- and object-code portability, and the separation of functionality from (automatable) performance tuning and accuracy-aware approximation tuning. Spatial³ is based on a valuable parallelism abstraction, but emphasizes hardware design. TVM⁴ supports diverse hardware targets but is narrowly focused on machine learning. Perhaps the best alternative is multilevel intermediate representation (MLIR), which is flexible enough to support a wide range of compiler goals and all missing features could be added in the future. At present, it lacks a virtual object-code format essential for object-code portability in accelerator-based systems. It also lacks an approximation framework (including diverse approximations, support for approximation tuning, and dynamic adaptation of approximations), which is crucial for many emerging applications on heterogeneous systems. We believe it would be valuable to add new "dialects" to MLIR based on the techniques described here in order to provide these missing capabilities.

HPVM OVERVIEW

HPVM is an open-source (<https://gitlab.engr.illinois.edu/llvm/hpvm-release/>) retargetable compiler infrastructure that uses a common abstraction of parallelism to define the compiler intermediate representation (IR), a virtual instruction set architecture (ISA), and a runtime system. Figure 1 shows the HPVM compiler stack. HPVM's parallel abstraction is designed to capture the multiple forms of parallelism available on heterogeneous systems, in a hardware-agnostic manner.

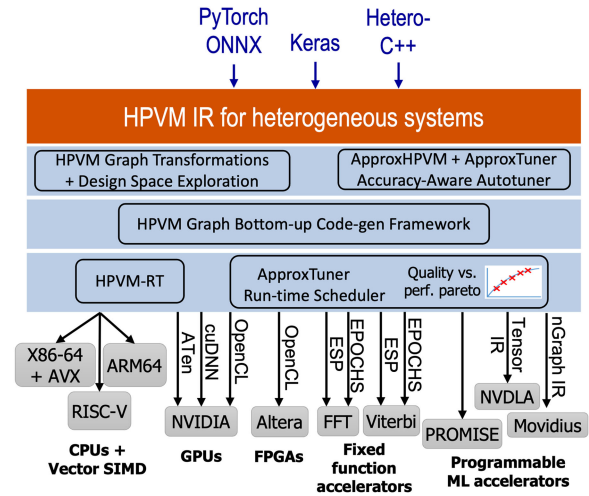


FIGURE 1. HPVM compiler infrastructure. General-purpose applications use HeteroC++ and are compiled using graph transformations and autotuning, and run using HPVM-RT. Tensor-domain applications use Keras, PyTorch, or ONNX, compiled to the ApproxHPVM tensor IR extensions, and optimized by ApproxTuner using static and dynamic approximation tuning.

The HPVM IR is currently built on top of LLVM, i.e., it uses LLVM IR to represent (scalar and vector) computations. HPVM is a modular framework that facilitates compiler optimizations and benefits from the advanced optimization and code-generation capabilities offered by LLVM. Our ApproxHPVM⁵ and ApproxTuner⁶ extensions provide a tensor-domain programming model and enable tradeoffs between accuracy and performance/energy for tensor applications, which we describe later. With that, HPVM provides a framework that can seamlessly integrate domain-specific tensor code with general-purpose code.

Starting with hardware-agnostic code written in a supported language—at present, these include a parallel dialect of C++ (HeteroC++) for general-purpose parallel code, and Keras, PyTorch, and ONNX for deep learning—HPVM's front ends translate the parallelism in the application into equivalent code using HPVM's parallel IR abstractions. HPVM's optimization frameworks, which support the optimization needs of different domains, tune the code for the target system. Finally, code is generated for individual (specified or default) target devices and for the host central processing unit (CPU) to generate executable binaries, leveraging existing LLVM code generators for each device wherever possible. Optional features enable design space exploration for FPGAs⁷ and GPUs, and an approximation autotuning and dynamic adaptation framework for tensor programs.

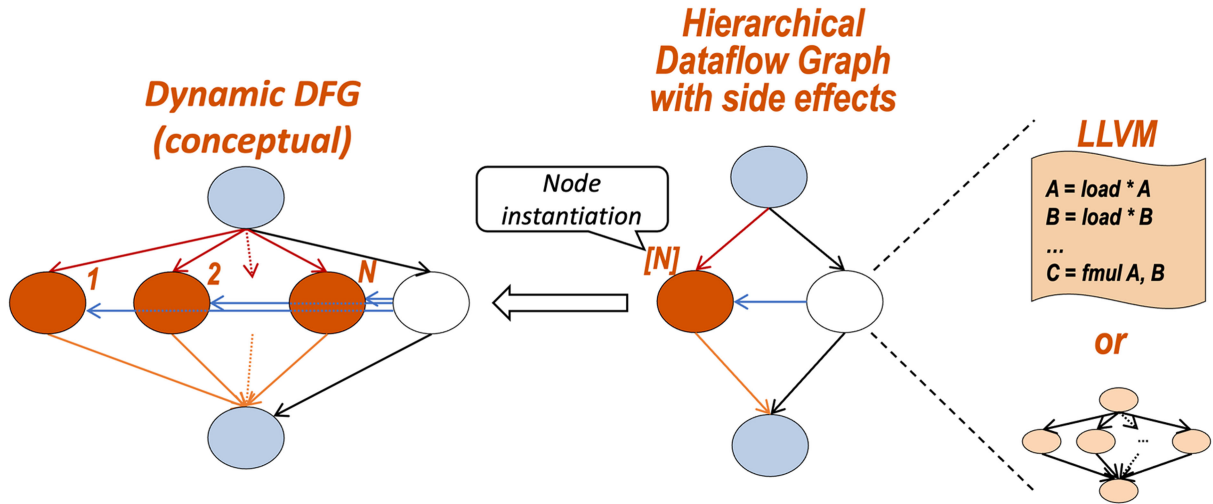


FIGURE 2. HPVM dataflow graph representation.

HPVM PARALLEL ABSTRACTION AND IR

A parallel program is represented in HPVM as a host program plus one or more acyclic dataflow graphs (DFGs; see Figure 2). Each DFG is hierarchical: a node is either a *leaf node* representing a unit of computation, or an *internal node* containing an entire “child” DFG. Edges in the DFG represent explicit, logical data transfer between nodes, with blocking semantics, i.e., the sink node of an edge will block if it attempts to access the data coming in via that edge. Leaf nodes can include references to global shared memory. HPVM programs are *assumed* to be data race free, i.e., explicit ordering through DFG edges or fine-grain synchronization operations must be used to enforce ordering or mutual exclusion between conflicting accesses to the same shared memory location.

Each node in a DFG has one or more dynamic instances, represented using a dynamic replication factor (e.g., for the iterations of a parallel loop). Node instances are assumed to be fully independent, i.e., safe to run in parallel. At the IR level, hints from the source program or front end are used to specify, which device a node must be targeted to, although an automatic partitioning system built on the autotuner described here will allow automatic assignment of nodes to devices.

HPVM Compiler Internal Representation

The HPVM compiler IR is a direct implementation of the hierarchical graph representation, using LLVM IR for both the host program and the node computations. Every HPVM node is represented by an ordinary LLVM function. The HPVM IR defines a set of LLVM “intrinsic functions” to describe the details of the DFGs, as well

as host operations to launch a graph computation asynchronously, wait for completion, and to transfer inputs and outputs to and from a graph (either as one-time transfers or as streams). An LLVM-to-DFG translation pass converts the intrinsics-based description into explicit graph data structures, before graph transformations and back-end code generation occurs.

HPVM leaf nodes use LLVM scalar and vector instructions for individual computational tasks and can use LLVM atomics and other synchronization operations to ensure data race freedom. Just like with LLVM, an HPVM IR program is a fully self-contained and executable “virtual object code” program, which can be used to achieve object code portability (via install-time translation) for any targets that support portable LLVM IR.

ApproxHPVM and Tensor Extensions

We extend the LLVM IR used by HPVM with higher level tensor intrinsics, such as matrix multiplication and convolution. This enables HPVM to retain domain-specific semantics for tensor application domains, such as deep learning and image processing. We leverage these intrinsics for three purposes as follows:

- 1) domain-specific optimizations, such as operator fusion;
- 2) code generation targeting machine learning accelerators, such as Movidius and NVDLA;
- 3) accuracy-aware optimizations for tensor operations, which apply and tune software and hardware approximation techniques on the tensor operators systematically such that user-specified quality of service (QoS) constraints are respected.

```

void* Section = __hetero_section_begin();
for(int i = 0; i < left_dim ; i++){
    for(int j = 0; j < right_dim; j++){
        __hetero_parallel_loop(
            /* Num Parallel Enclosing Loops */ 2,
            /* Num Input Objects */ 6,
            Res, Res_Size, V1, V1_Size, V2, V2_Size,
            left_dim, right_dim, common_dim,
            /* Num Output Objects */ 1,
            Res, Res_Size ,
            /* Optional Node Name */ "matmul_parallel_loop" );
        __hetero_hint(/* TARGET DEVICE */ Target::GPU);
        Res[i*right_dim + j] = 0;
        for(int k = 0 ; k < common_dim; k++){
            // Res[i,j] += V1[i,k] + V2[k,j]
            Res[i*right_dim + j] += V1[i*common_dim + k] * V2[k*right_dim + j];
        }
    }
}
__hetero_section_end(Section);

```

FIGURE 3. Matrix multiply written in HeteroC++. The outer two loops over i and j are parallel, whereas the innermost loop executes sequentially over k . Note that a pointer and the size of the allocated memory it points to together constitute a single input/output object.

HPVM FRONT ENDS

HeteroC++ Front End

HeteroC++ is an experimental parallel dialect of C/C++ that enables easy parallelization using HPVM. Computations in HeteroC++ are described as parallel tasks or parallel loops, which are lowered to nodes in the HPVM DFG using function outlining. The programmer annotates the incoming and outgoing pointer variables for each task or loop (see Figure 3), which are used to infer the edges between nodes. Importantly, HeteroC++ directives are *completely hardware-agnostic*, lacking any target-specific design or tuning parameters. Through this programming interface, the same high-level program in HeteroC++ can be compiled for CPUs, GPUs, FPGAs, and hardware accelerators by generating target-agnostic HPVM IR. HeteroC++ can be viewed as a small subset of OpenMP, supporting flexible nested loop and task parallelism and implicit offloading with “target hints.” We are currently extending our autotuning framework to automatically determine the target device for a node given a target heterogeneous system.

Deep Learning Front Ends

HPVM supports front ends for three popular neural network frameworks: Keras, PyTorch, and ONNX (a neural network exchange format). These front ends compile high-level code written in these frameworks to HPVM IR DFGs with tensor operations (convolutions and matrix multiplications) in the leaf nodes, and tensor data on the DFG edges.

HPVM BACK-END CODE GENERATION

To generate code for the different target devices, a bottom-up traversal of the DFG is performed by each device’s back end, translating every leaf node to device-specific code for one or more devices using the following device-specific translators, and generating host code as needed. Standard LLVM CPU back ends for the x86, ARM, and RISC-V families are used to generate CPU code, both for the host program and to compile DFG nodes targeted for the CPU.

GPU Back End

For every leaf node targeted to the GPU, the GPU back-end generates an OpenCL kernel in C, replacing HPVM’s thread-indexing intrinsics with OpenCL API calls. In addition, the necessary code to launch the kernel, set up its arguments, and copy arguments to device memory is generated for the host program. The dynamic replication factor of the kernel’s leaf node determines the OpenCL workgroup size.

FPGA Back End

The FPGA back end also performs OpenCL code generation for the corresponding leaf nodes in a similar manner to the GPU back end. The OpenCL kernels are compiled using Intel FPGA SDK for OpenCL to generate an FPGA bitstream for Altera FPGAs. Intel recommends using Single Work Item Kernels, so we added a node sequentialization transformation to generate sequential loops from the dynamic replication factors of nodes. Finally, the back-end uses separate OpenCL command

queues for each kernel to support concurrent execution on the FPGA, and OpenCL events are used to synchronize accordingly. This is done by issuing the corresponding calls to the HPVM Runtime (HPVM-RT).

Compiling to Fixed Function Accelerators

The HPVM fixed-function accelerator back-end requires LLVM functions (e.g., generated from C source code) representing the functional semantics of each of the target accelerator's operations. It looks for matches between application leaf node functions and accelerator operations by using a semantic matching scheme for LLVM functions described in Dasgupta *et al.*'s⁸ work based on program dependence graph isomorphism. When a match is found, it replaces the leaf node function body with code to launch the corresponding accelerator construct. This back end has been used to generate code for fast Fourier transform (FFT) and Viterbi fixed-function accelerators.

Back Ends for Deep Learning Accelerators

NVDLA is NVIDIA's programmable accelerator for energy-efficient tensor operations commonly used in deep learning (convolutions, matrix multiplication, relu, max pooling, etc.). HPVM directly maps HPVM tensor IR operations to NVDLA's tensor IR constructs for operations, such as relu and max pooling. It fuses HPVM IR operations, such as convolution and tensor, add before mapping to NVDLA constructs. The NVDLA IR is then compiled to object code by the vendor NVDLA compiler. We support two NVDLA modes: FP16 and INT8. For INT8, we perform floating-point to integer quantization using Distiller, a third-party tool for neural network quantization.⁹

Intel Movidius vision processing unit (VPU) is a deep learning accelerator used on edge devices, such as the Neural Compute Stick. Similar to the NVDLA workflow, the HPVM-Movidius back-end translates HPVM's DFG of tensor operations to the Intel nGraph IR—a DFG-based IR with tensor operations specific to the Movidius back end. The HPVM-Movidius back-end directly invokes nGraph compiler interfaces (linked with the HPVM toolchain), which apply hardware-specific optimizations and generate object code. These interfaces also insert code for offloading the compute kernels to the Movidius accelerator.

ATen Back End

HPVM tensor operations can be translated to high-performance libraries. We support compilation to the ATen

back end—the tensor library used by PyTorch for compiling to GPUs (using cuDNN) and CPUs (using MKL-DNN). This back-end enables HPVM to map to efficient approximate constructs supported by ATen, including support for sparse tensor operations (used in pruned neural networks) and quantized tensor operations.

HPVM OPTIMIZATION FRAMEWORKS

HPVM includes two code optimization frameworks:

- a graph optimization framework that automatically tunes an HPVM program for a specific accelerator target;
- an accuracy-aware optimization framework, ApproxTuner, that uses approximation techniques to gain performance and energy improvements.

Optimizations and Autotuning

The HPVM graph optimization framework includes both HPVM DFG graph optimizations and regular LLVM optimizations on the node functions. These optimizations are applied to HPVM leaf nodes (i.e., kernels) singly or in pairs. They include the following.

- *Argument privatization* finds pointer arguments that are marked as thread-private and creates a local (private) copy of them.
- *Loop unrolling* unrolls loops using a specified unroll factor.
- *Greedy loop fusion* considers fusing all pairs of LLVM loops in a single leaf-node function that are legal to fuse, from the outermost nesting level to the innermost.
- *Node fusion* fuses DFG nodes that are connected with an edge and have the same dynamic replication factor and no fusion-preventing dependencies.

HPVM's optimization framework incorporates autotuning using the HyperMapper design space exploration framework¹⁰ to automatically tune hardware-agnostic programs for FPGA⁷ and GPU (separately for now) using the abovementioned optimizations. A performance model is used for FPGA tuning to avoid long synthesis times, while we use direct execution on the GPU for GPU tuning. The autotuner selects which optimizations will be applied and, for loop unrolling, what unroll factor to use for each loop.

ApproxTuner

ApproxTuner is an end-to-end accuracy-aware optimization framework for tensor-based components of programs, such as deep neural networks and image

processing pipelines. ApproxTuner takes a program compiled to HPVM IR and a desired end-to-end quality (QoS) threshold, and automatically maps tensor operations to different approximations to maximize performance and/or energy benefits while ensuring that the QoS is achieved.

ApproxTuner uses a *novel three-phase, predictive tuning* strategy to map approximations on diverse hardware and maintain object-code portability. To enable efficient tuning, ApproxTuner uses a *predictive tuning approach*, which uses accuracy and performance prediction models instead of expensive empirical evaluations.

At development-time, ApproxTuner tunes the program with *hardware-independent* approximations, finding a number of configurations. A *configuration* is a mapping from each tensor operator to one or more approximations and parameter settings for those approximations. Predictive tuning is used to compute a Pareto-optimal frontier of these configurations in the performance-accuracy tradeoff space, and this Pareto curve is shipped with the HPVM IR. At install time, ApproxTuner retunes these configurations with any hardware-specific knobs that exist on the target hardware, and refines the Pareto curve by measuring real performance on the target hardware. The final Pareto curve is shipped with the application binary and used by a dynamic tuner at runtime (described in the next section).

ApproxTuner supports software approximations, such as reduction sampling (using a subset of inputs in the reduction), perforated convolutions, and sampled convolutions. At the hardware level, ApproxTuner supports reduced precision using FP16 or INT8, and mapping to low-voltage knobs on PROMISE, an experimental analog ML accelerator.¹¹

RUNTIME SYSTEMS AND SCHEDULERS

HPVM Runtime

The HPVM-RT enables programs compiled in HPVM to efficiently execute on a diverse range of hardware platforms. HPVM-RT provides a memory tracker, which maintains the location of the most recent *dirty* copy of memory objects used across HPVM DFGs, and uses that to determine when a memory copy between host and device is necessary. Also, HPVM-RT communicates with the corresponding device runtimes (OpenCL runtime for GPU and FPGA), creating the necessary OpenCL objects (Platform, Context, Kernels, and Command Queues), copying memory back and forth, setting kernel arguments, and launching kernels on the corresponding devices.

ApproxTuner Runtime Approximation Tuning

ApproxTuner's runtime approximation tuner enables applications to maintain performance goals under changing system or application conditions, e.g., low-power modes. The tuner adapts per-operation approximation knobs (described in the previous section) to adapt the accuracy-performance tradeoff while using a system monitor to detect system slowdowns. The tuner uses the Pareto curves shipped with the HPVM program binary to choose the most accurate approximation parameter settings that satisfy desired quality metrics. These parameter settings are used as arguments for the tensor operations on each invocation, making it easy to change configurations quickly.⁶

Third-Party SoC Schedulers

HPVM includes support for two third-party SoC schedulers as part of a collaborative project. The first is the ESP system¹² from Columbia, a hardware-design framework that enables easy integration of new accelerators into SoCs. ESP is being used as part of the EPOCHS project (led by IBM) to design an SoC specialized for autonomous vehicles with several accelerators, such as FFT, Viterbi, and NVDLA, and a RISC-V host. This approach enables hardware-agnostic programmability as well as potential accuracy-aware optimization for SoCs designed using ESP.

The second is the novel EPOCHS scheduler to schedule different application "tasks" onto the available accelerators in an SoC, including static and dynamic mappings. Our compiler targets the scheduler library API to launch the "tasks" and specify the possible target devices for each. These HPVM back ends are easily extensible to other similar SoC design and task scheduling frameworks.

EXPERIMENTAL EVALUATION

Because of space constraints, we focus on results from recent work. In an early HPVM publication, we reported results for seven Parboil benchmarks, showing that the HPVM infrastructure can compile the *same* hardware-agnostic code for NVIDIA GPUs and Intel vector instructions (AVX) and achieve performance competitive with that of separately hand-tuned OpenCL code for each target.²

In the following, we show experimentally that: 1) our GPU autotuner achieves excellent speedups; 2) a *single, hardware-agnostic* program can be partitioned for GPU and FPGA, achieving much higher speedups than on GPU alone; 3) when small reductions in end-to-end DNN inference accuracy are tolerable, ApproxTuner

can roughly double the performance of a wide range of DNNs⁶; and 4) dynamic approximation tuning enables a DNN to maintain image classification throughput despite a large reduction in GPU clock frequency, with only a small loss in inference accuracy.⁶ The first two experiments are new for this work.

Optimizing Applications for GPU

To evaluate the GPU autotuner, we conducted an experiment on seven benchmarks: a five-stage camera image processing pipeline (CAVA), (<https://github.com/yaoyuannnn/cava/>) an Edge Detection program for gray-scale images from² whose DFG is a six-node DAG, and five multikernel benchmarks from Rodinia¹³: breadth-first search (BFS), backpropagation (Backprop), speckle reducing anisotropic diffusion (SRAD), and both the “Euler” and “Pre-Euler” implementations of the computational fluid dynamics solver. The hardware was an Intel Xeon 4216 CPU and **NVIDIA RTX 2080 Ti GPU**. Each data point is the average of five runs, with error bars showing the range.

Figure 4(a) shows speedups compared to baselines compiled from the programs using HPVM’s single-threaded CPU back end, without applying our optimizations. Orange and gray bars show speedups achieved by HPVM without and with our optimizations, respectively, including autotuning in the latter. For Euler, Pre-Euler, and BFS, autotuning does not add much benefit over the GPU version. For the other four benchmarks, the widely varying configurations selected by autotuning demonstrate the importance of autotuning to achieve source code portability and hardware-agnostic programming. For example, Edge and Backprop had loops that were fully unrolled, Euler and SRAD had partially unrolled loops, while CAVA had only one loop unrolled. Making similar design choices manually requires trial and error, and achieving sufficient coverage of the search space through manual exploration is often impractical, and most seriously, can lead to hardware-specific tuned code, which compromises source-code portability.

SRAD suffers a slowdown on the GPU, even with autotuning, due to a sum-reduction kernel that is not parallelized by HPVM. Through autotuning, loop unrolling was able to reduce the slowdown from $5.5\times$ to $1.7\times$. Automating the parallelization of reductions for GPUs is planned, along with support for other important GPU optimizations like tiling for GPU registers and scratchpad memory.

Including initial random sampling to initialize HyperMapper, 212 designs were evaluated for CAVA and Backprop, while 400 were evaluated for Edge

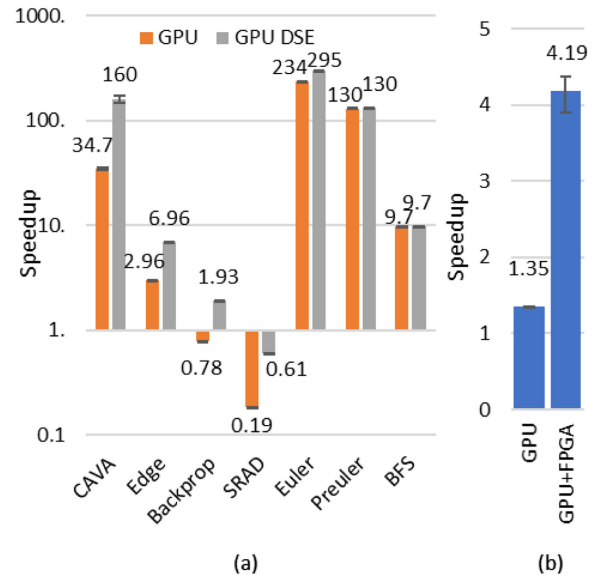


FIGURE 4. (a) Results for GPU autotuning. The figure shows a speedup of hardware-agnostic code running unoptimized on GPU (orange bar) and optimized using autotuning (gray bar) compared to CPU. GPU is NVIDIA 2080 Ti. (b) Speedups compared with CPU for two device partitionings of Edge Detection. GPU is NVIDIA Quadro P1000 and FPGA is Arria 10 GX.

Detection, proportional to their parameter counts. The remaining benchmarks had so few parameters that all designs in the search space were evaluated before reaching their set iteration counts, producing between 2 and 48 samples. This contributed to a large range of autotuning times, taking between 1 and 445 min for a full run (average 146). Time per sample ranged between 23 and 388 s (average 122) since each sample is executed on the GPU during autotuning.

Partitioning Applications on Multiple Devices

Partitioning applications across multiple devices (e.g., GPU and FPGA) poses two main challenges: 1) programmability since different devices tend to have different programming models and languages; and 2) partitioning decisions, which requires an intimate knowledge of the performance tradeoffs for each target device. HPVM allows us to overcome the first challenge by supporting a unified hardware-agnostic programming language and IR that can be targeted to multiple different devices. As a proof-of-concept, we manually partitioned the Edge Detection benchmark of the previous section using device hints on the nodes, putting the reduction kernel (which dominates the execution time)

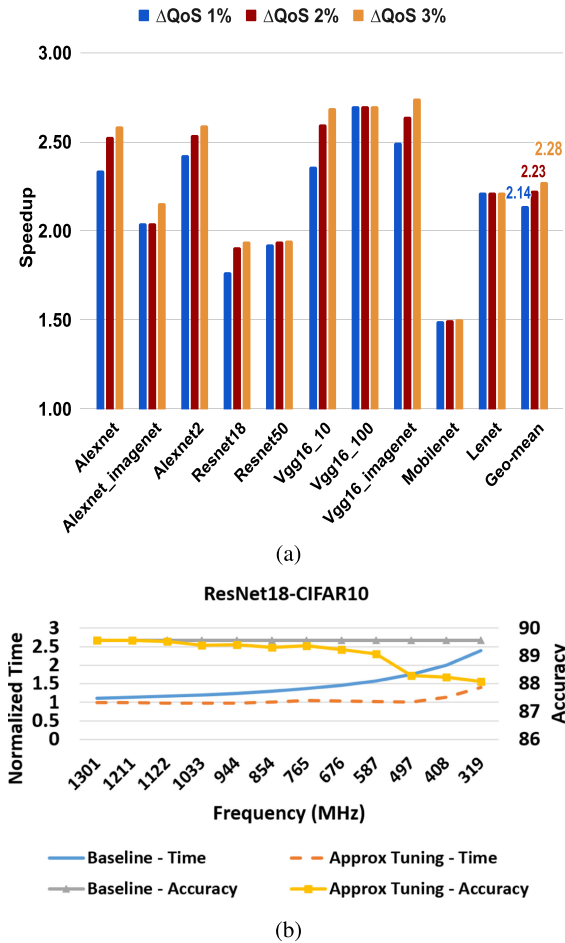


FIGURE 5. (a) Speedups achieved on GPU using approximations for ΔQoS 1%, ΔQoS 2%, and ΔQoS 3%. (b) Runtime approximation tuning maintains stable responsiveness (red line) with little loss of accuracy in most of the range (yellow), when GPU frequency is reduced. Time on the y -axis is relative to that at the highest GPU frequency (1.3 GHz). Without dynamic approximations, the application slows down (blue line). (From ApproxTuner.⁶)

on the FPGA and leaving the rest on the GPU. Our target GPU+FPGA system uses an Intel Xeon W-2275 CPU, an NVIDIA Quadro P1000 GPU, and an Intel Arria 10 GX FPGA. Figure 4(b) shows the performance results averaged over ten runs. The speedup increases by around $3\times$ by moving the reduction kernel to the FPGA, demonstrating the strong benefits of such a partitioning. (The base GPU speedup is lower than in Figure 4(a) because the Quadro P1000 is much slower than the RTX 2080 used there.) As in the previous experiment, the GPU performance of the reduction could be improved by parallelizing it, but the nature of pipeline parallelism in an FPGA is better suited for reductions and can handle

nonassociative reductions, which cannot be parallelized on a GPU.

This experiment shows that HPVM is able to overcome the first challenge described previously. In addition, once our autotuner is extended to support multiple target devices at the same time, the partitioning decisions would be automated and optimized as well, thus solving the second challenge of multidevice partitioning.

Performance Improvements Using ApproxTuner

Figure 5(a) shows our evaluation results⁶ using ApproxTuner on 10 popular CNN benchmarks measured on the Jetson Tegra Tx2 GPU. The graph shows speedups obtained compared to a baseline that does all operations in FP32 on the GPU, with no approximations. We configured ApproxTuner to choose from three different hardware-independent approximations per tensor operation: FP16, perforated convolutions, and sampled convolutions.

Across benchmarks, when allowed merely one percentage point drop in accuracy, ApproxTuner achieves a mean $2.1\times$ speedup compared with the baseline. Relaxing the accuracy threshold to 2 and 3 percentage points (red and orange bars, respectively) provides only a small increase in speedups, to $2.2\times$ and $2.3\times$, showing that most of the benefits are gained simply by allowing any approximation at all.

*AS HETEROGENEOUS SYSTEMS
ADOPT MORE DIVERSE
ACCELERATORS, WE WILL CONTINUE
EXPANDING HPVM WITH MORE
DEVICE BACK ENDS, AND MORE
ADVANCED OPTIMIZATION AND
TUNING TECHNIQUES.*

As expected, each network is amenable to a different set of approximations; there does not exist an approximation that provides the best accuracy-performance tradeoff on all networks. For example, AlexNet is amenable to perforated convolution and more sensitive to sampled convolution, while it is the opposite for the VGG networks, which ApproxTuner discovers through tuning. In addition, ApproxTuner finds that the first few and last convolution layers in a network cause the highest errors due to approximations, and it approximates these layers conservatively.

Figure 5(b) shows that ApproxTuner can counteract system slowdowns induced by low-frequency

KEY TAKEAWAYS

- ▶ HPVM enables hardware-agnostic programming of heterogeneous systems via a hierarchical DFG abstraction of parallelism that supports retargetable compilation to diverse hardware, such as CPUs, GPUs, FPGAs, and domain-specific accelerators.
- ▶ HeteroC++, PyTorch, and other hardware-agnostic front ends simplify programming heterogeneous systems and facilitate off-loading different application components to different devices while preserving source-code and (optionally) object-code portability.
- ▶ Sophisticated HPVM and LLVM optimizations, together with target-specific autotuning, deliver significant performance improvements without manual tuning, which greatly improves source-level portability and maintainability.
- ▶ The ApproxTuner automated approximation tuning framework for tensor operations supports powerful accuracy-aware optimizations and runtime adaptation while preserving hardware-agnostic programming and object-code portability.

modes on the GPU. As frequency lowers from left to right (x -axis), the normalized batch processing time (y -axis) shown by the blue line increases. ApproxTuner uses the shipped Pareto curve to pick configurations that increase speedups in order to counteract these slowdowns while sacrificing small amounts of accuracy. The red dotted line shows batch processing times stabilize when ApproxTuner dynamic tuning is enabled. The yellow line shows that as frequency decreases, the neural network accuracy gradually decreases since higher approximation levels are needed to counter greater slowdowns.

DIRECTIONS FOR FURTHER WORK

As heterogeneous systems adopt more diverse accelerators, we will continue expanding HPVM with more device back ends, and more advanced optimization and tuning techniques. This includes extending our autotuner to automatically partition programs across PEs, while also optimizing them for each target.

We are working on making HPVM even more retargetable for a wide range of emerging tensor

architectures, including CPU ISA extensions (Intel AMX and Power MMA), GPU extensions (NVIDIA's Tensor Cores and AMD's Matrix Cores), and custom ML accelerators (Amazon Inferentia/Trainium, Google TPU, and NVDLA).

We aim to leverage HPVM to greatly simplify DSL design and implementation for high-level applications that benefit from heterogeneous systems.

We are also interested in extending approximation tuning to emerging application domains, particularly edge computing domains such as mobile robotics, AR/VR, and video analytics.

ACKNOWLEDGMENT

This work was supported in part by NSF under Grants CCF 13-02641 and CCF 16-19245, in part by the Semiconductor Research Corporation and DARPA through the Center for Future Architectures Research (C-FAR) and the Applications Driving Architectures (ADA) center, in part by Intel Corp. through DARPA DSSoC Program, and in part by the Amazon AWS Machine Learning Research Awards and Amazon Research Awards Programs.

REFERENCES

1. V. Adve et al., "Virtual instruction set computing for heterogeneous systems," in *Proc. USENIX Workshop Hot Topics Parallelism*, 2012.
2. M. Kotsifakou et al., "HPVM: Heterogeneous parallel virtual machine," in *Proc. ACM SIGPLAN Symp. Princ. Pract. Parallel Program.*, 2018, pp. 68–80.
3. D. Koeplinger et al., "Spatial: A language and compiler for application accelerators," in *Proc. ACM SIGPLAN Conf. Program. Lang. Des. Implementation*, 2018, pp. 296–311.
4. T. Chen et al., "TVM: An automated end-to-end optimizing compiler for deep learning," in *Proc. USENIX Symp. Oper. Syst. Des. Implementation*, 2018, pp. 578–594.
5. H. Sharif et al., "ApproxHPVM: A portable compiler IR for accuracy-aware optimizations," in *Proc. ACM Program. Lang.*, 2019, Art. no. 186.
6. H. Sharif et al., "ApproxTuner: A compiler and runtime system for adaptive approximations," in *Proc. ACM SIGPLAN Symp. Princ. Pract. Parallel Program.*, 2021, pp. 262–277.
7. A. Ejeh et al., "HPVM2FPGA: Enabling true hardware-agnostic FPGA programming," in *Proc. IEEE Int. Conf. Appl.-Specific Syst., Archit., Processors*, 2022.
8. S. Dasgupta et al., "Scalable validation of binary lifters," in *Proc. ACM SIGPLAN Conf. Program. Lang. Des. Implementation*, 2020, pp. 655–671.

9. N. Zmora *et al.*, "Neural network distiller: A python package for DNN compression research," 2019, *arXiv:1910.12232*.
10. L. Nardi *et al.*, "Practical design space exploration," in *Proc. IEEE Int. Symp. Model., Anal., Simul. Comput. Telecommun. Syst.*, 2019, pp. 347–358.
11. P. Srivastava *et al.*, "PROMISE: An end-to-end design of a programmable mixed-signal accelerator for machine-learning algorithms," in *Proc. ACM/IEEE Int. Symp. Comput. Archit.*, 2018, pp. 43–56.
12. P. Mantovani *et al.*, "Agile SoC development with open ESP," in *Proc. IEEE/ACM Int. Conf. Comput. Aided Des.*, 2020, pp. 1–9.
13. S. Che *et al.*, "Rodinia: A benchmark suite for heterogeneous computing," in *Proc. IEEE Int. Symp. Workload Characterization*, 2009, pp. 44–54.

ADEL EJJEH is a Ph.D. student at the University of Illinois at Urbana-Champaign, Urbana, IL, 61801, USA. Contact him at aejjeh@illinois.edu.

AARON COUNCILMAN is a Ph.D. student at the University of Illinois at Urbana-Champaign, Urbana, IL, 61801, USA. Contact them at aaronjc4@illinois.edu.

AKASH KOTHARI is a Ph.D. student at the University of Illinois at Urbana-Champaign, Urbana, IL, 61801, USA. Contact him at akashk4@illinois.edu.

MARIA KOTSIFAKOU is a software engineer at Runtime Verification, Inc., Champaign, IL, 61820, USA. Contact her at maria.kotsifakou@runtimeverification.com.

LEON MEDVINSKY is a Ph.D. student at the University of Illinois at Urbana-Champaign, Urbana, IL, 61801, USA. Contact him at leonkm2@illinois.edu.

ABDUL RAFAE NOOR is a Ph.D. student at the University of Illinois at Urbana-Champaign, Urbana, IL, 61801, USA. Contact him at arnoor2@illinois.edu.

HASHIM SHARIF is a postdoctoral researcher at the University of Illinois at Urbana-Champaign, Urbana, IL, 61801, USA. Contact him at hsharif3@illinois.edu.

YIFAN ZHAO is a Ph.D. student at the University of Illinois at Urbana-Champaign, Urbana, IL, 61801, USA. Contact him at yifanz16@illinois.edu.

SARITA ADVE is a professor at the University of Illinois at Urbana-Champaign, Urbana, IL, 61801, USA. She is a Fellow of the IEEE. Contact her at sadve@illinois.edu.

SASA MISAILOVIC is an associate professor at the University of Illinois at Urbana-Champaign, Urbana, IL, 61801, USA. Contact him at misailo@illinois.edu.

VIKRAM ADVE is a professor at the University of Illinois at Urbana-Champaign, Urbana, IL, 61801, USA. He is a Member of the IEEE. Contact him at vadve@illinois.edu.