# Early-stage Automated Identification Tool for Shared Accelerators

Parnian Mokri
Tufts University
Medford, Massachusetts 02155
Email: parnian.mokri@tufts.edu

Mark Hempstead
Tufts University
Medford, Massachusetts 02155
Email: mark.hempstead@tufts.edu

The use of application-specific accelerators to improve systems' energy-efficiency and performance is becoming more prevalent. To overcome the tight area budget on embedded systems we propose an early detection tool that complements existing High-level Synthesis tools by identifying computationally similar synthesizable kernels that are used to build Shared Accelerators (SAs). SAs are specialized hardware accelerators that execute very different software kernels but share the common hardware functions between them. SAs can provide increased coverage if similarities between the dataflow and control flow of seemingly very different workloads are detected. Existing methods use either dynamic traces or analyze register transfer level (RTL) implementations to find these similarities which requires deep knowledge of RTL and the time-consuming RTL design process.
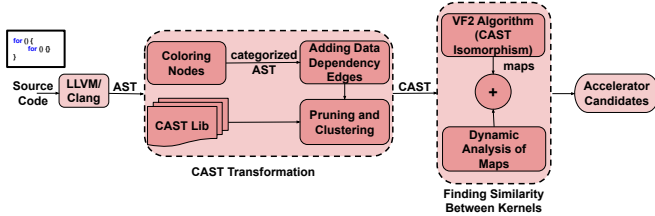


Fig. 1. Block diagram of the ReconfAST methodology.

This work leverages abstract syntax trees (ASTs) generated from LLVM's-clang to discover similar kernels among workloads. ASTs are compact, unlike control and dataflow representations, but contain extra syntax and variable node ordering that complicates workload comparison. As shown in Figure 1 our methodology, ReconfAST, transforms the AST into a new clustered AST (CAST) representation that further removes unneeded nodes and uses a flexible tree-traversal and regular expression matching scheme to detect and group common node patterns. ReconfAST transforms ASTs into a hardware implementable tree by removing whitespace nodes to remove differences that are resulted from coding style. We run a dynamic analysis of these static structural similarities, to further refine SA candidates by making sure these maps represent hot code. Finally, we prune the candidates based on their static data dependency class. This step will remove cases when a variable inside the acceleration candidate depends
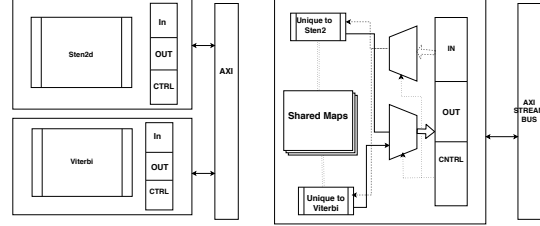


Fig. 2. The two Dedicated Accelerators (on the left) are replaced by one Shared Accelerator (on the right), increasing accelerator's coverage and saving area.



Fig. 3. Maximum Dynamic Coverage (percentage of total execution time) measured of the matching (isomorphic) sub-graphs found between the CASTs of each workload.

on outside variables. In addition, the tool warns the user in cases where many data dependencies are found inside the acceleration candidate since that would limit the ability of HLS tools to use common hardware optimizations to improve performance and energy efficiency. Figure 2 shows a simplified example of a system with accelerators for two MachSuite benchmarks, *Stencil2D* and *Viterbi*. Figure 3 shows that the common source code between stncil2d and vterbi was 94% of stencil2D hot-code. Therefore, a common accelerator can accelerate both workloads.

The presence of data dependencies, the cost of reconfiguration, and the difference between the size of accelerators affect the efficiency of SAs. We have designed over 700 of these accelerators using Vivado_HLS. A good Shared Accelerator, on FPGAs has comparable speedup to dedicated accelerators and reduces the resources required for FPGA implementations: 37% FFs, 16% DSPs, and 10% on LUTs on average.