

Software Synthesis for High Productivity ExaScale Computing

Armando Solar-Lezama (MIT), Rastislav Bodik and James Demel (UC Berkeley)

This project builds on recent results on software synthesis to develop a programming model that integrates validation, synthesis and autotuning. The programming model is based on the concept of programmer guided synthesis; the idea is to leverage synthesis and validation technology while leaving programmers in control of the high-level implementation decisions. The goal is to simplify the process of programming on high performance heterogeneous architectures.

Key components of the programming model

Program Refinement: The programming model supports stepwise development through refinement. A refinement of a program P is program Q that computes the same task as P. Often Q is obtained from P by optimizing a specific aspect of P's computation. For example, program P may assume a simple shared memory model while Q implements a more scalable strategy for partitioning data over distributed memory.

Many high-performance applications are already developed this way; developers start with a simple reference implementation, and optimize different aspects of the program one at a time, checking their work after each step by testing each new version of the program against the previous one. Our programming model provides direct support for this development strategy by allowing programmers to describe these refinement relationships, and leverages this information to help with the more difficult aspects of implementing high-performance code.

Functional Equivalence Checking: One of the ways we support the refinement model is by helping programmers check that each new iteration of the system is indeed a refinement of the previous version. This validation is performed using model checking algorithms implemented on top of SAT/SMT solvers, as well as automated testing techniques that have proven effective for both hardware and software verification. The key research challenge is the application of these techniques to exascale programs with billion-way parallelism. The key to scalability in this setting lies in symmetry reduction, exploiting the fact that most threads in such a program are doing more or less the same thing.

Synthesis with partial programs: When performing a refinement, programmers are expected to have a high-level idea of the overall goal of the refinement—e.g. replace accesses to a global array with local accesses followed by a global exchange step. However, the low-level details of such refinement are often difficult to derive by hand. Synthesis can automatically derive these low-level details when given adequate information about the high-level structure of the solution. The synthesizer can also be forced to focus on solutions that satisfy important resource constraints, such as avoiding bank conflicts or ensuring a balanced partitioning of a data-structure. In this way, programmers can focus on the high-level implementation strategies without having to worry about the low-level details.

Autotuning: When multiple correct completions to a template exist, they will be performance-evaluated with an autotuner through empirical evaluation. In this way, autotuning and synthesis collaborate with each other to enumerate and explore large spaces of correct implementations to arrive at an efficient implementation of the high-level idea provided by the programmer.

The Language Infrastructure

Our language is based on the Scala programming language, making it easy to define high-level programming abstractions for the synthesis constructs. Our infrastructure translates the Scala programs written by the user into the Sketch intermediate language, allowing the system to leverage the Sketch synthesis and analysis infrastructure, and providing a clean representation from which to generate efficient code.

Our first milestone is to demonstrate the approach by showing the use of synthesis to implement efficient kernels in CUDA. We have already completed an initial translation from Scala to the Sketch intermediate language, and are experimenting with new algorithms to allow for efficient analysis of data parallel kernels.