

Luthier: A Dynamic Binary Instrumentation Framework Targeting AMD GPUs

Matin Raayai-Ardakani, Norman Rubin (Advisor), and David Kaeli (Advisor)

Department of Electrical and Computer Engineering, Northeastern University, Boston MA USA

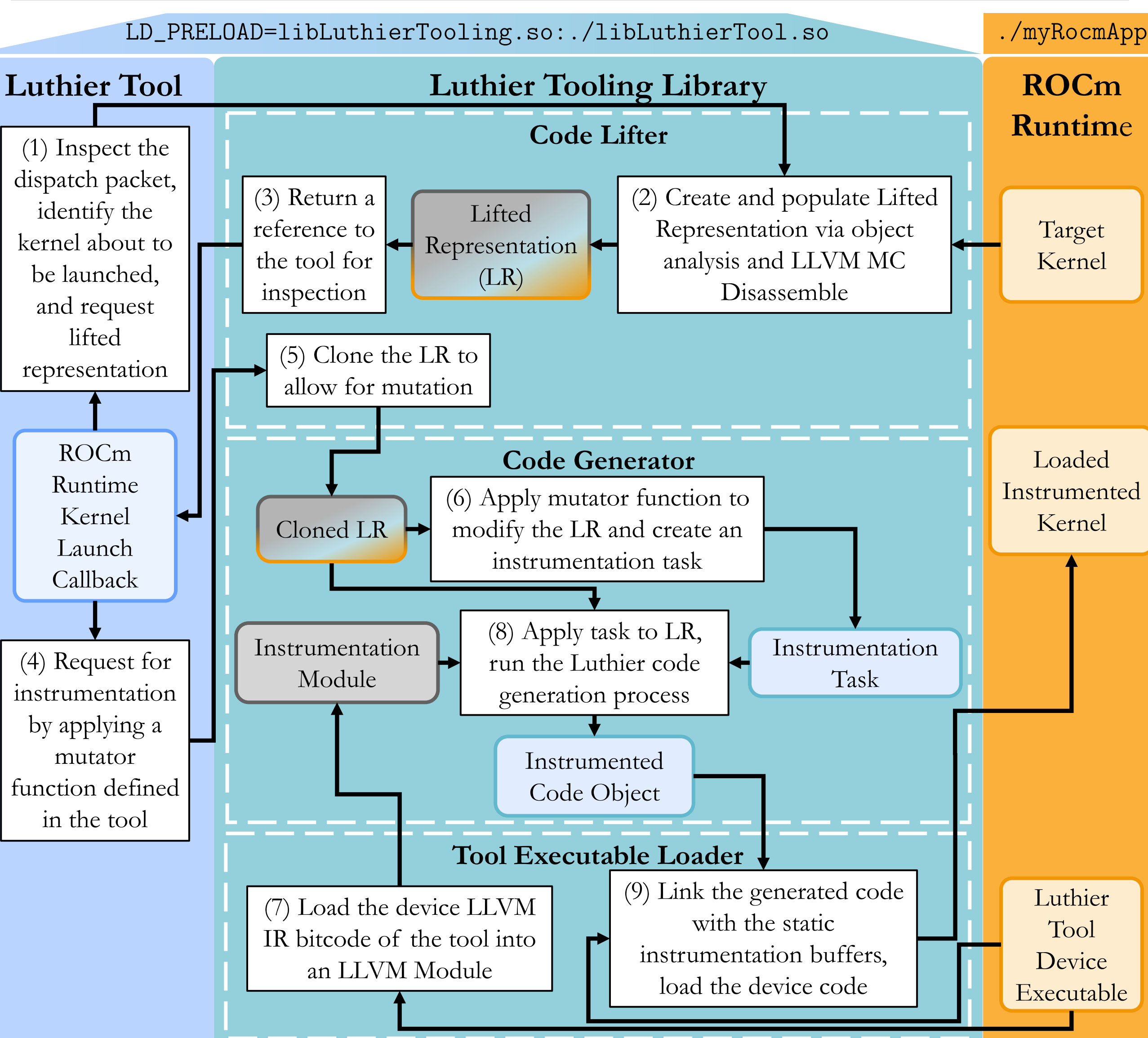


Motivation

- AMD GPUs have experienced growing adoption in HPC for accelerating large-scale scientific computing and AI applications
- Current AMD ROCm profiling tools rely on reading hardware counters and Program Counter (PC) sampling provided by ROCm's ROCprofiler-SDK
- Instrumentation of device code can be used to collect significantly more detailed info regarding regions of interest at the cost of increased execution overhead
- Current efforts only focus on offline instrumentation of device code in ROCm
- Dynamic binary instrumentation offers multiple benefits:
 - Does not require application recompilation from source
 - Works with kernels written in assembly
 - Can be selectively applied to manage instrumentation overhead

We present Luthier, the first open-source dynamic binary instrumentation framework for AMD GPUs on the ROCm stack

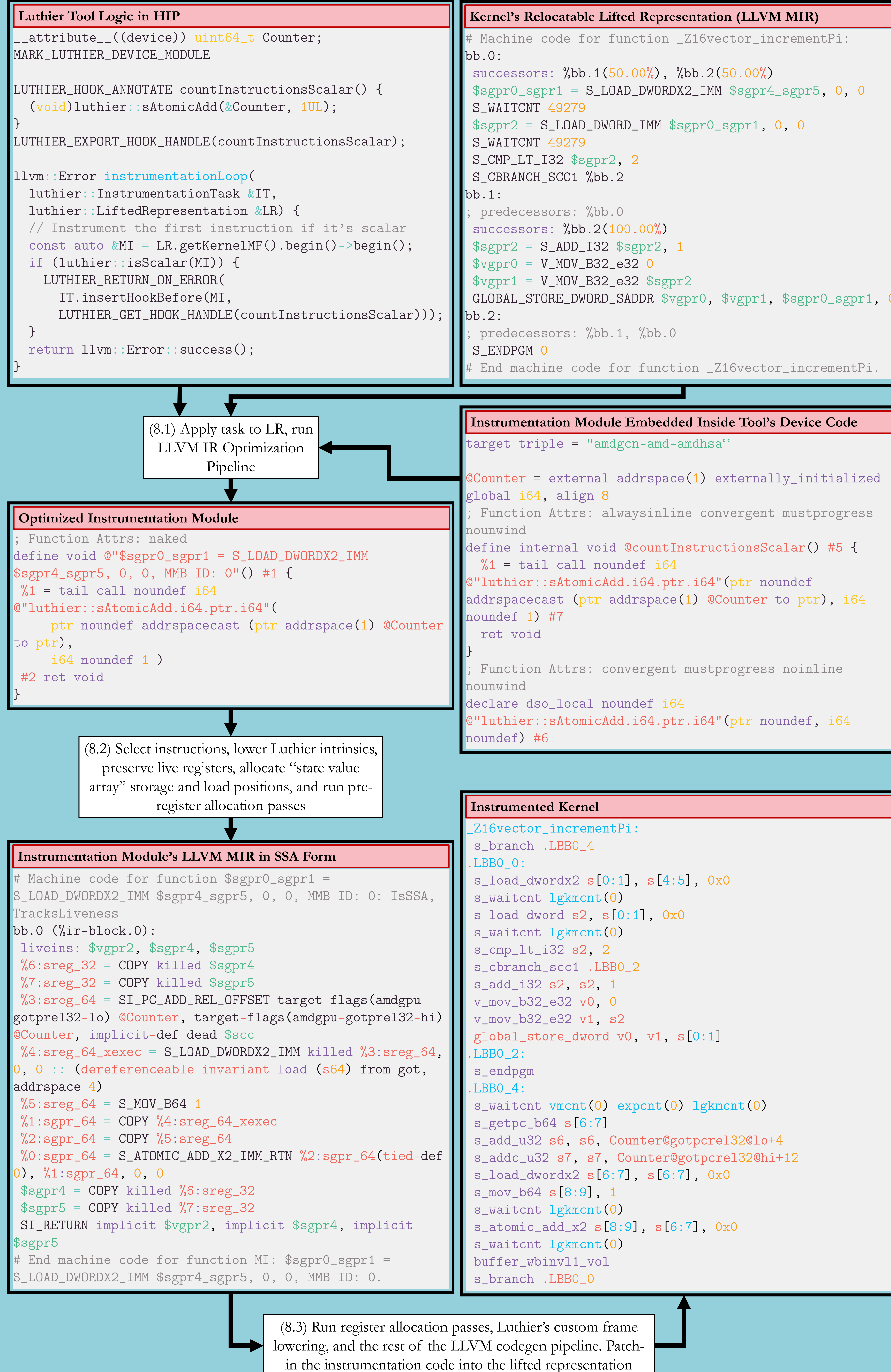
Luthier's Instrumentation Process



Key Design Features

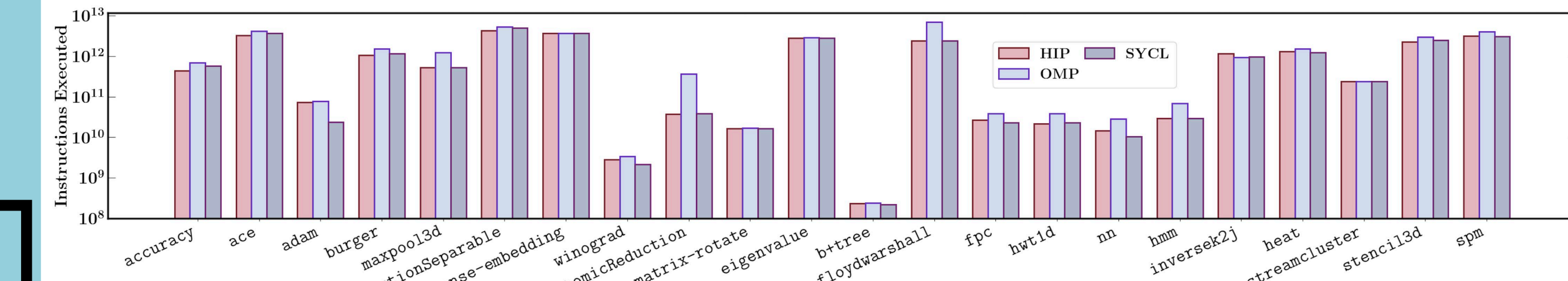
- Instrumentation is performed on the relocatable LLVM MIR representation of the target kernel and loaded as a separate copy
 - In-place instrumentation is not feasible due to variable-length instructions and limited range of direct jumps on AMD GPUs.
- Instrumentation code is Just-In-Time (JIT) compiled to:
 - Inline the instrumentation logic in the target kernel's code
 - Reuse available free registers of the target kernel in the instrumentation code
 - Only generate spill code if no free registers are available
 - Leverage additional static analysis on the target kernel to fully optimize instrumentation code
- Instead of using inline assembly, a custom intrinsic binding and lowering mechanism is used (Luthier::sAtomicAdd in the example code):
 - Allows for more aggressive register usage optimization
 - Can define custom lowering logic for existing intrinsics according to the instrumentation process
- Instrumentation code maintains its own stack and kernel argument values (\$vgpr2 in the example code):
 - Instrumentation stack is recovered/spilled in its custom emitted frame
 - Guarantees correct instrumentation of kernels written in assembly
 - Can be used to fetch thread identifier values and issue calls to printf

Luthier's Code Generation: A Closer Look

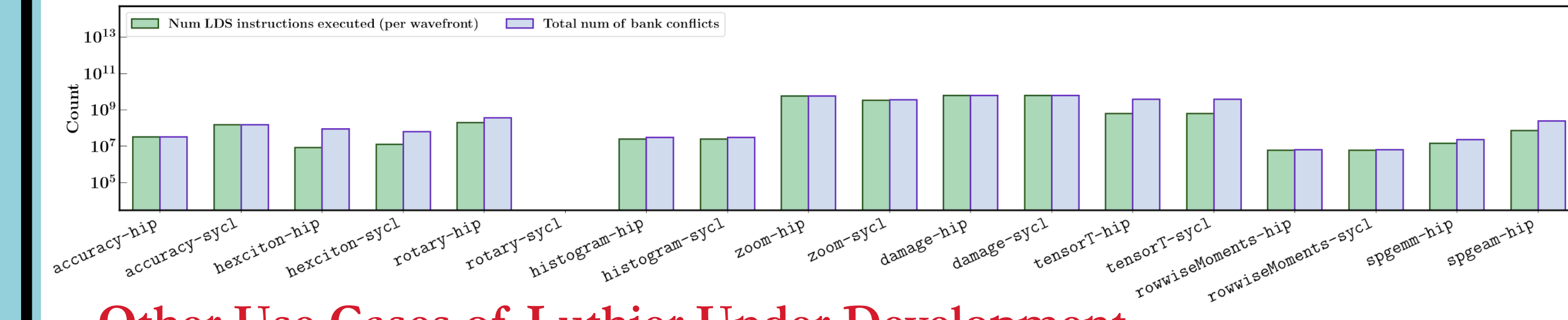


Example Luthier Use Cases

- Luthier can be used to count the total number of instructions executed on each GPU throughout the program as well as a breakdown of the executed instruction opcodes.
- Luthier's instruction count tool was applied to entries in HeC Bench^[1] across HIP, OpenMP, and SYCL running on an AMD Instinct MI100.



- Luthier can also be used to detect bank conflicts inside kernels.
- It can point out the exact instructions the bank conflict occurred, which can then be correlated back into the source code if debug information is present.
- The Luthier bank conflict tool was applied to entries in HeC Bench^[1] across HIP and SYCL and were run on an AMD Instinct MI100.

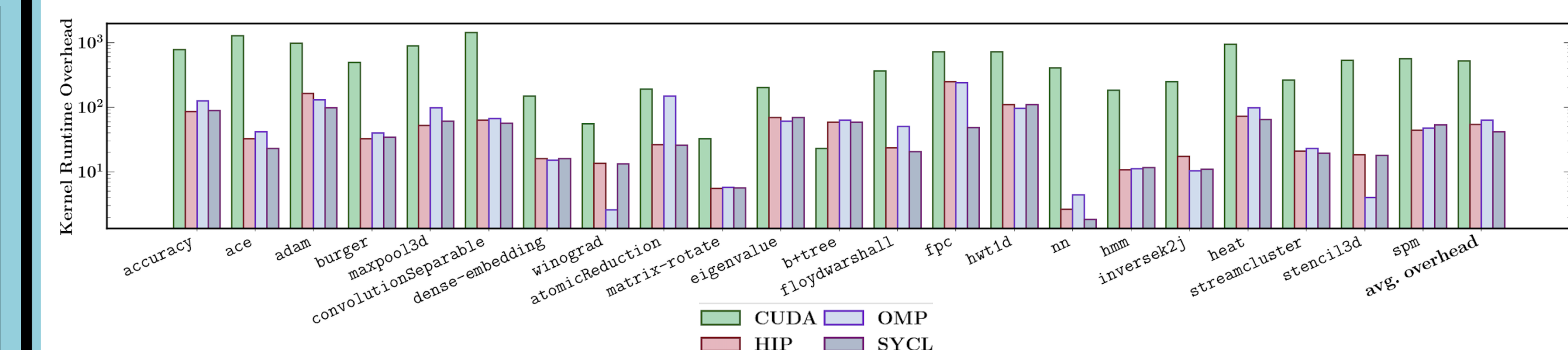


Other Use Cases of Luthier Under Development

- Causal Profiling^[2]: Identify most impactful optimization opportunities in exact regions of the source code
- Call Graph Construction: Identify functions called and their callers
- Reliability: E.g., inject software faults into the device code

Instrumented Kernel Runtime Overhead

- We compared the overhead of running Luthier's instruction count tool against NVBit's instruction counter tool
- Luthier and AMD benchmarks were run on an AMD MI100, while NVBit and the NVIDIA equivalent benchmarks were run on an NVIDIA V100 attached to the same machine
- On average, Luthier has a 50X runtime overhead, 10 times lower than NVBit



Current Focus of Development

- Instrumentation of very large kernels with stripped device function symbols
- Instrumentation of assembly kernels that do not conform to LLVM MIR
- Support for more AMD GPU targets

For more information, please refer to our paper: M. Raayai-Ardakani et al., "Luthier: A Dynamic Binary Instrumentation Framework Targeting AMD GPUs," 2025 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS), Ghent, Belgium, 2025, pp. 137-149, doi: 10.1109/ISPASS64960.2025.00022.

Acknowledgements

- Timour Paltashev, Benjamin Welton, Konstantin Zhuravlyov, Mathew Arsenaull, Jacob Lambert, Shilei Tian, Tony Tye, Brian Sumner (AMD)
- Rene Van Oostrum (AMD research)
- Omar Shair, Maya De Los Santos, Yuhui Bao, Zlatan Feric, Qishi Wang (Northeastern University)
- Luthier was supported in part by AMD

References

- Z. Jin and J. S. Vetter, "A Benchmark Suite for Improving Performance Portability of the SYCL Programming Model," 2023 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS), Raleigh, NC, USA, 2023, pp. 325-327, doi: 10.1109/ISPASS57527.2023.00041. (<https://ieeexplore.ieee.org/document/10158214>)
- Charlie Curtisinger and Emery D Berger. 2015. Coz: Finding code that counts with causal profiling. In Proceedings of the 25th Symposium on Operating Systems Principles. 184-197. (<https://dl.acm.org/doi/abs/10.1145/2815400.2815409>)