



Extending High-Level Synthesis for Task-Parallel Programs

Yuze Chi*, Licheng Guo*, Jason Lau*, Young-kyu Choi*+, Jie Wang*, Jason Cong*

*University of California, Los Angeles, †Inha University

{chiyuze,cong}@cs.ucla.edu



Widespread Adoption of High-Level Synthesis



SDC scheduling ^[DAC'06, TCAD'11], frequency optimizations ^[DAC'20, FPGA'21]

[DAC'06]: Cong and Zhang. An Efficient and Versatile Scheduling Algorithm Based On SDC Formulation. In DAC, 2006.
 [TCAD'11]: Cong et al. High-Level Synthesis for FPGAs: From Prototyping to Deployment. TCAD, 2011.
 [DAC'20]: Guo et al. Analysis and Optimization of the Implicit Broadcasts in FPGA HLS to Improve Maximum Frequency. In DAC, 2020.
 [FPGA'21]: Guo et al. AutoBridge: Coupling Coarse-Grained Floorplanning and Pipelining for High-Frequency HLS Design on Multi-Die FPGAs. In FPGA, 2021.

Limitations of Current High-Level

•••

void Gemm(float a[N][N], float b[N][N], float c[N][N]) { for (int i = 0; i < N; ++i) { for (int k = 0; k < N; ++k) { #pragma HLS pipeline II = 1 for (int j = 0; j < N; ++j) { #pragma HLS unroll c[i][k] += a[i][j] * b[j][k];

void Gemm(float a[N][N], float b[N][N], float c[N][N]) { hls::stream<float> a_00, a_01, ... hls::stream<float> b_00, b_01, ... hls::stream<float> c_00, c_01, ... hls::stream<float> fifo_00_10, ... #pragma HLS dataflow Load(a, a_00, a_01, a_10, a_11); Load(b, b_00, b_01, b_10, b_11); Compute(a_00, b_00, c_00, ...); Compute(a_01, b_01, c_01, ...); Compute(a_10, b_10, c_10, ...); Compute(a_11, b_11, c_11, ...); Store(c, c_00, c_01, c_10, c_11); }

Why Go Task-Parallel?

- Data parallelism may be limited
 - PageRank / graph convolutional networks / on-chip network switching
- Task parallelism can better accommodate variable DRAM latency
- Task parallelism shows better scalability
 - Data-parallel applications can be implemented as task-parallel programs
 - Systolic arrays [ICCAD'18, PolySA]
 - Stencil [ICCAD'18, SODA]

Benchmark	Application	#Tasks	#Task Instances	#Channels
cannon	Cannon's matrix multiplication algorithm	5	91	344
cnn	Systolic array generated by PolySA [ICCAD'18]	14	209	366
gaussian	Stencil generated by SODA [ICCAD'18]	15	564	1602
gcn	Graph convolutional network [ICLR'17]	5	12	25
gemm	Systolic array generated by PolySA [ICCAD'18]	14	207	364
network	8×8 Omega switching network	3	14	32
page_rank	Edge-centric PageRank accelerator	4	18	89

[ICLR'17]: Kipf and Welling. Semi-Supervised Classification with Graph Convolutional Networks. In ICLR, 2017. [ICCAD'18]: Cong and Wang. PolySA: Polyhedral-Based Systolic Array Auto-Compilation. In ICCAD, 2018. [ICCAD'18]: Chi et al. SODA: Stencil with Optimized Dataflow Architecture. In ICCAD, 2018.

A Motivating Example



- ◆ 4 PEs interconnected w/ a ring
 - PEs send packets to ring nodes & receive packets from ring nodes
 - Ring nodes forward packets
 conditionally based on the
 destination PE specified in the
 packet header

Extending Kernel Programming Interfaces





Equivalent kernel code w/ TAPA

For kernel code, TAPA is much shorter to write with peeking and transaction API support.

Simplifying Host-Kernel Communication Interfaces



Host code for Vivado HLS code as an OpenCL kernel

For host code, TAPA is also much shorter to write by allowing programmers call the kernel using a single function call.

Software Simulation







TAPA's coroutine-based simulator



[1]: <u>https://eli.thegreenplace.net/2018/measuring-context-switching-and-memory-overheads-for-linux-threads/</u>
 [2]: <u>https://www.boost.org/doc/libs/1_73_0/libs/coroutine2/doc/html/coroutine2/performance.html</u>

RTL Code Generation



All PEs are instances of the same task (C++ function). All ring nodes are also instances of the same task.



TAPA's hierarchical approach Run HLS Run HLS for for PE RingNode Û ſ Instantiate RTL modules Create FSM that controls them No redundant work!

...

Overall TAPA Workflow



Kernel LoC Reduction



Host LoC Reduction



Simulation Time Reduction



Code Generation Time Reduction



Resource Usage Comparison



Beyond Productivity Enhancement

- Extended AXI4 memory-mapped interface
 - Asynchronous access w/ separate AR/R/AW/W/B channels
 - Shared access among multiple task instances
- Detached tasks/smarter FIFO templates/<u>AutoBridge</u> integration/...
- For more details
 - TAPA is open-source: <u>https://github.com/UCLA-VAST/tapa</u>





Q&A

This work is partially supported by a Google Faculty Award, the NSF RTML program (CCF1937599), NIH Brain Initiative (U01MH117079), the Xilinx Adaptive Compute Clusters (XACC) program, and CRISP, one of six JUMP centers.

