Paper accepted to HPCA 2022 No final paper released yet

ScaleHLS: Scalable High-Level Synthesis through MLIR

Hanchen Ye, Cong Hao, Jianyi Cheng, Hyunmin Jeong, Jack Huang, Stephen Neuendorffer, Deming Chen Presenters: [Team 2] Insu Jang, Jae-Won Chung



COLLEGE OF ENGINEERING COMPUTER SCIENCE & ENGINEERING UNIVERSITY OF MICHIGAN

Agenda

- Background
 - High-level synthesis (HLS)
 - Multi-level IR (MLIR)
- ScaleHLS
 - Goal and architecture
 - Dataflow optimization
 - Loop optimization
- Evaluation
- Conclusion



Background: High-level Synthesis

- Programming FPGA is hard, requiring circuit knowledge
- HLS: A way of implementing HW design with high level language



Background: Multi-level IR (MLIR)

• ScaleHLS built on top of MLIR

Before MLIR

• Each language implements its own high-level IR for its own optimization



UNIVERSITY OF MICHIGAN

MLIR is

 easily customizable to add MLIR dialect or use common MLIR dialects



Goal of ScaleHLS

- End-to-end automation
- Capture different levels of optimization
 - Graph-level
 - Loop-level
 - Directive-level



ScaleHLS Architecture



Optimizations

Level	Pass
Graph	Legalize dataflow
	Split function
Loop	Loop perfectization
	Loop order optimization
	Variable bound removal
	Loop tiling
	Loop unrolling

Level	Pass
Directive	Loop pipelining
	Function pipelining
	Array partitioning

Graph-level Dataflow Optimization



A bypassing dataflow edge precludes finer-grained pipelining

Graph-level Dataflow Optimization



Loop Flattening and Perfectization

• Loop flattening: flatten multi-dimension nested loops into one loop



• Problem: Imperfect loops cannot be flattened



Loop Unrolling and Pipelining

- Add a pipeline directives in d
- k-loop is fully unrolled
- Pipelining is applied to j-loop
- Transformation is applied by Vivac
 * i-j nested loops can be automatically
 flattened since this is perfect loop
- Requirements: all sub-loops unrolled / sub-functions
 - Legalize the loop by unrolling sub-loop to be fully pipelined

```
for (int i=0; i<16; i++) {
   for (int j=0; j<i; j++) {
     for (int k=0; k<8; k++) {
        a[i][j] += b[i][k] * b[j][k];
}}</pre>
```

for (int i=0; i<16; i++) {
 for (int j=0; j<i; j++) {
 #pragma HLS pipeline II=3
 a[i][j] += b[i][0] * b[j][0];
 ...
 a[i][j] += b[i][7] * b[j][7];
}</pre>

Array Partitioning

• Loop unroll may not be effective due to **limited FPGA memory ports**



• **Automatically** detect memory access pattern and apply suitable array partition directives: #pragma HLS array_partition variable=a ...

Evaluation: ResNet-18

- Translate PyTorch model definitions to Graph-level IR
- Results
 - **3.9x more resource efficient** versus TVM-VTA¹
 - **3825.0x speedup** versus itself without any optimizations
 - **automatic optimization** versus nothing (This is the first work.)

[1] "A Hardware-Software Blueprint for Flexible Deep Learning Specialization," Moreau et al., arXiv:1807.04188

Evaluation: ResNet-18



More graph optimizations (Gn: finer dataflow granularity)



Speedup over Baseline

More loop optimizations (Ln: larger loop unroll factor)

Conclusion and Commentary

- ScaleHLS: Optimizations for high performance HLS
 - Well-architected system with extensible design
 - Hardware-aware optimization passes and modelling
 - A lot of engineering effort invested
- Some limitations, though.
 - Weak evaluation baseline
 - Not many new things
 - Not sure why this is scalable

Paper accepted to HPCA 2022 No final paper released yet

ScaleHLS: Scalable High-Level Synthesis through MLIR

Questions? Comments?

Presenters: [Team 2] Insu Jang, Jae-Won Chung



COLLEGE OF ENGINEERING COMPUTER SCIENCE & ENGINEERING UNIVERSITY OF MICHIGAN