Estimating Circuit Quality Directly from HLS Source Code

Lana Josipović

June 2025





Hardware acceleration for high parallelism and energy efficiency



How to perform hardware design?

... circuit design is often considered a "black art", restricted to only those with years of training in electrical engineering... [cacm.acm.org/magazines/2023/1/] ... chips take years to design, resulting in the need to speculate about how to optimize the next generation of chips... [ai.googleblog.com/2020/04]

High-Level Synthesis: From Programs to Circuits



High-Level Synthesis: From Programs to Circuits



SW

HW

How to make hardware design accessible to non-experts?

High-level synthesis (HLS): from software descriptions to hardware design

```
1 void(int* mem){
2 mem[512] = 0;
3 for(int i=0; i<512; i++)
4 mem[512] += mem[i];</pre>
```

5 }

Unoptimized software program, execution time = 27,236 clock cycles

\mathbf{v}

- 1 // Width of MPort = 16 * sizeof(int)
- 2 #define ChunkSize (sizeof(MPort)/sizeof(int))
- 3 #define LoopCount (512/ChunkSize)
- 4 // Maximize data width from memory
- 5 **void**(MPort * mem) {
- 6 // Use a local buffer and burst access
- 7 MPort buff[LoopCount];
- 8 memcpy(buff, mem, LoopCount);
- 9 // Use a local variable for accumulation
- 10 int sum=0;
- 11 for(int i=1; i<LoopCount; i++) {</pre>
- 12 // Use additional directives where useful
- 13 // e.g. pipeline and unroll for parallel exec.
 - **#**pragma PIPELINE
 - **for(int** j=0; j<ChunkSize; j++){
 - **#**pragma UNROLL

```
sum+=(int) (buff[i]>>j*sizeof(int)*8);}
```

```
18 mem[512]=sum;
```

```
19 }
```

14

15

16

17

Optimized software program, execution time = 302 clock cycles

George et al. FPL 2014.

SW

How to make hardware design accessible to non-experts?

How to automatically extract parallelism from software code?

Typical software features prevent hardware parallelism

Sparse-matrix dense-vector multiplication (SpMV)

Functional verification of circuits using hardware simulation → inefficient, limited, non-exhaustive





SW

How to make hardware design accessible to non-experts?

How to automatically extract parallelism from software code?

How to avoid hardware-level functional verification?

HW

SW

How to make hardware design accessible to non-experts?

How to automatically extract parallelism from software code?

How to avoid hardware-level functional verification?

How to understand and account for hardware implementation details?

FPGA technology mapping, placement, and routing → impact on circuit performance and power



Langhammer et al. ARITH 2015.

HW

SW

How to make hardware design accessible to non-experts?

How to automatically extract parallelism from software code?

How to avoid hardware-level functional verification?

How to understand and account for hardware implementation details?

How to perform efficient HLS design space exploration?

Design Space Exploration for HLS



(e.g., this kernel \rightarrow hundreds of points)

How to identify Pareto-optimal code configurations?

```
void(int* mem) {
    mem[512] = 0;
    for(int i=0; i<512; i++)
        mem[512] += mem[i];
    }
</pre>
```

execution time = 27,236 clock cycles

\mathbf{P}

// Width of MPort = 16 * sizeof(int) #define ChunkSize (sizeof(MPort)/sizeof(int)) #define LoopCount (512/ChunkSize) // Maximize data width from memory void(MPort* mem) { // Use a local buffer and burst access 6 MPort buff[LoopCount]; 7 memcpy(buff, mem, LoopCount); 8 // Use a local variable for accumulation 9 int sum=0; 10 for(int i=1; i<LoopCount; i++) {</pre> 11 // Use additional directives where useful 12 // e.g. pipeline and unroll for parallel exec. 13 #pragma PIPELINE 14 for(int j=0; j<ChunkSize; j++) {</pre> 15 #pragma UNROLL 16 sum+=(int) (buff[i]>>j*sizeof(int)*8);}} 17 mem[512] = sum;18 19

Design Space Exploration for HLS



Design Space Exploration for HLS





How to quickly and accurately estimate the quality of results (QoR)?



GNNs fully encode the program graph's topological information

Sohrabizadeh, Bai, Sun, and Cong. Robust GNN-Based Representation Learning for HLS. ICCAD 2023. Gao, Zhao, Lin, and Guo. Hierarchical Source-to-Post-Route QoR Prediction in High-Level Synthesis with GNNs. DATE 2024.



Sohrabizadeh, Bai, Sun, and Cong. Robust GNN-Based Representation Learning for HLS. ICCAD 2023. Gao, Zhao, Lin, and Guo. Hierarchical Source-to-Post-Route QoR Prediction in High-Level Synthesis with GNNs. DATE 2024.



Balor, our HLS QoR estimator, increases estimation accuracy and reduces computational cost

A new HLS-tailored graph compiler

• HLS-specific, information-dense graph

Revisited estimation stack

- Local GNNs localize cross-graph interactions
- Hierarchical GNNs increase information density



Balor, our HLS QoR estimator, increases estimation accuracy and reduces computational cost

DB4HLS LUT mean percentage error: 4% Latency mean percentage error: 6%

ML Contest for Chip Design with HLS 2024: 1st place in QoR estimation

Graph Representation

- Graph representation directly obtained from compiler intermediate representation (IR)
 - ProGraML, used by existing estimation approaches, builds graph directly from the LLVM IR







Cummins et. al. ProGraML: A Graph-based Program Representation for Data Flow Analysis and Compiler Optimizations. ICML 2021. Sohrabizadeh, Bai, Sun, and Cong. Robust GNN-based representation learning for HLS. ICCAD 2023. Gao, Zhao, Lin, and Guo. Hierarchical Source-to-Post-Route QoR Prediction in High-Level Synthesis with GNNs. DATE 2024.

- Previously: long-range program graph interactions forced the use of large receptive fields
 - **Receptive field:** a window into the graph



- Previously: long-range program graph interactions forced the use of large receptive fields
 - **Receptive field:** a window into the graph





- Previously: long-range program graph interactions forced the use of large receptive fields
 - Receptive field: a window into the graph





Receptive field radius = 1



- Previously: long-range program graph interactions forced the use of large receptive fields
 - Receptive field: a window into the graph





Receptive field radius = 1



Receptive field radius = 2



- Previously: long-range program graph interactions forced the use of large receptive fields
 - Receptive field: a window into the graph
- Problem 1: poor GNN scaling
 - Large receptive fields cause redundant processing
 - Redundant processing reduces estimation accuracy





Receptive field radius = 1



Receptive field radius = 2



- Previously: long-range program graph interactions forced the use of large receptive fields
 - Receptive field: a window into the graph
- Problem 1: poor GNN scaling
 - Large receptive fields cause redundant processing
 - Redundant processing reduces estimation accuracy
- Problem 2: the **receptive field is not even**
 - Long-range interactions unevenly accounted for





Arrives early: high impact



Arrives late: low impact

Receptive Field Size

"Wide" GNN
(Large Receptive Field)Local GNN
(Small Receptive Field)Poor GNN scaling
causes strongly redundant computationLow redundancy
enables increased accuracyLong-range interactions:
Unevenly accounted for during estimationLong-range interactions:
Not accounted for during estimation

due to late arrival

Information Density:

No problems

Information Density: Small receptive field requires "information-dense" graphs

How to devise information-dense graphs suitable for local GNNs?



Balor, our HLS QoR estimator, increases estimation accuracy and reduces computational cost

A new HLS-tailored graph compiler

• HLS-specific, information-dense graph

Revisited estimation stack

- Local GNNs localize cross-graph interactions
- Hierarchical GNNs increase information density

1. Create node for each code statement

1. Create node for each code statement

2. Break nodes into indivisible hardware actions

- 1. Create node for each code statement
- 2. Break nodes into indivisible hardware actions
- 3. Label each node with data type

- 1. Create node for each code statement
- 2. Break nodes into indivisible hardware actions
- 3. Label each node with data type
- 4. Add variable declaration nodes & connect them to reads/writes

- 1. Create node for each code statement
- 2. Break nodes into indivisible hardware actions
- 3. Label each node with data type
- 4. Add variable declaration nodes & connect them to reads/writes
- 5. Add dataflow edges

- Architecture: GNN DSE
 - The core of state-of-the-art architectures
- Dataset: DB4HLS
 - 36,296 data points, synthesized with Vitis HLS
 - 29 Machsuite kernels
- Estimate:
 - Latency (clock cycles)
 - Clock period (ns)
 - Resource usage (LUTs/FFs/BRAMs/DSPs)

Sohrabizadeh, Bai, Sun, and Cong. Automated Accelerator Optimization Aided by Graph Neural Networks. DAC 2022. Ferretti, Kwon, Ansaloni, Di Guglielmo, Carloni, and Pozzi. DB4HLS: A Database of High-Level Synthesis. ESL 2021

- Architecture: GNN DSE
 - The core of state-of-the-art architectures
- Dataset: DB4HLS
 - 36,296 data points, synthesized with Vitis HLS
 - 29 Machsuite kernels
- Estimate:
 - Latency (clock cycles)
 - Clock period (ns)
 - Resource usage (LUTs/FFs/BRAMs/DSPs)
- Percentage error:

 $\frac{\rm abs(Real \ Value - Inferred \ Value)}{\rm Real \ Value} * 100$

Example scatter plot, 5,673 test points

Sohrabizadeh, Bai, Sun, and Cong. Automated Accelerator Optimization Aided by Graph Neural Networks. DAC 2022. Ferretti, Kwon, Ansaloni, Di Guglielmo, Carloni, and Pozzi. DB4HLS: A Database of High-Level Synthesis. ESL 2021

Information-dense, HLS-tailored graphs reduce both error and cost

Sohrabizadeh, Bai, Sun, and Cong. Automated Accelerator Optimization Aided by Graph Neural Networks. DAC 2022.

Balor, our QoR estimator, increases estimation accuracy and reduces computational cost

A new HLS-tailored graph compiler

• HLS-specific, information-dense graph

Revisited estimation stack

• Local GNNs localize cross-graph interactions

Hierarchical GNNs increase information density

Receptive Field Size

"Wide" GNN
(Large Receptive Field)Local GNN
(Small Receptive Field)Poor GNN scaling
causes strongly redundant computationLow redundancy
enables increased accuracyLong-range interactions:
Unevenly accounted for during estimation
due to late arrivalLong-range interactions:
Not accounted for during estimation

Information Density:

No problems

Information Density: Small receptive field requires **"information-dense" graphs**

Are our custom graphs suitable for local GNNs?

Cross-Graph Interactions: The Pragma Challenge

Sohrabizadeh, Bai, Sun, and Cong. Automated Accelerator Optimization Aided by Graph Neural Networks. DAC 2022. Sohrabizadeh, Bai, Sun, and Cong. Robust GNN-Based Representation Learning for HLS. ICCAD 2023.

Cross-Graph Interactions: The Pragma Challenge

Sohrabizadeh, Bai, Sun, and Cong. Automated Accelerator Optimization Aided by Graph Neural Networks. DAC 2022. Sohrabizadeh, Bai, Sun, and Cong. Robust GNN-Based Representation Learning for HLS. ICCAD 2023.

Move to a Local Architecture

GNN-DSE (DAC '22)

Sohrabizadeh, Bai, Sun, and Cong. Automated Accelerator Optimization Aided by Graph Neural Networks. DAC 2022.

Local-focused approaches reduce redundant processing

Sohrabizadeh, Bai, Sun, and Cong. Automated Accelerator Optimization Aided by Graph Neural Networks. DAC 2022.

Balor, our QoR estimator, increases estimation accuracy and reduces computational cost

A new HLS-tailored graph compiler

• HLS-specific, information-dense graph

Revisited estimation stack

- Local GNNs localize cross-graph interactions
- Hierarchical GNNs increase information density

What are Hierarchical GNNs?

Flat GNNs use 1 network and 1 reduction:

Hierarchical GNNs use multiple networks and multiple reductions to increase information density:

How to cluster the graph to reduce its size?

How to Cluster?

How to Cluster?

Hierarchical Architecture

Local-Focused Architecture:

Sohrabizadeh, Bai, Sun, and Cong. Automated Accelerator Optimization Aided by Graph Neural Networks. DAC 2022.

1st Place at ML Contest for Chip Design with HLS

ML Contest for Chip Design with HLS

Machine Learning Contest for Chip Design with High-Level Synthesis

$$\mathrm{Minimize}
ightarrow rac{\sum_t \mathrm{RMSE}(t)}{\mathrm{F1} + 10^{-6}}$$

Rules Overview Data Discussion Leaderboard

Overview

1

This competition is hosted by UCLA Vastlab, UCLA DM lab, and AMD.

Please sign-up on our official sign-up sheet to be eligible to receive prizes

High-level synthesis (HLS) aims to raise the abstraction level in hardware design, enabling the design of domain-specific accelerators (DSAs) like field-programmable gate arrays (FPGAs) using C/C++ instead of hardware description languages (HDLs). The figure on the right shows an example of HLS designs, where the C++ code defines the functionality and compiler directives in the form of pragmas play a crucial role in modifying the microarchitecture within the HLS framework which determine the performance (e.g., latency and resources utilization rate) of the hardware.

Competition Host

Prizes & Awards Kudos Does not award Points or Medals

Participation

170 Entrants 51 Participants 27 Teams 608 Submissions

1.	1.26 → Balor
2.	1.72 (+37%)
3.	1.99 (+58%)

Overview Discussion Data Leaderboard Rules Δ # Team Members

Emmet Murphy

Score

1.26033

Zongyue Qin

Where Do We Go from Here?

- Balor: an HLS source code evaluator
 - Information-dense, HLS-tailored graphs
 - Local and hierarchical GNNs
 - Increased estimation accuracy and reduced computational cost

GNNs for accurate quality of results estimation directly from the source code

How to **generalize** to different HLS tools, FPGAs, unseen kernels,...?

Where Do We Go from Here?

- Balor: an HLS source code evaluator
 - Information-dense, HLS-tailored graphs
 - Local and hierarchical GNNs
 - Increased estimation accuracy and reduced computational cost

How to **provide insights** that guide design space exploration?

GNNs for accurate quality of results estimation directly from the source code

How to **generalize** to different HLS tools, FPGAs, unseen kernels,...?

SW

How to make hardware design accessible to non-experts?

How to automatically extract parallelism from software code?

How to avoid hardware-level functional verification?

How to understand and account for hardware implementation details?

HW

GNNs for accurate quality of results estimation directly from the source code

Typical software features prevent hardware parallelism

for (i=0; i<N; i++) { Irregular memory access patterns
 hist[x[i]] = hist[x[i]] + weight[i];</pre>

Features of OoO superscalar processors (e.g., speculation) → up to 14.9X faster than Vivado HLS circuits

Dynamatic: an open-source HLS compiler that generates dynamically scheduled circuits from C/C++

SW

How to make hardware design accessible to non-experts?

How to automatically extract parallelism from software code?

How to avoid hardware-level functional verification?

How to understand and account for hardware implementation details?

HW

Josipović, Ghosal, and Ienne. Dynamically Scheduled High-Level Synthesis. FPGA 2018 Best Paper Award Nominee Josipović, Sheikhha, Guerrieri, Ienne, and Cortadella. Buffer Placement and Sizing for High-Performance Dataflow Circuits. FPGA 2020 Best paper award

Functional verification of circuits using hardware simulation → inefficient, limited, non-exhaustive

Formal verification for HLS: prove HLS correctness and improve HLS-generated circuits

HW

hardware implementation details?

SW

Xu, Murphy, Cortadella, and Josipović. Eliminating excessive dynamism of dataflow circuits using model checking. FPGA 2023. Elakhras, Guerrieri, Erhart, Ienne, and Josipović. ElasticMiter: Formally Verified Dataflow Circuit Rewrites. ASPLOS 2025

Standard pipelining (register placement) is **unaware of circuit transformations** during logic synthesis and technology mapping

SW

How to make hardware design accessible to non-experts?

How to automatically extract parallelism from software code?

How to avoid hardware-level functional verification?

How to understand and account for hardware implementation details?

Total delay pre-synthesis = 3 ns

HW

Rizzi, Guerrieri, and Josipović. An Iterative Method for Mapping-Aware Frequency Regulation in Dataflow Circuits. DAC 2023 Wang, Rizzi, and Josipović. MapBuf: Simultaneous Technology Mapping and Buffer Insertion for HLS Performance Optimization. ICCAD 2023 Best Paper Award Nominee

Standard pipelining (register placement) is **unaware of circuit** transformations during logic synthesis and technology mapping

valid. LUT delav = 0.5 ns FORK valid_{out0} ready_{in0} **Total delay** NIOD post-place-and-route **Ins** Wire delay = 0 ns = 1 ns readv valid LUT delay = 0.5 ns FORK valid_{out0} ready_{in0} valid_{out1} ready_{in1}

Implementation-aware compiler optimizations for fast and small circuits

55

Rizzi, Guerrieri, and Josipović. An Iterative Method for Mapping-Aware Frequency Regulation in Dataflow Circuits. DAC 2023

Wang, Rizzi, and Josipović. MapBuf: Simultaneous Technology Mapping and Buffer Insertion for HLS Performance Optimization. ICCAD 2023 Best Paper Award Nominee

SW

HW

How to make hardware design accessible to non-experts?

How to automatically extract parallelism from software code?

How to avoid hardware-level functional verification?

How to understand and account for hardware implementation details?

do x++

FPGA implementation (synthesis, placement, and routing) is orders of magnitude slower than software compilation

SW

How to make hardware design accessible to non-experts?

How to automatically extract parallelism from software code?

How to avoid hardware-level functional verification?

How to understand and account for hardware implementation details?

while (x<N) return x X_{start} Merge Buff Add Fork Cmp Branch

To place: 151 LUTs, 124 FFs To route: 1299 wires (bits)

Slow (mins-hours) 😕

HW

Guerrieri, Guha, Lavin, Hung, Josipović, and Ienne. DynaRapid: Fast-Tracking From C to Routed Circuits. FPL'24. Best Paper Award

FPGA implementation (synthesis, placement, and routing) is orders of magnitude slower than software compilation

SW

HW

How to make hardware design accessible to non-experts?

How to automatically extract parallelism from software code?

How to avoid hardware-level functional verification?

How to understand and account for hardware implementation details?

20x implementation speedup 🙂

Make hardware design broadly accessible, fast, and reliable

DYNAMO: Digital Systems and Design Automation Group

Research group:

dynamo.ethz.ch

Balor QoR estimator

Thanks! 🙂

github.com/ETHZ-DYNAMO/balor