

# A Novel Approach for Workload-Based GPU Datapath Power Optimization

Deepayan Dasgupta, Tanuj Sengupta  
Samsung Austin Research Center



SPONSORED BY

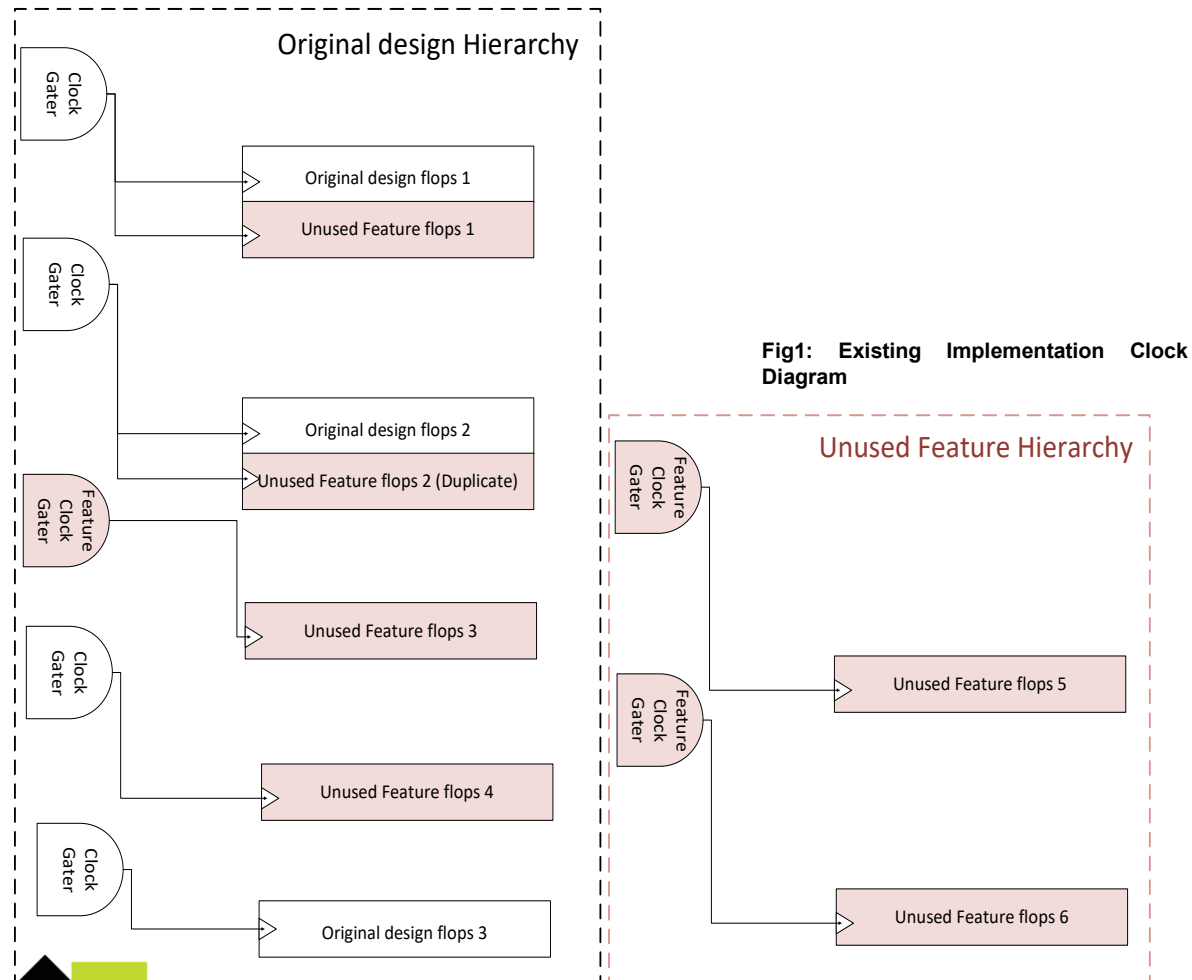


# Motivation

- A GPU is designed with support for multiple “features”. Features can be used to improve the quality of the rendered image or improve the speed at which the image is rendered.
- Although the GPU may have support for multiple features, each workload may use only a subset of the supported features, depending on the GPU compiler.
- One key GPU workload was analyzed to observe which features supported by the design are used by the workload.
- Some features can add a power penalty on the design despite not being used, as they require some structures in frequently used datapaths to be widened. This can include muxes that need to be widened and new flopbits are added to support the feature.
- These changes lead to a 2.64% power penalty on the design, worsening the overall performance per watt of the design.



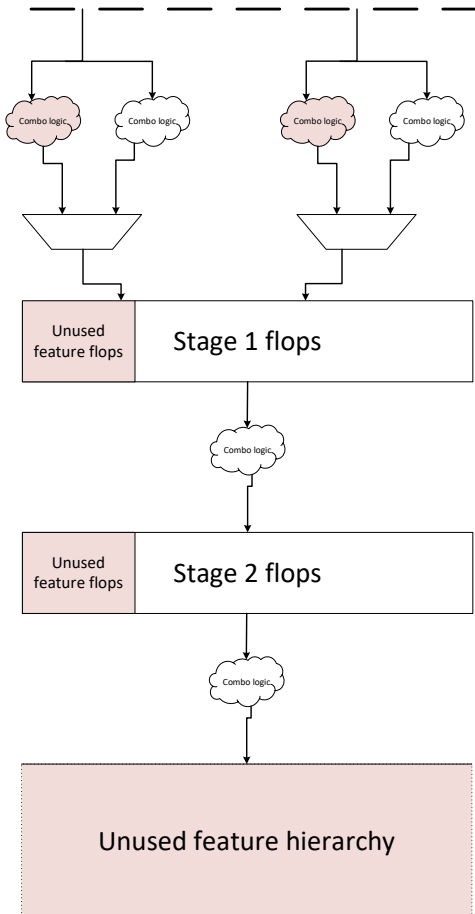
# Problem: Clock Power Penalty



- To understand what causes the power to increase despite the feature not being used, EDA tools related to synthesis and power are used to breakdown the delta.
  - Clock reports are generated to measure the flop-to-Integrated Clock Gating (ICG) ratio of existing structures before and after the feature is added. Further analysis revealed some of the flops are storing duplicate and it is possible to reuse existing flops instead.
  - Clock reports also explain which registers had to grow in size to accommodate the feature.
  - Some structures, despite not receiving any input data, burn power due to the module receiving ungated clock input signals. This leads to the flops in these modules burning sequential power, contributing to the overhead.

# Problem: Combo Power Penalty

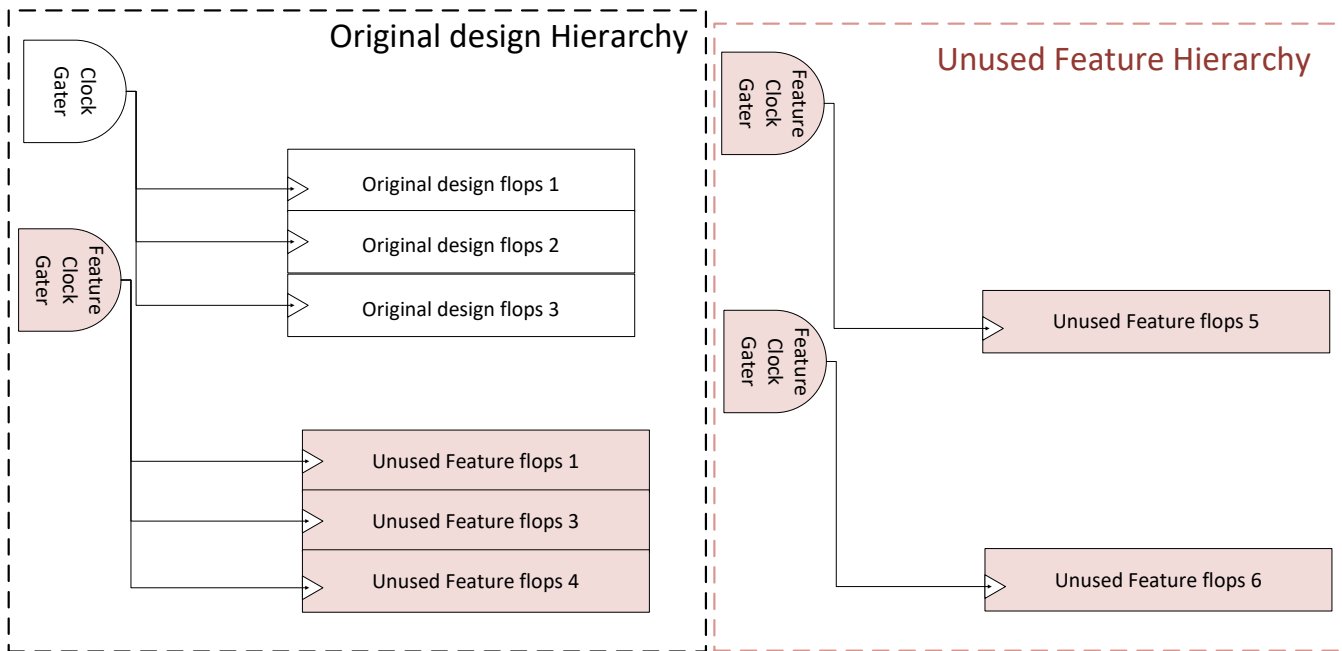
Fig2: Existing Implementation Datapath Diagram



- Muxes grew in size to arbitrate between existing data paths and new data paths added due to the feature.
- Fanout reports of existing mux structures were analyzed to understand why some muxes were reporting increased power despite no change to the mux tree in the RTL. Reports revealed that the fanout of these muxes had increased downstream due to the logic added for the feature, leading to these muxes being upsized and burning extra power.

# Implementation: Optimizing the Datapath for Power (Clock Gating)

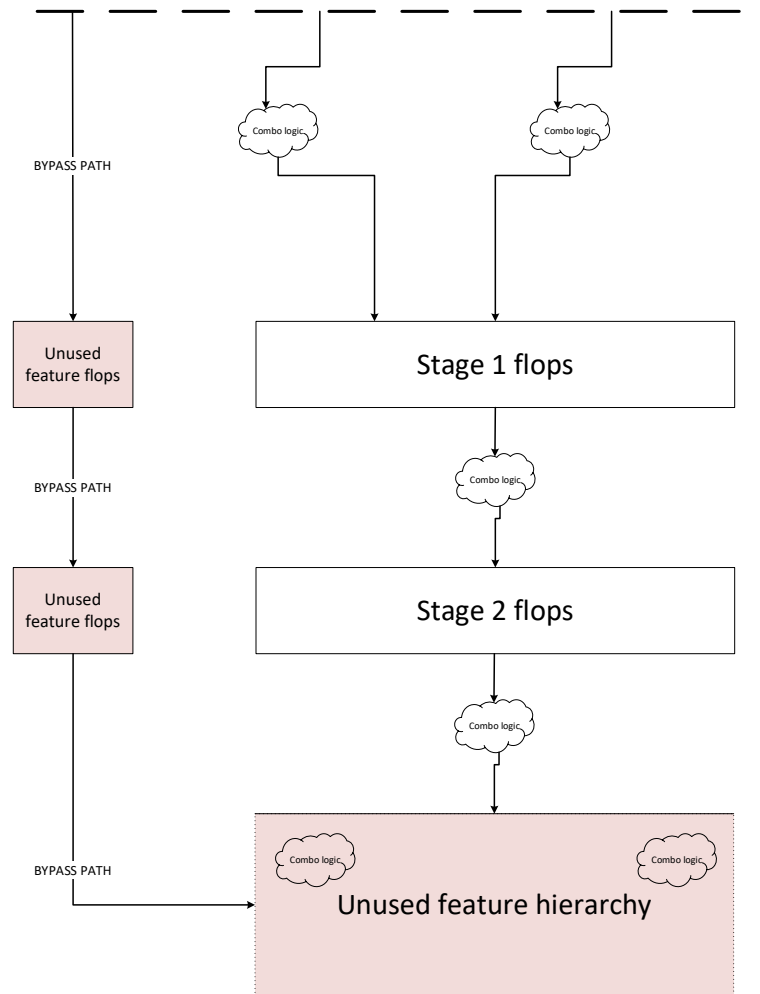
Fig3: New Implementation Clock Diagram



- To reduce this power overhead, the following design changes were made:
  - **CHANGE1:** To balance the flop-to-ICG ratio of existing structures, the RTL was recoded to remove redundant flops. The enable conditions of the flops were also rewritten to group more flops under the same enable, leading to fewer ICGs being instantiated. This led to an overall reduction in dynamic and leakage power (as the design now had fewer ICGs and flops).
  - **CHANGE2:** Additional qualifying signals were added to gate the registers that had to grow in size to accommodate the feature, reducing the clock and sequential power.
  - **CHANGE3:** A tighter clock signal was fed to the modules that were receiving a generic clock input, leading to reduced clock and sequential power.

# Implementation: Optimizing the Datapath for Power (Bypass Path)

Fig2: New Implementation Datapath Diagram (CHANGE 4+5)



- **CHANGE4:** Some floating point calculations were being done outside of hierarchy defined for the feature. This logic was moved into the hierarchy created for the feature.
- **CHANGE5:** This change allowed us to move most of the computation from the primary datapath to a bypass path which is used less often in the workload. This reduced the width of muxes in the primary datapath leading to reduction the combinational logic power penalty.

# Evidence

## ICG and Flop reduction (CHANGE1):

|                     |
|---------------------|
| ICG count % change  |
| -61.54%             |
| Flop count % change |
| -19.14%             |

## Power reduction due to ICG and flop reduction (CHANGE 1):

| Instance Name | Total Change (mW) | Dynamic Change (mW) | Leakage Change @NN75C (mW) | Clk Dyn. Change (mW) | Mem Dyn. Change (mW) | Logic Dyn. Change (mW) | Reg Dyn. Change (mW) |
|---------------|-------------------|---------------------|----------------------------|----------------------|----------------------|------------------------|----------------------|
| TOP           | -0.20%            | -0.09%              | -0.11%                     | -0.05%               | 0.00%                | -0.01%                 | -0.03%               |

## Power reduction due to improved clock gating (CHANGE 2+3):

| Instance Name | Total Change (mW) | Dynamic Change (mW) | Leakage Change @NN75C (mW) | Clk Dyn. Change (mW) | Mem Dyn. Change (mW) | Logic Dyn. Change (mW) | Reg Dyn. Change (mW) |
|---------------|-------------------|---------------------|----------------------------|----------------------|----------------------|------------------------|----------------------|
| TOP           | -0.35%            | -0.34%              | -0.01%                     | -0.07%               | 0.00%                | -0.19%                 | -0.09%               |

- Early power estimation flow did not report a lot of power savings from CHANGE1, as most of the redundant flops were already gated by the unused feature enabled, thus they were not burning any dynamic power.
- However, on the synthesized netlist, dynamic power savings were observed due to the clock networking shrinking, since it now has fewer buffers, ICGs and flops to drive. Leakage savings are observed from CHANGE1 due to the reduction of cells.
- CHANGE2+3 has dynamic power savings due to unused feature flops being clock gated correctly.



# Evidence

Power reduction due to compute being moved to bypass path (CHANGE 4+5):

| Instance Name | Total Change (mW) | Dynamic Change (mW) | Leakage Change @NN75C (mW) | Clk Dyn. Change (mW) | Mem Dyn. Change (mW) | Logic Dyn. Change (mW) | Reg Dyn. Change (mW) |
|---------------|-------------------|---------------------|----------------------------|----------------------|----------------------|------------------------|----------------------|
| TOP           | -1.52%            | -1.50%              | -0.01%                     | -0.12%               | 0.00%                | -1.21%                 | -0.18%               |

Std. Cell area savings due to all changes observed in a synthesized netlist:

|                                       | Flop Area | Comb area | Total Area |
|---------------------------------------|-----------|-----------|------------|
| Area Change (as % TOP Std. Cell Area) | -0.06%    | -0.28%    | -0.35%     |

- Removal of some of the muxes due to the bypass path led to significant dynamic power savings.
- There are area savings due to the following reasons:
  - Flop and ICG reduction
  - Removal of redundant combo logic due to the creation of a bypass path.





# Summary

- Motivation:
  - Some features supported by GPU hardware are not utilized by key workloads.
  - These features can add a power overhead depending on the implementation and type of feature.
- Optimization:
  - To reduce the power overhead from one unused feature, the following changes were made:
    - **CHANGE1:** Redundant flops were removed, and enable conditions were rewritten to group flops under fewer ICGs, reducing dynamic and leakage power.
    - **CHANGE2:** Added qualifying signals to gate larger registers, reducing clock and sequential power.
    - **CHANGE3:** Tighter clock signals replaced generic clock inputs, reducing clock and sequential power.
    - **CHANGE4:** Logic for floating-point calculations was moved into the appropriate feature hierarchy.
    - **CHANGE5:** Computation shifted to a less-used bypass path, reducing primary datapath mux width and combinational logic power.
- Conclusion:
  - The above power optimizations resulted in:
    - **61.5% ICG count reduction** in feature hierarchy.
    - **19.1% flopbits reduced** in feature hierarchy.
    - **2.07% total power reduction** at top level.
    - **0.35% std. cell area reduction** at top level.





**Deepayan  
Dasgupta**

Senior Engineer  
Samsung Austin  
Research Center



**Tanuj  
Sengupta**

Senior Engineer  
Samsung Austin  
Research Center