# DrGPU: A Top-Down Profiler for GPU
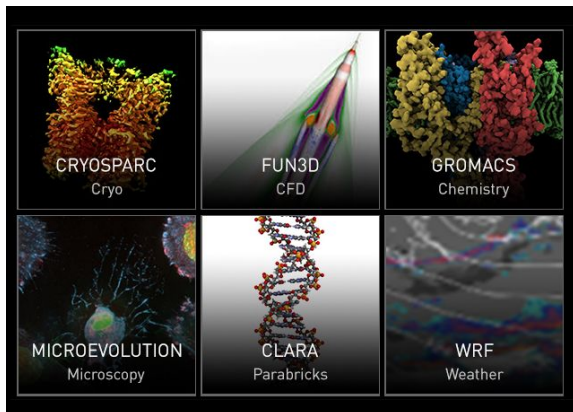
Yueming Hao[1], Nikhil Jain[2], Rob Van Der Wijngaart[2],
Nirmal Saxena[2], Yuanbo Fan[3], Xu Liu[1]

[1]North Carolina State University
[2]NVIDIA Corporation
[3]Tenstorrent Incorporated

# GPUs are Broadly Used for Acceleration

# GPU Programming

```
// Kernel definition
__global__ void VecAdd(float* A, float* B, float* C)
{
    int i = threadIdx.x;
    C[i] = A[i] + B[i];
}

int main()
{
    ...
    // Kernel invocation with N threads
    N = 1024
    VecAdd<<<1, N>>>(A, B, C);
    ...
}
```

→ GPU Kernel

# of GPU threads

1 warp = 32 threads

# Existing GPU Performance Tools



NVIDIA Nsight Compute

HPCToolkit

Existing tools apply high-level *hotspot* analysis

# DrGPU Contribution

DrGPU

➢ tells you where GPUs waste on stalling by a top-down tree

➢ provides analysis and optimization guidance for non-experts

# What are Bottlenecks of GPU?



SM

1 instruction per cycle (**IPC**) per warp scheduler per SM =>

Ideal instruction per cycle is **4**

# What are Bottlenecks of GPU?

Gap is large between achieved IPC and ideal IPC!

# Categories of Stall Reasons

# Instruction Related Stalls

Instruction Related Stalls

| Delay due to pending global stores before exit | Delay due to instruction dependency | Delay due to pipe contention |
|---|---|---|

```
__global__ kernelA(){
   ...
   arrayA[thread_id] += 10;
}
```

```
__global__ kernelB(){
   ...
   c = log10(a)
   f = c + e
   ...
}
```

```
__global__ kernelC(){
   ...
   for (int i = 0, i<1000,i++){
      c = a * b
      f = c * e
   }
   ...
}
```

# DrGPU Overview



Online Data Collection

Fully Optimized Binary

NCU Profiling

GPU0

Offline Data Processing

Profile Coalesce

Tree Construction

Optimization Guidance

DrGPU utilizes NCU(NVIDIA Nsight Compute) to profile GPU applications with customized GPU hardware counter list

DrGPU gives optimization suggestions based on collected hardware counters.

# An example of Analysis Trees



No-issue cycles
97.81%
Util/SOL: 48.54%(SM)
Issue IPC: 0.09

High level overview of a kernel

Delay due to device memory accesses
95.50% of no-issue cycles

Delay due to instruction dependency
3.92% of no-issue cycles

Main stall reasons

Max active warps: 32
Theoretical active warps: 8
Achieved active warps: 7.98
Register per thread: 254
Block size: 256
Limited by : Register, Blocksize

FP64
62% of all inst

integer
14% of all inst

Further breakdown for details

Try to reduce block size to 128 to increase active warps.

Try to change functions to their fast version to minimize the stalls.
Note: There may be an accuracy loss.

Optimization suggestions

chemistry_file.H:
2425 logFcent = log10(    15.90%
2490 qr[1] *= Corr * k_f / (exp(-g_RT[5] - g_RT[5] + g_RT[7]) * refC);
9.71%

Related code lines and contributions to stalls

# Evaluation Platforms

- GPU
  - V100  16GB
  - GTX 1650 4GB
- Applications
  - Rodinia benchmarks
  - YOLOv4 (Darknet)
  - LULESH2
  - PeleC
  - Castro

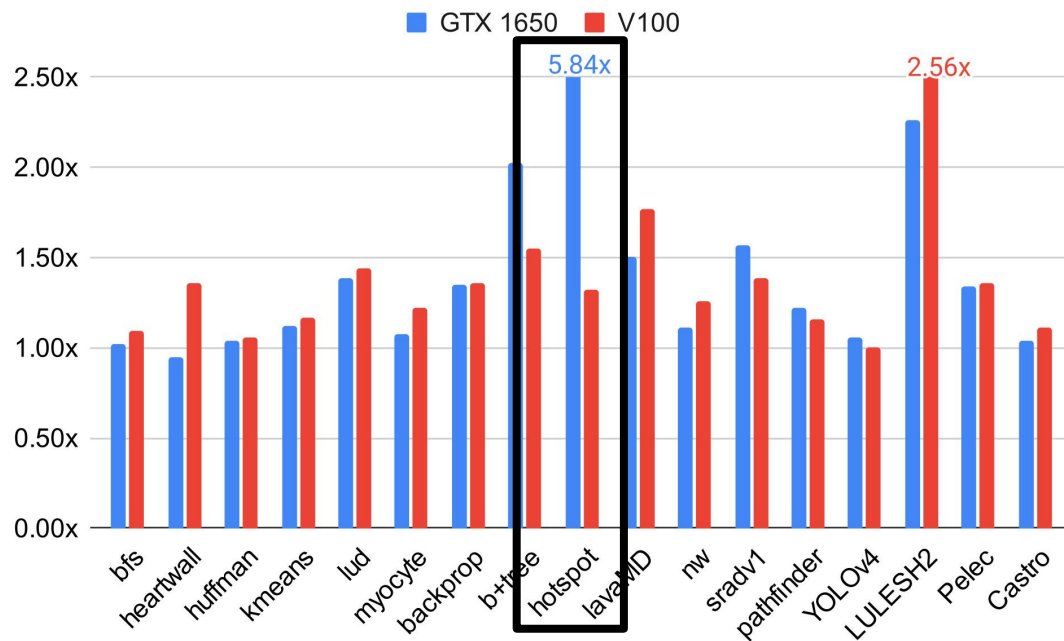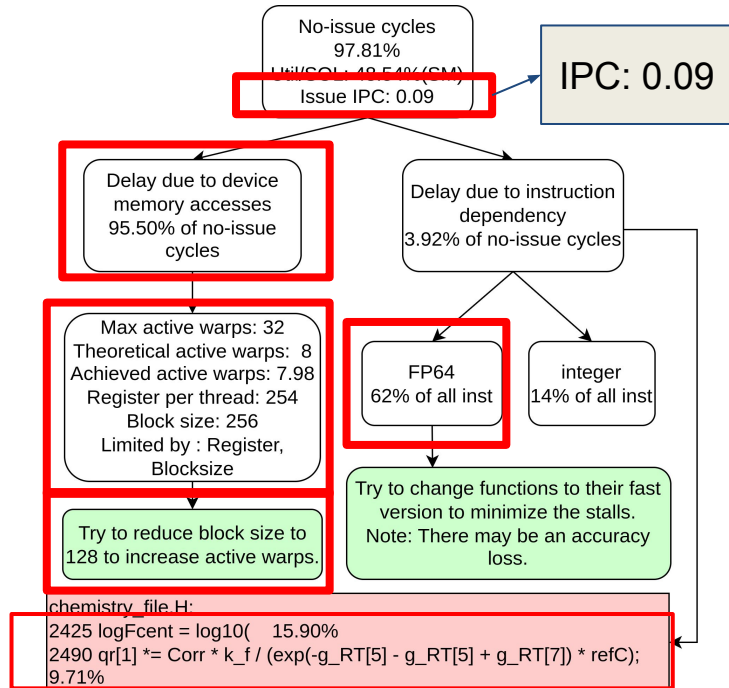| Application | Kernel | State | Optimization |
|---|---|---|---|
| bfs | Kernel | Long Scoreboard | Loop unrolling |
| heartwall | kernel | Wait | Loop unrolling |
| huffman | vlc_encode_kernel_sm64huff | Barriers | Restruct code |
| kmeans | kmeansPoint | Wait | Loop unrolling |
| lud | lud_diagonal | Wait/ Short Scoreboard/ No instruction | Restruct code |
| myocyte | solver_2 | Short Scoreboard | Function spliting |
| | | Math Pipe Throttle | Add use_fast_math |
| backprop | bpnn_layerforward_CUDA | Barrier | Remove unnecessary barriers |
| | | Wait | Restruct code |
| b+tree | findRangeK | Long Scoreboard | Restruct code |
| | | Barrier | Reduce blocksize |
| hotspot | calculate_temp | Wait | Remove inappropriate FP convertion |
| | | | Add use_fast_math |
| lavaMD | kernel_gpu_cuda | Long Scoreboard/Wait | Loop unrolling |
| | | Wait | Replace speical FP functions |
| nw | needle_cuda_shared_1 | Barriers | Remove unnecessary barriers |
| | | | Replace syncthreads with sync warp safely |
| sradv1 | reduce | Short Scoreboard | Loop unrolling |
| | | Barrier | Reduce blocksize |
| pathfinder | dynproc_kernel | Short Scoreboard | Replace shared memory with variables |
| | | Wait | Remove unnecessary iterations |
| Darknet | im2col_gpu_kernel_ext | Wait | Loop unrolling |
| LULESH2 | ApplyMaterialProperties AndUpdateVolume_kernel | Wait | Add use_fast_math |
| Pelec | react_state | Long Scoreboard | Increase occupancy |
| | | Wait | Replace speical FP functions |
| Castro | trace_ppm | long scoreboard | Increase occupancy |

# Speedups with Optimization Guided by DrGPU
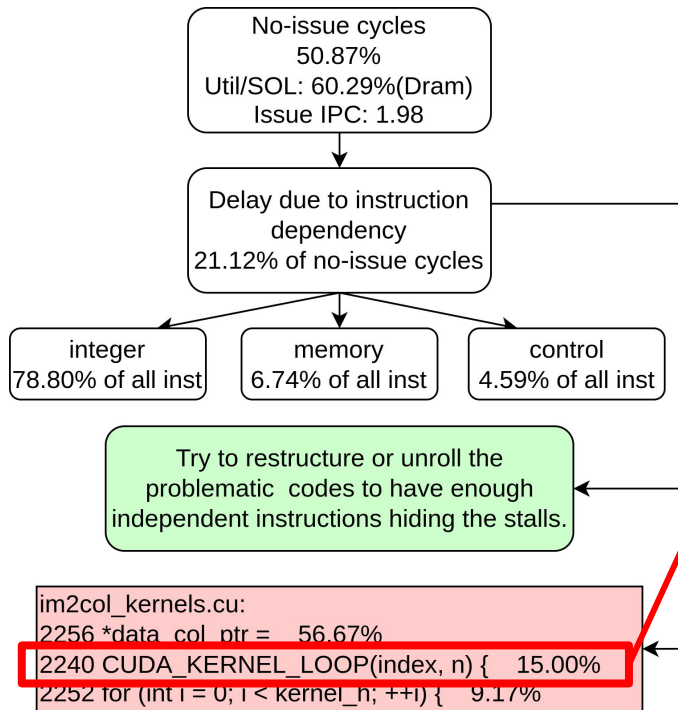


1.58X on GTX 1650

1.36X on V100

# **PeleC**



A portion of the analysis tree on GTX 1650

Optimizations

- Set blocksize to 128.
- Replace functions to their faster version. e.g., log10 -> log10f (0.1% precision loss)

1.34X speedup on GTX 1650
1.36X speedup on V100

# YOLOv4

Optimization
- Loop unrolling

1.06X speedup on GTX 1650

A portion of the analysis Tree on GTX 1650

# **Conclusions**

We propose DrGPU, a novel top-down profiler for GPU kernels.

➢ DrGPU quantifies stall cycles and decomposes them according to various hardware events for root causes.

➢ DrGPU generates performance analysis trees including source code location, root causes, and actionable guidance

➢ We optimized a number of applications with the insights provided by DrGPU with nontrivial speedups on both desktop and Summit NVIDIA GPUs

➢ Some of optimization suggestions proposed by DrGPU have been integrated to NVIDIA Nsight Compute

Code is available at: https://github.com/FindHao/drgpu

# Thanks!