# Summary & Background
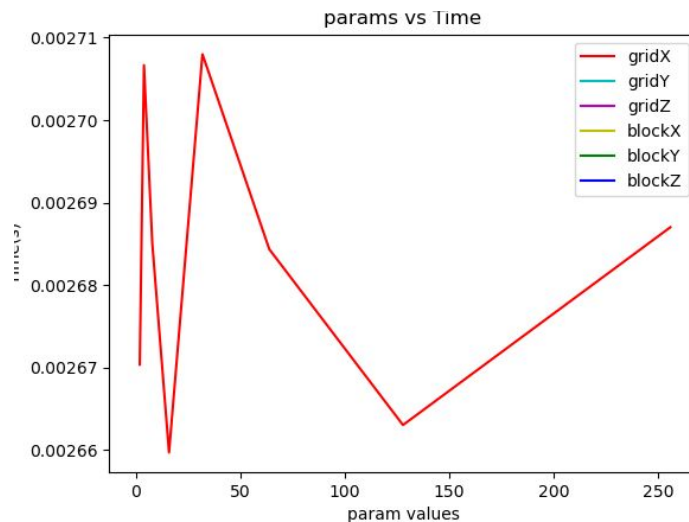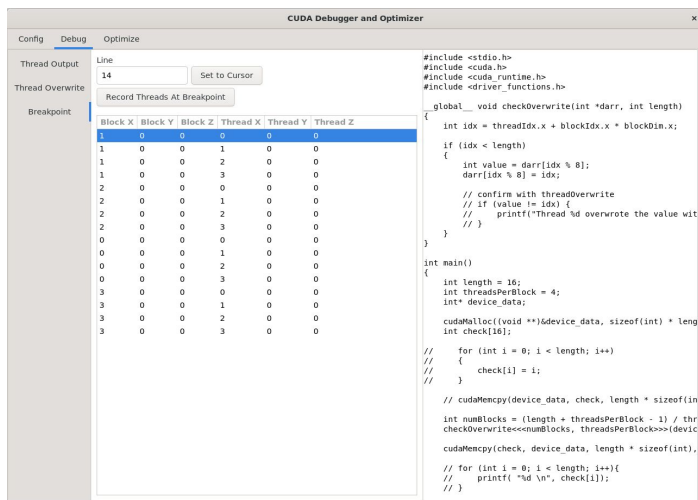
➔ Project: CUDA debugger/optimizer

  ◆ thread output, thread overwriting, breakpoint

  ◆ time/memory bottleneck, optimal macro config, speedup

➔ Produce values and visuals to simplify the debugging/optimizing process

# threadOutput & threadOverwrite

Output: blockIdx, threadIdx,

and updated variable value



Output: address, blockIdx, threadIdx, (if array, include

index), and updated variable value

# Time Bottleneck: gprof/GPU/API



**Function Execution Times**

main(66.67)

_init(33.33)

gprof: function runtime as a % of total runtime

GPU

API

# Memory Bottleneck: nvprof

staticMem



dynamicMem

Mem size

throughput

# Optimal Config



➤ User selects parameters they want to change, and gives a range to help generate all possible configs

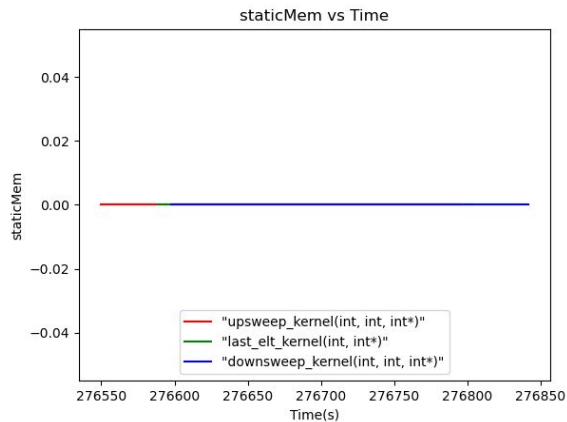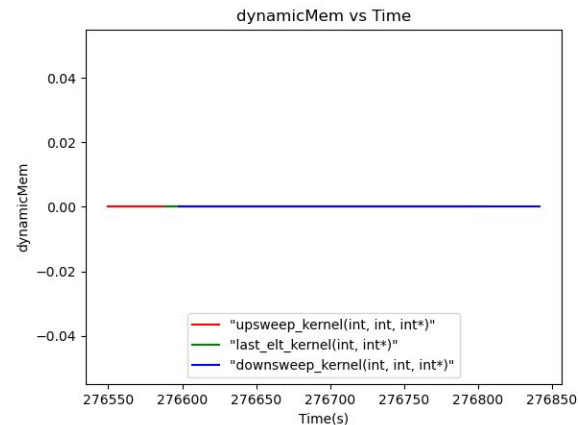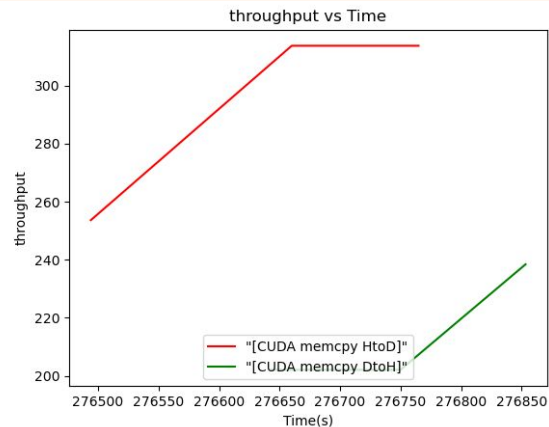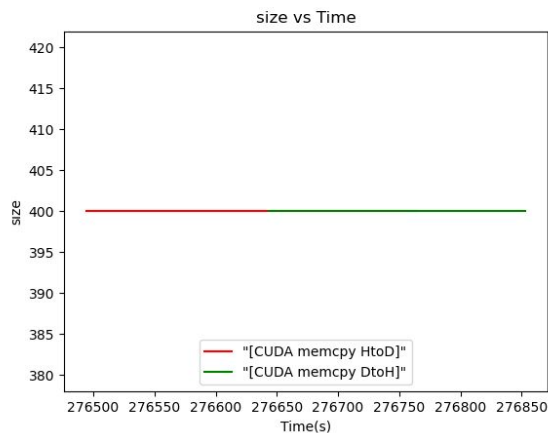➤ Run code with different params

➤ Add runtime of code to params line

➤ Note: Output of optimal config is used to generate parameter graphs

➤ Requires user to have definitions of macros in the code

# Parameter Graphs

# Breakpoint & Speedup

➔ Output: sorted list of threads (blockIdx, threadIdx)

◆ Purpose of sort: easier to identify missing threads

# CUDA Debugger & Optimizer

kflorend & huiningl

## Interactive Test Cases!

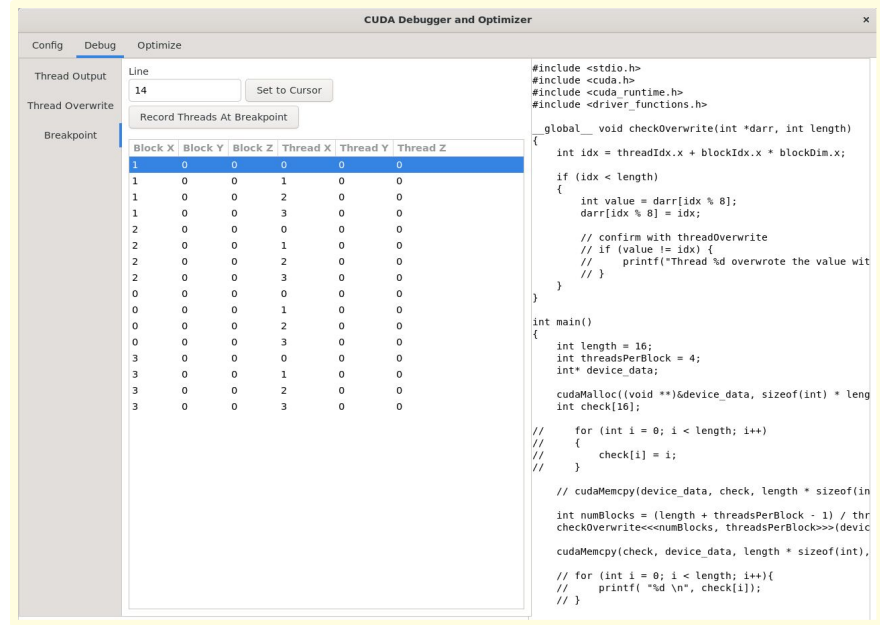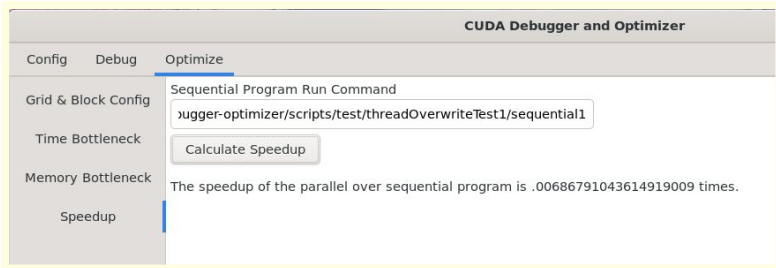| Flags | scan | test1 (sequential1) | test2 |
|---|---|---|---|
| -m | scan | test/threadOverwriteTest1 | test/threadOverwriteTest2 |
| -r | scan/cudaScan -m scan -i random -n 100 | test/threadOverwriteTest1/test1 (sequential1) | test/threadOverwriteTest2/test2 |
| -c | scan.cu | test1.cu (sequential1.cu) | test2.cu |
| -v | device_data[i+twod1-1] | darr[idx % 8] (check[i % 8]) | darr[0] |
| -t | int | int | int |
| -l | 63 | 16 (14) | 14 |
| -a | y | y | n |
| -b | BLOCK_DIM_X,,, | N/A | N/A |
| -g | GRID_DIM_X,,, | N/A | N/A |
| -v (opt) | 2 0 0 0 0 0,4 0 0 0 0 0 | N/A | N/A |
| -f | output/optimizeConfig.txt | N/A | N/A |
| -s | N/A | test/threadOverwriteTest1/sequential1 | N/A |
| -p | N/A | test/threadOverwriteTest1/test1 | N/A |