

Building Custom GPU Power Models With AccelWattch Project Report

Huy Dinh Tran
htran197@ucr.edu

University of California, Riverside

Daniel Wong
danwong@ucr.edu

University of California, Riverside

Abstract

Graphics Processing Units (GPUs) are the choice of accelerator that is widely utilized in data analytics and machine learning fields. Performance per watt is one of the most crucial evaluation metrics that determine the efficiency of a GPU. Performance and power models are equally important when researching and designing new technologies for GPU architects. However, while GPU performance modeling has seen great progression, power modeling has not been focused on to meet today's standards. Over the past years, researchers from Intel, Northwestern University, Purdue University, and University of British Columbia had developed and published an open-sourced power modeling framework named AccelWattch[4]. AccelWattch is an analytic model for modern GPUs that accurately estimated the constant, static, and dynamic power consumption. In this paper, we would be experimenting with creating our own power model of an NVIDIA Pascal GPU, the GTX 1050 Ti, using AccelWattch. The following sections will give a run-down procedure of the experiment which includes setting up validation benchmarks, hardware profiling for validation, running simulation with power model, and validating simulated estimated results with measured results from real GPU.

1 Introduction

1.1 Motivation

GPUs are widely used for many applications in data analytics and machine learning. The efficiency of a GPU is determined by the performance per watt. In designs, both performance and power models are important when researching and testing new technologies for new GPUs. While GPU performance modeling has progressed in great stride, power modeling has lacked behind and is not up-to-date with modern architecture. The previous state-of-art power modeling tool, GPUWattch, is outdated and does not accurately model modern GPU architecture. It is lack of cycle-level power model and is not able to capture constant and static power. To further aid the development of GPUs, researchers from Intel, Northwestern University, Purdue University, and University of British Columbia had developed and published an open-sourced power modeling framework named AccelWattch. AccelWattch is an analytic model for modern GPUs that accurately estimated the constant, static, and dynamic power consumption. AccelWattch was able to deliver promises that

GPUWattch wasn't able to do with better accuracy and model modern technology better.

1.2 AccelWattch

AccelWattch is an analytic model that accurately captures constant, static, and dynamic power consumption in the presence of DVFS, thread divergence, intra-warp functional unit overlap, variable SM occupancy, and power gating. AccelWattch was specially designed for the NVIDIA Volta GPU architecture. AccelWattch resolves long-standing needs for modern GPU architectures: the lack of a cycle-level power model, and the inability to capture the constant and static power with existing methodologies. AccelWattch is integrated with GPGPU-Sim[1] and Accel-Sim[5] which enable it to be directed by emulation (PTX[6]) or trace-driven (SASS[7]) software performance models, by hardware performance counters, or by a combination of the above. AccelWattch can enable reliable design space exploration. Directly applying the Volta power model on a GPU configuration resembling the Pascal and Turing architectures result in accurate power models for these architectures without tuning specifically for them

1.3 Flowchart

Figure 1 shows the overall flowchart of how AccelWattch create a power model to estimate the power consumption of a GPU. There are three sections of modeling which are constant, static, and dynamic power modeling. Each of these sections follows the methodology to accurately model the power consumption. These will then create the power model which can be used for validation against real hardware power measurements.

1.4 Methodology

The methodology of AccelWattch consists of three modeling section which serves to estimate the constant, static, and dynamic power consumption. It takes the presence of DVFS, thread divergence, intra-warp functional unit overlap, variable SM occupancy, and power gating into account to accurately build a power model.

1.4.1 Constant Power Modeling. Figure 2 shows a high-level equation of the total power consumption of a GPU. Figure 3 shows a breakdown equation of Figure 2 with capacitance, voltage, and frequency. m and n are constant. However, while this equation can model older generations of

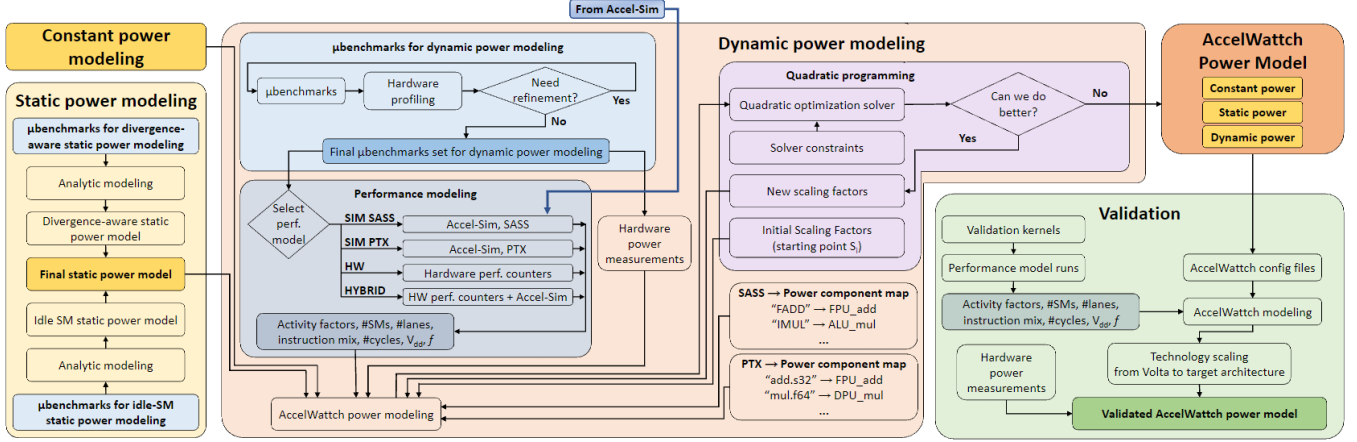


Figure 1. AccelWatch power modeling flowchart

$$P_{total} = P_{proc, dyn} + P_{mem, dyn} + P_{proc, static} + P_{mem, static} + P_{const}$$

Figure 2. High-level total power equation

$$P_{total} = \underbrace{mCV^2f}_{Dynamic} + \underbrace{nV}_{Static} + P_{const}$$

Figure 3. Older GPU high-level power equation

GPUs, recent technologies are built upon DVFS which will cause negative constant and static power estimations using this equation. To take DVFS into account, a new equation

$$P_{total} = \underbrace{\beta Cf^3}_{Dynamic} + \underbrace{\tau f}_{Static} + P_{const}$$

Figure 4. Modern GPU high-level power equation

using recently-published data for fully-realized processors which better models GPUs with DVFS are shown in Figure 4. β and τ are constant.

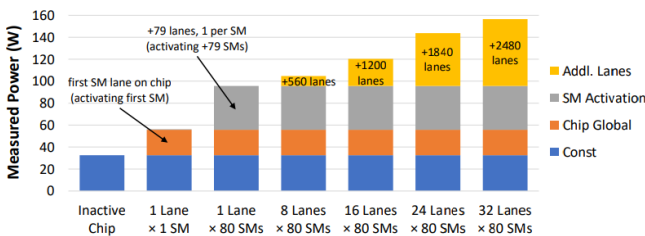


Figure 5. Inferring the power consumption of activating power-gated chip-wide and SM-wide components

1.4.2 Static Power Modeling. The static power modeling consists of three main awareness which are power gating, divergence, and idle SMs. For power gating, the model considers the impact on the power consumption of re-activating components on real hardware. It accurately captures the power consumption in the presence of power gating. For example, Figure 5 shows a chart that displays the power consumption if we just activate a local component (an SM with 1 lane), it will activate a shared global component (L2 cache) which results in a large amount of static power consumption.

$$P_{static, addLane} = \frac{P_{static, 32Lanes} - P_{static, firstLane}}{31}$$

$$P_{static, yLanes} = \begin{cases} P_{static, firstLane} + P_{static, addLane} \cdot (y - 1), & y \leq 16 \\ P_{static, firstLane} + \frac{1}{2} P_{static, addLane} \cdot 15 + \frac{1}{2} P_{static, addLane} \cdot (y - 17), & y > 16 \end{cases}$$

Figure 6. Warp-divergence based on number of active threads equation

For divergence, the model considers how a warp operates. In modern GPUs, a warp executes by running two 16-thread half-warps, one after the other. Figure 6 shows an equation that models the power consumption of divergence within a warp depending on the number of active threads. y is the number of active threads.

$$P_{perIdleSM} = \sqrt[n]{\prod_{i=1}^n \frac{P_{total, i} - P_{const} - \frac{P_{total, 80SMs, i} - P_{const}}{80} \cdot N_{activeSM}}{N_{idleSMs}}}$$

Figure 7. Power consumption of per idle SM equation

For idle SMs, Figure 7 shows the power consumption per idle SM. The equation considers the dynamic plus static

power consumption per active SM when running microbenchmark i . Based on that, it will then derive the power consumption for idle SMs by eliminating all other powers and repeating it for n microbenchmarks to find the power per idle SM. $n = \#$ microbenchmarks, $i =$ specific microbenchmark (microbenchmark i), 80 = NVIDIA Volta GV100 has 80 SMs, $N_{\text{activeSM}} = \#$ active SMs, $N_{\text{idleSM}} = \#$ idle SMs.

$$P_{est.} = \underbrace{\sum_{i=1}^N \frac{a_i \cdot \hat{E}_i}{T_{elapsedTime}} \cdot x_i}_{Dynamic} + \underbrace{\frac{P_{static, yLanes}}{80} \cdot k + P_{perIdleSM} \cdot (80 - k)}_{Static} + \underbrace{P_{const}}_{Constant}$$

Figure 8. Total estimated power consumption with dynamic power equation

1.4.3 Dynamic Power Modeling. Figure 8 shows an equation that overall estimates the power consumption of a GPU where it consists of dynamic, static, and constant power. The equation to estimate the dynamic power is the sum of each initial estimate of component i's energy per access \hat{E}_i , its activity factor a_i (the number of accesses to it during execution), and the run time $T_{elapsedTime}$ for N microarchitectural components. x_i is the unknown variable caused by inaccuracies that depend on some microarchitectural components. a_i = activity factor, k = # SMs, \hat{E}_i = initial estimate of component i's energy consumption per access, x_i = unknown variable (caused by inaccuracy), i = specific component (component i), N = # microarchitectural components.

2 Implementation

For this project, we followed the "AccelWattch MICRO'21 Artifact Appendix Manual" which is a procedure of steps to produce the results similar to the original AccelWattch's paper. Our GPU for the implementation is the NVIDIA GeForce GTX 1050 Ti which is based on the Pascal architecture. There are three main steps which consist of setting up validation benchmarks, hardware profiling to create a power model, and running simulations with the power model to estimate the power consumption of benchmarks.

2.1 Setting up validation benchmarks for AccelWattch

Firstly, we compiled the microbenchmarks, validation set benchmarks for simulator runs, and power profiling that were provided. The datasets of the benchmarks were also provided which needed to be extracted from a compressed file for later use. Table 1 shows the list of benchmarks used for AccelWattch. Figure 9 shows a screenshot of all the microbenchmark binaries that were compiled and ready to use for hardware profiling and simulations.

| List of benchmarks | |
|------------------------|----------------------|
| backprop-rodinia-3.1 | binomialOptions |
| b+tree-rodinia-3.1 | dct8x8 |
| fastWalshTransform | histogram |
| hotspot-rodinia-3.1 | kmeans-rodinia-3.1 |
| mergeSort | parboil-mri-q |
| parboil-sad | parboil-sgemm |
| pathfinder-rodinia-3.1 | quasirandomGenerator |
| sobelQRNG | srad_v1-rodinia-3.1 |

Table 1. List of benchmarks for AccelWattch[illegible]

Figure 9. All microbenchmark binaries

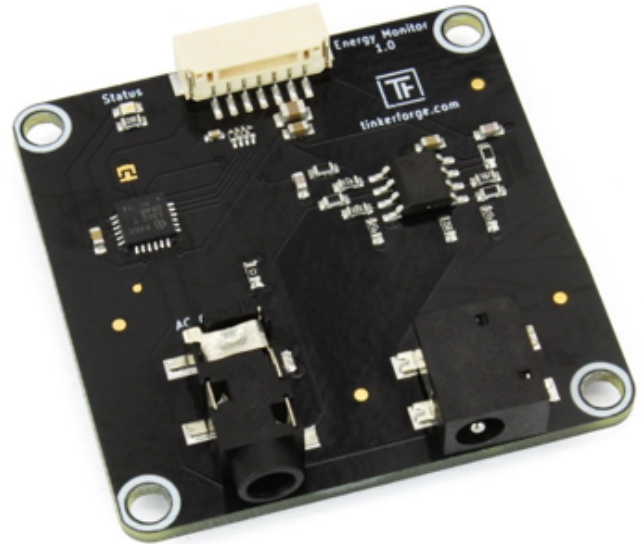


Figure 10. TinkerForge Energy Monitor Bricklet

2.2 Hardware profiling for AccelWattch validation

Unfortunately, the NVIDIA GTX 1050 Ti was not equipped with an internal power consumption monitor out of the box. Alternatively, we bought an Energy Monitor Bricklet[9] from TinkerForge, shown in Figure 10, which is a small current

sensor that we connected to the GPU to measure the total power consumption. We then integrated the Bricklet into AccelWatch to measure the power. Figure 12 shows a screenshot of the source code files of the Bricklet. Figure 13 shows the screenshot of AccelWatch’s power measuring source code which we integrated the Bricklet code into it to measure the power consumption for the hardware profiling. After setting up the Bricklet, we then performed hardware profiling for validation. This step measured power five times and calculated the average for the runs for each validation set kernel on real hardware. It will first heat up the GPU chip to $\geq 65^{\circ}\text{C}$ before taking power measurements. Temperature variability affects static power exponentially. Keeping a constant temperature during hardware measurements eliminates this noise. Figure 11 shows a screenshot of the hardware profiling runs.

Figure 11. Hardware profiling for validation runs

Figure 12. Bricklet source files

Figure 13. Integrating Bricklet into AccelWatch power measuring source code

2.3 Running AccelWattch and collecting power model results

```
b+tree-rodinia-3.1.log
1 kernel_name = findRangeK
2 kernel_launch_uid = 1
3
4 Kernel Average Power Data:
5 kernel_avg_power = 74.1722
6 gpu_avg_IBP, = 0.130758
7 gpu_avg_ICP, = 0
8 gpu_avg_DCP, = 0
9 gpu_avg_TCP, = 0
10 gpu_avg_CCP, = 0
11 gpu_avg_SHRDP, = 0
12 gpu_avg_RFP, = 3.64172
13 gpu_avg_INTP, = 2.26018
14 gpu_avg_FPUP, = 0
15 gpu_avg_DPUP, = 0
16 gpu_avg_INT_MUL24P, = 0
17 gpu_avg_INT_MUL32P, = 0
18 gpu_avg_INT_MULP, = 0.66215
19 gpu_avg_INT_DIVP, = 0
20 gpu_avg_FP_MULP, = 0
21 gpu_avg_FP_DIVP, = 0
22 gpu_avg_FP_SQRTP, = 0
23 gpu_avg_FP_LGP, = 0
24 gpu_avg_FP_SINP, = 0
25 gpu_avg_FP_EXPP, = 0
```

3 Result


```
[4]: ptDir = "/home/bram377/accol-sim-framwork/accoluthat/power_reports/pastcal_gfx_xin"
ptCol = ["benchmark", "kernel 1 power", "kernel 2 power"]
ptFilePower = pd.DataFrame(columns=ptCol)
for filename in os.listdir(ptDir):
    f = os.path.join(ptDir, filename)

    line5 = linecache.getline(f, 5)
    line299 = linecache.getline(f, 299)

    for t in line5.split():
        try:
            linePower = float(t)
        except ValueError:
            pass
    for t in line299.split():
        try:
            line299Power = float(t)
        except ValueError:
            pass

    ptFilePower = ptFilePower.append({"benchmark": filename[0:-4], "kernel 1 power": linePower, "kernel 2 power": line299Power}, ignore_index=True)

print("PTX")
print(ptFilePower)
```

```
PTX
kernel 1 power    benchmark kernel 2 power
0    10.4937      parboil-wf1q      81.08075
1    83.8674  backprop-rodinia-3.1    82.5867
2    83.8674  knnns-rodinia-3.1      82.5867
3    90.4211      fastclusterform     91.1011
4    90.4211      parboil-wd         91.1011
5    83.1567  quasicrandGenerator     83.8564
6    83.1467      cobs1000c           83.8564
7    83.1467      merge3d             83.8564
8    83.1467      parboil-gemm         83.8564
9    74.4558  stad-wf-rodinia-3.1     71.5451
10  72.7933  bf-tree-rodinia-3.1      71.5451
11  72.7933  histogramOptim          71.5451
12  72.7933      dct8b8              71.5451
```

Figure 15. Collecting data

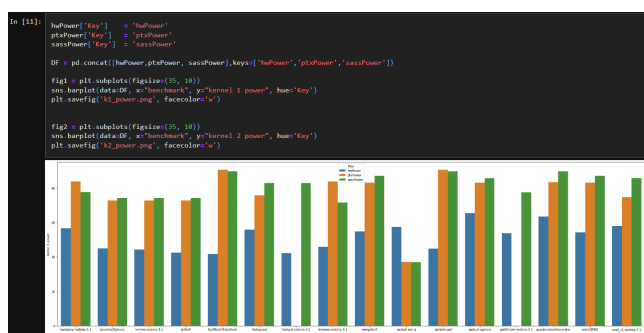


Figure 16. Plotting

```

In [51]:
sasIMPk1 = []
ptbIMPk1 = []

for i, row in scatDFk1.iterrows():
    if math.isnan(scatDFk1['halfway'][i]) or math.isnan(scatDFk1['ptbPower'][i]) is True:
        sasIMPk1.append(np.NaN)
        scatDFk1['halfway'][i], scatDFk1['sasPower'][i] = None

for i, row in scatDFk1.iterrows():
    if math.isnan(scatDFk1['halfway'][i]) is True and math.isnan(scatDFk1['ptbPower'][i]) is True:
        ptbIMPk1.append(np.NaN)
        scatDFk1['halfway'][i], scatDFk1['ptbPower'][i] = None

print("sasIMPk1 k1:", np.average(sasIMPk1))
print("ptbIMPk1 k1:", np.average(ptbIMPk1))

sasIMPk2 = []
ptbIMPk2 = []

for i, row in scatDFk2.iterrows():
    if math.isnan(scatDFk2['halfway'][i]) is True and scatDFk2['halfway'][i] > 0.8 and math.isnan(scatDFk2['sasPower'][i]) is True:
        sasIMPk2.append(np.NaN)
        scatDFk2['halfway'][i], scatDFk2['sasPower'][i] = None

for i, row in scatDFk2.iterrows():
    if math.isnan(scatDFk2['halfway'][i]) is True and scatDFk2['halfway'][i] > 0.8 and math.isnan(scatDFk2['ptbPower'][i]) is True:
        ptbIMPk2.append(np.NaN)
        scatDFk2['halfway'][i], scatDFk2['ptbPower'][i] = None

print("sasIMPk2 k2:", np.average(sasIMPk2))
print("ptbIMPk2 k2:", np.average(ptbIMPk2))

sasIMPk1 k1: 62.28806615918332
ptbIMPk1 k1: 62.3441317657629
sasIMPk2 k2: 60.14936245781224

```

Figure 17. Calculating MAPE

3.1 Kernel 1 results

Figure 18 shows a bar plot of all results from running kernel 1 of each benchmark for real, PTX, and SASS power.

3.2 Kernel 2 results

Figure 19 shows a bar plot of all results from running kernel 2 of each benchmark for real, PTX, and SASS power.

3.3 Real Power vs Simulated Power Kernel 1

Figure 20 shows the correlation plot of real vs PTX and real vs SASS in kernel 1 of benchmarks. We got the average MAPE for PTX of 62.94% and the average MAPE for SASS of 62.28%.

3.4 Real Power vs Simulated Power Kernel 2

Figure 21 shows the correlation plot of real vs PTX and real vs SASS in kernel 2 of benchmarks. We got the average MAPE for PTX of 68.34% and the average MAPE for SASS of 60.15%.

4 Analysis

After getting all results together, we concluded that the model we created and tested using AccelWatch was inaccurate compared to what was advertised in the AccelWatch published paper. The average MAPE for both kernels on PTX and SASS for our experiment was 63.42%. In the paper, the MAPE of the Pascal TITAN X they got was $11 \pm 3.8\%$. We can clearly see that our model is far from accurate compared to their results. There are several possible causes for the inaccuracy.

- Since AccelWattch is driven by a performance model, our performance model was never validated at any stage of the procedure which by itself can be the only reason for this result. The performance model is the core part of AccelWattch which feeds execution stats to AccelWattch for estimating power consumption.
- There is a difference in binary code between running on real hardware vs simulation. Since PTX is a virtual ISA that will get translated to native SASS, there will be optimizations that will change the code. In addition, GPGPU-Sim does not simulate actual SASS instructions. It actually simulates PTXPlus[2] instructions which claimed to have a one-to-one mapping with SASS instructions. In general, the code simulated in GPGPU-Sim is not identical to the one running in real GPU.
- Another major cause is the inaccuracies within the GPGPU-Sim configuration file. Since internal microarchitecture information about the 1050 Ti is limited, we were not able to configure GPGPU-Sim to fit exactly our GPU configuration. Some of the configurations were not found so we used the provided configuration of the TITAN X.

5 Challenges

Overall, we were able to achieve the main goal of building the power model using AccelWattch and validating it against real power consumption readings. However, the process was not smooth sailing and we did encounter several challenges.

- Since AccelWattch was implemented by a small group of researchers, there were many small inconsistencies within the scripts. Source and destination directories are either hard-coded, inconsistent, or not updated which results in some directories not being found. Some Makefiles are also not fully working because of included flags that caused errors when compiling files. Since AccelWattch was specially designed for Volta architecture, there are TensorCore-related benchmarks

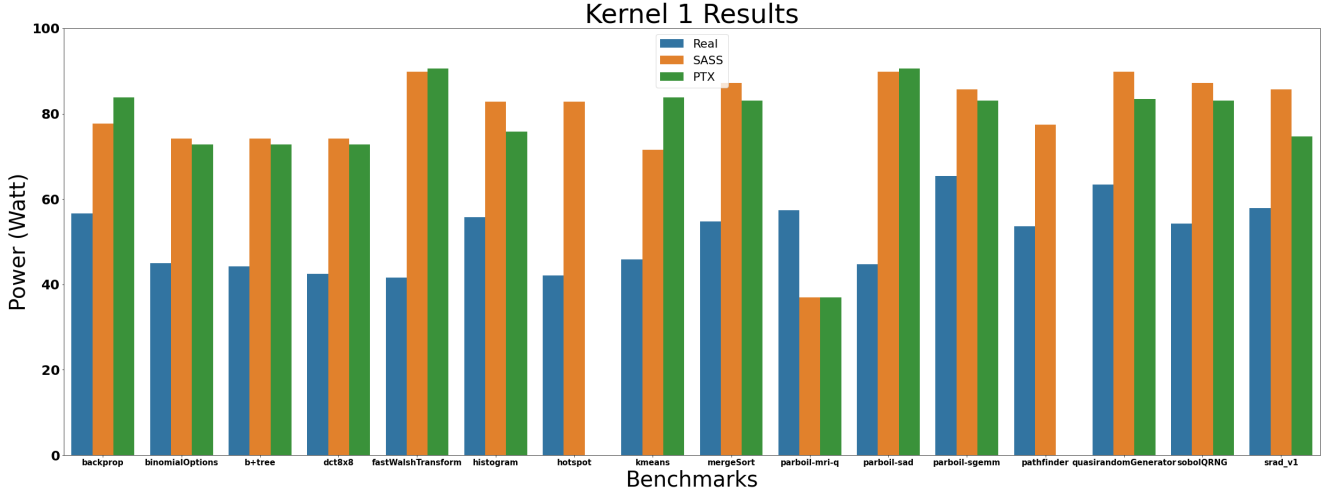


Figure 18. Real, PTX, SASS power results over kernel 1 of each benchmark

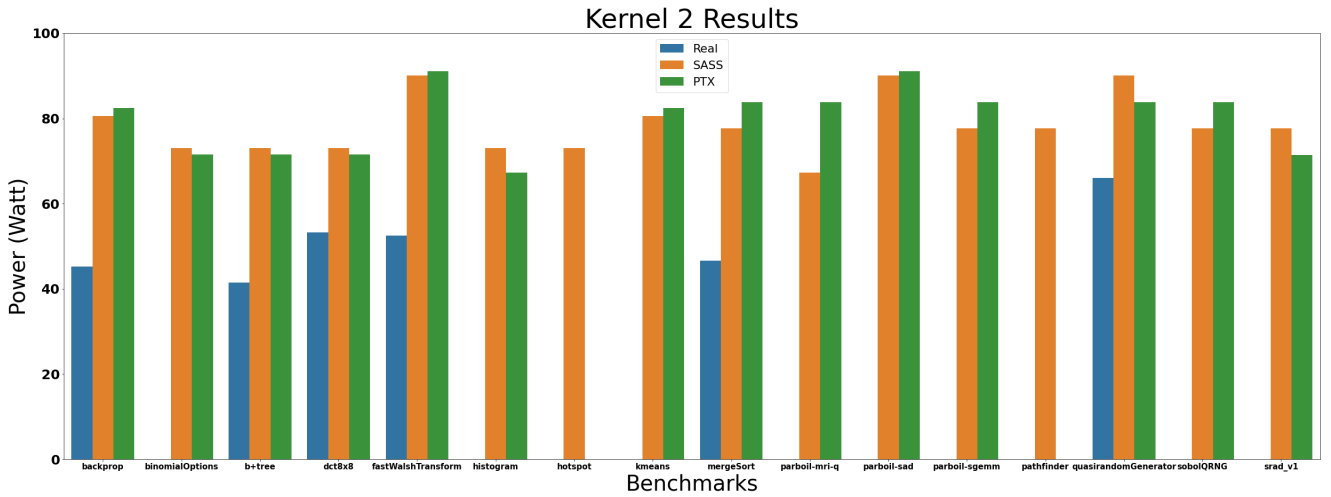


Figure 19. Real, PTX, SASS power results over kernel 2 of each benchmark

that were forced to compile or run which is not supported for Pascal since it doesn't have Tensor cores.

- We also faced problems with running some simulations of benchmarks that failed without much feedback from the runs. There are several causes that we thought prevented these benchmarks to run. According to AccelWattch's researchers, setup environment scripts were not fully successfully ran which causes CUDA's libraries not to be found by the benchmarks which we checked but didn't work. It might be our system is not fully set up correctly or maybe some programs are not compatible or outdated. There were instances where NVIDIA's software was not fully functioning which causes some CUDA-related errors.

References

- [1] Ali Bakhoda, George L. Yuan, Wilson W. L. Fung, Henry Wong, and Tor M. Aamodt. 2009. Analyzing CUDA workloads using a detailed GPU simulator. In *2009 IEEE International Symposium on Performance Analysis of Systems and Software*. 163–174. <https://doi.org/10.1109/ISPASS.2009.4919648>
- [2] GPGPU-Sim. 2009. PTXPlus. http://gpgpu-sim.org/manual/index.php/Main_Page#PTX_and_SASS
- [3] Jupyter. 2015. Jupyter Energy Monitor Bricklet. <https://jupyter.org/>
- [4] Vijay Kandiah, Scott Peverelle, Mahmoud Khairy, Junrui Pan, Amogh Manjunath, Timothy G. Rogers, Tor M. Aamodt, and Nikos Hardavellas. 2021. AccelWattch: A Power Modeling Framework for Modern GPUs. In *MICRO-54: 54th Annual IEEE/ACM International Symposium on Microarchitecture (Virtual Event, Greece) (MICRO '21)*. Association for Computing Machinery, New York, NY, USA, 738–753. <https://doi.org/10.1145/3466752.3480063>

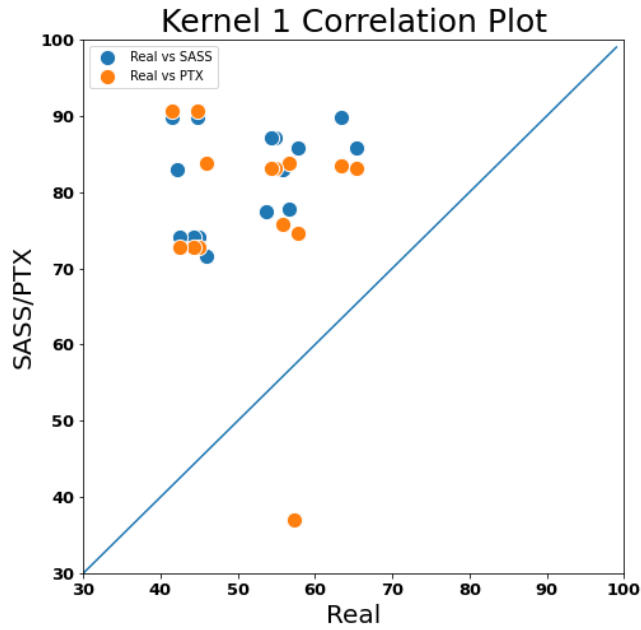


Figure 20. Correlation plot of kernel 1

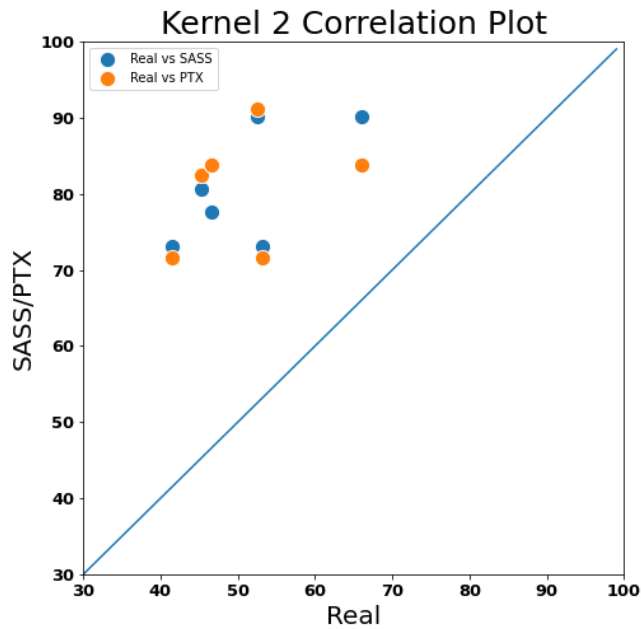


Figure 21. Correlation plot of kernel 1

- [7] NVIDIA. 2007. PTX and SASS Assembly Debugging. https://docs.nvidia.com/gameworks/content/developertools/desktop/ptx_sass_assembly_debugging.htm
- [8] SchedMD. 2003. Slurm, Workload Manager. <https://slurm.schedmd.com/>
- [9] TinkerForge. 2011. TinkerForge Energy Monitor Bricklet. <https://www.tinkerforge.com/en/shop/energy-monitor-bricklet.html>
- [10] Michael Waskom. 2012. seaborn: statistical data visualization. <https://seaborn.pydata.org/>

- [5] Mahmoud Khairy, Zhesheng Shen, Tor M. Aamodt, and Timothy G. Rogers. 2020. Accel-Sim: An Extensible Simulation Framework for Validated GPU Modeling. In *2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA)*. 473–486. <https://doi.org/10.1109/ISCA45697.2020.00047>
- [6] NVIDIA. 2007. Parallel Thread Execution ISA Version 8.1. <https://docs.nvidia.com/cuda/parallel-thread-execution/index.html>