


Iris Bench Report

GPU Benchmarking for Performance and NetZero at IRIS

1. Introduction

The IRIS GPU Bench project, led by the Science and Technology Facilities Council (STFC), aims to create an open-source GPU benchmarking suite tailored to the specific scientific needs of the IRIS community. Supporting a range of scientific fields from particle physics to astrophysics, IRIS seeks to optimize its computational resources and reduce its carbon footprint. The IRIS GPU Bench will address the gaps in current tools by providing performance insights and CO2 emission data relevant to scientific applications such as neutron spectroscopy, tomography, and machine learning for astrophysics, thus promoting both efficiency and sustainability.

2. Accomplishments Of Project So Far

- Code Repo with the necessary to deploy and run IRIS Bench: [IRIS GPU Bench Repo](#)
- Accompanying Documentation necessary to deploy, run, maintain, and develop IRIS Bench:  [IRIS GPU Bench](#)
 - The documentation contains the following pages an Overview, Installation Steps, Docker Image Management, Command-Line Interface Usage, Example Commands, Result Collection, Live Monitoring of Benchmarks, Accuracy Considerations, and Future Development Tasks. This report supports the Documentation, however, the documentation is worth reviewing as well.
- Accompanying Grafana Dashboards to visualise data produced by IRIS Bench: [Dashboard](#)
- Functional Scientific/User Benchmarks (Mantid/SciML) integrated with IRIS Bench

3. Evaluating Accuracy

3.1 Carbon Metrics

- **Data Source and Frequency:** Carbon data is sourced in real-time from the [National Grid ESO Regional Carbon Intensity API](#), with updates every approximately 30 minutes. The monitor averages values at the start and end of each period, which may limit accuracy for containers running longer than 30 minutes.
- **Variability:** Carbon forecasts fluctuate based on weather, time, and energy demand, affecting total emissions estimates. To provide a broader context, total energy can be multiplied by the UK's average carbon emission rate of [162 gCO2/kWh in 2023](#).

3.2 GPU Metrics

- **Source and Error Margin:** GPU metrics are obtained from `pynvml` and `nvidia-smi`, with a reported power draw error margin of $\pm 5\%$ (up to $\pm 30W$) ([source](#)).
- **Energy Calculation:** Energy use is estimated using trapezoidal integration of power readings. Smaller monitoring intervals improve accuracy but increase overhead, which may affect results.

3.3 Benchmark Results: Integration with Meerkat (HIGH PRIORITY)

To enhance reliability, the benchmarks should be integrated with **Meerkat** for periodic execution. This will allow consistent testing over time, enabling calculation of medians, standard deviations, and error bars added to the benchmark results, which can be calculated through **Grafana** dashboards. Regular testing helps identify trends, maintain performance consistency, and provide statistically significant results.

Being able to run the Benchmark Periodically will also show how other external factors may affect the performance of the GPUs, such as hardware location, weather temperature (ie heatwaves vs winter) throughout the year.

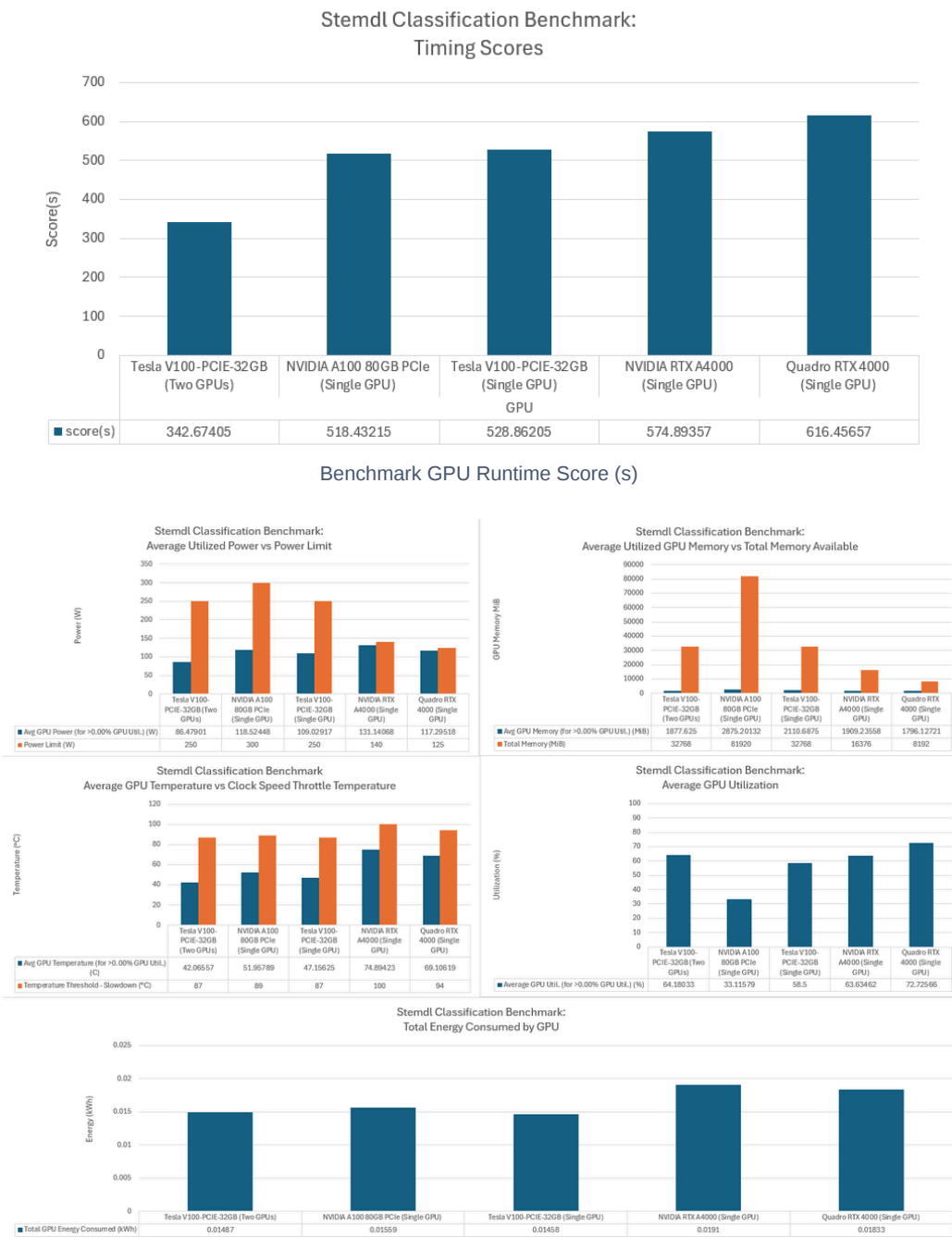
4. Results and Findings

4.1 Benchmark Results

The results below are from single runs of the benchmarks using IRIS Bench and containers. Hence, integration with Meerkat will allow for these tests to be repeated periodically (see 3.3) better validity.

Future work: review the results below to develop an understanding of GPU performance for various IRIS Workloads.

4.1.1 Sciml Bench: Stemdl Classification (Multi GPU Benchmark)



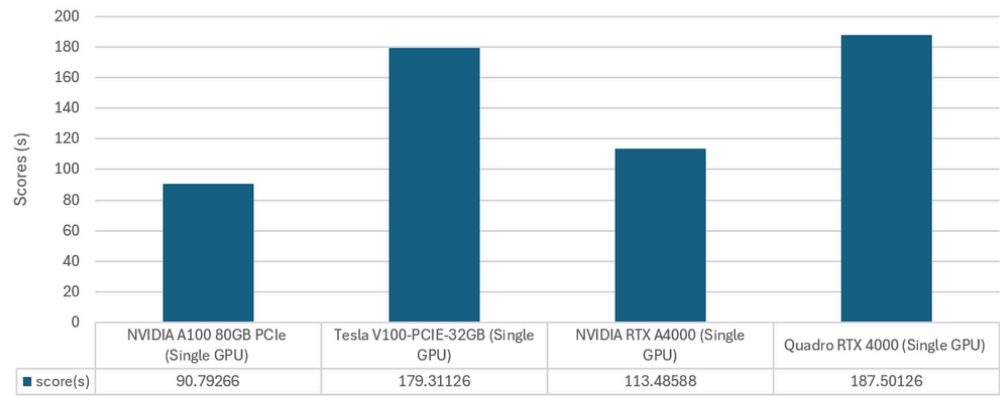
4.1.2 Sciml Bench: Synthetic Regression (Single GPU Benchmark)



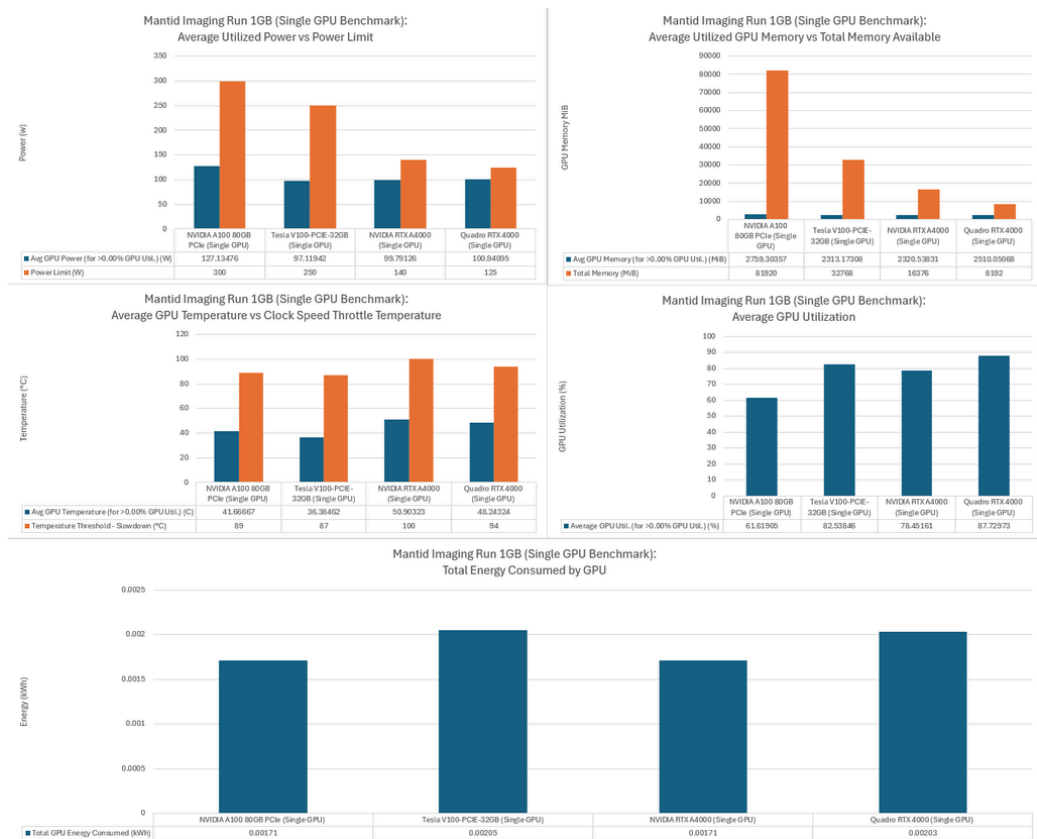
More Details on Performance

4.1.3 MANTID: Mantid Imaging Reconstruction Benchmark 1GB (Single GPU Benchmark)

Mantid Imaging Run 1GB (Single GPU Benchmark):
Timing Scores



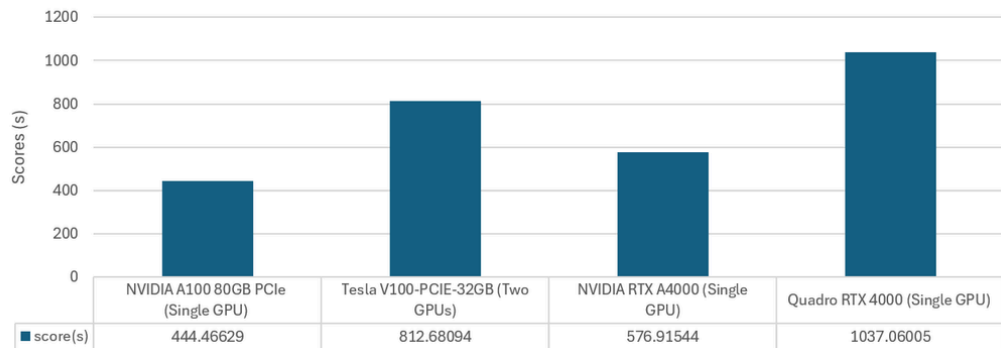
Benchmark GPU Runtime Score (s)



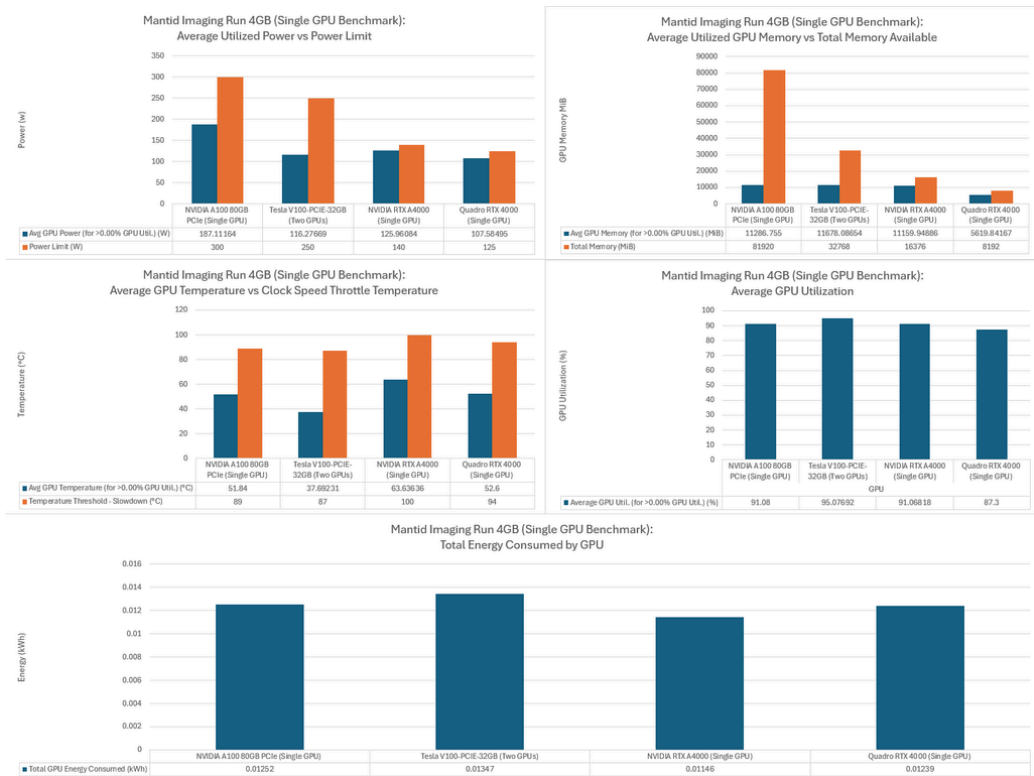
More Details on Performance

4.1.4 MANTID: Mantid Imaging Reconstruction Benchmark 4GB (Single GPU Benchmark)

Mantid Imaging Run 4GB (Single GPU Benchmark):
Timing Scores



Benchmark GPU Runtime Score (s)



More Details on Performance

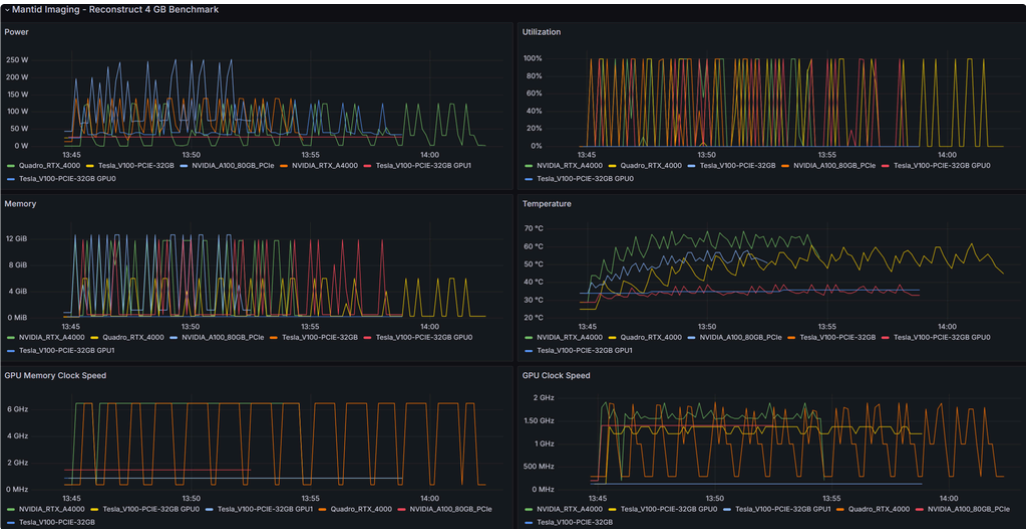
For this benchmark, a single V100 VM lacked sufficient memory, so the results are from a VM with two V100 GPUs. Consequently, the energy consumption for the dual-GPU setup is slightly higher than that of a single GPU due to the additional idle power.

4.2 Grafana Results Example

Example of Grafana Results for Mantid Imaging Benchmark Run with IRIS GPU Bench on a range of GPU.
From 2024-09-11 13:43:16 to 2024-09-11 14:03:04



Collected Carbon Forecast For South Of England During Benchmark Runs



Benchmark Run GPU Metrics

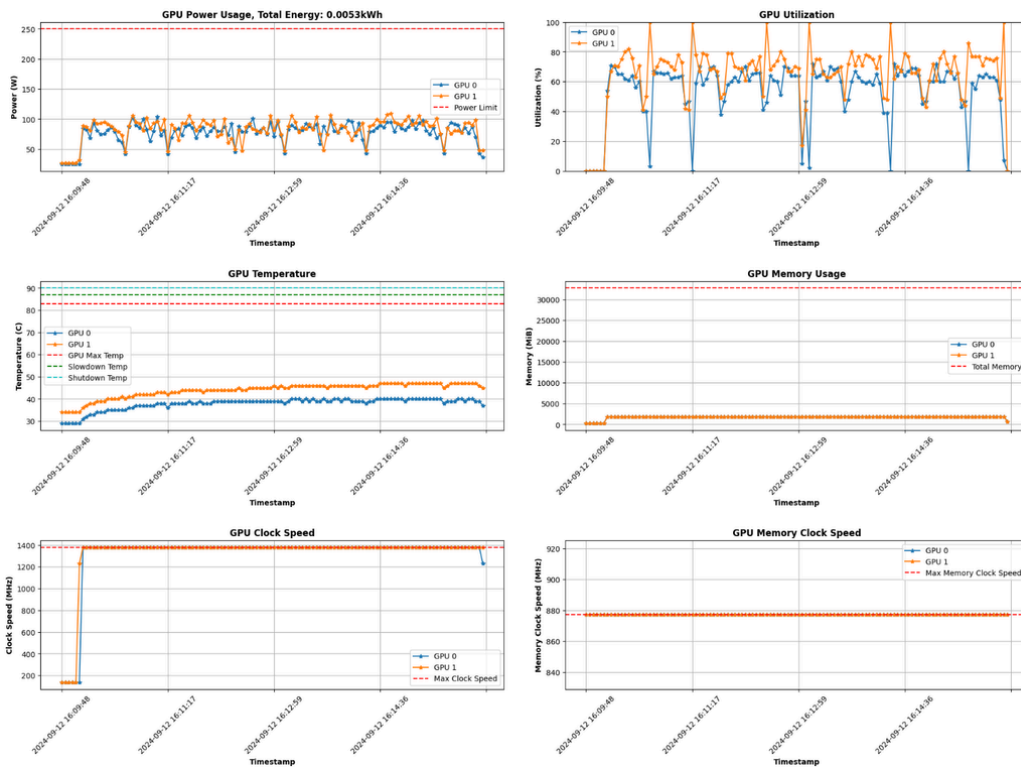
| Benchmark Results | | | | | | | |
|---------------------|-----------------------|----------------|---------------------|---------------------------------|-----------------------------------|--------------------------------|--------------|
| Time | GPU Type | Number of GPUs | Benchmark Score (s) | Total GPU Energy Consumed (kWh) | Total GPU Carbon Emissions (gCO2) | Avg Carbon Forecast (gCO2/kWh) | Avg GPU Memo |
| 2024-09-11 14:02:30 | Quadro_RTX_A4000 | 1 | 1037.06005 | 0.01239 | 1.1706 | 94.50000 | |
| 2024-09-11 13:45:00 | Tesla_V100-PCIE-32GB | 1 | | 0.00003 | 0.00340 | 98.00000 | |
| 2024-09-11 13:59:00 | Tesla_V100-PCIE-32GB | 2 | 812.68084 | 0.01347 | 1.32051 | 98.00000 | |
| 2024-09-11 13:52:30 | NVIDIA_A100-80GB-PCle | 1 | 444.46629 | 0.01252 | 1.22660 | 98.00000 | |
| 2024-09-11 13:54:50 | NVIDIA_RTX_A4000 | 1 | 576.91544 | 0.01146 | 1.12322 | 98.00000 | |

Grafana Table Containing Final Benchmark Results

4.3 IRIS Bench Results

When IRIS Bench is runs a benchmark on a VM the following results are generated locally (see more on this in the documentation). There are results for Benchmark Score, Energy Consumed, Carbon Emissions, Averaged GPU Metrics as well as a png figure providing more detail about the benchmark than is collected by Grafana.

The locally saved timeseries plot includes shutdown and slowdown temperatures, allowing us to better understand any performance drops if temperatures exceed these thresholds. The plot also features horizontal lines indicating total memory usage (though memory utilization might be a more relevant metric), power limits, maximum clock speed, and maximum memory clock speed. Additionally, while Grafana's minimum scrape interval is 10 seconds, shorter intervals can be selected for these locally saved plots to capture sharper results, such as when utilization fluctuates from 100% to 0% multiple times within a 10-second period.



timeseries_plot.png from stemdl_classification with two V100 GPUs

GPU and Carbon Performance Results

| Metric | Value |
|-----------------------------------|-----------------------|
| Benchmark: | stemdl_classification |
| Benchmark Score (s) | 1053.29596 |
| Elapsed Monitor Time (s) | 1202.78311 |
| Total GPU Energy Consumed (kWh) | 0.00561 |
| Total GPU Carbon Emissions (gCO2) | 0.54695 |

Carbon Information

| Metric | Value |
|------------------------------------|---------------------|
| Average Carbon Forecast (gCO2/kWh) | 97.5 |
| Carbon Forecast Start Time | 2024-09-11 14:25:55 |
| Carbon Forecast End Time | 2024-09-11 14:45:58 |

GPU Information

| Metric | Value |
|---|------------------------------------|
| GPU Type | Tesla V100-PCIE-32GB |
| No. of GPUs | 2 |
| Average GPU Utilization (for >0.00% GPU Util.) (%) | 43.03604 |
| Average GPU Power (for >0.00% GPU Util.) (W) | 65.31081 (Power Limit: 250) |
| Average GPU Temperature (for >0.00% GPU Util.) (°C) | 38.58559 |
| Temperature Threshold - Slowdown (°C) | 87.00 |
| Average GPU Memory (for >0.00% GPU Util.) (MiB) | 1877.58896 (Total Memory: 32768.0) |
| Average Clock Speed (MHz) | 1340.87 (Max: 1380.00) |
| Average Memory Clock Speed (MHz) | 877.00 (Max: 877.00) |

formatted_metrics.txt

5. Future Work and Expansion

5.1 Additional Benchmarks to Integrate

- McStas Neutron Monte Carlo Ray-Tracing
- **RFI GPU Benchmarks:** Explore and integrate relevant Radio Frequency Interference (RFI) GPU benchmarks to cover a wider range of scientific performance evaluations.
- **Highly Optimized NVIDIA HPL Benchmarks:** Add the [NVIDIA HPC Benchmarks](#), focusing on the HPL benchmarks designed for Grace Hopper tests. This integration will require referencing the [NVIDIA HPL Benchmark documentation](#) to ensure proper implementation and optimization.

5.2 Enhancing Benchmark Results

- **Integration into Meerkat (HIGH PRIORITY):** Complete the integration with Meerkat to facilitate regular benchmark testing. This will allow the generation of more robust statistical data, such as error bars, by calculating mean and standard deviation across multiple tests.
- **Normalization of Results:** Implement normalization techniques to standardize benchmark results across different GPU models and workloads, facilitating meaningful comparisons.
- **FLOP Estimations and Efficiency Metrics:** Calculate Floating Point Operations per Second (FLOPs) to determine performance per watt, providing insights into computational efficiency and energy consumption.

5.2.1 Estimating Floating Point Operations Per Second (FLOPs)

Comparisons based solely on clock speed can be misleading due to differences in GPU architectures, CUDA cores, tensor cores, and operations per clock cycle. Note clock speeds vary for various reason such as size of workload, temperature, power supply.

To estimate the FLOPs of a GPU:

`FLOPs = CUDA Cores × Clock Speed (Hz) × Operations per Clock Cycle`

Steps to calculate FLOPs:

- Determine the number and type of cores (e.g., CUDA cores, tensor cores) and if they are being utilized.
- Clock speed can be found using IRIS BENCH.
- Identify the operations per clock cycle, specific to the GPU's architecture (ie core generation). ["The A100 SM includes new third-generation Tensor Cores that each perform 256 FP16/FP32 FMA operations per clock"](#)
- Investigate how `nvidia-smi` calculates utilization and whether more granular utilization data can be obtained.
- Review how tools like `sciml-bench` and `pytorch` report core usage and incorporate these insights into metrics collection.

| Data Sheets Contain Information About the Number of Cores |
|---|
| V100 datasheet |
| A4000 Datasheet |
| <input type="checkbox"/> NVIDIA A100 |
| RTX4000 DataSheet |

5.3 Explore NVIDIA Nsight Systems

- Utilize NVIDIA Nsight Systems for detailed performance profiling, identifying optimization opportunities, and potentially getting insight into the activated cores for FLOP calculations. This tool offers in-depth insights into GPU performance across various workloads and configurations.

5.4 Additional Metrics and Areas for Measurement in IRIS Bench

- **Utilization Time:** Measure the total time the GPU is actively utilized, which can provide insights into idle periods and workload efficiency.
- **Stability: Crash Frequency:** Track and report any GPU crashes or visual artifacts during benchmarks to assess stability.
- **Throttling Events:** Monitor instances of clock speed reductions due to high temperatures or power constraints.
- **Memory Bandwidth:** Measure the data transfer rate between the GPU and system memory to identify potential bottlenecks and optimize performance

5.5 Other Ideas for Future Implementation

- **Consistent Hardware Configurations:** Ensure that all GPUs being tested use the same hardware configurations (e.g., memory, CPUs) to eliminate variability and produce consistent results.

- **Continuous Integration for Performance Testing:** Encourage IRIS users to integrate GPU benchmarks into their CI workflows. Implement automated performance tests on every pull request; if performance drops by a specified percentage, the pull request would fail, ensuring that code changes do not degrade performance.
- **Experimenting with Precision to Utilize Tensor Cores fully:** For GPUs equipped with tensor cores, utilize lower precisions (e.g., FP16) for matrix operations where feasible. This can lead to significant performance gains, depending on the workload's precision requirements.

5.6 Improvements for GitHub Repository

- **Continuous Integration (CI) Tests:** Develop and integrate comprehensive CI tests using GitHub Actions to maintain code reliability, ensure consistent performance, and catch issues early in the development cycle.
- **Carbon Index Calculation:** Enhance the environmental impact analysis by calculating the carbon index throughout the entire benchmarking run, rather than at just the start and end, to provide a more accurate representation.