

Predicting the performance of a kernel

ANN XIAO,^{1,*} XINYU LI,^{2,*} AND YIDAN QIN^{3,*}

^{*}New York University, Courant

¹yx3493@nyu.edu

²xl5280@nyu.edu

³yq2421@nyu.edu

Abstract: Optimizing the performance of CUDA kernels is crucial for improving the efficiency of GPU-based applications, particularly for memory-bound operations. An accurate simulation method to predict the kernel speedup can save significant computational resources and lead the way for performance optimization. In this work, we propose a machine learning-based framework that utilizes static profiling metrics, such as the Compute-to-Global-Memory-Access (CGMA) ratio and shared memory utilization, combined with static analysis to predict kernel performance. With our Random Forest regression model, We explain over 92% of the variance in kernel speedup, achieving a Mean Absolute Percentage Error (MAPE) of 55.15%. Feature importance analysis highlights CGMA as the most impactful metric for predicting performance, while metrics like occupancy showed minimal influence. Our findings highlight the practicality and scalability of machine learning for GPU performance prediction, offering insights for kernel optimization.

1. Introduction

CUDA kernel performance prediction is a critical area in GPU computing, enabling developers to optimize workloads efficiently without incurring significant computational overhead. Accurate performance estimation is essential for reducing the iterative profiling and tuning cycles required for kernel optimization, especially for memory-bound operations that dominate many real-world applications. Existing methods, such as static analysis [1], dynamic profiling [2], and machine learning-based approaches [3], have addressed challenges such as warp divergence, memory-computation overlap, and execution time prediction. While static analysis excels at providing worst-case bounds, it often lacks the flexibility to adapt to diverse workloads. Dynamic profiling, on the other hand, can provide accurate runtime information but requires kernel execution, leading to significant overhead. Machine learning approaches, while accurate, often rely heavily on runtime profiling data and large, diverse training datasets, which can limit their scalability and generalizability.

In this project, we address the challenging task of predicting the speedup of memory-bound CUDA kernels relative to single-thread performance, without requiring kernel execution. This problem is particularly important because it provides a practical means to assess kernel performance in scenarios where running the kernel on various configurations is computationally prohibitive. Our approach utilizes a combination of static analysis and runtime profiling to identify critical metrics that influence kernel performance. Specifically, we designed 11 distinct memory-bound kernels to capture diverse performance patterns, including scenarios involving shared memory usage and global memory access.

From these kernels, we identified seven critical metrics—problem size, grid configuration, occupancy, compute-to-global memory access ratio (CGMA), global memory access reduction per thread, branch divergence, and percentage of shared memory utilization—that are instrumental in determining performance. These metrics were used to create a dataset of nearly 700 samples, collected through experimentation on the CIMS CUDA2 GPU platform with an NVIDIA GeForce RTX 2080 Ti. This dataset was then used to train a Random Forest regression model, chosen for its robustness and ability to handle nonlinear interactions between features.

Our proposed approach combines the efficiency of static analysis with the adaptability of

48 machine learning to predict kernel speedup in memory-bound scenarios. By focusing on shared
49 memory usage, global memory access, and critical performance metrics, our framework provides
50 a scalable and practical solution for performance prediction, significantly reducing the reliance
51 on runtime profiling. Furthermore, the results demonstrate the effectiveness of our methodology
52 in delivering accurate predictions while providing interpretable insights into the factors that
53 impact kernel performance.

54 **2. Literature Survey**

55 CUDA kernel performance prediction is a critical area of research, with significant implications
56 for optimizing GPU workloads. Existing approaches to address this problem can broadly be
57 categorized into static analysis, dynamic profiling, and machine learning-based methods, each
58 offering unique advantages and limitations.

59 *2.1. Static Analysis Methods*

60 Static analysis methods aim to provide worst-case or upper-bound estimates of kernel performance
61 without executing the kernel. Recent advances include methods like the Many-BSP analytical
62 model, which achieves a 12.33% prediction error through static analysis [4], compute-memory
63 bound performance estimation across devices [5], and static analysis-guided optimization
64 frameworks that predict execution costs for informed kernel improvements [6].

65 Muller et al. [1] utilize a static analysis framework to identify performance bottlenecks
66 such as warp divergence and memory access conflicts. While effective at estimating resource
67 requirements, this method lacks predictive modeling capabilities, limiting its utility for nuanced
68 performance forecasting.

69 Other studies have explored static analysis models for execution time prediction. For instance,
70 one such model uses PTX code and GPU scheduling simulation to achieve accurate predictions
71 for compute-intensive kernels, though it struggles with memory-bound workloads due to limited
72 memory hierarchy modeling [7]. Hong and Kim [8] propose an analytical model that incorporates
73 memory-level and thread-level parallelism to estimate kernel performance. This approach is
74 valuable for identifying bottlenecks but struggles with accurately modeling memory-bound
75 workloads due to simplified assumptions about memory hierarchies. Similarly, the Many-BSP
76 model [4] achieves a prediction error of 12.33% using static analysis but requires manual tuning
77 for diverse workloads.

78 *2.2. Dynamic profiling Approaches*

79 Dynamic profiling methods utilize runtime information to analyze kernel behavior. GPUPerf,
80 introduced by Jaewoong et al. [2], combines static and dynamic profiling to identify bottlenecks
81 and assess memory-computation overlap. While the improved MWP-CWP model provides
82 interpretable metrics for optimization, its reliance on runtime profiling introduces overhead and
83 limits its applicability to real-time scenarios.

84 Other dynamic profiling frameworks, such as those proposed by Wang et al. [9], use runtime
85 data to balance performance and energy efficiency through techniques like dynamic voltage and
86 frequency scaling (DVFS). However, these methods require extensive profiling data, which can
87 be computationally expensive to collect for large-scale applications.

88 *2.3. Machine Learning Approaches*

89 Machine learning-based methods have gained popularity for their ability to model complex,
90 non-linear relationships in GPU performance. Tiwari et al. [3] utilize supervised learning to
91 predict execution time, leveraging runtime profiling data to train the model. While effective, these
92 methods often require large, diverse datasets and may struggle to generalize across architectures.

Hybrid approaches that combine analytical modeling with machine learning, such as those by Meng et al. [10], have demonstrated improved prediction accuracy for diverse workloads. These methods incorporate architectural details into the model but require careful feature selection and tuning, increasing development complexity.

2.4. Comparison and Limitations

The existing literature demonstrates significant progress in CUDA kernel performance prediction but highlights trade-offs between prediction accuracy, computational overhead, and generalizability. Static analysis excels at providing low-overhead, architecture-independent estimates but lacks adaptability to diverse workloads. Dynamic profiling delivers high accuracy but incurs runtime costs and is not scalable for real-time prediction. Machine learning methods offer flexibility and accuracy but depend on large, representative datasets and careful feature engineering.

2.5. Motivation for This Work

Despite these advances, a gap remains in developing a framework that balances efficiency, scalability, and accuracy for memory-bound kernels. Our work addresses this gap by integrating static analysis and machine learning, employing critical metrics such as CGMA, branch divergence, and shared memory utilization. By focusing on memory-bound workloads and minimizing profiling requirements, we aim to offer a scalable and practical solution that bridges the strengths of existing methods.

3. Proposed Idea

Our approach combines GPU-friendly kernel design with critical performance metrics to develop a predictive model for CUDA kernel speedup. We designed 11 kernels, identified 7 key metrics, and used a Random Forest model for accurate regression. Below are the details for microbenchmarks, metric definitions, and machine learning implementation.

3.1. Microbenchmarks

We implemented 11 kernels as microbenchmarks to simulate different GPU workloads, representing computational and memory access patterns. They are extracted from various real-time scenarios and are all GPU-friendly. These benchmarks demonstrate various challenges and opportunities in GPU programming.

Elementwise-Add and Elementwise-Multiply These 2 kernels are the fundamental parts in many algorithms, they are easy to implement and with scalability to add more blocks to fit for larger problem size. Elementwise multiply is a nature progression from elementwise addition, emphasizing the arithmetic operations.

Heat Diffusion/Distribution Modeling scientific simulations, this kernel evaluates memory-bound performance by computing values based on neighboring elements. It highlights the challenges of high memory access and computational dependencies.

Count Number The histogram computation involves counting the occurrences of each element in a given dataset. With the implementation on GPU, the access pattern of parallel updates to the output array can lead to a large ratio of memory conflicts, so we used atomicAdd to ensure correctness updates to global memory. While this eliminates race conditions, it introduces fully sequential updates when accessing the global memory.

133 **Count numbers with shared memory** To optimize the count number application, we utilized
134 shared memory as an intermediate buffer to reduce global memory latency. Each thread updates
135 shared memory using `atomicAdd`, and the results are then aggregated into global memory.
136 Although there is also a partially sequential access pattern, it significantly reduced contention
137 and improved performance compared to the previous implementation.

138 **2D Convolutional Layer** As a fundamental component of convolutional neural networks
139 (CNNs), 2D convolutional layer is widely used in tasks like image recognition, segmentation, and
140 object detection. This kernel calculates each output element as a weighted sum of overlapping
141 input elements covered by a filter. While the current implementation accesses global memory
142 directly, introducing redundant reads, future optimizations could use shared memory to cache
143 overlapping regions, reducing latency. Despite lacking shared memory optimization, the kernel
144 remains computationally demanding, making it an effective GPU performance benchmark.

145 **Activation Layer (Sigmoid Function)** The sigmoid function introduces non-linearity in
146 neural networks, making it compute-intensive. This kernel demonstrates efficient thread-level
147 parallelism while maintaining high computational demands.

148 **Maximum Distance** This kernel is inspired by real-world applications like computational
149 geometry and spatial analysis in urban planning. It calculates the maximum distance from each
150 point in the input array to predefined "landmark" positions. Each thread independently calculates
151 and updates the maximum distance without requiring synchronization. Parallelism is achieved
152 using grid-stride loops.

153 **Array Reversal** Simplified from matrix transposition, this kernel performs in-place 1D array
154 reversal. Each thread processes elements independently, ensuring conflict-free memory updates
155 and scalability.

156 **Softmax Normalization** Softmax normalization converts raw scores into probabilities. The
157 implementation uses shared memory for intermediate sums, reducing global memory latency and
158 improving performance for larger datasets.

159 **Binary Search** This kernel employs a fundamental algorithm with $O(\log n)$ complexity to
160 efficiently locate a target value in a sorted array. Our implementation dynamically allocated
161 shared memory for intermediate results. Each thread in the multi-threaded setup searches a
162 segment of the array independently, followed by a reduction step to consolidate results.

163 3.2. Metrics

164 We identified 7 key metrics that capture performance characteristics, enabling comprehensive un-
165 derstanding and accurate prediction. Below, we provide calculation formulas and implementation
166 details for each metric.

167 3.2.1. Kernel Input Parameters

168 Three Parameters:

169 **Problem Size.** 10000, 50000, 100000, 500000, 1000000, 5000000, 10000000

170 **Number of Blocks.** Calculated based on specific problem size and number of threads per block

171 **Number of Threads per Block.** 32, 64, 128, 200, 256, 400, 500, 512, 1024

172 3.2.2. Maximum Percentage of Blocks

173 **Definition:** Reflects the fraction of total blocks that can run concurrently, measuring GPU
174 resource utilization.

175 **Calculation Formula.**

$$\text{Max Percentage of Blocks} = \frac{\text{Concurrent Blocks}}{\text{Total Blocks}} \times 100$$

176 Where:

$$\text{Concurrent Blocks} = \text{Blocks per SM} \times \text{Number of SMs}$$

177

$$\text{Blocks per SM} = \min \left(\frac{\text{Total Registers}}{\text{Registers per Block}}, \frac{\text{Total Shared Memory}}{\text{Shared Memory per Block}}, \frac{\text{Maximum Threads per SM}}{\text{Threads per Block}} \right)$$

178 **Implementation.** The custom script `maxBlock_occupancy.py` automates this calculation, processing
179 kernel configurations and GPU constraints.

180 3.2.3. Occupancy

181 **Definition:** Measures GPU efficiency by calculating the fraction of active warps per SM.

182 **Calculation Formula.**

$$\text{Occupancy} = \frac{\text{Active Warps per SM}}{\text{Maximum Warps per SM}}$$

183 Where:

$$\text{Active Warps per SM} = \text{Blocks per SM} \times \text{Warps per Block}$$

184

$$\text{Blocks per SM} = \min \left(\frac{\text{Total Registers}}{\text{Registers per Block}}, \frac{\text{Total Shared Memory}}{\text{Shared Memory per Block}}, \frac{\text{Maximum Threads per SM}}{\text{Threads per Block}} \right)$$

185 **Implementation.** The script `maxBlock_occupancy.py` computes occupancy values based on kernel
186 configurations and hardware constraints.

187 3.2.4. Branch Divergence

188 **Definition:** Quantifies serialization within warps caused by divergent thread execution paths.

189 **Calculation Formula.**

$$\text{Divergence Percentage} = \frac{\text{Divergent Executions}}{\text{Total Executions}} \times 100$$

190 Where:

- 191 • **Divergent Executions:** Difference in iterations among threads in the same warp.
192 • **Total Executions:** Total iterations executed by all threads in the warp.

193 **Implementation.** The script `divergence.py` simulates warp execution, identifying divergence by
194 calculating maximum and minimum iterations within each warp.

195 3.2.5. CGMA (Compute-to-Global Memory Access)

196 **Definition:** Measures computational intensity relative to memory access.

197 **Calculation Formula.**

$$\text{CGMA} = \frac{\text{Number of Computations}}{\text{Number of Global Memory Accesses}}$$

198 **Implementation.** The script `autoprofiler.py` uses NVIDIA Nsight Compute (`ncu`) to extract
199 metrics like `inst_executed` and `dram_bytes`.

200 3.2.6. Shared Memory Utilization

201 **Definition:** Evaluates shared memory usage and its impact on performance.

202 **Calculation Formula.**

$$\text{Shared Memory Utilization} = \frac{1024 \times \text{Number of Blocks}}{\text{Problem Size}}$$

203 **Implementation.** Metrics were calculated using kernel input parameters for each experiment.

204 3.2.7. Global Memory Access

205 **Definition:** Measures efficiency of global memory usage under multi-threaded execution.

206 **Calculation Formula.**

$$\text{Reduction} = (\text{dram_read} + \text{dram_write})_{\text{one_thread}} - \frac{(\text{dram_read} + \text{dram_write})_{\text{multi_thread}}}{\text{Number of Blocks} \times \text{Threads per Block}}$$

207 **Implementation.** Metrics `dram__bytes_read.sum` and `dram__bytes_write.sum` were col-
208 lected using `ncu` for different execution scenarios.

209 3.3. Machine Learning Model

210 3.3.1. Model Selection

211 Our task is to address a regression problem by constructing a predictive model that maps the
212 selected metrics to the corresponding speedup values and also provides some insights for the
213 inner relationships across different metrics. Employing a machine learning model removes the
214 need for manual formula derivation and enables efficient and accurate performance estimation.

215 Some simple regression model, such as linear and ridge regression model assume linear
216 dependencies and treat metrics independently, making them less suitable for capturing the
217 complexity of CUDA kernel performance.

218 We employed Random Forest for its ability to capture non-linear dependencies and inter-metric
219 relationships. Hyperparameter tuning was performed using Optuna, optimizing for Mean
220 Absolute Percentage Error (MAPE).

221 3.3.2. Implementation

222 The Random Forest model was trained on a dataset of 639 samples, extracted from kernel
223 executions on the NVIDIA GeForce RTX 2080 Ti. The final model achieved high accuracy, with
224 metrics indicating strong performance prediction capabilities.

225 4. Experimental Setup

226 All experiments were conducted on a single GPU environment.

227 4.1. Microbenchmarks

228 The microbenchmarks were grouped in a single `run.cu` file, which contains the entry of the
229 comparison functions and definitions of all kernels. To reproduce the experiments, follow the
230 steps outlined in `readme.md`, where detailed instructions are provided for compiling the source
231 code, running the executable files, and collecting results.

232 4.2. Data Collection

233 Speedup values for each kernel configuration were directly output during execution. Additional
234 metrics, such as memory access patterns and performance profiles, were collected using NVIDIA
235 Nsight Compute (`ncu`) and custom scripts. These metrics were recorded across varying problem
236 sizes and grid configurations to build a comprehensive dataset. Manual validation was conducted
237 to ensure accuracy in the reported results.

238 4.3. Scripts and Tools

239 Several script files were written to automate metric collection, reducing human error and ensuring
240 consistency. Table 1 summarizes the methods and tools used to calculate each performance
241 metric.

Table 1. Methods for Calculating Metrics

Metric	Calculation Method
Occupancy	Computed with the command <code>python3 maxBlock_occupancy.py</code> .
Branch Divergence	Simulated and calculated with the command <code>python3 divergence.py</code> .
CGMA (Compute-to-Global Memory Access)	Computed with the command <code>python3 autoprofiler.py</code> .
Shared Memory Utilization	Calculated using kernel input parameters for each experiment.
Maximum Percentage of Blocks	Derived using <code>maxBlock_occupancy.py</code> with appropriate kernel configurations.
Global Memory Access	Retrieved via <code>ncu</code> . Data collected and calculated manually.
Problem Size, Number of Blocks and Number of Threads Per Block	Defined explicitly in the kernel inputs to cover varying thread and block configurations.

242 4.4. Hardware Specifications

- 243 • **GPU Model:** NVIDIA GeForce RTX 2080 Ti
- 244 • **CUDA Version:** 12.4
- 245 • **System Environment:** CIMS CUDA2 Server

246 5. Results and Analysis

247 5.1. Measures of Success

248 The primary measures of success for this project are:

- 249 • **Prediction Accuracy:** Measured using Mean Absolute Percentage Error (MAPE), Mean
250 Absolute Error (MAE), Root Mean Squared Error (RMSE), and R^2 scores. A lower MAPE
251 and higher R^2 indicate better predictive performance.
- 252 • **Feature Interpretability:** Understanding which features contribute most to kernel speedup
253 prediction using feature importance analysis.
- 254 • **Scalability:** The ability of the model to handle diverse kernels with varying configurations
255 and problem sizes.

256 We expect the model to achieve a MAPE of less than 60% and an R^2 score above 0.90,
257 indicating a high degree of prediction accuracy.

258 5.2. Experimental Procedure

259 The experiments were conducted on an NVIDIA GeForce RTX 2080 Ti GPU. The collected
260 dataset of nearly 700 samples was then preprocessed by:

- 261 • Removing rows with missing values.
- 262 • Transforming numerical columns with non-standard formats (e.g., commas) into usable
263 numeric data.
- 264 • Splitting the data into an 80:20 training and testing set.

265 A RandomForestRegressor model was trained using Optuna for hyperparameter optimization.
266 50 trials were conducted. The evaluation metrics were calculated on the test set for the
267 best-performing model.

268 5.3. Results and Key Metrics

269 Table 2 summarizes the key metrics achieved by the best model.

Table 2. Performance Metrics for the Best Model

Metric	Value
Mean Absolute Percentage Error (MAPE)	55.15%
Mean Absolute Error (MAE)	442.78
Root Mean Squared Error (RMSE)	3767.63
R^2 Score	0.9287
Explained Variance Score (EVS)	0.9291

270 The results demonstrate that the model achieves high accuracy, with R^2 and EVS values
271 exceeding 0.92, indicating that over 92% of the variance in kernel speedup is explained by the
272 model.

273 5.4. Feature Importance

274 Feature importance analysis reveals which metrics contribute most to the model's predictions.
275 Table 3 lists the importance of each feature.

276 As shown in Table 3, we can conclude:

- 277 • **CGMA Dominates:** CGMA has the highest importance (0.908) as it directly impacts
278 performance in memory-bound kernels by balancing computation and global memory
279 access.
- 280 • **Low Occupancy Impact:** Occupancy has minimal importance (0.0005) because memory
281 latency, not thread scheduling, dominates performance in memory-bound scenarios.
- 282 • **Other Features:** Metrics like shared memory utilization and branch divergence show
283 limited influence, reflecting their minor role in the tested kernels.

Table 3. Feature Importance for Kernel Speedup Prediction

Feature	Importance
CGMA	0.908
Global Memory Access	0.0487
Maximum Percentage of Blocks	0.0183
Num of Bblocks	0.0143
Shared Memory Utilization	0.0054
Problem Size	0.0030
Threads per Block	0.0024
Branch Divergence	0.0010
Occupancy	0.0005

284 5.5. Visual Analysis

285 Figure 1 visualizes the feature importance, highlighting that CGMA is the most influential metric for
286 predicting kernel speedup, while other metrics like occupancy and branch_divergence
287 have minimal impact.

288 Additionally, the visualization of the trained Random Forest model is presented in Appendix
289 A.

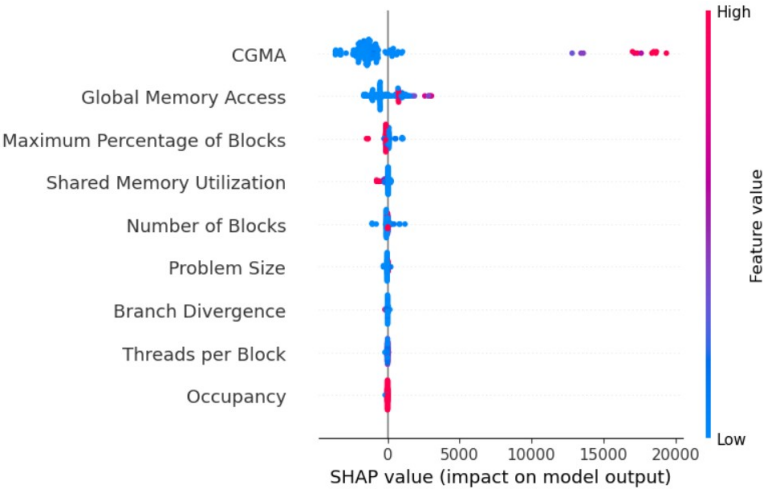


Fig. 1. Feature Importance Visualization

290 5.6. Discussion

- 291 • **Strengths:** The high R^2 score and low MAE demonstrate that the model accurately
292 predicts kernel speedup. The feature importance analysis provides interpretable insights
293 into which metrics significantly influence performance.
- 294 • **Weaknesses:** The MAPE value of 55.15% suggests that prediction errors are higher for

295 some kernels, particularly those with extreme configurations. This could be attributed to
296 insufficient representation of edge cases in the dataset.

297 Overall, the model performs well for the majority of configurations and highlights opportunities
298 for improving dataset diversity to handle outliers effectively.

299 To further enhance the results, future work could involve:

- 300 • Testing on additional datasets and architectures.
- 301 • Incorporating advanced models like Gradient Boosting or Neural Networks for comparison.
- 302 • Exploring additional metrics or transformations to improve accuracy.

303 6. Conclusion

304 In this project, we developed a machine learning-based approach to predict the speedup of
305 different CUDA kernels using a set of microbenchmarks and performance metrics. By combining
306 insights from static analysis with runtime profiling, we demonstrated the ability to model and
307 predict kernel performance without requiring actual kernel execution. Our approach takes
308 advantage of the flexibility of machine learning, providing a scalable and efficient solution for
309 performance prediction.

310 The key findings of this project include:

- 311 • **Accurate Performance Prediction:** The Random Forest model achieved a high R^2 score
312 of 0.9287 and an explained variance score of 0.9291, successfully explaining over 92% of
313 the variance in CUDA kernel speedup. The Mean Absolute Percentage Error (MAPE) of
314 55.15% further demonstrates its reliability across diverse kernels.
- 315 • **Importance of CGMA:** The Compute-to-Global-Memory-Access (CGMA) ratio emerged
316 as the most influential feature, underscoring its critical role in kernel performance. Other
317 metrics, like memory reduction per thread, also contributed meaningfully.
- 318 • **Model Limitations:** Edge cases with high MAE highlighted areas where dataset diversity
319 and feature representation could be improved, particularly for metrics like occupancy and
320 branch divergence, which showed minimal impact.

321 Despite the model's strengths of robust predictions and interpretable feature importance metrics,
322 its limitations also suggest opportunities for further improvement. The relatively high MAE for
323 certain kernels may reflect difficulties in modeling irregular performance patterns or inherent
324 variability in kernel configurations. Additionally, the current dataset may not comprehensively
325 represent all possible configurations or GPU architectures, potentially limiting the model's
326 generalizability.

327 Future work could address these limitations by:

- 328 • Expanding the dataset with diverse kernel configurations and testing across different GPU
329 architectures to improve generalization.
- 330 • Refining profiling techniques for metrics like branch divergence and shared memory
331 utilization to better represent their impact on performance.
- 332 • Including additional GPU-specific features, such as warp scheduling and memory band-
333 width, to enhance predictive accuracy.

334 In conclusion, this project provides a practical and scalable framework for CUDA kernel
335 performance prediction, reducing the reliance on runtime profiling. The results demonstrate the
336 model's potential for optimizing and understanding GPU workloads, while its limitations outline
337 clear directions for future advancements.

References

1. S. K. Muller and J. Hoffmann, "Modeling and analyzing evaluation cost of cuda kernels," *ACM Trans. Parallel Comput.* **11** (2024).
2. J. Sim, A. Dasgupta, H. Kim, and R. Vuduc, "A performance analysis framework for identifying potential benefits in gpgpu applications," *SIGPLAN Not.* **47**, 11–22 (2012).
3. R. Membarth, O. Reiche, F. Hannig, *et al.*, "Hipa^{cc}: A domain-specific language and compiler for image processing," *IEEE Trans. Parallel Distrib. Syst.* **27**, 210–224 (2016).
4. A. Riahi, A. Savadi, and M. Naghibzadeh, "Many-bsp: an analytical performance model for cuda kernels," *Computing* **106**, 1519–1555 (2024).
5. E. Konstantinidis and Y. Cotronis, "A practical performance model for compute and memory bound gpu kernels," in *2015 23rd Euromicro International Conference on Parallel, Distributed, and Network-Based Processing*, (2015), pp. 651–658.
6. M. Lou and S. K. Muller, "Automatic static analysis-guided optimization of cuda kernels," in *Proceedings of the 15th International Workshop on Programming Models and Applications for Multicores and Manycores*, (Association for Computing Machinery, New York, NY, USA, 2024), PMAM '24, p. 11–21.
7. G. Alavani, K. Varma, and S. Sarkar, "Predicting execution time of cuda kernel using static analysis," in *2018 IEEE Intl Conf on Parallel & Distributed Processing with Applications, Ubiquitous Computing & Communications, Big Data & Cloud Computing, Social Computing & Networking, Sustainable Computing & Communications (ISPA/IUCC/BDCLOUD/SocialCom/SustainCom)*, (2018), pp. 948–955.
8. S. Hong and H. Kim, "An analytical model for a gpu architecture with memory-level and thread-level parallelism awareness," in *Proceedings of the 36th Annual International Symposium on Computer Architecture*, (Association for Computing Machinery, New York, NY, USA, 2009), ISCA '09, p. 152–163.
9. Q. Wang, C. Liu, and X. Chu, "Gpgpu performance estimation for frequency scaling using cross-benchmarking," in *Proceedings of the 13th Annual Workshop on General Purpose Processing Using Graphics Processing Unit*, (Association for Computing Machinery, New York, NY, USA, 2020), GPGPU '20, p. 31–40.
10. D. Lymberopoulos, O. Riva, K. Strauss, *et al.*, "Pocketweb: instant web browsing for mobile devices," in *Proceedings of the Seventeenth International Conference on Architectural Support for Programming Languages and Operating Systems*, (Association for Computing Machinery, New York, NY, USA, 2012), ASPLOS XVII, p. 1–12.

A. Additional Materials

[Click here to access the tree model diagram](#)