



An Overview of Theoretical Performance Prediction Models for GPUs

Taabish Jeshani



Agenda

Introduction

Background

A Performance Prediction Model for the CUDA GPGPU Platform

An Adaptive Performance Modeling Tool for GPU Architectures

Input-Aware Auto-Tuning of Compute-Bound HPC Kernels

GPGPU Performance and Power Estimation Using Machine Learning

A Simplified and Accurate Model of Power-Performance Efficiency on Emergent GPU Architecture

MWP-CWP

Future Work & Conclusion



Introduction

Problem

How it relates to my research?

Why?



Background





A Performance Prediction Model for the CUDA GPGPU Platform

Combination of the BSP model of Valiant, the PRAM model of Fortune and Wylie, and the QRQW model of Gibbons, Matias, and Ramachandran

The model separates memory accesses and computations and accounts for the cost of both separately.

Also takes into account the effects of scheduling, pipelining, and memory hierarchy on performance

The model can be used to analyze pseudo-code for a CUDA kernel and predict its performance by estimating the number of cycles required for computation and memory accesses



How it works?

The cost of computation is estimated by considering the cycle requirement of each operation and adding them up.

The cost of memory accesses is estimated by considering the deep memory hierarchy in the GPU and the large variation in access time for each level of the hierarchy.

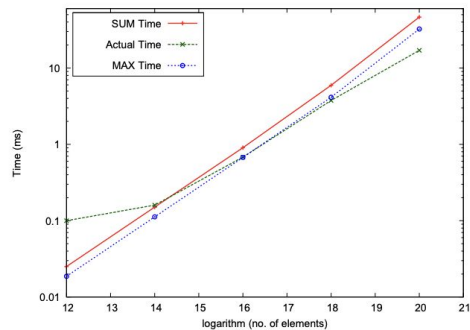
Parameter	Definition
N_w	Number of warps per block
N_t	Number of threads per warp = 32
D	Pipeline depth of a core
N_c	Number of cores per SM
K_i	Kernel i on the GPU
$C_t(K)$	Max. number of cycles required by any thread in kernel K
R	Clock rate of GPU
$T(K)$	Time taken by a kernel K

Table 1. List of Parameters in our Model

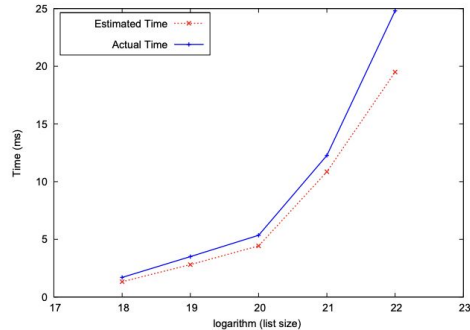
Results

The model assumes that scheduling is preemptive in nature, but the actual scheduling performed inside the GPU is not public knowledge.

The model does not take into account the effect of intra-kernel synchronization steps such as `__syncthreads()` on overall runtime.



(a) Matrix Multiplication



(b) List Ranking



An Adaptive Performance Modeling Tool for GPU Architectures

The model is designed to provide performance information to an auto-tuning compiler and assist it in narrowing down the search to the more promising implementations.

It can also be incorporated into a tool to help programmers better assess the performance bottlenecks in their code.

The authors analyze each GPU kernel and identify how the kernel exercises major GPU microarchitecture features.

They introduce an abstract interpretation of a GPU kernel, work flow graph, based on which they estimate the execution time of a GPU kernel.

How it works?

The model takes into account factors such as warp-level parallelism (WLP), data-level parallelism (DLP), and instruction-level parallelism (ILP).

The model introduces the concept of a work flow graph (WFG) to represent the computation of a kernel and combine the effects of different performance factors.

The model estimates the execution latency of a kernel by combining the effects of available concurrency (WLP, DLP, ILP) and latencies of the SIMD pipeline and memory system.

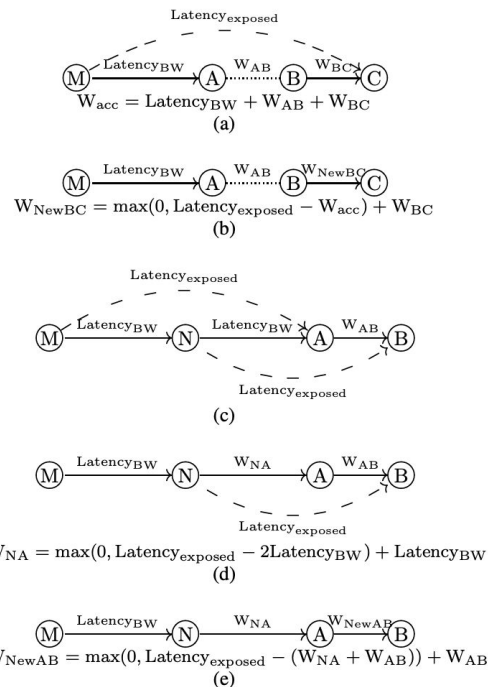


Figure 4. Reducing data-dependence arcs. (a) A sequential WFG subgraph with the long-latency data-dependence arc dashed. (b) Data-dependence arc removed after the untolateral memory latency is incorporated into W_{NewBC} . (c) Two interleaved long latency memory operations. (d) Data-dependence arcs are removed in the order that their corresponding uses appear in WFG. (e) W_{NA} that includes the untolateral memory latency of M is incorporated in covering the memory latency for N; the overlapped latency is counted once toward the warp execution latency.

Results

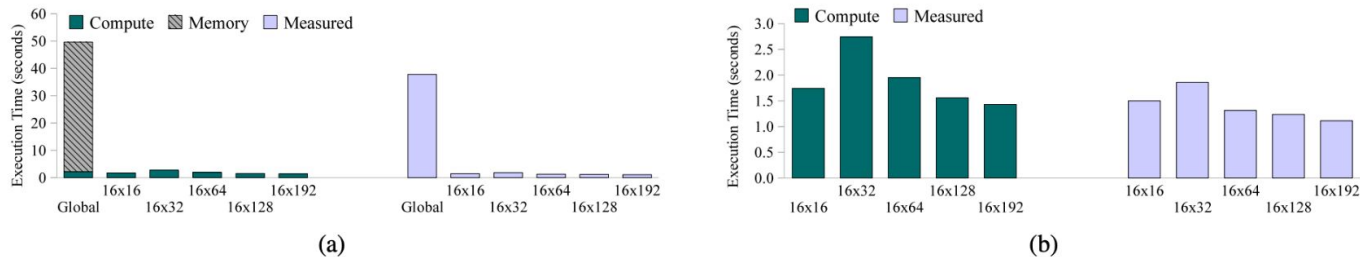


Figure 7. Matrix multiply kernels. (a) Initial and tiled kernels: a breakdown of the predicted time for the global version shows the portion of memory stalls versus compute time. (b) Zoomed for tiled kernels.

Input-Aware Auto-Tuning of Compute-Bound HPC Kernels

The paper presents an input-aware auto-tuning framework for matrix multiplications and convolutions, called ISAAC.

ISAAC uses predictive modeling techniques to drive highly parameterized PTX code templates towards not only hardware-, but also application-specific kernels.

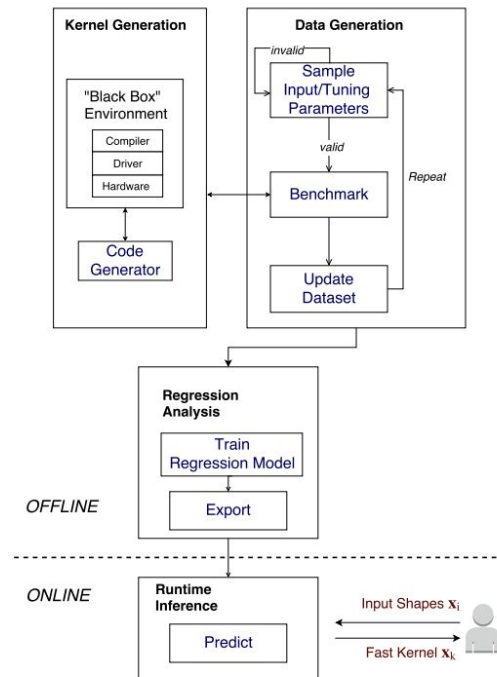


Figure 1: Overview of ISAAC



How it works?

The system is composed of four major components:

1. code generation/parameterization techniques for matrix multiplication and convolution
2. a process for generating training data for the input-aware predictive model
3. a multi-layer perceptron used as the input-aware predictive model
4. and a method for using this model at runtime to quickly infer globally optimal kernels given any input configuration.

Results

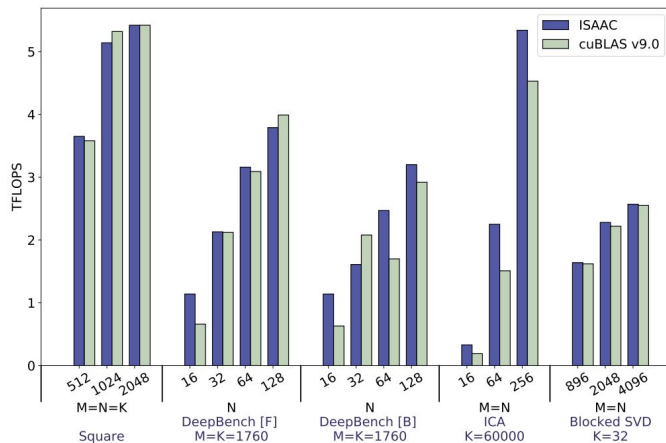


Figure 6: SGEMM performance on the GTX 980 TI

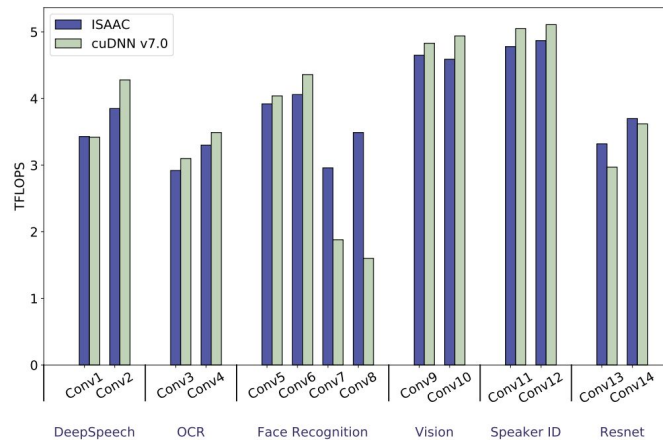


Figure 9: SCONV performance on the GTX 980 TI



GPGPU Performance and Power Estimation Using Machine Learning

The paper presents a model for estimating the performance and power of GPUs using machine learning techniques.

The model is trained on measurements from real GPU hardware, allowing it to accurately predict how applications scale as the GPU's configuration is changed.

The model uses hardware performance counter values to predict which scaling curve from the training data best represents a new application, allowing it to estimate the performance and power of the new application at different GPU configurations.

How it Works?

The model is trained on a collection of applications that are run at numerous different hardware configurations. From the measured performance and power data, the model learns how applications scale as the GPU's configuration is changed.

The model uses hardware performance counter values to predict which scaling curve from the training data best represents a new application. These dynamic counter values are fed into a neural network that predicts which training kernel is most like this new kernel.

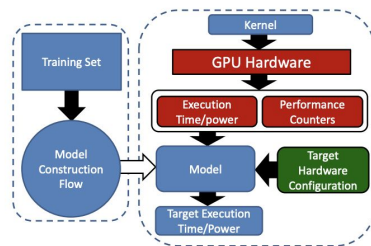


Fig. 2: The model construction and usage flow. Training is done on many configurations, while predictions require measurements from only one.

Kernel name	CU count, Engine freq., Mem. freq.		Perf. Count. 1.		Perf. Count. 2	...
	4,300,375	...	32,1000,1375	
Kernel 1						
.....						
Kernel N						

Execution Times/Power

Performance Counter Values gathered on base hardware configuration

Fig. 3: The model's training set, which contains the performance or power of each training kernel for a range of hardware configurations.

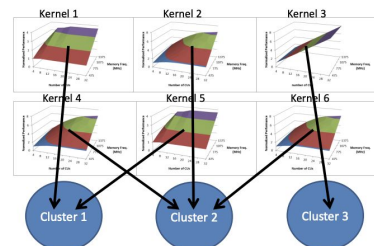


Fig. 4: Forming clusters of kernel scaling behaviors. Kernels that scale in similar ways are clustered together.

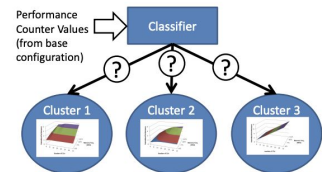
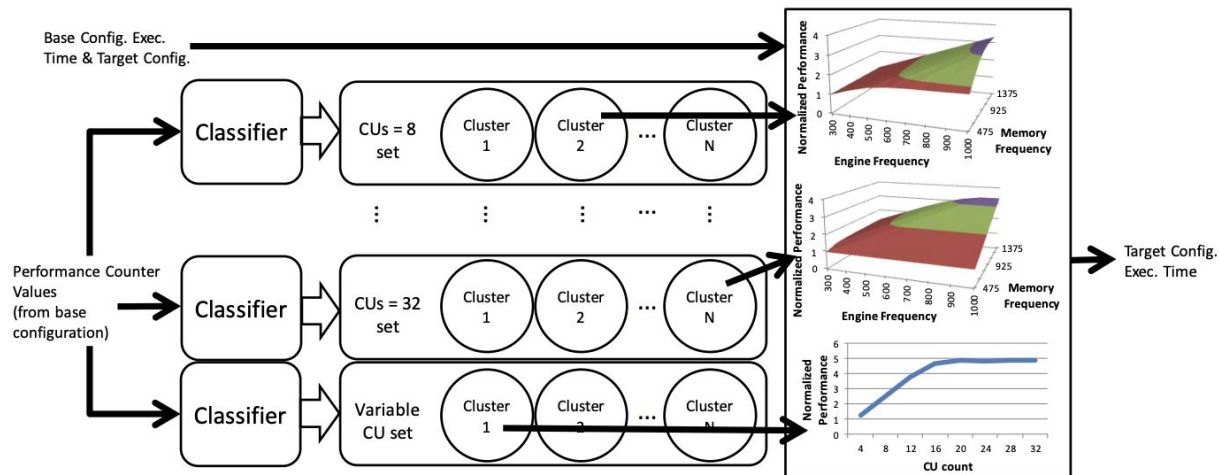


Fig. 5: Building a classifier to map from performance counter values to clusters.

How it Works? Cont.



Results

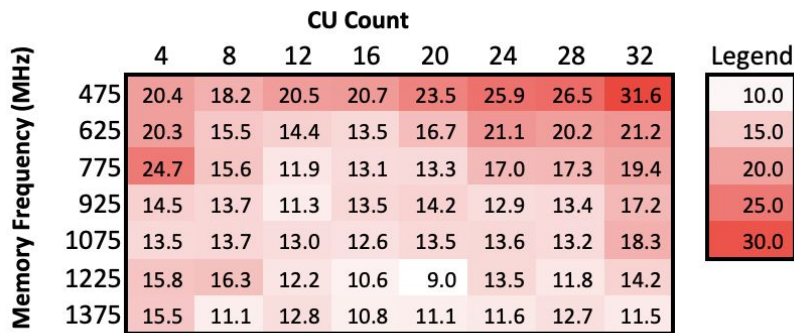


Fig. 10: Validation set error heat map at 1000 MHz core frequency. Each point represents the average error of estimating from that point's base configuration to all other configurations (including all other frequencies).

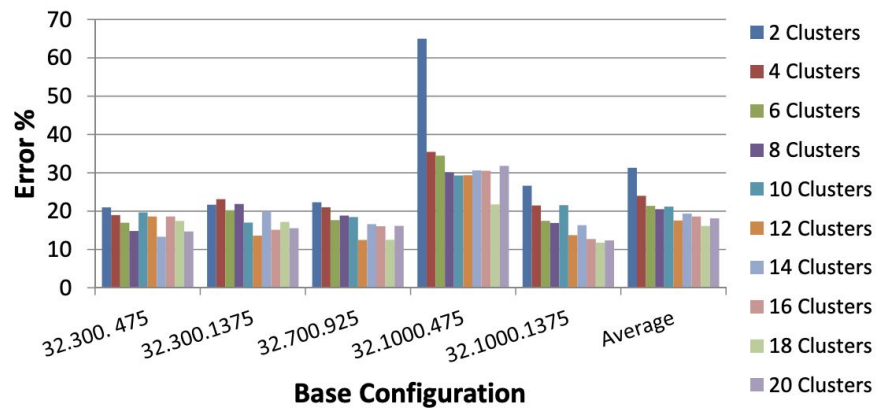


Fig. 13: Validation set error variation across cluster count.

A Simplified and Accurate Model of Power-Performance Efficiency on Emergent GPU Architecture

The paper proposes a model that combines hardware performance counter data with machine learning and advanced analytics to model power-performance efficiency for modern GPU-based systems.

The resulting model is accurate for predicting power and performance for application kernels on modern GPUs and can identify power-performance bottlenecks and their root causes for various complex computation and memory access patterns.

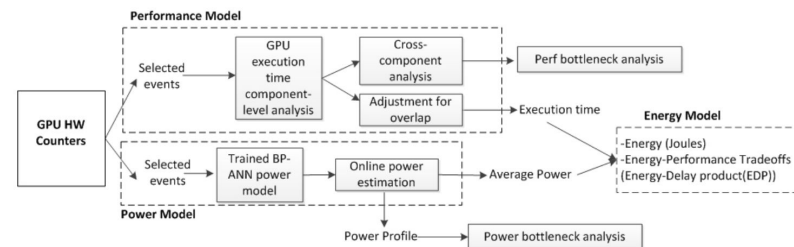


Fig. 3: Detailed view of our GPU performance counter based approach to estimate energy use on a real system.

How it works?

TABLE II: Performance Events For Training

events	composition	category
G_{load}	gld_request+l1_global_load_hit +l1_global_load_miss	global memory
G_{store}	global_store_transaction	global memory
L_{local}	local_load	local memory
L_{store}	local_store	local memory
$exec_inst$	inst_executed	instruction pipeline
$branch$	branch	instruction pipeline
$divergent_branch$	divergent_branch	instruction pipeline
S_{Access}	Shared_load+l1_shared_bank_conflict	shared memory
T_{hits}	tex0_cache_sector_queries	texture memory
T_{miss}	tex0_cache_sector_misses	texture memory

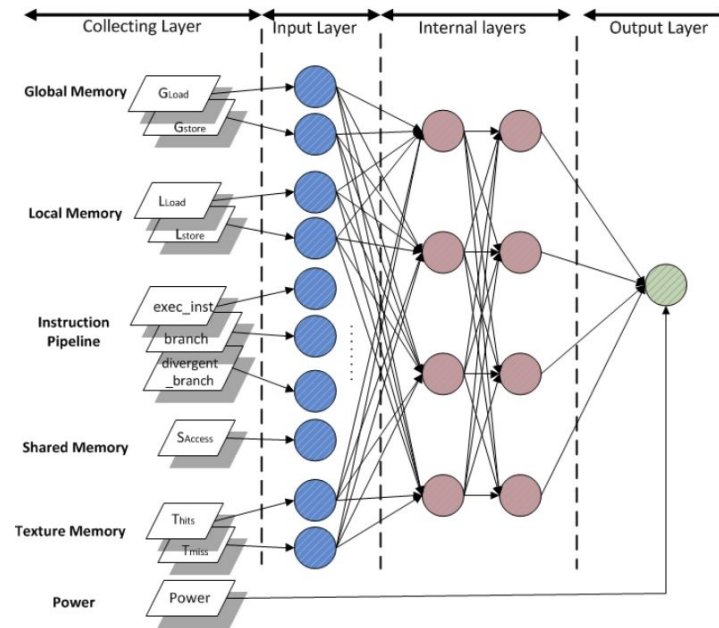


Fig. 6: Design diagram for the BP-ANN based GPU power model. Circles represent neurons.

Results

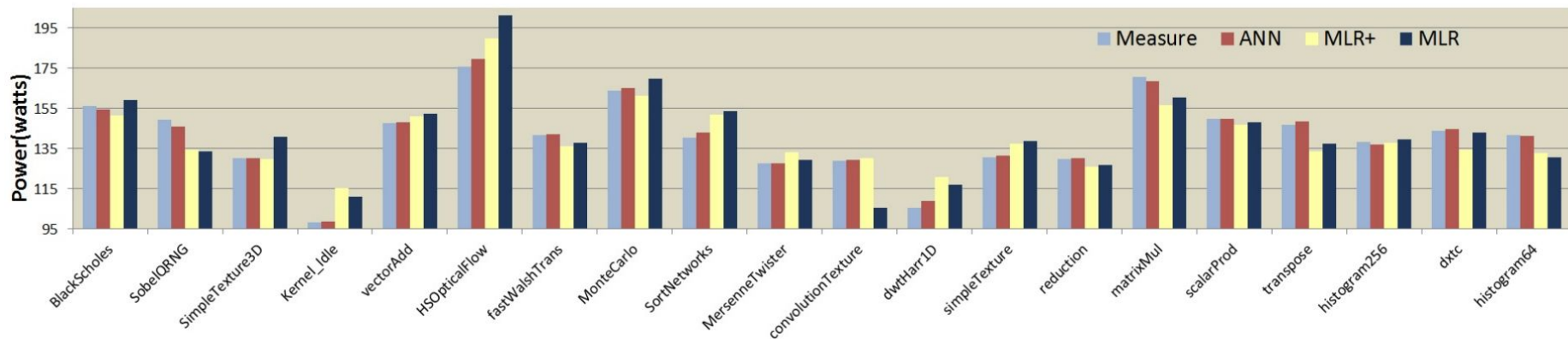


Fig. 9: Prediction accuracy comparisons of BP-ANN, MLR and MLR+ for 20 CUDA kernels.

An Analytical Model for a GPU Architecture with Memory-level and Thread-level Parallelism Awareness

The model takes into account the degree of memory warp parallelism and computation warp parallelism to estimate the cost of memory requests and overall execution time of a program

The model can provide insights into the performance bottlenecks of parallel applications on GPU architectures and help programmers tune their applications for better performance.

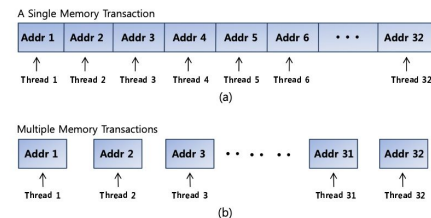


Figure 2: Memory requests from a single warp. (a) coalesced memory access (b) uncoalesced memory access



How it Works?

The model estimates the execution time of massively parallel programs on GPU architectures by taking into account the degree of memory warp parallelism (MWP) and computation warp parallelism (CWP).

MWP represents the maximum number of warps in each SM that can access memory simultaneously during one memory warp waiting period.

CWP represents how much computation can be done by other warps while one warp is waiting for memory values.

The model estimates the cost of memory requests based on MWP and CWP, and uses this information to estimate the overall execution time of a program.

The model takes into account factors such as the number of active warps per SM, memory access patterns (coalesced/uncoalesced), and synchronization effects to accurately estimate execution time.

Results

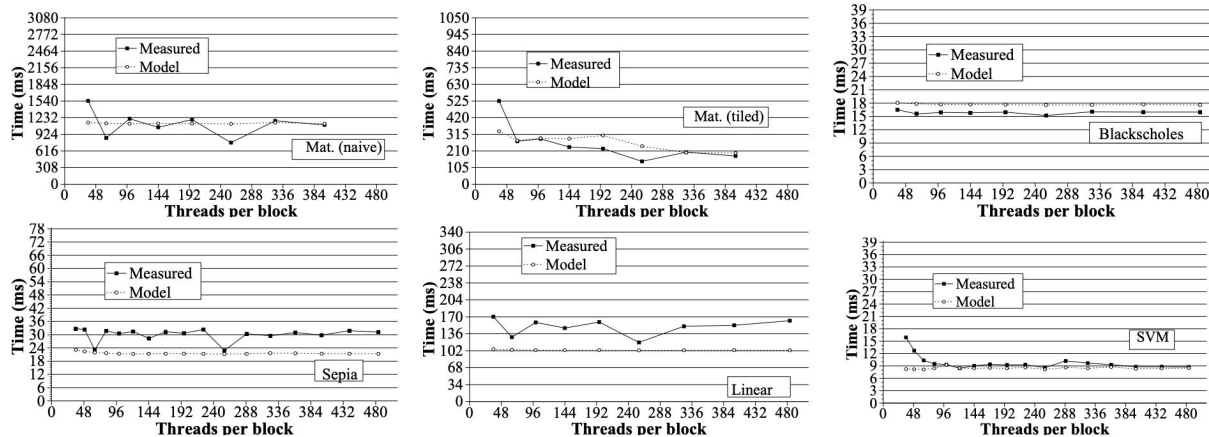


Figure 16: The total execution time of the Merge benchmarks on GTX280



Future Work & Conclusion

Machine learning based Models

More performance metrics

Dynamic adaptations to newer GPU architectures

KLARAPTOR