

Performance Prediction Modeling of GPU Applications: Analytical Modeling and Machine Learning

Marcos Amarís González*

Advisor: Alfredo Goldman*

co-advisor: Raphael Y. de Camargo†

University of São Paulo*
Federal University of ABC†
São Paulo, Brasil



SC17 - November 2017

Timeline

- 1 Introduction
- 2 BSP model Vs. ML techniques
 - Methods
 - Results
- 3 ML techniques with feature extraction
 - Methods
 - Results
- 4 Conclusions

- 1 Introduction
- 2 BSP model Vs. ML techniques
- 3 ML techniques with feature extraction
- 4 Conclusions

Introduction

The prediction of application execution times over GPUs is a great challenge and is essential for efficient Job Management Systems.

I am going to present two works

- 1 A Comparison of GPU Execution Time Prediction using Machine Learning and Analytical Modeling with regular applications
- 2 Feature Extraction and Machine Learning Techniques to Predict Execution Times of Irregular GPU Applications

3 Machine Learning Techniques

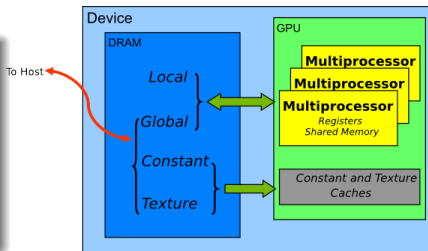
- **Linear Regression (LR)**
- **Support Vector Machines (SVM)**
- **Random Forest (RF)**

- 1 Introduction
- 2 **BSP model Vs. ML techniques**
- 3 ML techniques with feature extraction
- 4 Conclusions

CUDA, GPUs and Memory spaces

A GPU has many processors P , all processors have the same clock rate R and they are organized in Multiprocessors.

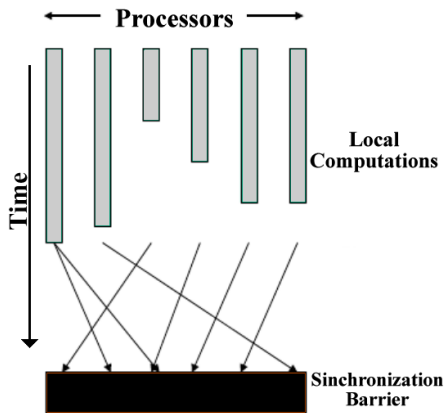
A CUDA **Kernel** can be composed of thousands or millions of threads t .



Type	On Chip	Cacheable	Instructions	Visibility	g Latency
Registers	Yes	No	Load/Store	Thread	1 cycle
Shared-L1	Yes	No	Load/Store	Block	5 cycles
Cache L2	No	Yes	Load/Store	Kernel	250 cycles
Global	No	Yes	Load/Store	Kernel	500 cycles

Table: Memory types in GPUs supported by CUDA

Bulk Synchronous Parallel Model



The cost to execute the i -th super-step is then given by:

$$w_i + gh_i + L \quad (1)$$

The total execution time of the application is given by:

$$T = W + gH + LS \quad (2)$$

Figure: Super-step in the BSP model

Proposed Analytical Model

$$T_k = \frac{t \cdot (\text{Comp} + \text{Comm}_{\text{SM}} + \text{Comm}_{\text{GM}})}{R \cdot P \cdot \lambda} \quad (3)$$

$$\text{Comm}_{\text{GM}} = (\text{ld}_1 + \text{st}_1) \cdot \text{g}_{\text{GM}} \quad (4)$$

$$\text{Comm}_{\text{SM}} = (\text{ld}_0 + \text{st}_0) \cdot \text{g}_{\text{SM}} \quad (5)$$

comp, **ld₀**, **st₀**, **ld₁** and **st₁** are obtained from the source code or profiling process.

Divergence, optimizations in the communication and differences between architecture or GPUs are adjusted by one parameter, λ .

¹M. Amaris, D. Cordeiro, A. Goldman, and R. Y. Camargo, "A simple bsp-based model to predict execution time in gpu applications," in *High Performance Computing (HiPC), 2015 IEEE 22nd Int'l Conference on*, December 2015, pp. 285–294.

GPUs of the Testbed

Model	C.C.	Memory	Bus	Bandwidth	L2	Cores/SM	Clock
GTX-680	3.0	2 GB	256-bit	192.2 GB/s	0.5 M	1536/8	1058 Mhz
Tesla-K40	3.5	12 GB	384-bit	276.5 GB/s	1.5 MB	2880/15	745 Mhz
Tesla-K20	3.5	4 GB	320-bit	200 GB/s	1 MB	2496/31	706 MHz
Titan	3.5	6 GB	384-bit	288.4 GB/s	1.5 MB	2688/14	876 Mhz
Quadro K5200	3.5	8 GB	256-bit	192.2 Gb/s	1 MB	2304/12	771 Mhz
Titan X	5.2	12 GB	384-bit	336.5 GB/s	3 MB	3072/24	1076 Mhz
GTX-970	5.2	4 GB	256-bit	224.3 GB/s	1.75 MB	1664/13	1279 Mhz
GTX-980	5.2	4 GB	256-bit	224.3 GB/s	2 MB	2048/16	1216 Mhz
Pascal-P100	6.0	12 GB	4096	224.3 GB/s	1.75 MB	1664/13	1279 Mhz

Table: Hardware specifications of the GPUs in the testbed

Algorithm Testbed

9 different regular CUDA kernels

- Matrix Multiplications in 4 different optimizations:
 - * Global Memory - MMGU
 - * Global Memory with coalesced accesses - MMGC
 - * Global and Shared Memory - MMSU
 - * Global and shared Memory with coalesced accesses - MMSC
- Matrix Addition in 2 different optimizations:
 - * Global Memory - MAU
 - * Global Memory with coalesced accesses - MAC
- Dot Product - dotP
- Vector Addition - vAdd
- Maximum Subarray Problem - MSA

Features of the Machine Learning Techniques

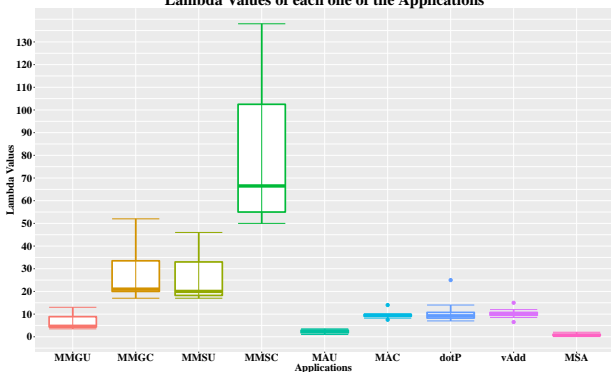
13 features were used to feed the Machine learning Techniques.

Feature	Description
num_of_cores	Number of cores per GPU
max_clock_rate	GPU Max Clock rate
Bandwidth	Theoretical Bandwidth
Input Size	Size of the problem
totalLoadGM	Load transaction in Global Memory
totalStoreGM	Store transaction in Global Memory
TotalLoadSM	Load transaction in Shared Memory
TotalStoreSM	Store transaction in Global Memory
FLOPS SP	Floating operation in Single Precision
BlockSize	Number of threads per blocks
GridSize	Number of blocks in the kernel
No. threads	Number of threads in the applications
Achieved Occupancy	Ratio of the average active warps per active cycle to the maximum number of warps ed on a multiprocessor.

Use Cases of the Analytical Model

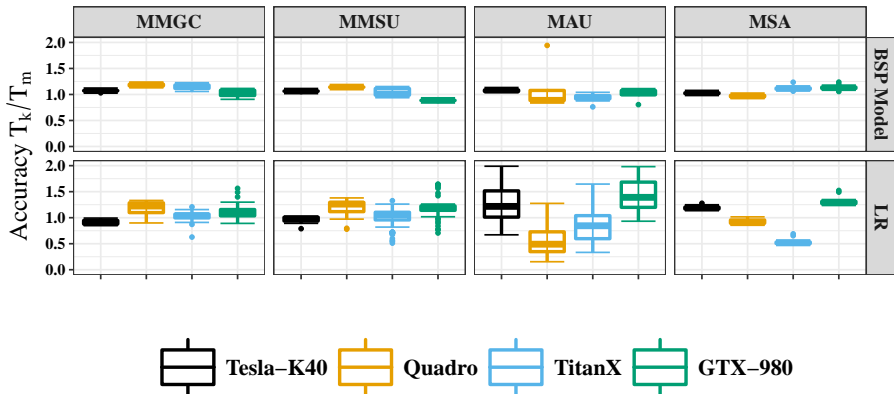
Par.	Matrix Multiplication				Matrix Addition		vAdd	dotP	MSA
	MMGU	MMGC	MMSU	MMSC	MAU	MAC			
comp	$N \cdot \text{FMA}$				$1 \cdot 24$			$1 \cdot 96$	$(N/t) \cdot 100$
ld ₁	$2 \cdot N$				2			2	N/t
st ₁	1				1			1	5
ld ₀	0				$2 \cdot N$			0	N/t
st ₀	0	1			0			$1 + \log(t)$	5

Lambda Values of each one of the Applications



Results Machine Learning VS Analytical Model

BSP Analytical Model VS Linear Regression with Regular Applications



- 1 Introduction
- 2 BSP model Vs. ML techniques
- 3 ML techniques with feature extraction**
- 4 Conclusions

Algorithm Testbed

- Correlation analysis
- Hierarchical clustering analysis.

Each cleaned sample of the dataset resulted in 85 kernel features plus 11 GPU architecture features.

Application	Berkeley Dwarf	Param.	Kernels	Samples
Back Propagation (BCK)	Unstructured Grid	1 - [57]	layerforward adjust-weights	57
Gaussian Elimination (GAU)	Dense Linear Algebra	1 - [32]	Fan1 Fan2	34800
Heart Wall (HWL)	Structured Grid	1 - [84]	heartWall	5270
Hot Spot (HOT)	Structured Grid	2 - [5,4]	calculate-temp	396288
Hot Spot 3D (H3D)	Structured Grid	2 - [3,10]	hotspotOpt1	150150
LU Decomposition (LUD)	Dense Linear Algebra	1 - [32]	diagonal	8448
			perimeter	8416
			internal	8416

Table: Rodinia applications used in the experiments

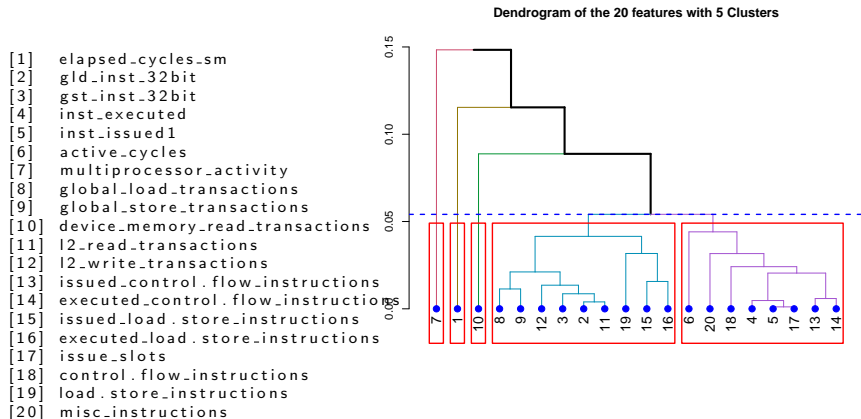


Figure: Dendrogram with 5 clusters to select 1 feature from each one

2 Contexts were tested with each ML

- 1 GPUs: we evaluated if the ML algorithm could predict the execution time over a previously unseen GPU
- 2 Kernels: in this context we predict the execution time of a previously unseen CUDA kernel.

Results

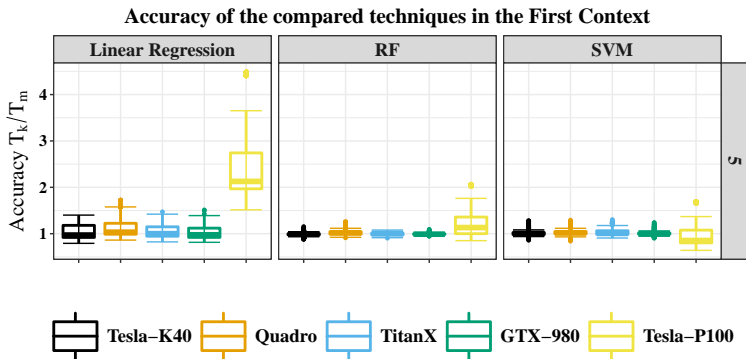


Figure: Accuracy of the first context with all the machine learning techniques used over 5 target GPUs

Results

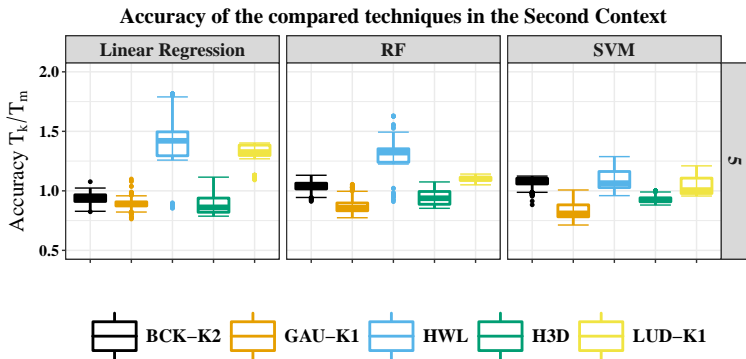


Figure: Accuracy of the second context with all the machine learning techniques used over 5 target irregular kernels

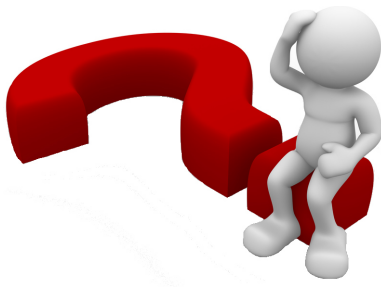
- 1 Introduction
- 2 BSP model Vs. ML techniques
- 3 ML techniques with feature extraction
- 4 Conclusions**

Conclusions

- Performance of regular GPU applications can be reasonably predicted with simple analytical models.
- Few parameters about communication and/or computation of the application are enough for these applications.
- Machine learning techniques worked better for Irregular Applications.
- Linear models are enough if the features are linearized.
- Predictions with high accuracy are possible with an optimum statistical process of feature extraction.



Thanks for your attention



Links of the work: <https://github.com/marcosamaris/>