Automatically Generating Performance Imperatives via Machine Learning

Biplab Kumar Saha Texas State University San Marcos, TX b_s183@txstate.edu

Tiffany Connors Texas State University San Marcos, TX tac117@txstate.edu

Saami Rahman **Texas State University** San Marcos, TX saami.rahman@txstate.edu

Apan Qasem **Texas State University** San Marcos, TX apan@txstate.edu

Recent interest in machine-learning based methods have produced many sophisticated models for performance modeling and optimization. These models tend to be sensitive to architectural parameters and are most effective when trained on the target platform. Training of these models, however, is a fairly involved process and requires knowledge of statistics and machine learning that the end users of such models may not possess. This poster presents a framework for automatically generating machine-learning based performance models. Leveraging existing open-source software, we develop a tool-chain that provides automated mechanisms for sample generation, dynamic feature extraction and selection, data labeling, validation and model selection. The resulting models yield high prediction accuracy and can deliver improved performance and energy efficiency in many different contexts. We present summary experimental results that demonstrate the utility and efficiency of the system.



Degree of practitioner involvement highlighted in each stage



Survey of 112 papers on performance modeling and tuning via machine learning (MLMT) 2006-2016

Scores assigned based on work described in studied paper only



—AOS —CA



Inclusion of working set in feature vector and its impact on learning effectiveness

recommendation w/out reuse distance recommendation with reuse distance block 1 block 1 block 2 block 2 erchange nterchange unroll 1 unroll 1 unroll 2 unroll 2 pad nad recommended lov

Differences in recommendation with and without composite features in feature vector

Bottleneck classification and performance insight



0.80

2 0.60

0.40

Pinpointing bandwidth saturation points for GPU kernels divergent memory access: (a) AMD 7850, (b) AMD8500 and (c) AMD Radeon Nano (CGO17)

Performance and energy gains



Can generate ML performance models with minimal user intervention

- Compiler optimizations sequence,
- specialized code transformations
- CPU, GPU and heterogeneous platforms

Auto-generated models exhibit comparable accuracy; can boost performance and improve energy efficiency

Exploring the feature space

Use of diagnostic and composite features



Context-specific scaling: program characterization with unscaled (left) and scaled features



Visualizing the feature space: PCA Varimax Rotation diagrams for feature space containing GPU performance events



Identifying performance patterns in the feature space of arithmetic intensity: (a) AMD 7850, (b) AMD8500 and (c) AMD Radeon Nano (CGO17)

Cross-platform performance prediction



Prediction accuracy of four classifiers on two generations of NVidia GPUs : Kepler and Maxwell. X-Y indicates model trained on X and invoked on Y

nread block size prediction in GPU kernels Hardware prefetching on Intel Xeon Phi (HPCC 2015) (SC17, under review)



Thread affinity prediction for multithreaded CPU workloads (IGSC15)



Register cap prediction for GPU kernels (under review)

Models provide new insight about known performance problems

Acknowledgements

This research funded by NSF award CNS 1253292

Work on data layout was done while Apan Qasem was a Visiting Scholar at AMD Research

We thank John Mellor-Crummey and Nathan Tallent for their help with the HPCToolkit interface