

Real World HPC Systems for Big Data/AI Research

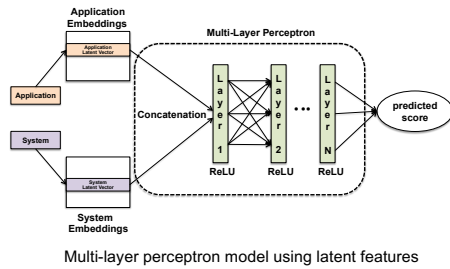
Learning Neural Representations for Predicting GPU Performance¹

Problem

- Modeling performance of applications across systems with different GPU microarchitectures

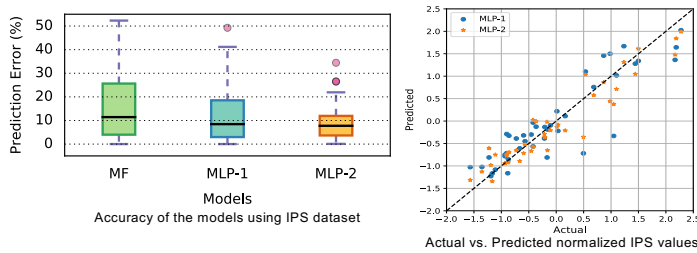
Proposal

- A collaborative filtering based matrix factorization (MF) model to automatically learn latent features describing performance of applications on systems
- A multi-layer perceptron (MLP) to model complex non-linear interactions between applications and systems



Results

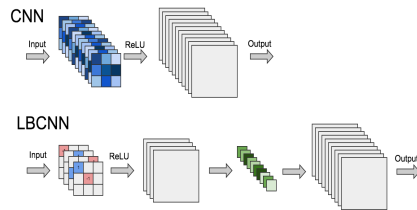
- MLP with 2 layers (MLP-2) achieves 90.6% accuracy when predicting instructions per second (IPS)



Can Local Binary Convolutions Make Neural Networks Models Smaller?³

Background

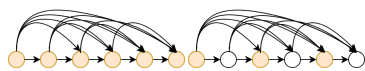
- Local Binary Convolution Neural Networks (LBCNN)
 - models use Local Binary Convolution (LBC) layers
 - (+) less learned parameters
 - (-) less accuracy



Problem

- Not applicable to Large Model
 - Too slow
 - Too much Acc. loss

● LBC ● LBC ○ Normal Conv

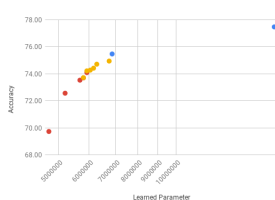


Proposed Method

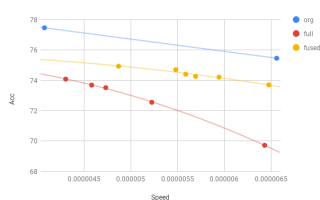
Only replace half of the normal convolution layer by LBC

Results

Results of models for different number of difference maps generated by LBC filters



Accuracy and Speed of different models



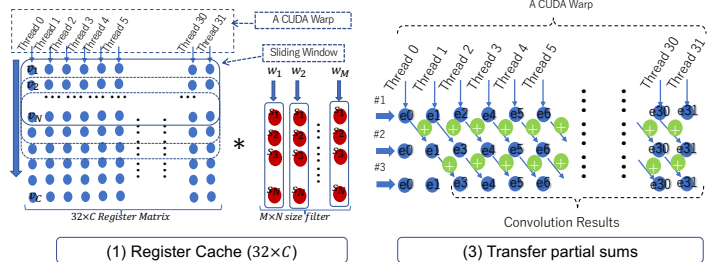
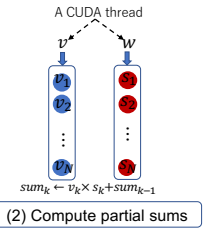
Efficient 2D Convolution on CUDA-enabled GPUs²

Problem

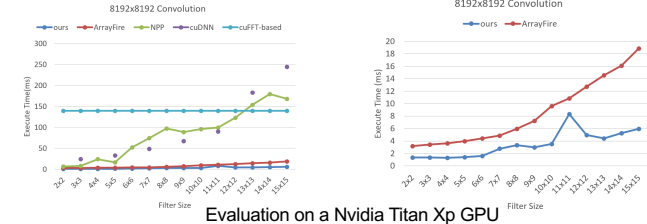
- 2D convolution computation on CUDA-enabled GPUs

Proposal

- Using registers as a cache
 - Lower latency, high throughput
 - High data reuse
- Computing partial sums in parallel
- Transferring partial sums in parallel
 - No stalling within thread
 - Efficient communication between threads by shuffle



Results



Framework for Transpilation Between Python and Fortran⁴

Problem

Fortran

- + top performance
- + HPC legacy
- hard to maintain

Python

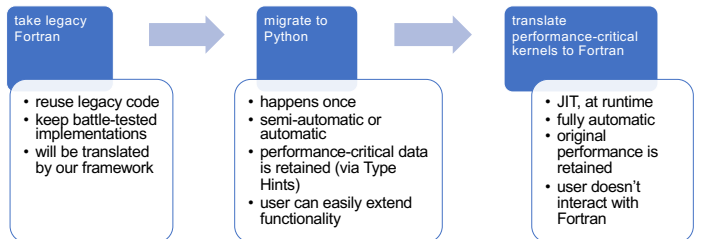
- + ease of programming
- + general-purpose tools
- big runtime overhead

Dilemma

- Performance or ease-of-programming?

Solution

- Python for development, Fortran at runtime.



Results

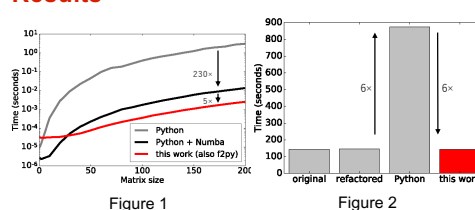


Figure 1: DGMEMM performance the same as Fortran. 5x better than Numba.

Figure 2: Migrated Miranda IO benchmark retains original performance.

[1] Shweta Salaria, Aleksandr Drozd, Artur Podobas, Satoshi Matsuoka, Learning Neural Representations for Predicting GPU Performance, ISC'19

[2] Peng Chen, Mohamed Wahib, Shin'ichiro Takizawa, Ryousei Takano, Satoshi Matsuoka, Efficient Algorithms for the Summed Area Tables Primitive on GPUs. IEEE CLUSTER'18

[3] Haoyu Zhang, Mohamed Wahib, Satoshi Matsuoka, Can Local Binary Convolutions Make Neural Networks Models Smaller?, Submitted to ICPP'19 (Poster)

[4] Mateusz Bysiek, Aleksandr Drozd, and Satoshi Matsuoka. Migrating Legacy Fortran to Python While Retaining Fortran-level Performance Through Transpilation and Type Hints, 6th Workshop on Python for High-Performance and Scientific Computing, co-held with SC16