

# Designing for Portable, Efficient and Explainable Performance

Tze Meng Low Assistant Research Professor Carnegie Mellon University

10 May 2022





## Electrical & Computer

4

#### **Performance Portability is Necessary**



- Manual heroic optimization effort is not scalable
- Repeated optimization with changes
  - Computing platform
  - Different applications
  - Different ML networks
- Blackbox tools and libraries provide little insights into why an implementation is bad/good



## **Key Questions**

How do we capture knowledge about performance?

How do we apply it to different applications and architectures?

**Explainable Designs though Analytical Models** 

# **Capture HW-SW interactions with Models**

- SW optimizations target specific HW features
- Analytical models map how SW is tuned based on available HW features
- Models inform changes in SW as HW and application requirements change



Electrical & Computer

# Hardware-driven constraints

- HW must be fully utilized for good performance
- Computational units characterized by
  - Number of Units,
  - Number of output computed
  - Time to compute each output (Latency)



https://en.wikichip.org/wiki/intel/microarchitectures/haswell

Core	Core	Core	Core	LD/ST	SFU
Core	Core	Core	Core	LD/ST	SFL
Core	Core	Core	Core	LD/ST	SFL
Core	Core	Core	Core	LD/ST	SFL
Core	Core	Core	Core	LD/ST	SFL
Core	Core	Core	Core	LD/ST	SFL
Core	Core	Core	Core	LD/ST	SFL
Core	Core	Core	Core	LD/ST	SFL

 $L_{\rm FMA}N_{\rm FMA}N_{vec}$ 

**Hardware Constraints** 

Tze Meng Low, Francisco D. Igual, Tyler M. Smith, and Enrique S. Quintana-Orti. 2016. Analytical Modeling Is Enough for High-Performance BLIS. ACM Trans. Math. Softw. 43, 2, Article 12 Electrical & Computer







## **Data Orchestration is Key**



- Data orchestration is about keeping functional units busy
  - Performing data movement
  - Managing the caches





# Intermediate Layouts Is Necessary



#### **Matrix Matrix Multiplication**

#### **Direct Convolution**



Electrical & Computer

#### **Fast Fourier Transforms**





#### **Further questions**

#### How general is the approach?

# How can models for one application be used for another?



## **Rethinking FFT algorithms**



- FFTs usually considered memorybound
- Caches are fast enough to sustain **FFT** computation
  - Use different layouts at different stages of the computation



# Sharing caches with SMT Threads

- SMT Threads share "private" caches
- Controlled data movement across different layers of the cache
  - Use of temporal loads/stores
  - Need locks to ensure correctness
  - Split caches by assigning ways to caches
- Sharing of functional units
  - NOPs introduced to allow both threads to proceed with their tasks





Elliott Binder, Tze Meng Low, and Doru Thom Popovici. "A Portable GPU Framework for SNP Comparisons." 2019 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)



## **High Performance Machine Learning**

**Compute across different inputs blocks** 

**Reuse Kernel across outputs** 

- Associativity & Cache Size

- Register Limit

- Cache Length Size

Satisfy ILP requirement

- Associativity & Cache Size

- L1 / L2 Size

#### **High Performance Direct Convolution**



Leverage model-based data layout

for ii = 1 to C i/C ib do

for kk = 1 to  $W \circ / W \circ b$  do

for m = 1 to W f do

for i = 1 to <u>C</u> ib do for k = 1 to W ob do

for j = 1 to C ob do

for n = 1 to H f do

for l = 1 to H o do

#### Performance normalized to OpenBLAS GEMM on AMD PileDriver 4.0 GHz, 4/4 cores/threads



Jiyuan Zhang, Franz Franchetti, Tze Meng Low, "High Performance Zero-Memory Overhead Direct Convolutions", 2018, International Conference of Machine Learning (ICML)

O[...] += I[...] \* F[...] - SIMD Length

- L2 / L3 Size





## **Towards Graphs & Sparse Linear Algebra**

- Graph algorithms in the language of linear algebra
- New algorithms for
  - finding patterns/subgraphs
  - propagating information
- Modeling to find out HW bottlenecks

IBM-GraphBLAS redisgraph redislabs







repository
Distributed 9btl
Lawrence Livermore
National Laboratory

**Galois GB** 

SuiteSparse

Mark P. Blanco, Tze Meng Low, and Kyungjoo Kim, "Exploration of Fine-Grained Parallelism for Load Balancing Eager K-truss on GPU and CPU", IEEE HPEC 2019 , **Graph Challenge Champion** Tze Meng Low, Daniele G. Spampinato, Anurag Kutuluru, Upasana Sridhar, Doru Thom Popovici, Franz Franchetti, Scott McMillan (CMU) ,

"Linear Algebraic Formulation of Edge-centric K-truss Algorithms with Adjacency Matrices ", IEEE HPEC 2018 Graph Challenge Finalist

Social Network

, Electrical & Computer



**Protein Interactions** 

21



#### Summary

- Analytical models capture key SW-HW interactions
  - Flexibility and portability across architectures and applications
  - Rethink of many current algorithms & implementations
  - Tools integration (e.g. Polly-LLVM) 
     programmer productivity



lowt@cmu.edu