FSA: Co-Designing Future High Performance Systems for Efficiency and Scalability

Shuaiwen Leon Song

Director, Future System Architecture Lab (FSA) University of Sydney CS and UW Seattle ECE

DSI presentation, August 18th, 2021







I acknowledge the tradition of custodianship and law of the Country on which the University of Sydney campuses stand. We pay our respects to those who have cared and continue to care for Country.

Our Research Overview





Our People and Collaborators



Google







Data-Centric Bespoke Design

- ➢ We constantly struggle to identify the best design and optimization strategies to consider application features, data characteristics, and the system constraints from software and hardware. We always work under constraints unless we customize our own hardware and software for everything.
- > Challenge: Hardware lottery between general-purpose design vs customized acceleration



Focus of HPC Research

- Scalability, efficiency and optimization of scientific and nonscientific applications under current hardware constraints;
- Effective big data reduction for scientific simulation and modelling;
- High-performance System Designs for emerging ML services and applications (SystemML)

HPC Architectures: Extreme Heterogeneity



- ECP RFI feedback from vendors suggests that exascale platforms and beyond will become increasingly heterogeneous.
- Reports also suggest a variety of novel technologies are likely to be integrated into a single extreme heterogeneous system.
- We already see sings of this in the planned system architectures for **CORAL and APEX** exascale supercomputer acquisition.

Co-Design Dilemma: Architecture Becomes More Complex



(B) V100-based DGX-1 with NVLink-V2

Showcase I: Efficiency and Scalability

S-BLAS: High performance parse linear algebra kernels for complex multi-accelerator based HPC architectures

https://github.com/pnnl/s-blas

- Sparse-Matrix-Vector-Multiplication (SpMV)
- Sparse-Triangular-Solve (SpTRSV)
- Sparse-Matrix-Transposition (SpTrans)
- Sparse-Matrix-Matrix-Multiplication (SpMM)



The SOTA SpTrsv on Multi-GPU Based HPC Architectures

- cuSparse lib: csrsv2()
 - Basic idea: some components are independent and can be processed simultaneously.



- A Synchronization-Free SpTrsv
 - Basic idea: Migrating the level-sets analysis at runtime
 - Components are scheduled by the hardware warp-switch of the GPU
 - Update the intermediate value (in_degree and left_sum) using GPU atomic operations

partial order.

SpTrsv with Unified Memory

- Convenient
 - Hide complexity from users
- Page fault mechanism
 - Coarse-grained data copy
- Data contentions
 - System-wide atomic update will access the data simultaneously
- Workload unbalance
 - Dependency are unidirectional



Require Low Overhead Fine-Grained Communications!!

GPU CUDA NVSHMEM

- OpenSHMEM-based PGAS programming interface for multi-GPUs
- GPU-side interface allows GPU threads to
 - 1. Access distributed memory via data movement API
 - 2. Direct load/store (LD/ST) where GPUs are P2P-accessible
 - 3. Highly-concurrent fine-grained messaging
 - 4. Asynchronously one-side data communication



Peer Direct LD/ST and Global Remote Access

Using NVSHMEM for Resolving Dependency

• For SpTrsv, we convert the system-wide atomic update from unified memory to NVSHMEM shared memory space



Workload Balance among GPUs



Our design can achieve an average of 3.53× speedup on a DGX-1 system and 3.66 × speedup on a DGX-2 system over the Unified-Memory design, respectively.

TASK DASEU WOLKIDAU DISTLIDUTION

- More tasks per GPU: workload becomes more balanced among GPUs
- Less task per GPU: can exploit in-task data locality for better performance



Showcase II: Big Data Reduction via Lossy Compression

Background ●○○○ Introduction

Design 000000 Evaluation

Conclusion

Trend of Supercomputing Systems Gap Between Compute and I/O

The compute capability is ever growing while storage capacity and bandwidth are developing **more slowly** and not matching the pace.

supercomputer	year	class	PF	MS	SB	MS/SB	PF/SB
Cray Jaguar	2008	1 pflops	1.75 pflops	360 тв	240 дв/з	1.5k	7.3k
Cray Blue Waters	2012	10 pflops	13.3 pflops	1.5 рв	1.1 тв/s	1.3k	13k
Cray CORI	2017	10 pflops	30 pflops	1.4 рв	1.7 тв/s [•]	0.8k	17k
IBM Summit	2018	100 pflops	200 pflops	>10 рв●●	2.5 тв/s	>4k	80k

PF: peak FLOPS MS: memory size SB: storage bandwidth

when using burst buffer
 counting only DDR4

Source: F. Cappello (ANL)

Table 1: Three classes of supercomputers showing their performance, MS and SB.



Scientific Simulations and Experiments: Flood of Data!

Today's scientific research using simulations or instruments produces extremely large datasets

 \circ Many datasets are in petabyte (PB = 10^{18} Bytes)!

Cosmology Simulation

- HACC: hardware/hybrid accelerated cosmology code (twice IEEE/ACM Gordon Bell Prize Finalists)
- A total 20 PB data when simulating one trillion of particles
- Peta-scale system's file system ~ 20 PB
- Mira supercomputer has 26 PB FS, 20 PB / 26 PB ~ 80%
- NSF Blue Waters (1TB/s I/O bandwidth), 20 x 10¹⁵ / 10¹² seconds (5h30m) to store the data
- Data reduction of about a factor of 10 is needed 2 currently drop 9 snapshots over 10 (called decimation in time)







How SZ works?



Issues with SZ and Its Current FPGA Implementation

Low throughput of SZ

encoding

predictors (2 bits)

- Iack of parallelism: SIMD and SIMT cannot apply
- Limitations in FPGA GhostSZ
 - totally performance-driven design
 - 3 predictors in use, need extra bits to encode
 - more "workflow pipelines" (more resource)
 - low compression ratio

GhostSZ encoding prediction error in quantized form (14 bits)

WAVESZ, SZ-1.4

encoding prediction error in quantized form (16 bits)

New use scenarios of adopting FPGA

- real-time processing; "inline processing" (Intel, 2018)
- ExaNet—an FPGA-based direct network architecture of the European exascale systems [Ammendola et al. 2018]



"hard"

-0.01

CLDLOW.

0.00

Figure 3: CESM-ATM

"easy"

outlier

1

0.01

Figure 1: Loop-carried dependencies due to writeback.

Temporal-Spatial Mapping



- FPGA + wavefront memory layout = more pipelining control
- Ideally, suppose Λ cycles to finish (prediction + quantization), no stall if
 II = 1, and ② (vertically) iterating over Λ points from (r, c) to (r, c+1)
- **b** BODY ("perfect loop") unrolled with factor Λ (= vertical dimension) and II = 1

Showcase III: System ML for Accelerated Scientific Discovery

- Innovative software-hardware co-design strategies for emerging ML model architectures (*with UW, MIT*)
- Large-scale data-center and HPC ML services design and optimization (*with Alibaba Research*)
- Noise and randomness control in ML training (*with Google*)
- Resource-constrained training and inference (*with UW*)
- Neural Architecture Search and Tiny ML (*with Brown*)



CNN Architecture Evolution



21

Linear vs. Non-Linear Networks





Challenge: Memory Shortage

23



Our Goal



Largest Accelerator Memory:trainingTraining Memory Request:Tens of GBshigh speedHundreds of GBs

State-of-the-art Frameworks







(1) No Liveness Analysis,(2) Limited Support in Non-linear Networks (1) Inefficient Memory Offloading
 (2) No performance-memory tradeoff
 (3) No minibatch-LR tuning (1) Inefficient Memory Offloading(2) Speed-Centric(3) No minibatch-LR tuning



- > Dynamic memory provisioning framework for deep non-linear networks
- Recognizing non-uniform memory distribution across layers
- Runtime bound the minimal peak memory usage (peak_m = max(l_i)), at layer-wise granularity

PPoPP 2018

Opportunity 1: Physical Memory Reuse



<u>Core Idea</u> Reuse the same physical memory at different time partitions

Liveness Analysis



Example: Liveness analysis on AlexNet



It reduces more than 1/3 of memory requirement for AlexNet

Opportunity 2: Computation Pattern



Computation and Memory intensity differ across the layers:

Layers	Comp (%)	Memory (%)	Checkpoint? (Y/N)
CONV	~50%	~50%	Y
POOL, ACT, BN, LRN	~30%	~50%	Y, but requires opt
DROPOUT, SOFTMAX, FC	~20%	~1%	Ν

Unified Tensor POOL (UTP)

Key Operations:



Extensible to Various Physical Memory Pools



Example: Liveness Analysis + UTP on AlexNet



It further reduces another 1/6 of memory requirement for AlexNet

Opportunity 3: Store vs. Recompute



Observation: POOL, ACT, LRN and BN forward computation only accounts for less than 10% of total time, but with 50% memory consumption.



Memory-Centric (recompute dependencies) vs. Speed-Centric(reuse)

Example: Liveness Analysis + UTP on AlexNet + Cost-Aware



Memory is bounded by the layer peak and peak is further reduced for AlexNet

Many other works in this direction

- Machine Learning Optimizing Compiler vs. XLA and TVM (ASPLOS'22, under review, A*)
- Bayesian accelerator design for cloud and autonomous driving (MICRO'20, MICRO'21, A*)
- Large-scale Bayesian inference system for Argonne's genome analysis application and cosmic tagger application (Supercomputing'21, A*)
- LSTM and Transformer's accelerator design (ISCA'21, HPCA'21, A*)
- Capsule network acceleration design (HPCA'2020, TPDS architecture research highlight 2020, A*)
- Reinforcement learner for interactive virtual reality system design (ASPLOS'21, A*)

Edge/IoT/Embedded Scenarios

- Real time constraint, low power , low storage require better quality of data;
- We do not have luxury of even powerful chips; everything is working under a tighter constraints;
- We have to overperform our design to match users' demand;
- We have to squeeze every bit of efficiency out of our design;
- Better optimization, better hardware utilization, better compatibility with cloud, better security mechanisms, etc.

Data-centric Design is very important !

Showcase IV: Lightweight Reinforcement Learning Based VR SoC

• Strong relationship between users' input and rendering workload





Q-learning Based Eccentricity Controller using User Input



• Index the user input to simplify the design complexity for mobile devices

Software-Only Prototyping on Unity



() unity



Contact me if you are interested in our projects!

HPC Memory and Data Centric Accelerator Design Quantum Compiler and Architecture Design

. . .

<u>Shuaiwen.song@Sydney.edu.au</u> https://shuaiwen-leon-song.github.io/





Australian Government Australian Research Council



C-) Alibaba Cloud