Graph Based Contractions for Calculating Correlation Functions

Redstar: Analysis Software Suite

Jie Chen, Eloy Romero Alcalde, Frank Winter and Robert Edwards







Talking Points

- What is graph contraction?
- Redstar.
 - Current optimizations.
 - Reduce number of contraction calculations and run-time memory footprint.
 - Hadron contractions.
 - Carry out real tensor contractions on CPUs and GPUs.
- Performance
 - Mostly on GPUs.
 - Observations.
- Challenges.







Correlation Functions and Graphs



Different Types of Graphs









Graph Nodes (V)

- Meson (2 quarks), Baryon (3 quarks)
- Meson
 - Each vertex $M_{\alpha\beta}^{ij}$ where $\alpha,\beta=0,1,2,3$ (spins) and $i,j \sim$ hundreds (distillation space).
 - e.g. each vertex occupies 384*384*16*16 = 37 MB
- Baryon
 - Each vertex $B_{\alpha\beta\gamma}^{ijk}$ where $\alpha, \beta, \gamma=0,1,2,3$ (spins) and i,j,k ~ 100s (distillation space)
 - e.g. each vertex occupies 128*128*128*16*64= 2GB, each vertex occupies about 128*128*128*16*8=268 MB if only two upper spin indices are used ($\alpha, \beta, \gamma=0, 1$).







Graph Edges (E)

- Meson-Meson contraction $A^{ij}_{\alpha\beta} B^{jk}_{\beta\gamma}$
 - Batched matrix multiplication with batch size <= 64
 - $O(N^3)$ for calculations, $O(N^2)$ for memory
- Meson-Bryon contraction $A^{ij}_{\alpha\beta} B^{jkl}_{\beta\gamma\delta}$
 - Batched tensor contraction with batch size <= 256
 - $O(N^4)$ calculation complexity, $O(N^3)$ memory
- Baryon-Baryon Contraction
 - One index contraction $A_{\alpha\beta\gamma}^{ijk} B_{\gamma\delta\varepsilon}^{klm}$
 - $O(N^5)$ calculation complexity, $O(N^4)$ memory
 - Huge memory (e.g 128⁴ * 16 * 16 = 68GB for two upper spins)
 - Batched tensor contraction with batch size up to 1024
 - Two-index contraction $A^{ijk}_{\alpha\beta\gamma}B^{jkl}_{\gamma\beta\delta}$
 - $O(N^4)$ calculation complexity, $O(N^3)$ memory.





Try to avoid if possible



How Many Graphs ?

- Number of Graphs ∝ N ! N is the number of freedom of degrees
- Recently approaching 1M graphs
- Calculations across many time slices



I=1/2 K*π arXiv:1406.4158

Name	a ₀	f ₀	a	roper	Deuterium	Tritium
# Graphs	19,041	27,999	385,512	84,894	119,191	6,208
Туре	MM	MM	MM	MB	BB	BB
Vector Size	256	256	128	64	64	32







Office of

Science

Graph Contractions



Removal of one edge after another until two hadron nodes are left

The sequence of the edge removal is important









Graph Contractions









Redstar Workflow









Graph Generation and Classification (*redstar_npt*)

- Parallel contraction diagram and graph generations (*Wick Contraction*)
 - Parallel graph classifications to identify unique graphs



GPU ???

Science

- Graph isomorphism problem.
- Perform canonical labeling of each graph in parallel.
- Identify unique graphs in O(n) time.
 - Brendan D. McKay (Nauty)



Optimal Graph Contraction Path and Memory Reduction



 Each edge is assigned a weight (calculation cost/the number of appearances in all the graphs, O(N))



Optimal Graph Contraction Path and Memory Reduction

- Minimum spanning tree algorithm is used to select an optimal path to reduce the number of contractions.
 - Prefers 2-index contractions (smaller calculation cost)
 - > For meson systems, the edges with most appearances are chosen first.

• Similar graphs are grouped together during the construction of a contraction queue to minimize the memory footprint of the intermediate contraction results.







Optimal Graph Contraction Path and Memory Reduction









Redstar and Hadron Contraction

• A gigantic *DAG* is created after the optimal contraction path is determined for each graph.



 A contraction queue is constructed from the DAG. Independent contractions can be grouped together to form arrays of contractions (*vector-form* contractions).





- Batched tensor contractions on CPUs or GPUs
 - Batched zgemm coupled with tensor permutation if needed.
 - Vector-form of hadron contraction to increase batched zgemm batch size on CPUs.
 - Vector-form of hadron contraction to utilize multiple GPUs each could use multiple streams.
 - Each element in a contraction queue (an array) has no information about on which GPU it executes.
 - Hadron library utilizes an algorithm that decides where to put a contraction based on where the contraction sources have been.
 - Keep data locality high and maintain good load-balance.







Contraction APIs

HadronDistOperatorRep hadronExteriorContract(const HadronDistOperatorRep& src1_rep, int ind_j, const HadronDistOperatorRep& src2_rep, int ind_k);

HadronDistOperatorRep hadronExteriorContract(const HadronDistOperatorRep& src1_rep, const std::vector<int>& ind_j, const HadronDistOperatorRep& src2_rep, const std::vector<int>& ind_k);









More about the DAG and the Contraction Queue

- Contraction queue is constructed from the DAG by breadth-first walk through of the DAG.
 - No gpu information assigned to any contraction.
- Independent contractions are grouped together to form an array of contractions.
 - NPT_BATCH_SIZE
 - Array of contractions
 - Increase batch size for batched zgemm on KNL (or many core CPU) to take advantage of computing power of many cores.
 - Works very well for KNL nodes.
 - Each contraction in the array is carried out by a dedicated GPU if multiple GPUs are available.
 - If the sources of all contractions in the array are different, we expect a linear speedup.
 - Unfortunately, we have situation like : extcontract(a, 0, b, 1);

extcontract(a, 1, b, 0);

extcontract(a, 0, c, 1);

extcontract(b, 0, c, 0);

- In this example, a, b, c have to be copied or moved to different GPUs.
- Multi-gpus may not provide any benefit.







Hadron Contraction on GPUs

Performance bottlenecks

- Data in and out of GPUs kill the performance.
 - Need to consolidate all the spin components on host side into a contiguous memory block to reduce the number of memory copies.
 - Need data reuse.
 - Data stay on a GPU as long as possible unless they are needed on the host.
 - Memory evictions due to the limited memory on a GPU.
 - Data movements from one GPU to another.
- Many individual matrix multiplications perform poorly
 - Data movements and not enough computation on a GPU.
- Naive tensor permutations perform terribly on GPUs
 - Stride data access using the global memory on a GPU.
- Some matrix matrix multiplications are actually matrix vector multiplications depending on how two tensors' indices are folded.
 - Needs to identify those cases and implement batched version of gemv.







Hadron Contraction on GPUs Current Implementation and Optimizations

Major components

- GPU memory managers
 - One manager for one GPU.
 - cuda or hip malloc and free are going through the manager.
 - Contiguous gpu/host memory block to reduce I/O
 - Manages data evictions (improved LRU) and data reuse.
 - Intermediate contraction values stay on GPUs.

GPU matrix/tensor contraction kernels

- Batched zgemm from cublas/hipblas.
- Customized batched zgemv, zdotu, accumulate, contractall.
- GPU tensor transpose library (GPUTT)
 - Fast tensor transpose using the best transpose algorithm according to the shape of a tensor and the permutation indices. Again based on *cutt* by *Dmitry I. Lyakh*.
 - Batched transpose.
 - Permuted data stay on a GPU.
- Smart multi-gpu data placement algorithm for vector-form contractions.



ACM TACO Journal paper



IPDPS 22



Performance of Hadron Contraction on GPUs









Performance of Hadron Contraction on GPUs



Jefferson Lab





Redstar_npt execution time for 1-timeslice (one GPU)

- Recent redstar_npt runs on our test configurations.
- KNL 7250 host, MI100, MI250 and A100.

#	Name	N _{ref}	T _{orig}	KNL	MI100	MI250	A100
1	a ₀	256	260.7	87.95	11.94	8.23	5.96
2	f _o	256	412.2	138.89	21.14	12.63	10.55
3	a1	128	1427.11	866.70	107.56	64.28	63.22
4	Roper	10 64	778.2	3488.05	1933.64	219.86	1238.44
5	Deuterium	10 64	1471.9	3637.48	517.21	286.58	191.21
6	Tritium	10 32	1735	8208.01	2970.04	1803.66	1269.93



Redstar Performance





Redstar Performance



Performance speed-up values of a GPU over a single KNL host for calculating Correlation functions of multi-hadron systems.



Redstar Performance (Figure of Merit)





Redstar Performance on Multiple GPUs



NPT_BATCH_SIZE = Number of GPUs 31



Redstar Performance on Multiple GPUs: Impact of data movements

- For meson systems, no gpu memory evictions. The performance degradation is caused by GPU to GPU/Host movements and relative small matrix size (N=256 for a_0 , f_0 and N=128 for a_1).
- Define T_s as time to calculate contractions on a single GPU, BW as the effective unidirectional bandwidth between GPUs/Host, N_{flops} number of flops for all contractions, CN_{gflops} as the gflops/s for one contraction for matrix size N, Nbytes number of bytes transferred, T_n as the execution time on n gpus.

	a₀ 2gpu	a₀ 4gpu	f₀ 2gpu	f ₀ 4gpu	a₁ 2gpu	a₁ 4gpu
T _n	13.75	15.78	27.21	29.01	189.28	207.14
Ts	11.94	11.94	21.14	21.14	107.56	107.56
N _{flops}	2.52*10 ¹³	2.52*10 ¹³	5.79*10 ¹³	5.79*10 ¹³	5.26*10 ¹³	5.26*10 ¹³
CNgflops	3991	3991	3991	3991	1699	1699
Nbytes(GB)	115.63	145.19	195.66	236.20	681.93	1051.94
BW (GB/s)	23.28	22.18	14.68	12.59	7.01	8.56

 $T_n = T_s - N_{flops} / (CN_{gflops} * 10^9) + N_{flops} / (CN_{gflops} * 10^9) / n + Nbytes/BW$







Office of Science

Redstar Performance on Multiple GPUs

• For Baryon systems, data transfers and gpu memory evictions affect the performance.

Name	# GPUs	Ν	# Flops	Evictions	Bytes (GB)	Time (sec)
	1			32356		1933.64
Roper	2	64	2.83x10 ¹⁴	15290	3200	1073.32
	4			1671	5725	429.06
	1			17767		517.21
Deuterium	2	64	1.42x10 ¹⁵	235	8986	522.56
	4			667	14552	425.58
	1			2749		2970.04
Tritium	2	32	7.69x10 ¹⁵	2671	753	1790.77
	4			2354	1113	1323.79







Office of Science

Reduce Data Movements on Multi-GPUs by Partitioning the DAG

• Divide the DAG into multiple similar sized partitions with minimal number of edges crossing the borders of the partitions.







Near Term Challenges

Redstar and hadron contractions

- Reduce data movements on multi-gpu.
 - Hopefully partitioning the DAG will reduce the data movements.
 - Each partition of contractions is handled by a GPU through a separate contraction queue.
- Baryon-Baryon contractions on GPUs
 - One single contraction cannot fit on a single GPU.
 - Using multiple GPUs need to move partial GPU contraction results back to the host.
 - Memory evictions becomes difficult because each object is not small.
 - When to split a contraction on multiple GPUs ?
- Graph Isomorphism on GPUs?
 - Can we beat Nauty ?
- Optimal graph contractions.
 - Is local optimal selection = global optimal ?
- Multi-host redtar_npt through MPIs.
 - Still need to do partitioning the DAG (which host to do which parts of the DAG).
- Intel CPU/GPUs.





