

LaRIS: Targeting Portability and Productivity for LAPACK Codes on Extreme Heterogeneous Systems by Using IRIS

Mohammad Alaul Haque Monil, Narasinga Rao Miniskar, Frank Y. Liu, Jeffrey S. Vetter, Pedro Valero-Lara Computer Science and Mathematics Division, Oak Ridge National Laboratory

Presenter: Mohammad Alaul Haque Monil RSDHA Workshop, 2022

Introduction

- Heterogeneous systems come in variety
 - Ensuring portability and utilization are challenging
- Tuned vendor libraries for common kernels are not portable
 - Linear algebra kernels
- Making tiled linear algebra algorithm portable is a challenge
- There are solutions that provides tight coupling between algorithm and kernels
 - Not productive; changing source for every library support
- A separation between algorithm development and vendor library kernels is needed

What's Missing?

- Portability
 - Raja
 - Kokkos
 - OpenMP offload
- Heterogeneity
 - StarPU
 - Does not support hip/rocBLAS
 - Tightly coupled
 - MAGMA (CPU + GPU*)
 - hipMAGMA (same)

- Linear Algebra
 - PLASMA (CPU)
 - libFLAME (CPU)
 - Intel MKL (CPU)
 - ATLAS (CPU)
 - cuBLAS/cuSolver (GPU)
 - Hip/roc BLAS/Solver (GPU)
 - Others (Chameleon, etc.)

- Missing?
 - A solution that combines portability and multi-device heterogeneity for LAPACK codes

Research Question

Can we separate linear algebra algorithms from kernels to make seamless portability and utilization possible for different heterogeneous systems?



This Research

- A proof of concept focusing on portability and heterogeneity
- LaRIS LAPACK on IRIS runtime
- A test case -- considering LU factorization



[20] P. Valero-Lara, S. Catalan, X. Martorell, T. Usui, and J. Labarta. slass: 'A fully automatic auto-tuned linear algebra library based on openmp extensions implemented in ompss (lass library). J. Parallel Distributed Comput., 138:153–171, 2020

IRIS Runtime

- IRIS Runtime
 - Heterogeneous runtime
 - Supports CPU, GPU, FPGA and DSP
 - Orchestrate computation dynamically
- IRIS-BLAS
 - BLAS library API for different vendor libraries
 - Supports
 - hip/rocBLAS
 - cuBLAS
 - OpenBLAS
 - MKL
 - Under development



LaRIS Software Stack

- LaRIS: Unified API for portable LaPACK codes
 - Targets Four capabilities
 - Algorithm links vendor tuned kernels at runtime
 - Facilitated by IRIS runtime and IRIS-BLAS



LaRIS: LU Factorization

- Three types of kernels
 - GETRF Red
 - TRSM Yellow
 - GEMM Blue
- Graph expressed with dependency
- Dynamic scheduling
 - IRIS_ANY
- Automatic tiling capability



Figure 4: DAG of an 8×8 LU factorization that creates 204 tasks and their dependencies. Red ellipses are GETRF, orange ellipses are TRSM, and blue ellipses are GEMM.

LaRIS: Systems Used

Table I: Heterogeneous sy	stems used	in	this	research
---------------------------	------------	----	------	----------

Nodes	NVIDIA DGX-1	Frontier-like node	CADES node
Facility	ExCL at ORNL	Crusher (early access) at ORNL	CADES at ORNL
	Total 4 GPUs	Total 8 GPUs	Total 8 GPUs
GPUs	4× NVIDIA V100	$4 \times$ AMD MI250X	$4 \times$ NVIDIA A100
		(each contains two)	$4 \times$ AMD MI100
CPU	Intel Xeon E5-2698, 20 cores	AMD EPYC 7A53, 64 cores	AMD EPYC 7763, 128 cores
Compiler	GNU-9.4.0	GNU-8.5.0	GNU-8.5.0
CUDA and ROCm versions	CUDA-11.7	CUDA-11.7 and ROCm-5.1.0	CUDA-11.7 and ROCm-5.1.2
Math libraries for GPUs	cuBLAS and cuSOLVER	hipBLAS	cuBLAS, cuSOLVER and hipBLAS
Math libraries	OpenBLAS-0.3.20 and	OpenBLAS-0.3.17 and	OpenBLAS-0.3.20 and
for CPU	MKL-2022.1.0	MKL-2020.4.304	MKL-2020.4.304

LaRIS: Portability

- 16384 x 16384 matrix
- 32 x 32 tile = 11440 tasks
- Portability
 - CPUs (AMD and Intel)
 - GPUs
 - NVIDIA (A100 and V100)
 - AMD (MI100 and MI250X)
- Colors
 - GETRF Red
 - TRSM Yellow
 - GEMM Blue



LaRIS: Zoomed-in Trace of the Cades Node



• Main Message: Trace shows the utilization of all the heterogeneous processors without much gap in the timeline

LaRIS: Strong Scaling





- Main Message: Increased number of processors provide better performance, however, reaches a saturation. Why?
 - Tradeoff between computation capability and overhead

LaRIS: Decomposition vs Matrix Sizes



Figure 7: Finding the best tile configuration for different matrix sizes when all the processing elements are used.

- Main Message: Different system has different best tile decomposition that depends on Matrix sizes.
- Trade-off between parallelism vs overhead

64 x 64 ~ Total 90,000 task and kernel calls Table II: Best tile configuration for different matrix sizes

Matrix	DGX-1	Crusher node	CADES node
8,192 imes 8,192	8×8	4×4	4×4
$16,384\times16,384$	8×8	8×8	8×8
$32,768\times32,768$	8×8	16×16	16×16
$65,536\times65,536$	NA	16×16	16×16

Discussion and Future Opportunities

- This work focuses on portability and heterogeneity
 - Kernel level performance is ensured by the tuned vendor library
 - Optimization opportunities at runtime level
- Reduction of runtime overhead
 - Memory transfer optimization
 - Scheduling algorithm improvement
- Fresh from the Oven
 - We are working on both of these optimization and preliminary result is promising (unpublished)
 - We are working towards a public release as MatRIS
 - Combines LaRIS and IRIS-BLAS on top of IRIS

Conclusion

- LaRIS separates LAPACK algorithms from vendor specific libraries
- Seamless portability is obtained in LaRIS through IRIS runtime and IRIS-BLAS
- IRIS runtime schedular can utilize all the processors in a heterogeneous system to compute LU factorization
- Portability and heterogeneity is demonstrated as a proof of concept, however, there is scope for performance improvement.

Acknowledgments

This research was supported by DOE BLUESTONE and ECP PROTEAS-TUNE project



Thank you. Question?



// Encapsulate trsm tasks into the graph
laris_trsm_top_graph(graph, trsm_tasks[(step*
TILE_NUM) +
tile_jj], step, step, tile_jj,
TILE_SIZE, TILE_SIZE,
IRIS_Ctile[step*TILE_NUM+step],
A tile[step][step], TILE SIZE,
TRIS Ctile[step*TILE NUM+tile ii].
A tile[step][tile_ii], TILE_SIZE):
)
for (int tile ii = 1 + step: tile ii < TILE NUM:
tile ii++) /
// trom tacko
iris task create/
<pre>stram tasks[(tile ii*TILE NUM)+step1);</pre>
// Dependencies for trem tasks
if (step == 0) /
irio tack trom depend tacke[] = (
<pre>getrf tasks[(step + TILE NUM) + step]).</pre>
iris task depend/ tram tasks[(tile ii +
TILE NUM) + step] 1 trem depend tasks):
l else /
inia task twom demond tasks[] = (
rris_cask crsm_depend_casks() = (
getri_tasks((step*liLE_NOM)+step),
demm_tasks(((step-1)*TILE_NOM*TILE_NOM)+
(crie_ri*rink_NOM)+scep) };
iris_task_depend(
trsm_tasks[(tile_11 * TiLE_NUM) + step],
<pre>2, trsm_depend_tasks);</pre>
// Warners later bracks into the same
// Encapsulate trsm tasks into the graph
laris_trsm_left_graph(graph,
trsm_tasks[(tile_11*TILE_NUM)+step],
step, tile_11, step, TILE_SIZE, TILE_SIZE,
<pre>IRIS_Ctile[step * TILE_NUM + step],</pre>
A_tile[step][step], TILE_SIZE,
<pre>IRIS_Ctile[tile_ii * TILE_NUM + step],</pre>
A_tile[tile_ii][step], TILE_SIZE);
}
<pre>for (int tile_i = step + 1; tile_i < TILE_NUM;</pre>
tile_i++) {

for (int tile_j = step + 1; tile_j < TILE_NUM : tile_j++) { // gemm tasks iris task create (&gemm tasks [(step * TILE NUM *TILE_NUM) + ((tile_i*TILE_NUM)+tile_j)]); // Dependencies for gemm tasks if (step == 0) { iris_task gemm_depend_tasks[] = { trsm_tasks[(tile_i*TILE_NUM)+step], trsm_tasks[(step*TILE_NUM)+tile_j] }; iris_task_depend(gemm_tasks[(step*TILE_NUM*TILE_NUM)+((tile_i* TILE_NUM)+tile_j)], 2, gemm_depend_tasks); else (iris_task gemm_depend_tasks[] = { trsm_tasks[(tile_i*TILE_NUM)+step], trsm_tasks[(step*TILE_NUM)+tile_j], gemm_tasks[((step-1)*TILE_NUM* TILE_NUM) + ((tile_i*TILE_NUM)+tile_j)) }; iris_task_depend(gemm_tasks[(step*TILE_NUM*TILE_NUM)+ ((tile_i*TILE_NUM)+tile_j)], 3, gemm_depend_tasks); // Encapsulate gemm tasks into the graph laris_gemm_graph(graph, gemm_tasks[(step*TILE_NUM*TILE_NUM)+ ((tile_i*TILE_NUM)+tile_j)], step, tile_i, tile_j, TILE_SIZE, TILE_SIZE, TILE_SIZE, -1.0, IRIS_Ctile[tile_i*TILE_NUM+step], A_tile[tile_i][step], TILE_SIZE, IRIS_Ctile[step*TILE_NUM+tile_j], 114 A tile[step][tile_j], TILE_SIZE, 1.0, IRIS_Ctile[tile_i*TILE_NUM+tile_j], A_tile[tile_i][tile_j], TILE_SIZE); 119 } 120 } iris_graph_submit(graph, iris_default, 1); Listing 1: LaRIS's tiled LU factorization (GETRF) code.

2, trsm_depend_tasks);

LaRIS: Crusher and DGX-1



LaRIS: Zoomed in Graph

