



Close Coupling of AOCL-BLAS in AOCL-LAPACK

Sridhar G
Pradeep Rao
Vasanth Kumar R

AMD 
together we advance_

AGENDA

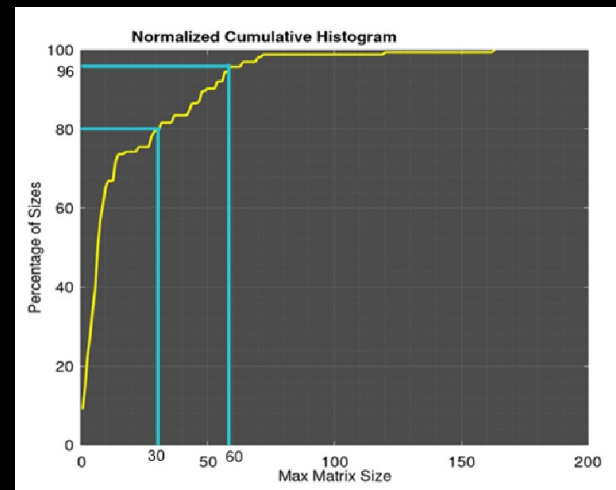
- Introduction
- Motivation
- Observation
- Proposed Solution
- Results
- Challenges
- Q & A

Introduction

- ❖ AOCL (AMD optimizing CPU libraries) is AMD's CPU Math Library, tuned for AMD processor
- ❖ AOCL-LAPACK, a fork of libFLAME library (repository from UT Austin), is optimized as part of AOCL
- ❖ AOCL-LAPACK depends on AOCL-BLAS (a fork of BLIS library optimized as part of AOCL) for various low-level vector and matrix operations
- ❖ Our last public release is AOCL 4.2:
 - Improved performance of GESVD, GETRF, GETRS and GESV
 - OpenMP parallelism enabled in {c,z}hetrd, {s,d}sytrd_sb2st and iparam2stage
 - Option to link with AOCL-BLAS during build to enable invoking AOCL_BLAS internal APIs
 - Cmake improvements
 - Test suite framework enhancement
- ❖ Latest AOCL-LAPACK source code is available on github: <https://github.com/amd/libflame>

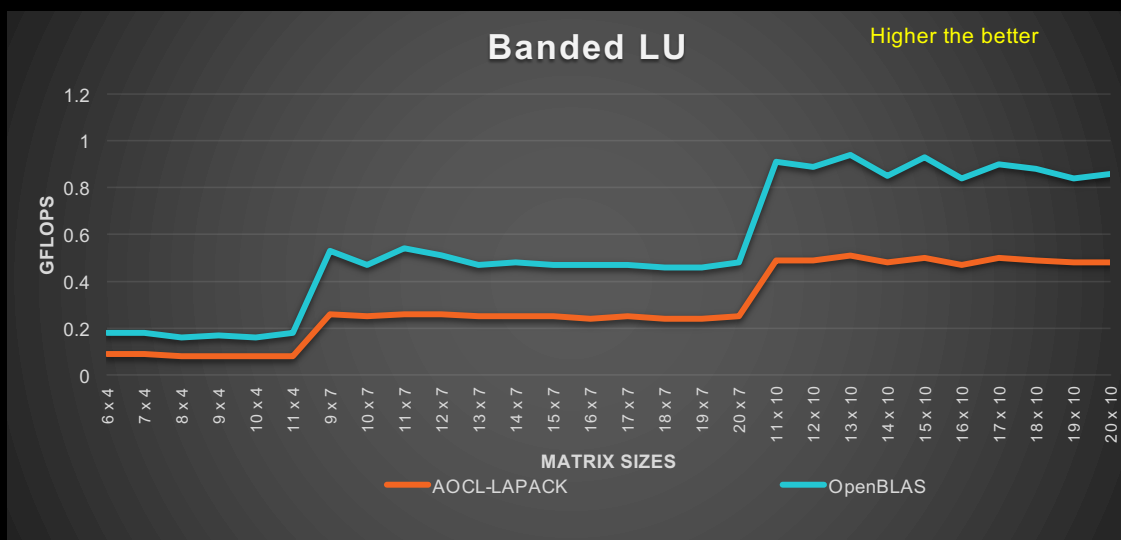
Motivation (1/2)

- ❖ Performance of AOCL-LAPACK is competitive with OpenBLAS for medium and large matrix inputs, especially with well tuned AOCL-BLAS library
- ❖ OpenBLAS is an optimized open-source implementation of BLAS and LAPACK libraries
- ❖ Problem sizes commonly encountered in many of our client benchmarks are small
- ❖ Graph depicts the problem sizes for DGBTRF API used in one of these benchmarks



Motivation (2/2)

How is the performance of DGBTRF for these sizes?



❖ We are, on an, average 50% behind openBLAS for small matrix sizes

LU Decomposition

❖ Definition of LU for matrix A:

$A = L * U$ where, L – Lower triangular matrix
 U – Upper triangular matrix

$$\begin{pmatrix} a_{00} & a_{01} & a_{02} \\ a_{10} & a_{11} & a_{12} \\ a_{20} & a_{21} & a_{22} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ l_{10} & 1 & 0 \\ l_{20} & l_{21} & 1 \end{pmatrix} \begin{pmatrix} u_{00} & u_{01} & u_{02} \\ 0 & u_{11} & u_{12} \\ 0 & 0 & u_{22} \end{pmatrix}$$

❖ Operations in LU for each column:

- (1) Partition input A as in top part of figure
- (2) Update $l_{21} = a_{21}/\alpha_{11}$ (*DSCAL*)
- (3) $u_{11} = \alpha_{11}$, $u_{12}^T = a_{12}^T$
- (4) Update $A'_{22} = A_{22} - l_{21} * u_{12}^T$ (*DGER*)
- (5) $A = A'_{22}$ and repeat steps (1) to (5)

$$A = \left(\begin{array}{c|c} \alpha_{11} & a_{12}^T \\ \hline a_{21} & A_{22} \end{array} \right)$$

$$L = \left(\begin{array}{c|c} 1 & 0 \\ \hline l_{21} & L_{22} \end{array} \right) \quad U = \left(\begin{array}{c|c} u_{11} & u_{12}^T \\ \hline 0 & U_{22} \end{array} \right)$$

Banded LU Decomposition (GBTRF)

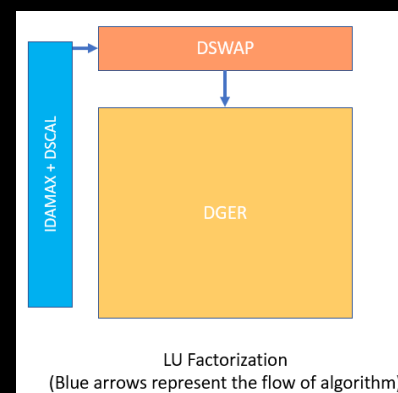
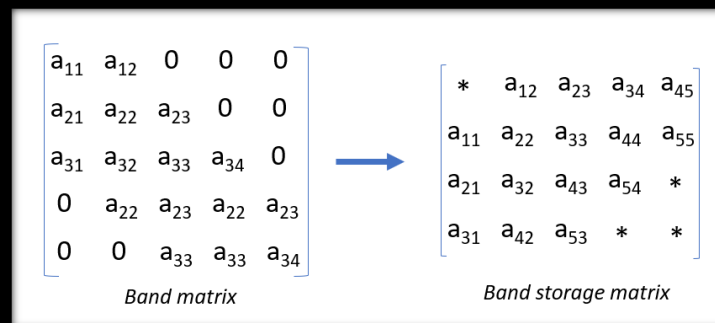
❖ GBTRF is a variant of LU factorization API in LAPACK, which performs factorization on banded storage matrix

❖ Advantages:

- **Space Efficiency:** Only the non-zero elements are stored, which significantly reduces memory usage for large matrices
- **Performance:** Accessing elements within the band is faster due to the compact storage

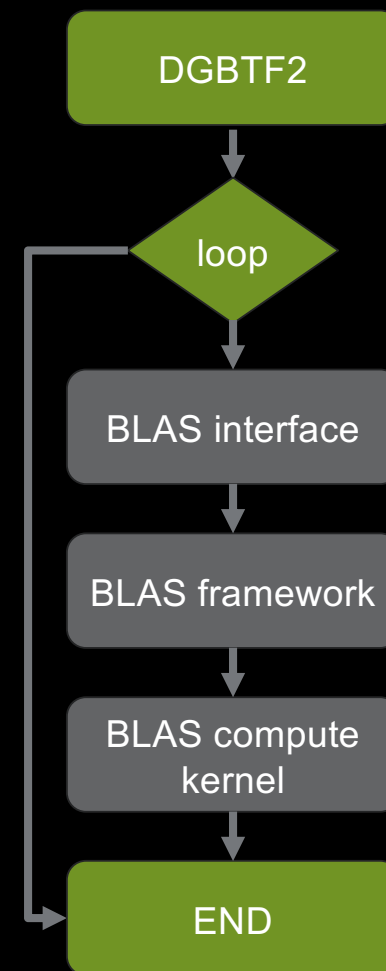
❖ Column wise computation involves:

- Pivoting and scaling: IDAMAX() + DSCAL()
- Row swapping: DSWAP()
- Rank-1 operation: DGER()



Observation (1/2)

- ❖ Existing design scales well for medium and large matrix inputs but shows significant performance degradation for small matrix inputs
- ❖ Although the relevant AOCL-BLAS APIs are optimized for small sizes, LAPACK layer was unable to fully leverage this optimization
- ❖ Profiling revealed around 30% to 40% of overhead was from AOCL-BLAS framework related functions
- ❖ When AOCL-BLAS API is called, a series of framework related functions are executed before the actual computational kernel runs
- ❖ These framework functions are essential for selecting the appropriate kernel based on factors such as machine instruction support and problem size

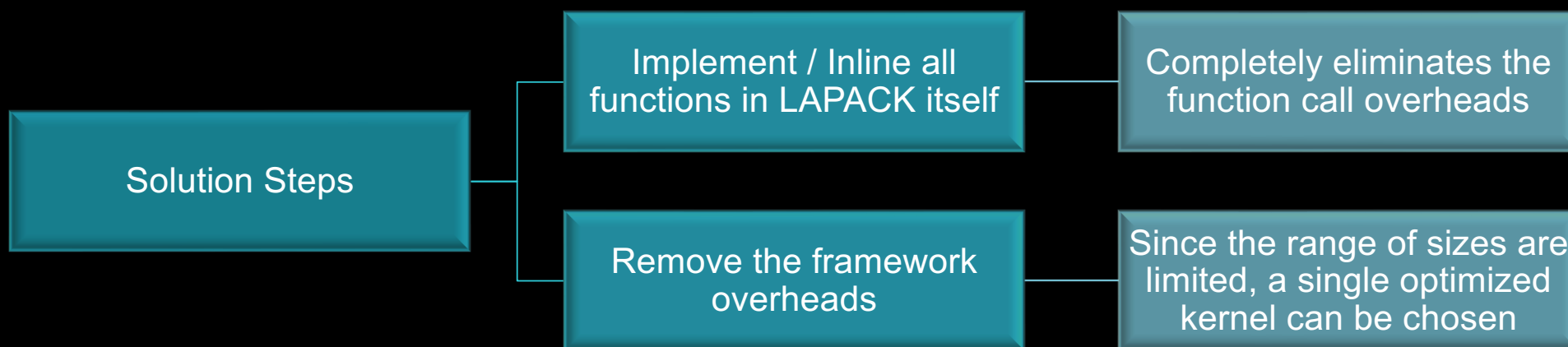


Observation (2/2)

- ❖ LU factorization is an iterative algorithm, and the computation of each column involves calling the same set of AOCL-BLAS APIs
- ❖ For small matrix inputs, the repeated overhead from framework functions associated with each AOCL-BLAS API call significantly affect the overall performance of the API
- ❖ Our goal was to minimize the overhead from the framework functions as much as possible to improve the performance

Proposed solution 1

Inline the BLAS compute kernels directly into AOCL-LAPACK



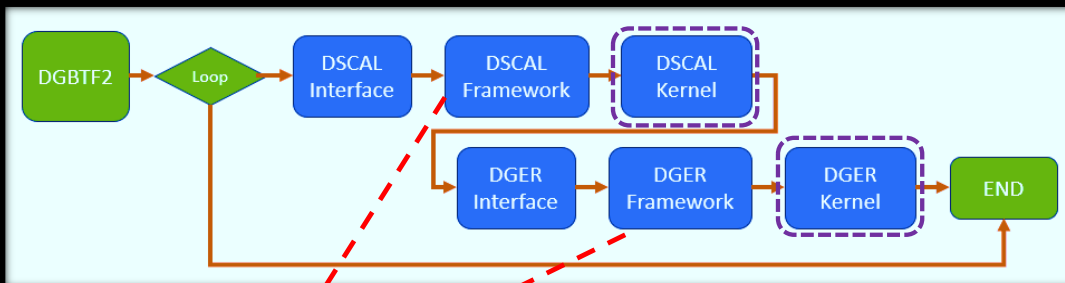
Drawbacks:

- ❖ It is a challenging task to apply across the library
- ❖ This solution may lead to code fragmentation – must be used only if required

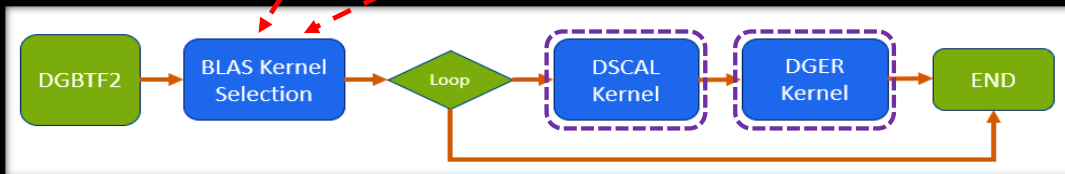
Proposed Solution 2 (1/2)

- ❖ Instead of calling the standard interface of AOCL-BLAS APIs, we propose invoking the corresponding compute kernels directly from AOCL-LAPACK
- ❖ Framework to choose the right AOCL-BLAS compute kernel will now be executed within AOCL-LAPACK layer
- ❖ AOCL-BLAS Framework code to select the appropriate compute kernels will now be executed only once per AOCL-LAPACK API call

Original Code Flow:



Optimized Code Flow:



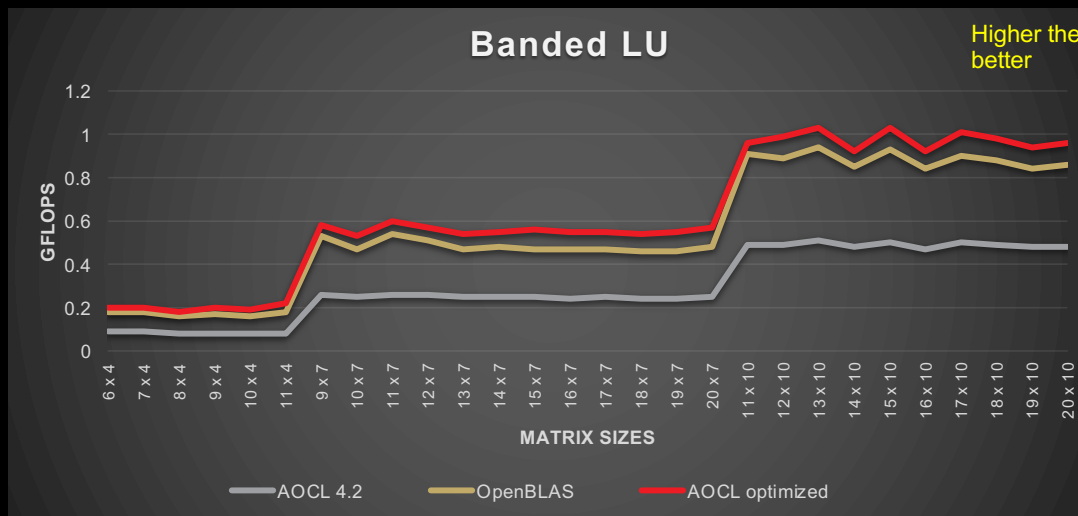
Proposed Solution 2 (2/2)

- ❖ We observed significant reduction in the number of instruction executed per AOCL-BLAS API call, with an average reduction of 60% to 80%

BLAS API	Instruction executed per call		Instruction difference
	Standard interface	Kernel Interface	
DGER	2056	306	-85.11%
IDAMAX	4734	579	-87.76%
DSCAL	480	165	-65.62%
DSWAP	385	71	-81.55%

- ❖ This solution allows LAPACK to fully leverage the computational power of AOCL-BLAS APIs, particularly for small input sizes
- ❖ From an Implementation perspective, it is relatively straightforward, involving only the retrieval and invocation of appropriate AOCL-BLAS compute kernels
- ❖ We were able to achieve better performance while preserving the modularity between the libraries

Results



- ❖ Performance improved on an average by 110%
- ❖ AOCL-LAPACK now outperforms OpenBLAS for small matrix sizes

Challenges

- ❖ Not all AOCL-BLAS compute kernels, and framework functions are exposed. We must ensure all the relevant code is open to AOCL-LAPACK
- ❖ Any kernel addition or modification must be appropriately updated within the AOCL-BLAS framework context, as AOCL-LAPACK is reliant on it for kernel selection

Questions

COPYRIGHT AND DISCLAIMER

2024 Advanced Micro Devices, Inc. All rights reserved.

AMD, the AMD Arrow logo, EPYC and combinations thereof are trademarks of Advanced Micro Devices, Inc. Other product names used in this publication are for identification purposes only and may be trademarks of their respective companies.

The information presented in this document is for informational purposes only and may contain technical inaccuracies, omissions, and typographical errors. The information contained herein is subject to change and may be rendered inaccurate releases, for many reasons, including but not limited to product and roadmap changes, component and motherboard version changes, new model and/or product differences between differing manufacturers, software changes, BIOS flashes, firmware upgrades, or the like. Any computer system has risks of security vulnerabilities that cannot be completely prevented or mitigated. AMD assumes no obligation to update or otherwise correct or revise this information. However, AMD reserves the right to revise this information and to make changes from time to time to the content hereof without obligation of AMD to notify any person of such revisions or changes.

THIS INFORMATION IS PROVIDED 'AS IS.' AMD MAKES NO REPRESENTATIONS OR WARRANTIES WITH RESPECT TO THE CONTENTS HEREOF AND ASSUMES NO RESPONSIBILITY FOR ANY INACCURACIES, ERRORS, OR OMISSIONS THAT MAY APPEAR IN THIS INFORMATION. AMD SPECIFICALLY DISCLAIMS ANY IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY, OR FITNESS FOR ANY PARTICULAR PURPOSE. IN NO EVENT WILL AMD BE LIABLE TO ANY PERSON FOR ANY RELIANCE, DIRECT, INDIRECT, SPECIAL, OR OTHER CONSEQUENTIAL DAMAGES ARISING FROM THE USE OF ANY INFORMATION CONTAINED HEREIN, EVEN IF AMD IS EXPRESSLY ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

