

KBLAS

KAUST BLAS (KBLAS) is a high performance CUDA library implementing a subset of **BLAS** as well as **LAPACK** routines on NVIDIA GPUs. Using recursive and batch algorithms, KBLAS maximizes the GPU bandwidth, reuses locally cached data and increases device occupancy. KBLAS supports operations for regular dense and hierarchical low-rank matrix workloads. KBLAS provides, therefore, the critical building blocks not only for the traditional high performance dense linear algebra libraries, but also for the emerging numerical libraries supporting hierarchical low-rank matrix computations. This is a major leap toward leveraging GPU hardware accelerators for high performance low-rank matrix approximations and computations, which currently drives the research agenda of the linear algebra scientific community.

RECURSIVE ALGORITHMS: GPU-Resident POTRF

```

Procedure POTRF_rec(uplo, n, A, lda, info)
    if n ≤ 32 then           // at the recursion stop size
        blockDim ← dim3(32, 32)
        POTRF_Kernel <<< 1, blockDim >>> (uplo, n, A, lda,
        info)
    else                      // invoke recursion
        n1 ← largest power of 2 less than n
        n2 ← n - n1
        POTRF_rec(uplo, n1, A, lda, info) // recursive POTRF
        TRSM( Right, Lower, Trans, NonUnit,
            n2, n1,
            1, A, lda,
            A + n1, lda)
        SYRK( Lower, NoTrans,
            n2, n1,
            -1, A + n1, lda,
            1, A + n1 * lda + n1, lda)
        POTRF_rec(uplo, n2, A + n1 * lda + n1, lda)
        // recursive POTRF
    return;

```

KBLAS 3.0

- ◆ Legacy Level-2 BLAS: [‡ ⊕ ⊖ ⊚] SYMV, GEMV, HEMV.
- ◆ Legacy Level-3 BLAS: [‡ ⊕ ⊖ ⊚] TRSM, TRMM, GEMM (⊖ only).
- ◆ Batch Level-3 BLAS: [‡ ⊕ ⊖ ⊚ = *] TRSM, TRMM, SYRK.
- ◆ Batch Triangular: [⊖ ⊕ ⊖ ⊚ = *] TRTRI, LAUUM.
- ◆ Batch Symmetric: [⊖ ⊕ ⊖ ⊚ = *] POTRF, POTRS, POSV, POTRI, POTI.
- ◆ Batch General: [⊖ ⊕ ⊖ ⊚ = *] GESVJ, GERSVD, GEQRF.
- ◆ **Batch Tile low-rank GEMM** [⊖ ⊕ ⊖ ⊚ =].
- ◆ **GPU-Resident POTRF kernel** [⊖ ⊕ ⊚].

‡ Standard precisions: s/d/c/z.

⊕ Real precisions: s/d.

⊖ Very small matrix sizes.

= Arbitrary sizes.

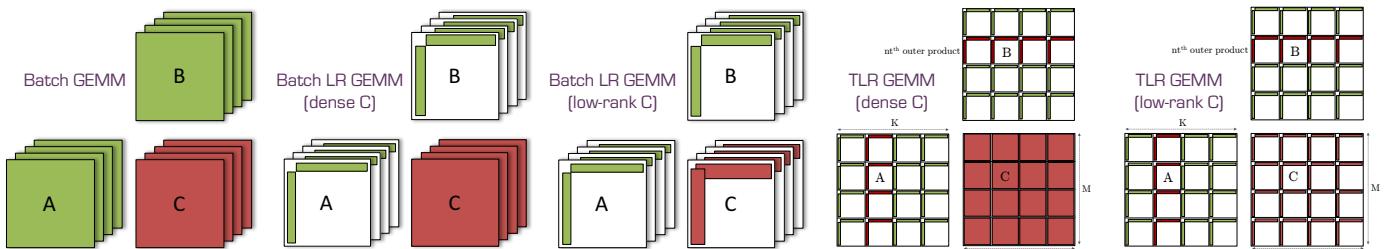
⊖ Single-GPU support.

⊕ Multi-GPU support.

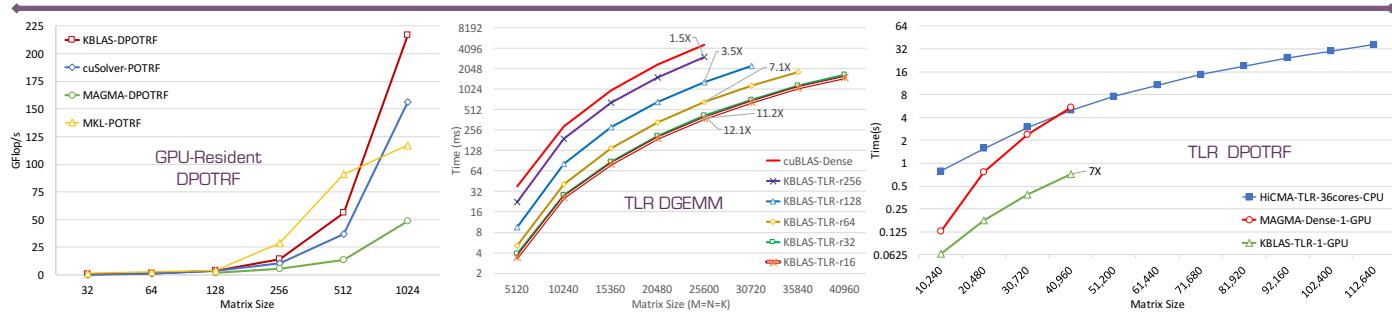
= Uniform batch sizes.

* Non-strided and strided variants.

BATCH ALGORITHMS: Tile Low-Rank GEMM – Bottom-Up Approach



PERFORMANCE RESULTS – NVIDIA V100 GPU



DOWNLOAD KBLAS AT: <https://github.com/ecrc/kblas>



A collaboration of



With support from



Sponsored by



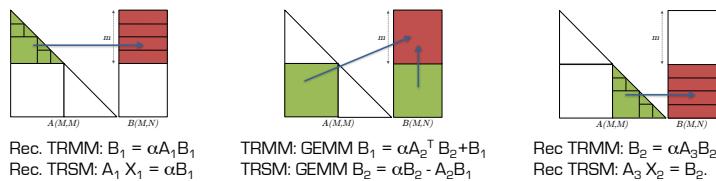
KBLAS



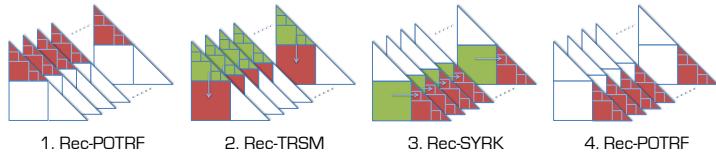
Extreme Computing
Research Center

KAUST BLAS (KBLAS) is a high performance CUDA library implementing a subset of BLAS as well as Linear Algebra PACKage (LAPACK) routines on NVIDIA GPUs. Using recursive and batch algorithms, KBLAS maximizes the GPU bandwidth, reuses locally cached data and increases device occupancy. KBLAS supports operations for regular dense and hierarchical low-rank matrix workloads. KBLAS provides

RECURSIVE ALGORITHMS: TRMM and TRSM



BATCH ALGORITHMS: Recursive Cholesky POTRF

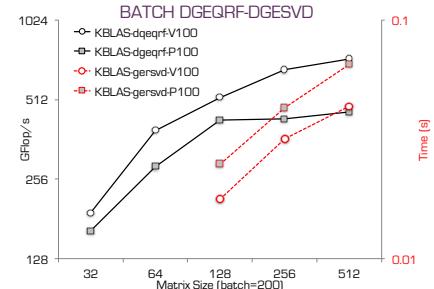
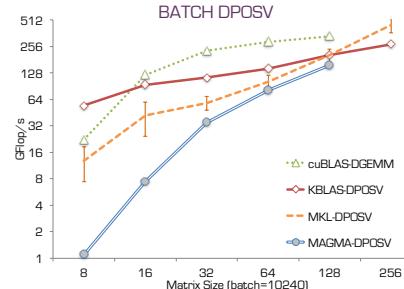
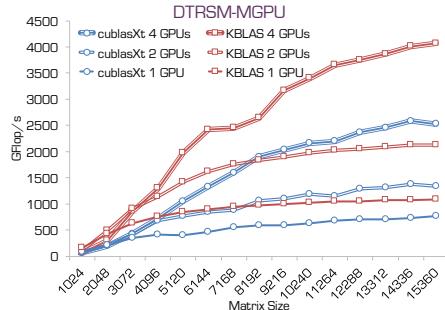
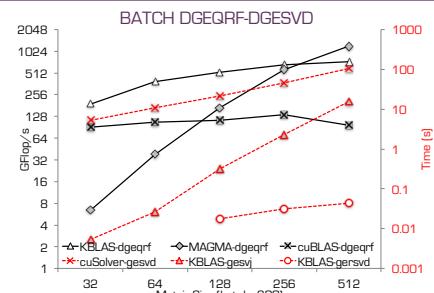
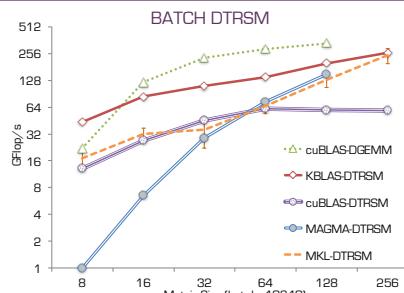
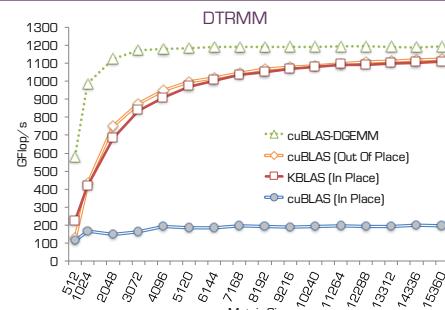


KBLAS HIGHLIGHTS

- ◆ KBLAS Level-2 (o): SYMV & HEMV
- ◆ KBLAS Level-3 (o): TRMM & TRSM



PERFORMANCE RESULTS



DOWNLOAD KBLAS AT: <https://github.com/ecrc/kblas>

KBLAS 2.0

- ◆ Legacy Level-2 BLAS: ($\ddagger \diamond \otimes o$) SYMV, GEMV, HEMV.
- ◆ Legacy Level-3 BLAS: ($\ddagger \diamond \otimes o$) TRSM, TRMM, GEMM (\otimes only).
- ◆ Batch Level-3 BLAS: ($\ddagger \diamond \otimes o = *$) TRSM, TRMM, SYRK.
- ◆ Batch Triangular: ($\diamond \ddagger \otimes o = *$) TRTRI, LAUUM.
- ◆ Batch Symmetric: ($\diamond \ddagger \otimes o = *$) POTRF, POTRS, POSV, POTRI, POTI.
- ◆ Batch General: ($\otimes \ddagger \otimes o = *$) GESVJ, GERSVD, GEQRF.

\ddagger Standard precisions: s/d/c/z. \diamond Single-GPU support.
 \ddagger Real precisions: s/d. \otimes Multi-GPU support.
 \diamond Very small matrix sizes. $=$ Uniform batch sizes.
 \diamond Arbitrary sizes. $*$ Non-Strided and Strided variants.

CURRENT RESEARCH

- ◆ Half Precision Legacy and Batch BLAS.
- ◆ Tile Low-Rank (TLR) BLAS on GPUs.
- ◆ Adaptive Cross Approximation (ACA) on GPUs.
- ◆ Vectorized Batch BLAS on x86.



A collaboration of



With support from



Sponsored by

