



# WORKSHOP ON BATCHED, REPRODUCIBLE, AND REDUCED PRECISION BLAS

May 18th - 19th, 2016

## REFERENCE IMPLEMENTATION AND TESTING OF BATCHED BLAS ROUTINES



The University of Manchester  
Innovative Computing Laboratory  
University of Tennessee  
and NLAFET

## Reference implementation

- Demonstrate that the specification is implementable
- Help to improve the specification
- Provide a reliable code to validate other implementations

## Testing

- Valid calling sequences
- Correctness of the computed results

## Remarks

The reference implementation does not aim to be efficient

This presentation focuses on Level 3 batched BLAS routines

1. Reference implementation
2. Accuracy testing
3. Experimental results
4. Concluding remarks

1. Reference implementation
2. Accuracy testing
3. Experimental results
4. Concluding remarks

# Overview of the reference implementation

## 1. Calling sequences with respect to the specification

```
dgemm_batch(const enum * transA, const enum * transB, const int * m,  
            const int * n, const int * k, const double * alpha,  
            const double * const double * arrayA, const int * lda,  
            const double * const double * arrayB, const int * ldb,  
            const double * beta, double **arrayC, const int * ldc,  
            const int batch_count, const enum batch_opts, int * info)
```

- **batch\_count**: style for the batched (**BATCH\_FIXED**, **BATCH\_VARIABLE**)
- **batch\_opts**: number of sub-problems to be processed
- **info**: error handling array

# Overview of the reference implementation

1. Calling sequences with respect to the specification
2. Arguments checking and error handling

## Critical arguments

- Fixed size case: **exit** if any argument is incorrect
- Variable size case:
  - ▶ **Exit** for whole batch error (e.g. `batch_count` )
  - ▶ **Continue** for sub-problem error (e.g. `matrix size`)
- **Xerbla** for error handling
- Flexible error handling strategy

# Overview of the reference implementation

1. Calling sequences with respect to the specification
2. Arguments checking and error handling
3. Call **reference** BLAS implementation to solve each sub-problem

```
for (iter = 0; iter < batch_count; iter++)  
{  
    Call a reference BLAS routine  
}
```

According to the specification, the reference implementation should provide results as accurate as reference BLAS routines

1. Reference implementation
2. Accuracy testing
3. Experimental results
4. Concluding remarks



## Basic idea

- The quantity of interest is the relative error  $err = \frac{\|C - \hat{C}\|}{\|C\|}$
- Typically we want  $err \leq \tau \varepsilon$  where  $\tau$  is some tolerance and  $\varepsilon$  is the machine epsilon e.g.  $\varepsilon = 10^{-16}$

## Remarks

- The tolerance  $\tau$  is selected independently for each problem, e.g.  $\tau = 10, 30, 100$  in MAGMA/LAPACK etc
- The choice of  $\tau = 10, 30 \dots$  may not be very rigorous
- Investigate a rigorous approach for small BLAS problems?

## Error analysis approach

- Perform error analysis of BLAS routines [Nicholas J. Higham, 2002]
- Use the **forward/backward error bound** for accuracy testing
- $\|Forward\_Error\| \leq \|Forward\_Error\_Bound\|$  or  
 $\|Residual\| \leq \|Backward\_Error\_Bound\|$

## Remarks

- No need of intensive experiments to set a tolerance  $\tau$
- A rigorous error bound can seem pessimistic in general but meaningful for small BLAS problems

## Forward error bound of GEMM

- Formula:  $C = \alpha AB + \beta C_{init}$
- Forward error  $\|C - \hat{C}\|_{\infty}$
- Forward error bound:  $(N|\alpha|\|A\|_{\infty}\|B\|_{\infty} + |\beta|\|C_{init}\|_{\infty})\epsilon$

## Remark

The computed solution  $\hat{C}$  is acceptable if

$$\frac{\|C - \hat{C}\|_{\infty}}{(N|\alpha|\|A\|_{\infty}\|B\|_{\infty} + |\beta|\|C_{init}\|_{\infty})\epsilon} \leq 1$$

## Backward error analysis of TRSM [Nicholas J. Higham, 2002]

**Theorem:** Let the triangular systems  $Tx = b$  where  $T \in \mathbb{R}^{n \times n}$  is nonsingular, be solved by substitution with any ordering. Then the computed solution  $\hat{x}$  satisfies  $(T + \Delta T)\hat{x} = b$ ,  $|\Delta T| \leq \gamma_n |T|$ , where  $\gamma_n = \frac{n\varepsilon}{1-n\varepsilon} \approx n\varepsilon$ .

## Backward error bound of TRSM

- Formula:  $AX = \alpha B$
- Backward error bound:  $N \|A\|_\infty \|X\|_\infty \varepsilon$
- Note: TRSM's backward error bound does depend on  $B$
- The solution  $\hat{X}$  is acceptable if  $\frac{\|A\hat{X} - \alpha B\|_\infty}{N \|A\|_\infty \|\hat{X}\|_\infty \varepsilon} \leq 1$

# Accuracy testing: BLAS3 error bound summary

Routine	Formula	Error bound based success criteria
GEMM SYMM HEMM	$C = \alpha AB + \beta C_{init}$	$\frac{\ C - \hat{C}\ _{\infty}}{(N \alpha \ A\ _{\infty}\ B\ _{\infty} +  \beta \ C_{init}\ _{\infty})\varepsilon} \leq 1$
SYRK HERK	$C = \alpha AA^{T/H} + \beta C_{init}$	$\frac{\ C - \hat{C}\ _{\infty}}{(N \alpha \ A\ _{\infty}\ A\ _1 +  \beta \ C_{init}\ _{\infty})\varepsilon} \leq 1$
SYR2K HER2K	$C = \alpha AB^{T/H} + \alpha BA^{T/H} + \beta C_{init}$	$\frac{\ C - \hat{C}\ _{\infty}}{(N \alpha (\ A\ _{\infty}\ B\ _1 + \ B\ _{\infty}\ A\ _1) +  \beta \ C_{init}\ _{\infty})\varepsilon} \leq 1$
TRSM	$AX = \alpha B$	$\frac{\ AX - \alpha B\ _{\infty}}{N\ A\ _{\infty}\ X\ _{\infty}\varepsilon} \leq 1$

# Accuracy testing: some remarks on the forward error bound

$$\frac{\|C - \hat{C}\|_{\infty}}{(N|\alpha|\|A\|_{\infty}\|B\|_{\infty} + |\beta|\|C_{init}\|_{\infty})\varepsilon} \leq 1?$$

In our accuracy checking, we assume that the **correct solution**  $C$  of GEMM exists, what is **not true** in real-world applications.

$$\frac{\|C - \hat{C}\|_{\infty}}{(N|\alpha|\|A\|_{\infty}\|B\|_{\infty} + |\beta|\|C_{init}\|_{\infty})\varepsilon} \leq 2?$$

- Our exact solutions: reference BLAS
- Problem: reference BLAS results have the same error bound
- Pessimistic approach: authorise twice the error bound

1. Reference implementation
2. Accuracy testing
3. Experimental results
4. Concluding remarks

# Some implementation details

## Batched BLAS implementations tested

- **Batched\_cuBLAS**: GEMM and TRSM
- **Batched\_MAGMA**: GEMM, HERK, SYRK and TRSM
- **Batched\_MKL**: only GEMM

## Alternative to Batched on multi-core processor?

- **OMP\_LOOP\_MKL**: OpenMP loop over sequential MKL
- **PARALLEL\_MKL**: single loop over parallel MKL



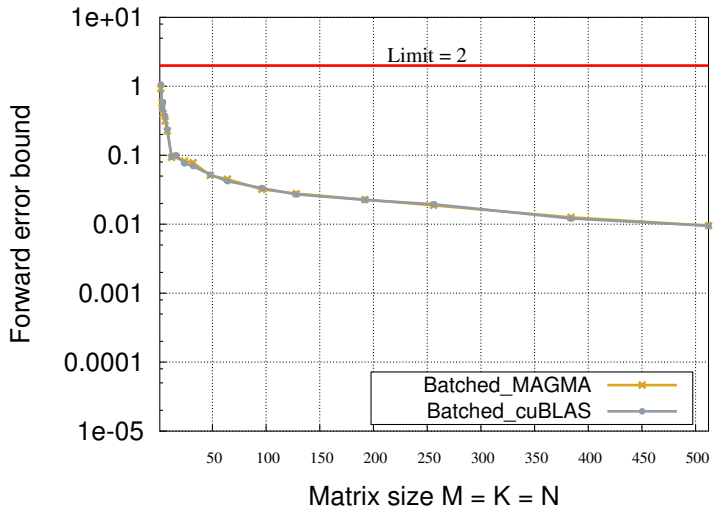
## About the source code

- Programming language: C
- Precision: single, double, single complex and double complex
- Documentation: well documented (doxygen)
- First release coming soon

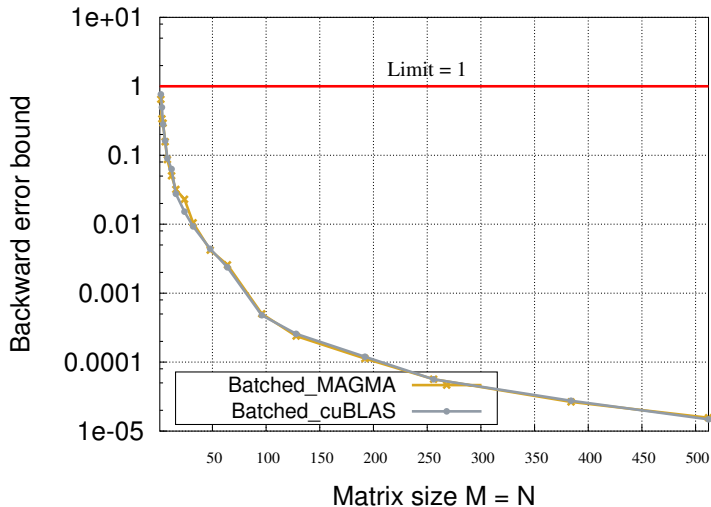
## Hardware

- CPU: 2x10 Intel Xeon E5-2650 v3, 25M Cache, 2.30 GHz
- GPU: Nvidia K40c GPU with 2, 880 CUDA cores

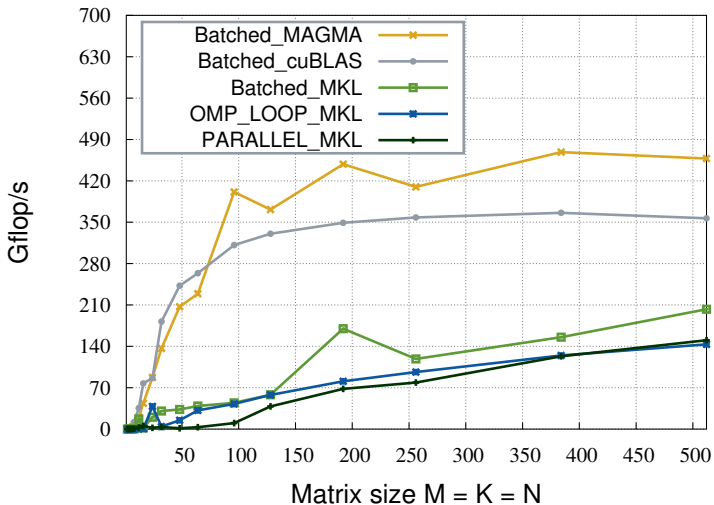
# Accuracy of 1K batched DGEMM



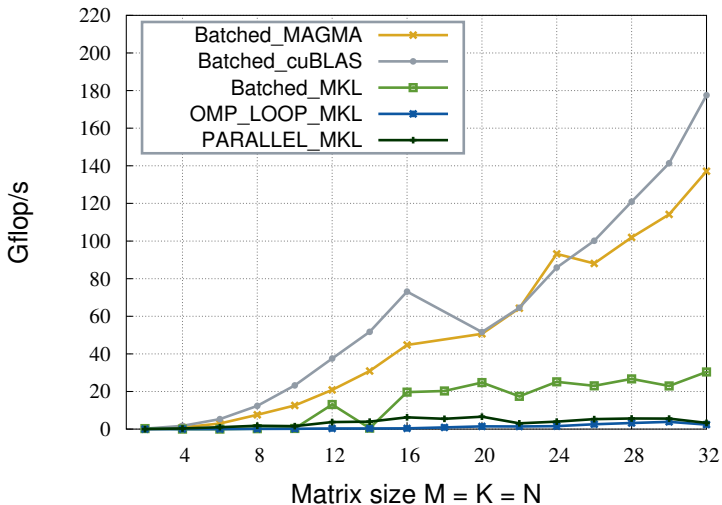
# Accuracy of 1K batched DTRSM



# Performance of 1K batched DGEMM



# Performance of 1K batched DGEMM



1. Reference implementation
2. Accuracy testing
3. Experimental results
4. Concluding remarks

# Concluding remarks and perspectives

## Conclusion remarks

- Investigation of new accuracy checking
- Competitive performance cuBLAS/MAGMA
- MKL batched BLAS seems promising

## Future work

- First release with respect to the final specification
- Extension to level 1 & 2 BLAS
- Provide fortran interface

# Concluding remarks and perspectives

## Conclusion remarks

- Investigation of new accuracy checking
- Competitive performance cuBLAS/MAGMA
- MKL batched BLAS seems promising

## Future work

- First release with respect to the final specification
- Extension to level 1 & 2 BLAS
- Provide fortran interface