# Accelerated Mixed-Precision Algorithms for Hermitian Eigenvalue Problem

Yaohung (Mike) TSAI **Piotr LUSZCZEK** Jack DONGARRA Innovative Computing Laboratory University of Tennessee Knoxville









FOR MORE INFORMATION



https://icl.utk.edu/~luszczek/sc22/pstr\_mxp\_eig\_sym/

#### ACKNOWLEDGEMENTS

the National Science Foundation under OAC Grant No.~2004541.

This material is also based upon work supported by the National Science Foundation under OAC Grant No.~2004541 and the University of Tennessee grant MSE E01-1315-038 as Interdisciplinary Seed funding. This research used the computational resources of the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No.~DE-AC05-000R22725 provided by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of the U.S. Department of Energy Office of Science and the National Nuclear Security Administration. The software library integration work was supported by

#### **ABSTRACT**

INTRODUCTION

scheduling

As the new hardware is being equipped with powerful low-precision capabilities driven primarily by the needs of the burgeoning field of Artificial Intelligence (AI), mixed-precision algorithms are now showing far greater potential and renewed interest in scientific computing community. The multi-precision methods commonly follow approximate-iterate scheme by first obtaining the approximate solution from a low-precision factorization and solve. Then, they iteratively refine the solution to the desired accuracy that is often as high as what is possible with traditional approaches. While targeting symmetric and Hermitian eigenvalue problems of the form  $Ax=\lambda x$ , we revisit the SICE algorithm by applying the Sherman-Morrison formula on the diagonally-shifted tridiagonal systems, we propose an updated SICE-SM algorithm. We exploited asynchronous scheduling techniques to take advantage of the new computational graph enabled by the use of mixed-precision in the eigensolver. This allowed us to hide some of the extra work required by mixed-precision algorithm behind the tasks required by the classic eigensolver and thus minimize the overheads resulting from introducing extra floating-point precision data and compute. This allowed us to maximize the efficiency on both multicore CPU as well as on GPU-accelerated platforms. By incorporating the latest two-stage algorithms from the PLASMA and MAGMA software libraries for numerical linear algebra, we achieved up to 3.6x speedup using the mixed-precision eigensolver with the blocked SICE-SM algorithm for iterative refinement when compared with full double complex precision solvers for the cases with a portion of eigenvalues and eigenvectors requested.

SICE (Subroutine for Improving Computed Eigevalues) algorithm [1]–[3]

improves accuracy of Hermitian (or symmetric) eigenvalue problem. To

improve its efficiency on modern platforms with hardware accelerators

we introduced two new parallel agorithms and exploited the

**ONE- & TWO-STAGE** 

EIGENVALUE SOLVER

The classic eigenvalue problem algorithm was one-stage (as shown

at the top): reduction to tridiagonal form was performed by a single

reducing bandwidth of the original matrix but those multiple stage

increased back-transformation cost and computing eigenvectors was

possible but not impractical. Hence, two stage algorithm (shown at

the bottom) struck a balance by only admitting two computational

stages: reduction from full to band form and reduction from band to

algorithm due to its parallel and compute intensive implementation

as large but it still could outperform the classic single-stage

DAG OF COMPUTE TASKS

The compute DAG (Direct Acyclic Graph) shown in the figure is a

simplified representation that underscores the main algorithmic

components of the methods we proposed. The asynchrony of

compute tasks is fully exploited and the hardware is presented with a

mix of workloads that vary in compute intensity and reliance on the

memory hierarchy and thus they use the hardware more efficiently.

not possible for single-stage algorithm.

tridiagonal. The computational overhead of transformation was twice

computational stage without intermediate steps. SBR (Successive

Band Reduction) Toolbox allowed multiple stage with each one

#### Algorithm 1 SICE algorithm : **Input:** Matrix $A \in \mathbb{R}^{n \times n}$ . An approximate eigenvalue $\lambda$ and the corresponding eigenvector x. $I_{max}$ denotes the maximum 2: Output: Refined eigenvalue $\lambda$ and its eigenvector x: function $[\lambda, x] \leftarrow SICE(A, \lambda, x, I_{max})$ obtain Schur decomposition $[Q,U] \leftarrow \operatorname{schur}(A)$ $A = QUQ^{\mathsf{T}}, \ QQ^{\mathsf{T}} = I.$ $[m,s] \leftarrow max(abs(x)); x \leftarrow x/m$ > Normalizing x so that $||x||_{\infty} = s_x = 1.$ for i in $1:I_{\max}$ do $r \leftarrow \lambda x - Ax$ $c \leftarrow -x - a_{\lambda s}$ $f^{\mathsf{T}} \leftarrow Q(s,:) = e_s^{\mathsf{T}} Q$ $\triangleright$ s-th row of Q. rotations $Q_1$ from $Q_1d = (P_2P_3...P_n)d = \gamma e_1$ where $\gamma = ||d||_2$ Solve the triangular system $\bar{U}_{\lambda}z = Q_2Q_1Q^{\mathsf{T}}r$ ▶ Update eigenvalue. $\lambda \leftarrow \lambda + y(s)$ $ightharpoonup \operatorname{Set} y(s) \text{ to } 0.$ $y(s) \leftarrow 0$ $x \leftarrow x + y$ end for 23: end function

### ORIGINAL SICE ALGORITHM

SICE is the classic refinement algorithm we implemented to profile it on modern hardware with contemporary software stack including two-stage eigenvalue solver that was not available at the time of SICE's introduction. The main performance downside of SICE algorithm is the fact that it is based on element-wise and vector-wise operations. In other words, it relies on BLAS Level 1 and 2 which is highly inefficient on current systems with hardware accelerators. The problem is exacerbated by the path that these low-performance increases with the increasing problem size.

• Tridiagonalization  $A = QTQ^{T}$ 

One stage via

First stage symmetric to band

Second stage band to tridiag

Householder transformations

operations are on the critical path of the algorithm and their overhead

Bulge chasing

### SICE-SM is our new parallel algorithm that overcomes the main

**NEW SICE ALGORITHM** 

bottlenecks of the original SICE algorithm by introducing Sherman-Morrison formula into the theoretical formulation. This allowed to lift some operations from lower-level BLAS and improve data reuse for better utilization of the memory hierarchy and accelerators' compute units.

Algorithm 2 SICE-SM algorithm: SICE algorithm with

: **Input:** Matrix  $A = A^{\mathsf{T}} \in \mathbb{R}^{n \times n}$ . An approximate eigenvalue  $\lambda$  and

the corresponding eigenvector x.  $I_{max}$  denotes the maximum

 $[m,s] \leftarrow max(abs(x)); x \leftarrow x/m \rightarrow Normalization of x so$ 

 $\triangleright$  s-th row of Q.

▶ Update eigenvalue.

Update eigenvector.

ightharpoonup Set y(s) to 0.

2: Output: Refined eigenvalue  $\lambda$  and eigenvector x.

3: **function**  $[\lambda, x] \leftarrow SICE SM(A, \lambda, x, I_{max})$ 

Sherman-Morrison formula

that  $||x||_{\infty} = s_x = 1$ .

25: end function

for i in  $1:I_{max}$  do

 $r \leftarrow \lambda x - Ax$ 

 $c \leftarrow -x - a_{\lambda s}$ 

 $d \leftarrow Q^{\mathsf{T}} c$  $f^{\mathsf{T}} \leftarrow Q(s,:) = e_s^{\mathsf{T}} Q$ 

 $\lambda \leftarrow \lambda + y(s)$ 

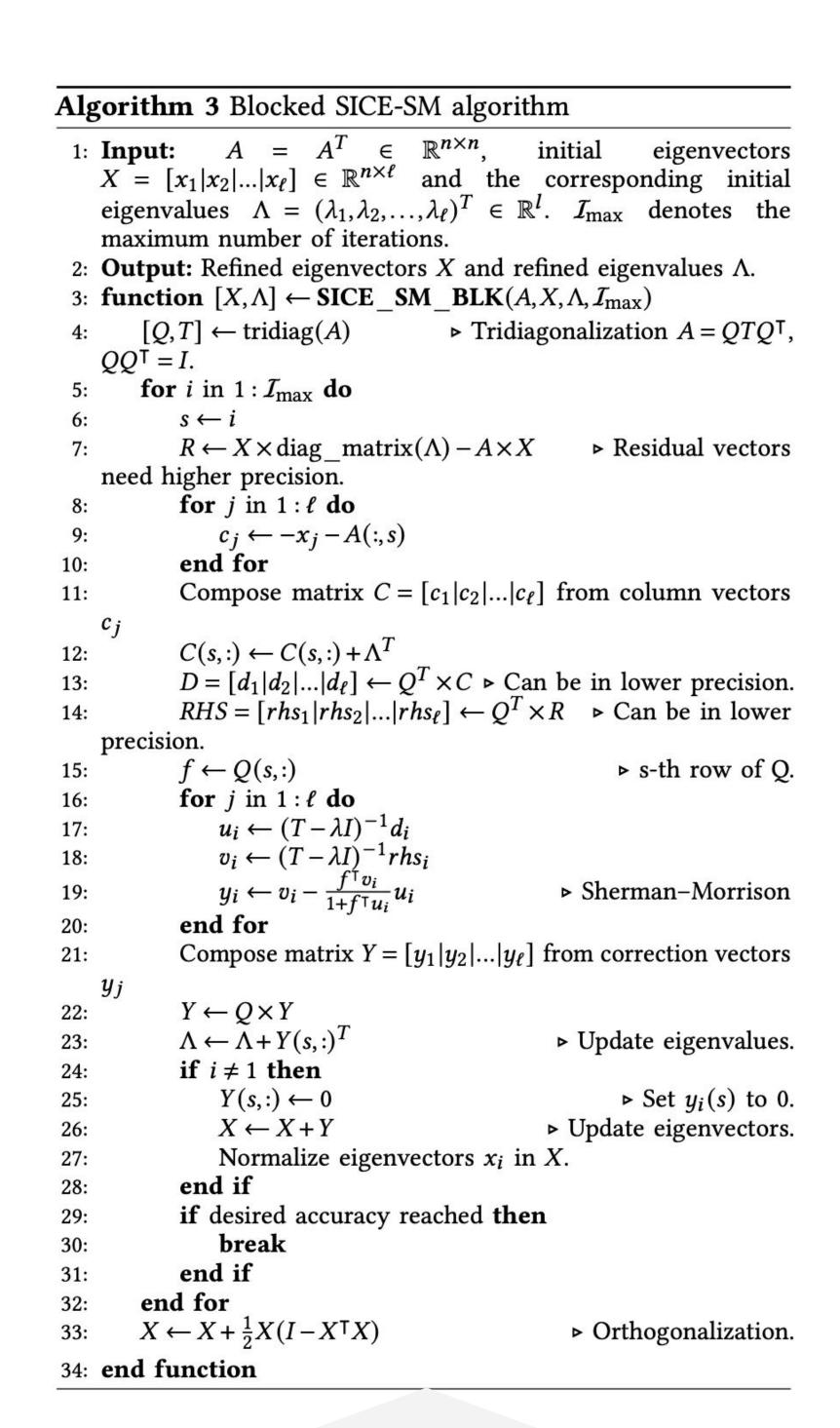
if  $i \neq 1$  then

 $y(s) \leftarrow 0$ 

 $x \leftarrow x + y$ 

if desired accuracy reached then

 $u \leftarrow (T - \lambda I)^{-1} d$  $v \leftarrow (T - \lambda I)^{-1} rhs$ 

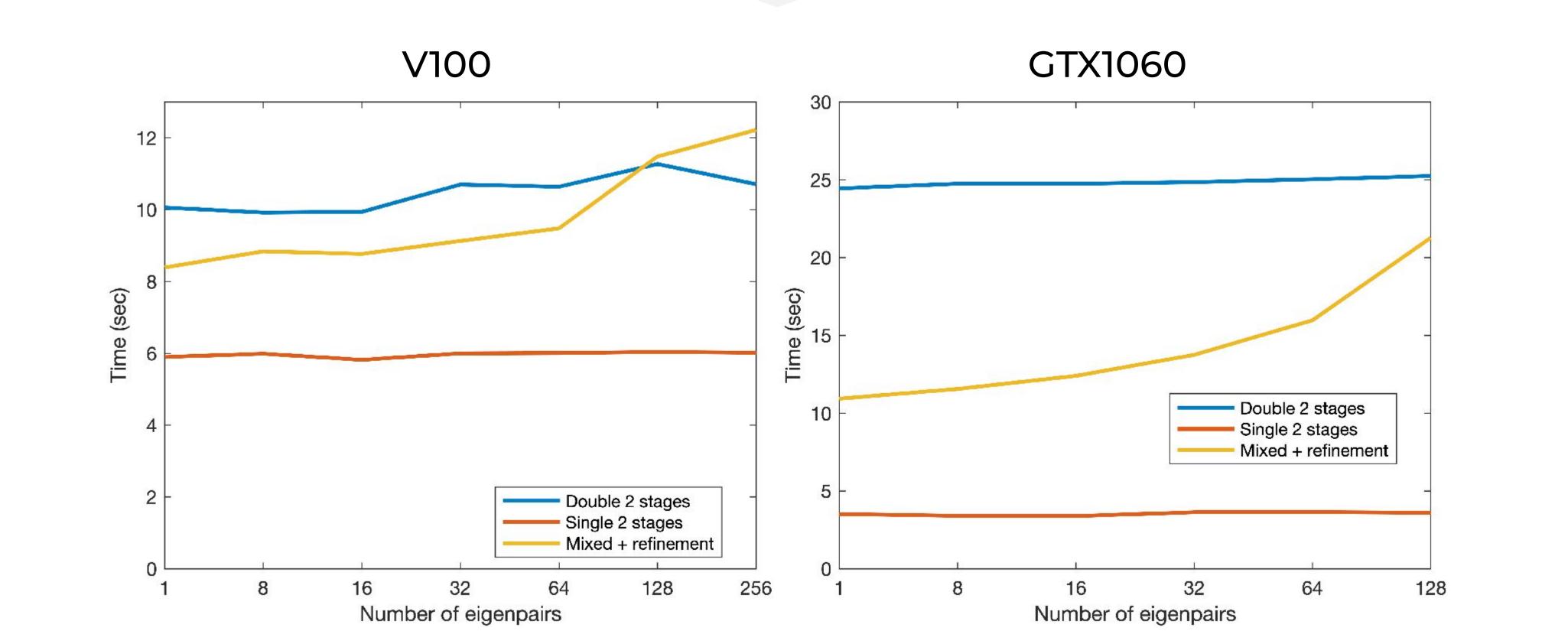


### NEW BLOCKED SICE-SM **ALGORITHM**

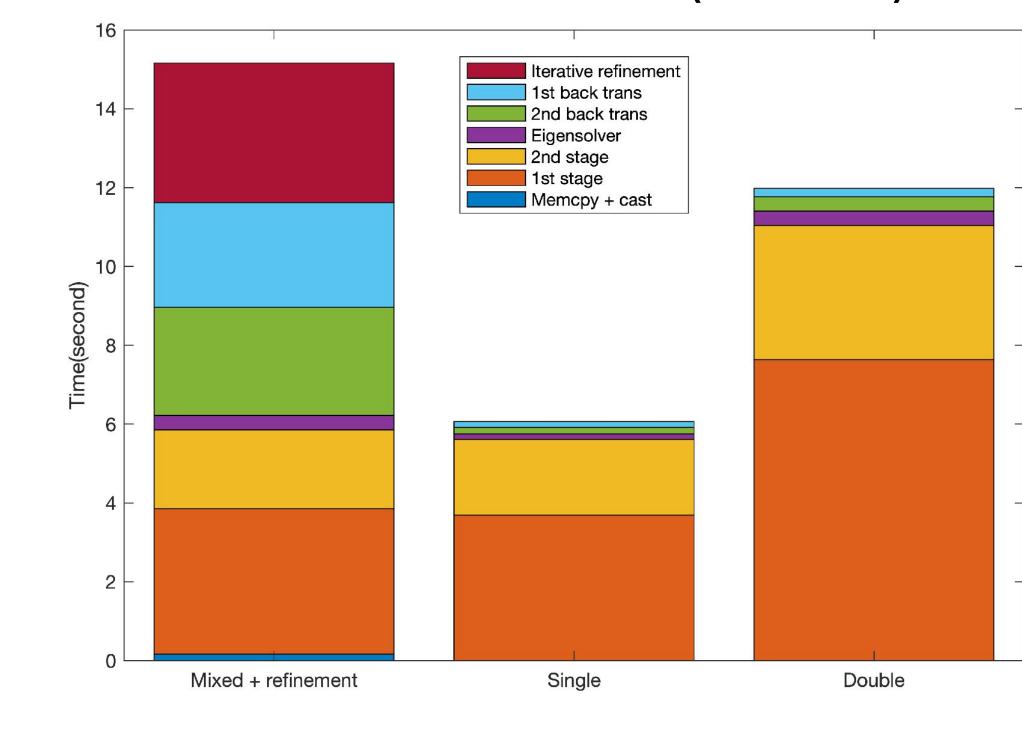
Blocked SICE-SM is our new parallel algorithm that further improves the cache-blocking opportunities that target the fastest levels of the memory hierarchy such as CPU caches, register files, and GPU shared memory. Better compute efficiency improves performance of the new mixed-precision tasks while the asynchronous scheduling continues to exploit independent operations for parallel execution and hiding less efficient tasks behind the compute-intensive ones for maximal overall

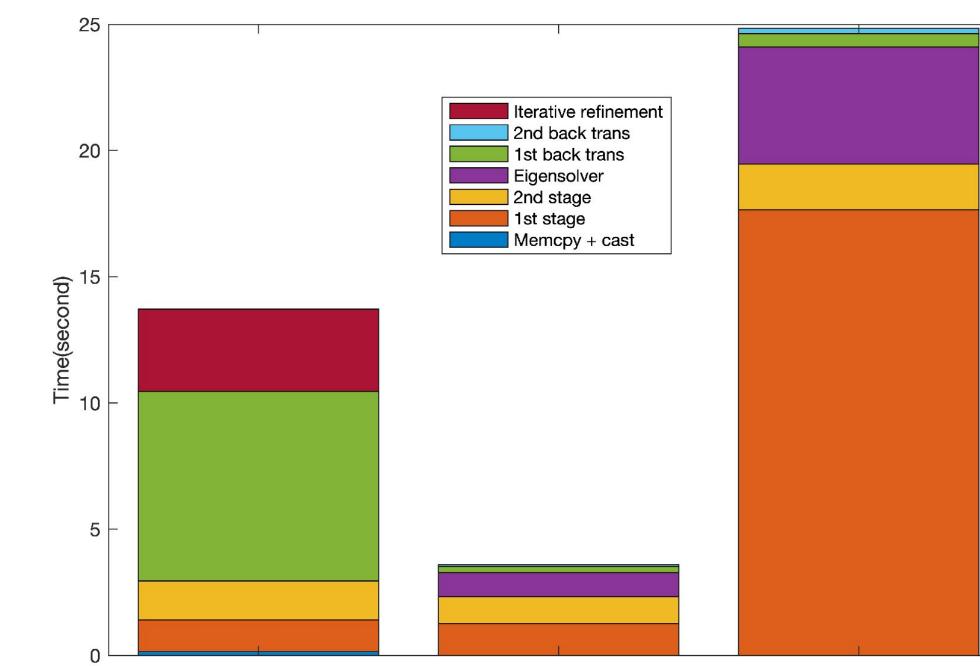
## PERFORMANCE FOR VARYING COUNTS OF EIGENPAIRS

The performance figures below show the behavior of our HPC implementations of two accelerated platforms that differ in their mixed-precision compute balance: 2x and 32x, respectively.



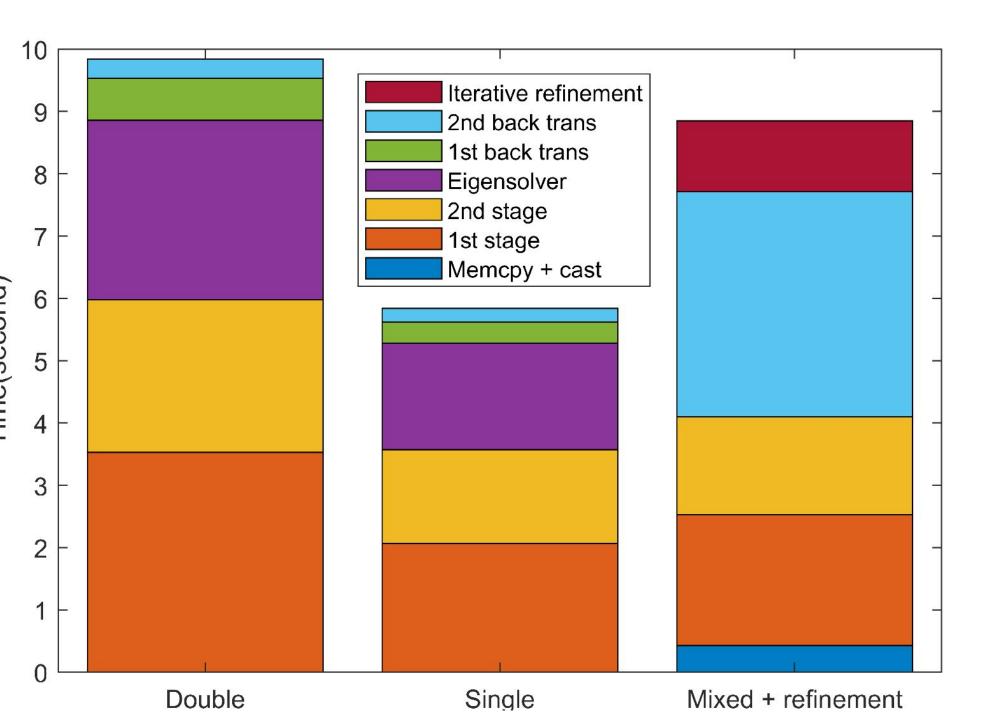




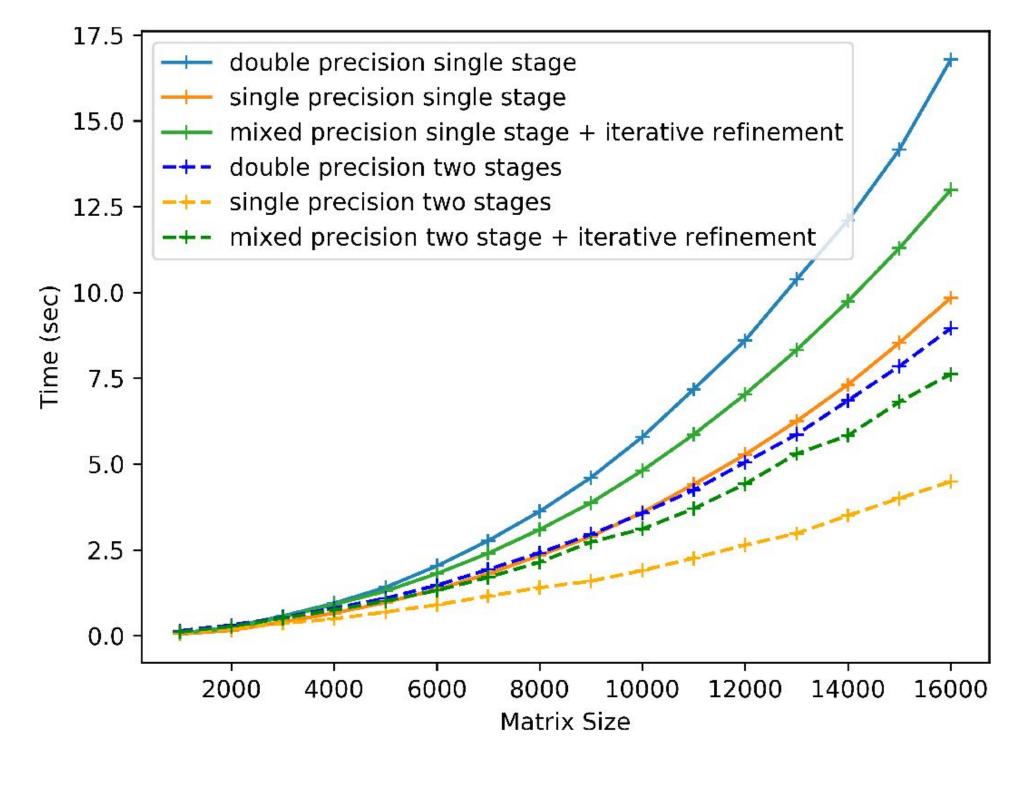


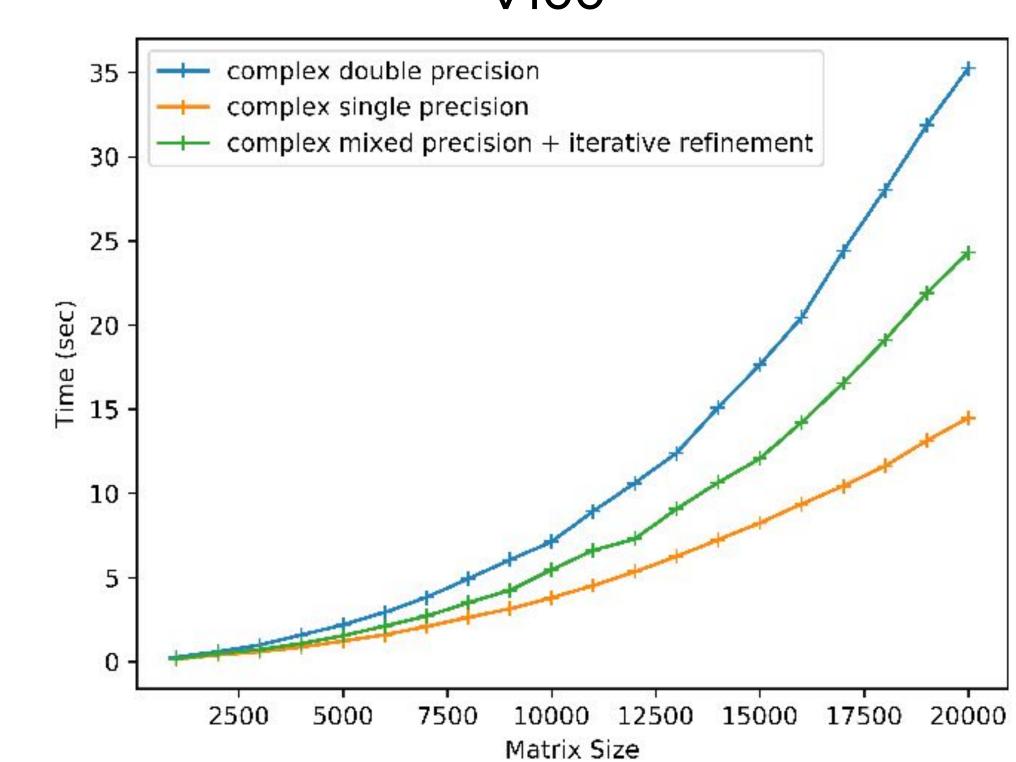
TIME BREAKDOWN (GTX1060)

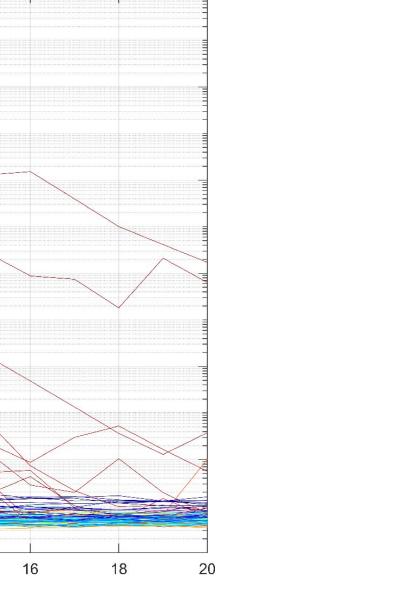
TIME BREAKDOWN (NVIDIA V100)



REAL MATRIX TIMINGS (V100)

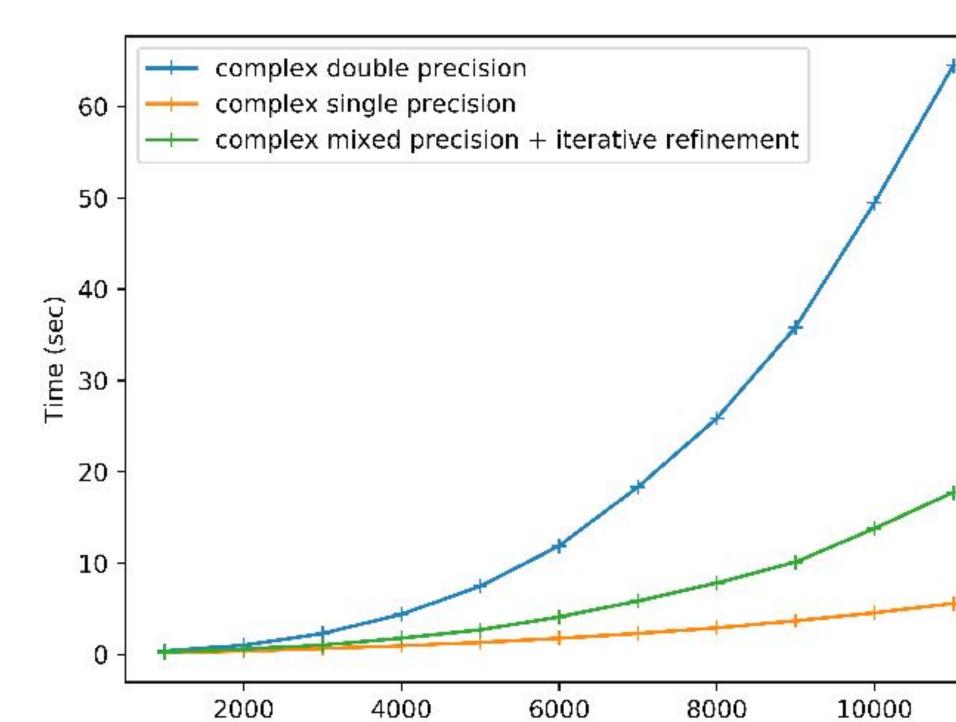






Convergence graphs represent the convergence rates across the eigen spectrum of a small matrix with entries from a uniform random distribution. The largest eigenvalues experience very rapid convergence, that deteriorates for the smaller ones further down the spectrum. This is the main potential for further extension to leverage eigen spectrum shifting to level up the convergence in any portion of





### **CONVERGENCE ACROSS THE EIGENSPECTRUM**

the eigenspectrum.