## Computing Accurate Eigenvalues of Symmetric Matrices With a Mixed Precision Jacobi Algorithm

## Nicholas J. Higham, Françoise Tisseur, Marcus Webb, Zhengbo Zhou

## Abstract

Modern hardware increasingly supports not only single and double precisions, but also half and quadruple precisions. These precisions provide new opportunities to considerably accelerate linear algebra computations while maintaining numerical stability and accuracy. Efforts on developing mixed precision algorithms in the numerical linear algebra and high performance computing communities have mainly focussed on linear systems and least squares problems. Eigenvalue problems are considerably more challenging to solve and have a larger solution space that cannot be computed in a finite number of steps [5].

There are two classes of algorithms for symmetric eigenproblems: (i) those that work directly on the matrix, such as the Jacobi algorithm and the QR-based Dynamically Weighted Halley (QDWHeig) algorithm and (ii) those that reduce the matrix to tridiagonal form in a finite number of steps and then employ an iterative scheme to compute all or just part of the eigenvalues and/or the eigenvectors, such as bisection and inverse iteration (BI), the QR algorithm, and the divide-andconquer algorithm (DC). All these algorithms have pros and cons. DC and the method of multiple relatively robust representations (MR), which is a sophisticated variant of inverse iteration, are generally much faster than QR and BI on large matrices, with MR performing the fewest floating point operations but at a lower MFLOPS rate than DC. The latter and QR are the most accurate algorithms with observed accuracy  $O(\sqrt{nu})$ , where u is the working precision, n the size of the matrix, and accuracy is measured in terms of scaled residual norms and loss of orthogonality for the eigenvectors [1]. None of these eigensolvers exploits the low precisions available in modern hardware.

A key question is how can we exploit access to multiple precisions arithmetic to accelerate symmetric eigensolvers while maintaining numerical stability and accuracy?

In terms of arithmetic cost, solving a symmetric eigenvalue problem is about 27 times more expensive than solving a symmetric positive definite linear system. Unlike for linear systems for which the  $O(n^3)$  part of the computation can be performed at low precision and the *n*-dimensional solution refined at working precision in  $O(n^2)$  operations, it can be shown that for the eigenvalue problem, some of the  $O(n^3)$  operations need to be performed in the working precision if one hopes to maintain numerical stability and achieve accuracy. So to gain any speedup, these should be BLAS 3 operations, i.e., highly optimized matrix-matrix multiplies. Modern architectures execute matrix multiplies of large size *n* at least 18 faster than symmetric eigensolvers on the same size matrices. Low precision arithmetic can be used to preprocess or to precondition the eigenproblem to allow for a faster solution.

In this talk we concentrate on symmetric positive definite matrices  $A \in \mathbb{R}^{n \times n}$  and consider a mixed precision preconditioned Jacobi algorithm that uses three precisions  $u_h < u < u_\ell$ . The preconditioner  $\tilde{Q}$  is an approximate eigenvector matrix that is efficiently computed using a combination of low and working precisions. Zhang and Bai [7] and Zhou [8] suggested to compute an eigenvector matrix at low precision and then orthogonalize it to working precision so that

$$\|\tilde{Q}^T \tilde{Q} - I\|_2 \le p_1 u < 1, \tag{1}$$

where  $p_1$  is a low degree polynomial in n. It is essential that the preconditioner  $\tilde{Q}$  satisfies (1) to ensure that the eigenvectors returned by the mixed precision preconditioned Jacobi algorithm are orthogonal to working precision u. We discuss several alternative efficient ways to construct such preconditioner and prove it reduces the off-diagonal entries of A to a level determined by the chosen low precision  $u_\ell$  so that the initial slow convergence phase of the Jacobi algorithm can be skipped.

Demmel and Veselič [2] showed that the eigenvalues computed by the Jacobi algorithm with stopping criterion  $|a_{ij}| \leq \sqrt{a_{ii}a_{jj}}$  for all i, j satisfy

$$\frac{|\lambda_i(A) - \lambda_i(A)|}{|\lambda_i(A)|} \le p(n) \, u \, \kappa_2^S(A),\tag{2}$$

where  $\lambda_i(A)$  and  $\lambda_i(A)$  denote the *i*th largest exact and computed eigenvalue of A, p(n) is a low degree polynomial and u is the working precision. Here  $\kappa_2^S(A)$  is the *scaled condition number* of A defined by

$$\kappa_2^S(A) = \kappa_2(DAD), \qquad D = \operatorname{diag}(a_{ii}^{-1/2}),$$

where  $\kappa_2(B) = \lambda_1(B)/\lambda_n(B)$ . For the QR and DC algorithms, the relative error is bounded by  $n^{1/2}p(n) u \kappa_2(A)$  so when  $\kappa_2(DAD) \ll \kappa_2(A)$ , the Jacobi algorithm can produce much more accurate approximations to the smaller eigenvalues than QR or DC algorithms.

Malyshev [6] and Drygalla [3, 4] suggest that preconditioning the matrix at a precision  $u_h$  higher than the working precision u improves the accuracy of the spectral decomposition computed by the preconditioned Jacobi algorithm. However, Malyshev only discuss the backward error and Drygalla only claims the high accuracy property without proving it. Let us denote by  $\widetilde{A}$  and  $\widetilde{A}_{comp}$  the product  $\widetilde{Q}^T A \widetilde{Q}$  computed in exact and floating point arithmetic, respectively. We prove under mild assumptions that the relative errors in the computed eigenvalues are proportional to  $u\kappa_2^S(\widetilde{A}_{comp})$  and  $u\kappa_2^S(\widetilde{A})$  instead of  $u\kappa_2^S(A)$  which appears in (2). Moreover, we prove that if  $\widetilde{A}$  is  $\theta$ -scaled diagonally dominant, i.e.,  $\theta = \|\widetilde{D}\widetilde{A}\widetilde{D}\|_2 < 1$  then the scaled condition numbers  $\kappa_2^S(\widetilde{A})$  and  $\kappa_2^S(\widetilde{A}_{comp})$  are of order 1. Hence, all the eigenvalues are computed to high relative accuracy. We remark that any preconditioner  $\widetilde{Q}$  such that off $(\widetilde{A})/\min_i(\widetilde{a}_{ii}) < 1$ , where off $(\widetilde{A}) = (\sum_{i\neq j} \widetilde{a}_{ij}^2)^{1/2}$ , yields an  $\widetilde{A}$  that is scaled diagonally dominant. For a preconditioned matrix  $\widetilde{A}$  that is not scaled diagonally dominant, we use a result by Demmel and Veselič [2, Prop. 6.2] to argue that if off $(\widetilde{A})$  is sufficiently small so that we can treat the diagonals of  $\widetilde{A}$  as its approximate eigenvalues, the scaled condition numbers  $\kappa_2^S(\widetilde{A}_{comp})$  and  $\kappa_2^S(\widetilde{A})$  are significantly smaller than  $\kappa_2^S(A)$ .

Finally, we present numerical results to support our theoretical analysis.

## References

- J. W. Demmel, O. A. Marques, B. N. Parlett, and C. Vömel. Performance and accuracy of LAPACK's symmetric tridiagonal eigensolvers. SIAM J. Sci. Comput., 30(3):1508–1526, 2008.
- [2] J. W. Demmel and K. Veselić. Jacobi's method is more accurate than QR. SIAM J. Matrix Anal. Appl., 13(4):1204–1245, 1992.
- [3] V. Drygalla. Extra precise preconditioning for non-Hermitian eigenvalue problems. Proc. Appl. Math. Mech., 6(1):713-714, Dec. 2006.
- [4] V. Drygalla. Exploiting mixed precision for computing eigenvalues of symmetric matrices and singular values. Proc. Appl. Math. Mech., 8(1):10809–10810, Dec. 2008.

- [5] N. J. Higham and T. Mary. Mixed precision algorithms in numerical linear algebra. Acta Numerica, 31:347–414, May 2022.
- [6] A. N. Malyshev. On iterative refinement for the spectral decomposition of symmetric matrices. Research Report 1651, INRIA/IRISA, Unité de Recherche, Rennes, France, 1992.
- [7] Z. Zhang and Z.-J. Bai. A mixed precision Jacobi method for the symmetric eigenvalue problem. Technical Report arXiv:2211.03339v1, 2022.
- [8] Z. Zhou. A mixed-precision eigensolver based on the Jacobi algorithm. M.Sc. Thesis, The University of Manchester, Manchester, UK, Sept. 2022.