THE LAPW METHOD WITH EIGENDECOMPOSITION BASED ON THE HARI–ZIMMERMANN GENERALIZED HYPERBOLIC SVD*

SANJA SINGER[†], EDOARDO DI NAPOLI[‡], VEDRAN NOVAKOVIĆ[§], AND GAYATRI ČAKLOVIĆ[¶]

Abstract. In this paper we propose an accurate, highly parallel algorithm for the generalized eigendecomposition of a matrix pair (H, S), given in a factored form (F^*JF, G^*G) . Matrices H and S are generally complex and Hermitian, and S is positive definite. This type of matrices emerges from the representation of the Hamiltonian of a quantum mechanical system in terms of an overcomplete set of basis functions. This expansion is part of a class of models within the broad field of Density Functional Theory, which is considered the golden standard in condensed matter physics. The overall algorithm consists of four phases, the second and the fourth being optional, where the two last phases are computation of the generalized hyperbolic SVD of a complex matrix pair (F, G), according to a given matrix J defining the hyperbolic scalar product. If J = I, then these two phases compute the GSVD in parallel very accurately and efficiently.

Key words. LAPW method, generalized eigendecomposition, generalized (hyperbolic) singular value decomposition, hyperbolic QR factorization

 $\mathbf{AMS} \ \mathbf{subject} \ \mathbf{classifications.} \ 65\mathrm{F}15, \ 65\mathrm{F}25, \ 65\mathrm{Y}05, \ 65\mathrm{Z}05$

1. Introduction. Density Functional Theory (DFT) is the Standard Model at the base of simulations in condensed matter physics. At the center of most DFT simulations lays the initialization of the Hamiltonian matrix H and its diagonalization. In many DFT methods the form and size of the Hamiltonian depends on the choice of the set of basis functions used to expand the atomic orbitals. When such a basis set is not orthonormal, a Hermitian positive definite overlap matrix S has to be computed and diagonalized simultaneously with H; this pair of matrices (H,S) define a generalized Hermitian eigenproblem (or eigenpencil, in short). In a subset of all DFT methods labeled as LAPW, the entries of both H and S are represented as multiple sums and products of smaller matrices with specific properties. We show how to exploit this peculiar representation to solve the generalized eigenvalue problem without explicitly assembling the H and S matrices. Our alternative method solves the eigenpencil using a cascade of phases ending with the Hari-Zimmermann algorithm for a generalized hyperbolic SVD. We demonstrate the scalability of a shared memory version of this method on a number of test cases extracted from concrete DFT simulations. If the matrix S is ill-conditioned, our method has the additional benefit of providing enhanced accuracy while avoiding the failure-prone Cholesky factorization.

The birth of DFT is marked by two fundamental articles by the Nobel prize winner Walter Kohn and his collaborator Lu J. Sham and Pierre Hohenberg [12, 15]. DFT

^{*}This work has been supported in part by Croatian Science Foundation under the project IP–2014–09–3670, and also in part by a bilateral research project "Optimization of material science algorithms on hybrid HPC platforms" funded by the Croatian Ministry of Science and Education (MZO) and the German Academic Exchange Service (DAAD).

[†]University of Zagreb, Faculty of Mechanical Engineering and Naval Architecture, I. Lučića 5, 10000 Zagreb, Croatia, (ssinger@fsb.hr).

[‡]Forschungszentrum Jülich, Jülich Supercomputing Centre, Wilhelm–Johnen-Straße, Jülich, 52425, Germany and RWTH Aachen University, AICES, Schinkelstraße 2, Aachen, 52062, Germany, (e.di.napoli@fz-juelich.de, dinapoli@aices.rwth-aachen.de).

[§]Completed a majority of his part of the research while being affiliated to Universidad Jaime I, Av. Vicent Sos Baynat, 12071 Castellón de la Plana, Spain, (novakoni@uji.es).

[¶]Ph.D. student, Forschungszentrum Jülich, Jülich Supercomputing Centre, Wilhelm-Johnen-Straße, Jülich, 52425, Germany, (g.caklovic@fz-juelich.de).

provides an approach to the theory of electronic structure that is alternative to the solution of the Schrödinger equation. While in the latter the emphasis is on a many-electron wave function describing the dynamics of electrons in a multi-atomic system, in DFT the electron density distribution $n(\mathbf{r})$ plays a central role. Besides providing a complementary perspective, DFT has made possible the simulation of much larger systems than the conventional multi-particle wave function methods. Depending on the specific DFT method, computing complexity scales at most with the cube of the number of atoms, with ongoing progress towards bringing it down to linear scaling.

Despite being a general theory, DFT can be realized in as many flavors as are the sets of basis functions one can choose from. Two widely spread classes of basis functions build on the simplicity of plane waves to build more complex and rich sets of basis functions, namely Projected Plane Waves (PAW) [23] and Linearized Augmented Plane Waves (LAPW) [26]. The complexity of these sets lays in that they are made up of non-orthogonal basis functions. In the case of LAPW, the set of functions is also overcomplete. The consequence of non-orthogonality is that the matrix S, whose entries are the scalar products among all the basis functions of any given finite size set, is usually dense. In the particular case of LAPW methods such matrix is positive definite but could have few singular values quite close to zero. This potential problem is due to the overcompleteness of the basis set and tends to worsen as the number of atoms increases since the number of basis functions grows linearly with it.

In DFT methods the dynamics of the quantum systems is described by a Hamiltonian operator. In practice, the Hamiltonian is translated into a Hermitian matrix H whose size and structure depends on the specific DFT method. This is because the matrix H is the result of the projection of the Hamiltonian operator over the finite set of basis functions of the given method. In the LAPW method, the mathematical form of the functions leads to an expression for both H and S in terms of a sum of smaller matrices over all possible atoms N_A ,

$$H = \sum_{a=1}^{N_A} (A_a^* T_a^{[AA]} A_a + A_a^* T_a^{[AB]} B_a + B_a^* T_a^{[BA]} A_a + B_a^* T_a^{[BB]} B_a),$$

$$(1.1)$$

$$S = \sum_{a=1}^{N_A} (A_a^* A_a + B_a^* U_a^* U_a B_a),$$

where $A_a, B_a \in \mathbb{C}^{N_L \times N_G}$, with N_G and N_L ($N_G \geq N_L$) being the size of the basis set and the total number of angular momentum states, respectively. The remaining matrices in (1.1) are complex, square, of order N_L , with some additional properties. Matrices U_a are real and diagonal, for all $T_a^{[AB]}$ holds $(T_a^{[AB]})^* = T_a^{[BA]}$, while $T_a^{[AA]}$ and $T_a^{[BB]}$ are Hermitian. Except for U_a , the other matrices are in general dense, and can have a range of sizes dictated by the constants N_A , N_G , and N_L (see section 2 for some their typical range). Despite the formulation above could lend itself to computation through specialized middleware libraries such as the Basic Linear Algebra Subprograms (BLAS), the standard approach followed by most code developers was one based on minimizing memory footprint and FLOP count [5, 16].

Recently, an alternative method for the assembly of the matrices H and S was presented in [8] and further developed in [6]. In their work [8], Di Napoli et al. consolidate the underlying matrix structure of the operations and proceed to encapsulate them in terms of the level 3 BLAS kernels. For instance, to maximize the arithmetic

intensity of the computation, matrix H is written as $H = H_{AA} + H_{AB+BA+BB}$, with

$$H_{AA} = \sum_{a=1}^{N_A} A_a^* T_a^{[AA]} A_a, \quad H_{AB+BA+BB} = \sum_{a=1}^{N_A} \left(B_a^* Z_a + Z_a^* B_a \right),$$

where $Z_a=T_a^{[BA]}A_a+1/2T_a^{[BB]}B_a$. Each of the Z_a and B_a matrices is then packed in memory in two consecutive 2-dimensional arrays Z_* and B_* , respectively. In the end, the sum $H_{AB+BA+BB}$ is computed by just two ZHER2K BLAS subroutines. A similar procedure holds for the matrix S. Once assembled the algebraic dense generalized eigenproblem is solved by standard methods. A Cholesky factorization $LL^* = S$ is used to reduce the problem to standard form $A \leftarrow L^{-1}AL^{-*}$. In turn, the standard problem is solved by a dense direct algorithm such as MRRR [7] provided by the LAPACK library [1], or an iterative eigensolver specialized for DFT computation (e.g., the ChASE library [29]). When the assembled S matrix is ill-conditioned, as it may happen for quantum systems with a large number of atoms (> 100), the Cholesky factorization may fail and makes it practically very hard to solve the corresponding generalized eigenproblem. This issue is typically solved by the practitioners by modifying the mathematical model so to avoid increasing the basis set, which is the source of an ill-condioned S, but thus compromising on the robustness of the DFT approach.

In this work, we propose an numerical method alternative to the physics-based approach for solving the eigenpencil (H, S) without forming the matrices explicitly. The core of the method is based on the generalized hyperbolic singular value decomposition (GHSVD) [2]. Not only such a method solves for the eigenproblem directly without assembling H and S, but also could give more accurate results when S is nearly singular. This is possible since the GHSVD decomposition acts directly on the multiplying factors making up S, conceivably reducing the singularity down to the square root of the condition number of S. As a surplus, if J, the matrix of the hyperbolic scalar product, is equal to the identity, the GHSVD reduces to the generalized SVD (GSVD), which is computed very efficiently in parallel.

The paper is subdivided into eight sections. In section 2, we present in more detail the physics of the problem and the mathematical model leading to the expression (1.1) for H and S. Section 3 is devoted to formulating the problem in precise algebraic terms, and provides an overview of our algorithm. The next four sections deal with the four phases of the algorithm, where the first three of them belong to the algorithm proper, and the fourth one completes the computation of the GHSVD and is unrelated to the underlying mathematical physics origin of the problem. Since each phase is an algorithm and a reusable software contribution in its own right, at the end of each section we present the numerical results and the parallelization techniques applied. The paper concludes with a note on the related future work in section 8.

2. The H and S matrices in LAPW methods. At the core of DFT are a set of equations, called Kohn-Sham equations, that have to be solved for each of the single particle wave function ψ_i

(2.1)
$$\hat{H}_{KS} \psi_i(\mathbf{r}) = \left[-\frac{\hbar^2}{2m_e} \nabla_{\mathbf{r}}^2 + V[n(\mathbf{r})](\mathbf{r}) \right] \psi_i(\mathbf{r}) = \epsilon_i \psi_i(\mathbf{r}), \quad i = 1, \dots, N_e.$$

The peculiarity of these equations is that the Hamiltonian operator \hat{H}_{KS} depends implicitly on all the ψ through the charge density function $n(\mathbf{r})$, which makes the entire set of Kohn–Sham equation strongly coupled and non-linear. In particular the function $n(\mathbf{r})$ is the sum of the squares of all ψ up to the total number of electrons N_e in any given quantum system

(2.2)
$$n(\mathbf{r}) = \sum_{i=1}^{N_e} |\psi_i(\mathbf{r})|^2.$$

Because the equations (2.1) are non-linearly coupled, they can be solved only self-consistently: one starts from a reasonable guess for the charge density $n(\mathbf{r})_{\text{start}}$, computes the potential $V[n(\mathbf{r})]$, and solves (2.1). The resulting functions ψ_i and values ϵ_i are then used to compute a new density as in (2.2), which is compared to the starting one. If the two densities do not match, the self-consistent loop is repeated with a new mixed charge density. The loop stops only when the new and the old density agree up to some defined constant.

So far we have described the general setup. There are many methods that translate this setup into an algorithm, and this is where the various "flavors" of DFT differ. The first difference is in the choice of the set of functions φ_t used to expand every one particle wave function ψ_i

(2.3)
$$\psi_i(\mathbf{r}) = \sum_{t=1}^{N_G} c_{t,i} \, \varphi_t(\mathbf{r}).$$

In the LAPW method [13, 28], the configuration space where the atomic cells are defined is divided in two disjoint areas where the wave functions have distinct symmetries: close to the atomic nuclei, solutions tend to be spherically symmetric and strongly varying, while further away from the nuclei, they can be approximated as uniformly oscillating. The qualitative structure of the solution leads to a space composed of non-overlapping spheres—called muffin tins (MT)—separated by interstitial (INT) areas. The complete set of basis functions φ_t are given by a piece-wise definition for each of the N_A atoms and relative surrounding regions.

(2.4)
$$\varphi_t(\mathbf{r}) = \begin{cases} \sum_{l=0}^{l_{\text{max}}} \sum_{m=-l}^{l} \left[A_{(l,m),a,t} u_{l,a}(r) + B_{(l,m),a,t} \dot{u}_{l,a}(r) \right] Y_{l,m}(\hat{\mathbf{r}}_a), & a^{\text{th}} \text{ MT} \\ \frac{1}{\sqrt{\Omega}} \exp(i\mathbf{k}_t \cdot \mathbf{r}), & \text{INT.} \end{cases}$$

In the MT spheres, each basis function depends on specialized radial functions $u_{l,a}$, their derivatives $\dot{u}_{l,a}$ and the spherical harmonics $Y_{l,m}$; the former only depend on the distance \mathbf{r} from the MT center, while the latter form a complete basis on the unit sphere defined by $\hat{\mathbf{r}} = \mathbf{r}/|\mathbf{r}|$ and so depends solely on the MT spherical angles. Despite being piece-wise functions, φ_t must be continuous and differentiable for each index t and each atomic index a. The coefficients $A_{l,m,a}, B_{l,m,a} \in \mathbb{C}$ are set to guarantee that $\varphi_t \in C^1$ for each of the values of the indices $L \equiv (l,m)$ and a. The variable t ranges over the size of the plane wave functions set in INT, and is used to label the vector \mathbf{k}_t living in the space reciprocal to \mathbf{r} . As such, the momentum \mathbf{k}_t characterizes the specific wave function entering in the basis set. The total size of the basis set is determined by setting a cutoff value $\mathbf{K}_{\text{max}} \geq \mathbf{k}_t$.

When one substitutes the expansion of ψ_i (2.3) in (2.1), the Kohn–Sham equations become an algebraic generalized eigenvalue problem that needs to be solved for the

 N_G -tuples of coefficients $c_i = (c_{1,i}, \dots, c_{N_G,i})^T$

$$\sum_{t=1}^{N_G} (H)_{t',t} c_{t,i} = \epsilon_i \sum_{t=1}^{N_G} (S)_{t',t} c_{t,i}.$$

The complexity of the LAPW basis set is transferred to the definition of the entries of the Hamiltonian and overlap matrices, respectively H and S, given by

$$(H)_{t',t} = \sum_{a} \iint \varphi_{t'}^*(\mathbf{r}) \hat{H}_{KS} \varphi_t(\mathbf{r}) \, d\mathbf{r}, \quad (S)_{t',t} = \sum_{a} \iint \varphi_{t'}^*(\mathbf{r}) \varphi_t(\mathbf{r}) \, d\mathbf{r}.$$

By substituting explicitly the functions φ_t of equation (2.4) and computing the integrals, one ends up with the following expressions for H and S:

$$(2.5) \qquad (H)_{t',t} = \sum_{a} \sum_{L',L} \left(A_{L',a,t'}^* T_{L',L;a}^{[AA]} A_{L,a,t} \right) + \left(A_{L',a,t'}^* T_{L',L;a}^{[AB]} B_{L,a,t} \right) \\ + \left(B_{L',a,t'}^* T_{L',L;a}^{[BA]} A_{L,a,t} \right) + \left(B_{L',a,t'}^* T_{L',L;a}^{[BB]} B_{L,a,t} \right),$$

$$(2.6) \qquad (S)_{t',t} = \sum_{a} \sum_{L=(l,m)} A_{L,a,t'}^* A_{L,a,t} + B_{L,a,t'}^* B_{L,a,t} \|\dot{u}_{l,a}\|^2.$$

The new matrices $T_{L',L;a}^{[...]} \in \mathbb{C}^{N_L \times N_L}$ are dense and their computation involves multiple integrals between the radial basis functions $u_{l,a}$ and the non-spherical part of the potential V multiplied by Gaunt coefficients (for details see [16][26, Ch. 5][8, App.]). As can be seen by simple inspection, equations (2.5)–(2.6) are equivalent to equations (1.1): while the former are written with all indices explicit, the latter have a subset of them implicit which highlights their matrix form.

We conclude with a small excursus on the structure of the self-consistent loop and its computational cost. In the first step, a starting charge density $n(\mathbf{r})_{\text{start}}$ is used to compute the Kohn–Sham Hamiltonian $H_{\rm KS}$. In a second step the set of basis functions is set up and the set of A, B coefficients is derived. Then, the Hamiltonian H and overlap S matrices are initialized, followed by the fourth step when the generalized eigenvalue problems $Hc = \epsilon Sc$ is solved numerically to return the eigenpairs $(C, \operatorname{diag}(\epsilon))$. Finally a new charge density $n(\mathbf{r})$ is computed and convergence is checked before starting a new loop. Out of all the steps above, initializing H and S and solving the eigenproblem accounts for more than 80% of CPU time. Having cubic complexity $\mathcal{O}(N_G^3)$, the eigenproblem solution is usually considered the most expensive of the two. It turns out that generating the matrices may be as expensive. If N_A and N_L , respectively, are the range of the summations \sum_a and \sum_L , then it can be shown that equations (2.6) and (2.5) have complexity $\mathcal{O}(N_A \cdot N_L \cdot N_G^2)$ and $\mathcal{O}(N_A \cdot N_L \cdot N_G \cdot (N_L + N_G))$. A typical simulation uses approximately N_G basis functions, with N_G ranging from about $50 \cdot N_A$ to about $80 \cdot N_A$, and an angular momentum $l_{\text{max}} \leq 10$, which results in $N_L = (l_{\text{max}} + 1)^2 \le 121$. It follows that the factor $N_A \cdot N_L$ is roughly of the same order of magnitude as N_G so that the generation of H and S also displays cubic complexity $\mathcal{O}(N_G^3)$. In practice, the constants above have values in the following orders of magnitude: $N_A = O(100)$, $N_G = O(1000) - O(10000)$, and $N_L = O(100)$.

3. Problem formulation. Our intention is to keep the matrices H and S in their factored form given in (1.1). The core of the process, Phase 3, is a one-sided Jacobi-like method for the implicit diagonalization, that computes a hyperbolic analog of the generalized SVD.

DEFINITION 3.1. For the given matrices $F \in \mathbb{C}^{m \times n}$, $m \geq n$, $J \in \mathbb{R}^{m \times m}$, $J = \operatorname{diag}(\pm 1)$, and $G \in \mathbb{C}^{p \times n}$, where G is of full column rank, there exist a J-unitary matrix $U \in \mathbb{C}^{m \times m}$ (i.e., $U^*JU = J$), a unitary matrix $V \in \mathbb{C}^{p \times p}$, and a nonsingular matrix $X \in \mathbb{C}^{n \times n}$, such that

$$(3.1) F = U\Sigma_F X, G = V\Sigma_G X, \Sigma_F \in \mathbb{R}^{m \times n}, \Sigma_G \in \mathbb{R}^{p \times n}.$$

The elements of Σ_F and Σ_G are zeros, except for the diagonal entries, which are real and nonnegative. Furthermore, Σ_F and Σ_G satisfy $\Sigma_F^T \Sigma_F + \Sigma_G^T \Sigma_G = I$. The ratios $\Sigma_{ii} := (\Sigma_F)_{ii}/(\Sigma_G)_{ii}$ are called the generalized hyperbolic singular values of the pair (F, G). If the pair (F, G) is real, then all matrices in (3.1) are real.

We choose to define the generalized hyperbolic SVD (GHSVD) only if the matrix G is of full column rank. This implies $p \geq n$, and there is no need to mention this in the definition. In the case of full column rank G, the matrix $S := G^*G$ is positive definite, and the matrix pair (H, S), where $H := F^*JF$, is Hermitian and definite, so it can be simultaneously diagonalized by congruences (see, for example, [22]).

If the GHSVD is computed as in (3.1), then the generalized eigenvalues and eigenvectors of (H, S) are easily retrieved, since

$$\begin{split} H &= F^*JF = X^*\Sigma_F^*U^*JU\Sigma_FX = X^*\Sigma_F^*J\Sigma_FX := X^*\Lambda_FX, \\ S &= G^*G = X^*\Sigma_G^*V^*V\Sigma_GX = X^*\Sigma_G^*\Sigma_GX := X^*\Lambda_GX. \end{split}$$

Substituting $X^* = SX^{-1}\Lambda_G^{-1}$ in the expression for H above, we get

$$(3.2) \hspace{1cm} HZ = SZ\Lambda; \quad Z := X^{-1}, \quad \Lambda := \Lambda_G^{-1}\Lambda_F.$$

Thus, from (3.2), the generalized eigenvalues $\operatorname{diag}(\Lambda)$ of the matrix pair are the squared generalized hyperbolic singular values, with the signs taken from the corresponding diagonal elements in J, i.e., $\operatorname{diag}(\Sigma^T J \Sigma)$, and the matrix of the generalized eigenvectors Z is the inverse of the matrix X of the right generalized singular vectors. For theoretical purposes it can be assumed that $\operatorname{diag}(\Lambda)$ is sorted descendingly, though for simplicity it is not the case in our implementation.

An approach that uses the SVD on a matrix factor, instead of the eigendecomposition on the multiplied factors, usually computes small eigenvalues more accurately.

In the first phase of the algorithm, we transform the initial problem by assembling the Hermitian matrices T_a ,

(3.3)
$$T_a = \begin{bmatrix} T_a^{[AA]} & T_a^{[AB]} \\ T_a^{[BA]} & T_a^{[BB]} \end{bmatrix}$$

and factoring them into a form suitable for the GHSVD computation.

After that, we are left with two tall matrices, which have, in our test examples, between 2 and 22 times more rows than columns. Since the Jacobi-like SVD algorithms are more efficient if the factors are square, in the second (optional) phase we could preprocess the factors: F by the hyperbolic QR factorization (see [24]), and G by the tall-and-skinny QR factorization (ZGEQR routine from LAPACK), into square ones.

The third phase is a complex version of the implicit Hari–Zimmermann method—a modification of the one-sided real method presented in [21]. The complex transformations, for the two-sided method, were derived by Vjeran Hari in his PhD thesis [9].

- **3.1.** Overview of the algorithm. The sequence of phases of our algorithm is:
- 1. The problem is expressed as $H = F_0^* \operatorname{diag}(T_1, \dots, T_{N_A}) F_0$, $S = \widetilde{G}^* \widetilde{G}$, the matrices F_0 and \widetilde{G} are assembled, and the matrices T_a , formed from $T_a^{[AA]}$, $T_a^{[AB]}$, $T_a^{[BA]}$, and $T_a^{[BB]}$, are simultaneously factored by the Hermitian indefinite factorization with complete pivoting, reformulating H as $H = \widetilde{F}^* \widetilde{J} \widetilde{F}$, with $\widetilde{J} = \operatorname{diag}(\pm 1)$.
- 2. Optionally, the tall-and-skinny matrices \widetilde{F} and \widetilde{G} are shortened: \widetilde{F} by the indefinite, \widetilde{J} -QR factorization to obtain the square factor F and a new signature matrix J, and \widetilde{G} by the QR factorization to obtain the square factor G.
- 3. The GEVD of (H, S) is computed by the *J*-GHSVD of (F, G) (or the \widetilde{J} -GHSVD of $(\widetilde{F}, \widetilde{G})$ if the Phase 2 is skipped) by the implicit Hari–Zimmermann method.
- 4. Optionally, the GHSVD process is formally completed by explicitly computing the right generalized singular vectors X from the generalized eigenvector matrix Z.
- **3.2. Testing environment and data.** The testing environment consists of a node with an Intel Xeon Phi 7210 CPU, running at 1.3 GHz with Turbo Boost turned off, in Quadrant cluster mode with 96 GiB of RAM and 16 GiB of flat-mode MC-DRAM, under 64-bit CentOS Linux 7 with the Intel compilers (Fortran, C) and Math Kernel Library (MKL) version 19.0.5.281, and GNU Fortran 8.3.1 for the error testing.

The software code, freely available in https://github.com/venovako/FLAPWxHZ repository, of all the phases presented in this paper is written mostly in Fortran, with some auxiliary parts in C, while the parallelization relies on the OpenMP constructs.

The phases are meant to be run in a sequence, where each phase is executed as a separate process with several OpenMP threads. Since the modern compute nodes generally have enough memory to hold all required data, the algorithms are implemented for the shared memory, but at least the algorithms for the Phases 1, 3, and 4 can be transformed into distributed-memory ones, should the volume of data so require.

In testing it was established that each thread should be bound to its own physical CPU core, with OMP_PROC_BIND=SPREAD placement policy. The double precision and the double-complex BLAS and LAPACK routines were provided by the thread-parallel MKL, but with only one (i.e., the calling) thread allowed per call, except for the ZSWAP, ZROT, and ZGEQR routines in Phase 2, ZGETC2 routine in Phase 4, and ZGEMM, ZHERK, ZHEGV, and ZHEGVD routines in subsection 6.5.4, where the MKL was allowed to use the test's upper limit on the number of threads. The nested parallelism is therefore possible but not required in our code. Hyper Threading was enabled but not explicitly utilized, though nothing precludes a possibility that on a different architecture the BLAS or LAPACK calls could benefit from some form of intra-core symmetric multi-threading. A refined thread placement policy OMP_PROC_BIND=SPREAD, CLOSE might then allow better reuse of data in the cache levels shared among the threads of a core.

Apart from the maximal number of threads set to the number of CPU cores in a node, the tests were also performed with half that number, to assess the effects on the computational time of the larger block sizes and the availability of the whole L2 data cache (1 MiB, shared among two cores) to a thread. The algorithms do not constrain the number of threads in principle, but are not intended to be used single-threadedly.

Our main test node has 64 cores, but a subset of the tests were repeated on a faster JUWELS [14] node, with two Intel Xeon Platinum 8168 CPUs, running at 2.7 GHz with 1 MiB L2 cache per each of $2 \times 24 = 48$ cores, with a similar software setup, for a comparison of the GHSVD and the generalized eigendecomposition approaches (see subsection 6.5.4). When the results obtained on JUWELS are shown, the test's number of threads is emphasized (e.g., 48), to distinguish them from the main results.

Another hardware feature targeted is the SIMD vectorization: each core of both

machines has a private L1 data cache of 32 kiB with a line size of 64 B, and equally wide (e.g., 8 double precision floating-point numbers) vector registers upon which a subset of AVX-512 instructions is capable to operate in the SIMD fashion. The vectorization is employed both implicitly, by aligning the data to the cache line size whenever possible and instructing the compiler to vectorize the loops, and semi-explicitly, as will be described in the following sections. The code is parametrized by the maximal SIMD length (i.e., the number of 8 B lanes in the widest vector register type) \mathbf{v} , and it vectorizes successfully on other architectures (e.g., on AVX2, with $\mathbf{v}=4$).

Under an assumption that the compiler-generated floating-point reductions (e.g., those of the SUM Fortran intrinsic) obey the same order of operations in each run, and due to the alignment enforced as above, the algorithms should be considered conditionally reproducible, in a sense that the multiple runs of the same executables on the same data in the same environment should produce bitwise-identical results.

3.2.1. Datasets. Each dataset under test contained all matrix inputs $(A_a, B_a, U_a, T_a^{[AA]}, T_a^{[BA]}, T_a^{[BB]})$ for a single problem instance. With 8 datasets from Table 3.1 we believe to have a representative coverage of the small-to-medium size problems from practice. As already mentioned in section 2, the maximum value of the momentum

Table 3.1 The datasets under test. For A datasets $N_L=121,\ N_A=108,\ and\ m=2N_LN_A=26136,$ while for B datasets $N_L=49,\ N_A=512,\ and\ m=2N_LN_A=50176.$ Also, $n=N_G.$

ID	AuAg	n	ID	AuAg	n	ID	NaCl	n	ID	NaCl	n
A1	2.5	3275	A3	3.5	8970	B1	2.5	2256	ВЗ	3.5	6217
A2	3.0	5638	A4	4.0	13379	B2	3.0	3893	B4	4.0	9273

 \mathbf{K}_{\max} which appears as an index to the dataset label (e.g., $\mathtt{AuAg_2.5}$) determines the size of the basis functions set N_G . This is why datasets with same label (e.g., \mathtt{AuAg}) but different index (e.g., 2.5 vs. 3.0) have differing values for N_G . In the following, the datasets are referred to by their IDs.

- 4. Phase 1 simultaneous factorizations of T_a matrices. The goal of this section is to rewrite the problem (1.1) in a form suitable for GHSVD computation.
 - **4.1. Problem reformulation.** The first step is to write (1.1) as

$$(4.1) \hspace{1cm} H = \sum_{a=1}^{N_A} H_a^* T_a H_a, \quad S = \sum_{a=1}^{N_A} S_a^* S_a, \quad H_a = \begin{bmatrix} A_a \\ B_a \end{bmatrix}, \quad S_a = \begin{bmatrix} A_a \\ U_a B_a \end{bmatrix}.$$

Furthermore, the matrices in (4.1) can be expressed as

$$(4.2) H = \begin{bmatrix} H_1^* & \cdots & H_{N_A}^* \end{bmatrix} \operatorname{diag}(T_1, \dots, T_{N_A}) \begin{bmatrix} H_1^* & \cdots & H_{N_A}^* \end{bmatrix}^* := F_0^* T F_0,$$

$$S = \begin{bmatrix} S_1^* & \cdots & S_{N_A}^* \end{bmatrix} \begin{bmatrix} S_1^* & \cdots & S_{N_A}^* \end{bmatrix}^* := \widetilde{G}^* \widetilde{G}.$$

In (4.2), diag (T_1,\ldots,T_{N_A}) stands for a block-diagonal matrix with the prescribed diagonal blocks $T_a,\ a=1,\ldots,N_A$ from (3.3). Newly defined matrices have the following dimensions: $H_a,S_a\in\mathbb{C}^{(2N_L)\times N_G},\ T_a\in\mathbb{C}^{(2N_L)\times (2N_L)},\ F_0,\widetilde{G}\in\mathbb{C}^{(2N_AN_L)\times N_G},$ and $T\in\mathbb{C}^{(2N_AN_L)\times (2N_AN_L)}$. From now on, let $m:=2N_AN_L$, and $n:=N_G$.

To efficiently exploit the structure of the problem, the matrix T needs to be diagonal, with its diagonal elements equal to either 1 or -1 (possibly with some zeros in the case of a singular T). There is no theoretical obstacle to apply the simultaneous

(J-)orthogonalization in the computation of the GHSVD on the matrices F_0 , G, and T implicitly, but the repeated multiplication (in each reduction step) by T is slow. Therefore, T should be either factored concurrently, by using a modified version of the Hermitian indefinite factorization of all T_a blocks, or diagonalized concurrently: all factorizations (or diagonalizations) are independent of each other and can proceed in parallel. Since the diagonalization, compared to the Hermitian indefinite factorization, is a slower process, our choice is to factor all the diagonal blocks T_a .

4.2. Hermitian indefinite factorization. Each T_a is factored by the algorithm described in [27]. The algorithm for each T_a consists of the Hermitian indefinite factorization with a suitable pivoting [4], followed by the transformation of the block-diagonal matrix. Such factorization has the following form

$$(4.3) T_a = P_a^T M_a^* D_a M_a P_a,$$

where P_a is a permutation (in the LAPACK sense), M_a is upper triangular, and D_a is block-diagonal, with diagonal blocks of order 1 or 2.

Then, D_a is transformed into $\widehat{J}_a = \operatorname{diag}(\pm 1)$. If D_a has a diagonal block of order 1 at position k, then \widehat{J}_a stores the sign of this block in its kth diagonal element, and the kth row of M_a is scaled by $|(D_a)_{kk}|^{1/2}$. In the case of a (Hermitian) pivot block of order 2, this block is diagonalized by a Jacobi rotation R_k , and the corresponding two rows of M_a in (4.3) are multiplied by R_k . Two transformations of the new diagonal elements of D_a are then performed, as above. To speed-up the process, the rotation and the scaling of two rows of M_a are combined and then applied as a single transformation.

The outer permutations P_a are generated starting from the identity, and stored as the partial permutations of the principal submatrices, as in LAPACK, according to the pivoting of choice. Since the matrices T_a are of a relatively small order, our choice is the complete pivoting from [4].

4.3. Postprocessing. After the factorization, with a postprocessing step we obtain $T_a = \widehat{M}_a^* \widehat{J}_a \widehat{M}_a$, where $\widehat{M}_a = M_a P_a$, and \widehat{M}_a does not need to remain triangular.

Finally, by applying an inner permutation \widehat{P}_a , \widehat{J}_a can be rearranged into a diagonal matrix \widetilde{J}_a , where the positive signs precede the negative ones on the diagonal. This property of \widetilde{J}_a matrices can be exploited to speed-up computation of the hyperbolic scalar products x^*Jx in the subsequent phases (see subsections 4.5.1, 5.1, and 6.4.2). Let the whole factorization routine described thus far be called ZHEBPJ. Then, $T_a = \widetilde{M}_a^*\widetilde{J}_a\widetilde{M}_a$, $\widetilde{M}_a = \widehat{P}_a\widehat{M}_a$, and H_a is multiplied by \widetilde{M}_a from the left as $\widetilde{H}_a = \widetilde{M}_aH_a$.

After such preprocessing, H from (4.2) is written as

$$H = \widetilde{F}^* \widetilde{J} \widetilde{F}, \quad \widetilde{F}^* = [\widetilde{H}_1^*, \dots, \widetilde{H}_{N_A}^*], \quad \widetilde{J} = \operatorname{diag}(\widetilde{J}_1, \dots, \widetilde{J}_{N_A}).$$

In datasets A, each \widetilde{J}_a has 3 positive and 239 negative signs. In datasets B, a non-consecutive half of \widetilde{J}_a matrices are positive definite, and the others are negative definite.

4.4. Implementation and testing. The computational tasks for different indices a are fully independent, and are performed in parallel such that each thread is responsible for one or more indices a, as indicated in the pseudocode of Algorithm 4.1. Each thread, in turn, performs the three steps for its index a sequentially, up to a possible usage of a parallel BLAS in the first two steps. The last step, i.e., computing U_aB_a to assemble \tilde{G} , is a loop with the independent iterations, and could be done in parallel, using the nested parallelism within each thread, should N_G be large enough

Algorithm 4.1 A pseudocode for the Phase 1 algorithm.

```
for all atoms a, 1 \leq a \leq N_A do {an OpenMP parallel do} factorize T_a = \widetilde{M}_a^* \widetilde{J}_a \widetilde{M}_a; {ZHEBPJ with BLAS level 1 and 2 routines} multiply \widetilde{H}_a = \widetilde{M}_a H_a; {1 ZGEMM, of a 2N_L \times 2N_L and a 2N_L \times N_G matrix} scale the rows of B_a as U_a B_a; {N_L ZDSCALs, each on a row of N_G elements} end for
```

and should also the newly spawned threads for that loop have enough computational resources available to warrant the overhead of the additional thread management.

Each thread is responsible for allocating (MCDRAM is not explicitly used) and accessing the memory for the data it processes, so the data locality is achievable whenever each NUMA node has enough storage. The algorithm is thus viable in the heavily non-uniform memory access settings, like the Intel Xeon Phi's SNC-4 mode.

In a distributed memory setting (e.g., using the MPI processes), the assembling of \widetilde{F} , \widetilde{J} , and \widetilde{G} can be done by assigning to each process a (not necessarily contiguous) subrange of the iteration range of the for-all loop from Algorithm 4.1, while inside the process all atoms assigned to it are processed exactly as above, within an OpenMP parallel-do loop. The matrices \widetilde{F} , \widetilde{J} , and \widetilde{G} would then end up being distributed in the chunks corresponding to the chosen subranges among the processes.

4.4.1. Testing. In Table 4.1 the average per-atom wall execution time of Phase 1 is shown. The results suggest that it is beneficial to have more L2 data cache available

Table 4.1

The average per-atom wall execution time (wtime) of Phase 1 with 32 and 64 threads. Since the routine weights are rounded to the nearest per mil, their sum may not yield 100%. The first weight corresponds to ZHEBPJ, the second one to ZGEMM, and the third one to ZDSCALs step of Algorithm 4.1.

ID	average wtim 32 threads	e [s] per atom 64 threads	routine weig 32 threads	hts %:%:% 64 threads
A1	0.243186	0.277876	71.9:26.2: 1.9	66.1 : 32.1 : 1.8
A2	0.279318	0.312057	59.6:36.6: 3.8	53.8:41.4: 4.9
A3	0.345113	0.409165	47.9:46.4:5.7	41.4 : 46.2 : 12.4
A4	0.436803	0.536776	37.8:53.5: 8.8	31.7 : 50.9 : 17.5
В1	0.023099	0.027365	56.3:38.4: 5.3	49.9:45.6: 4.5
B2	0.030430	0.033248	40.6:52.3: 7.2	38.6:54.8: 6.7
В3	0.045586	0.060505	25.6:65.1: 9.3	20.2:71.7: 8.1
B4	0.070669	0.139536	16.1:71.6:12.3	8.3:80.9:10.8

per thread, as is the case with 32 threads overall. In the breakdown of the weights (i.e., percentages of time taken) of each computational step it is confirmed that ZGEMM starts to dominate the other computational steps of Algorithm 4.1 as the ratio n/m increases. It is a strong indication that even a procedure more expensive than ZHEBPJ, such as a diagonalization of T_a , may be applied on the datasets having a square-like shape, without considerably degrading the relative performance of Phase 1.

For a fully vectorized, cache-friendly alternative to applying ZDSCAL with a non-unit stride in Algorithm 4.1 please refer to section S.1 of the supplementary material.

4.5. An alternative way forward. After this phase has completed, one can proceed as described in the rest of the paper, should the condition numbers of (the

yet unformed) matrices H and S be large enough to severely affect the accuracy of a direct solution of the generalized Hermitian eigenproblem with the pair (H, S).

An alternative and more time-efficient way to proceed would be to explicitly form H and S. For $S = \widetilde{G}^*\widetilde{G}$, one ZHERK call would suffice. For $H = \widetilde{F}^*\widetilde{J}\widetilde{F}$, a copy of \widetilde{F} should be made, and that copy's rows should be scaled in parallel by the diagonal elements of \widetilde{J} . One ZGEMM call on \widetilde{F}^* and $\widetilde{J}\widetilde{F}$ then completes the formation of H. After that, an efficient solver for the generalized Hermitian eigenproblem can be employed on (H,S), such as ZHEGV or ZHEGVD from LAPACK, as shown in subsetion 6.5.4.

4.5.1. Row scaling. A cache-friendly implementation of the row scaling by \widetilde{J} is to iterate sequentially over the rows of a fixed column j, and change the sign of each element \widetilde{F}_{ij} for which $\widetilde{J}_{ii} = -1$, while the outer parallel-do loop iterates over all column indices j. However, that implementation can be optimized further.

If \widetilde{J} has its diagonal partitioned into (regularly or irregularly sized) blocks of the same sign, then it can be compactly encoded as a sequence of pairs $(i_-,l)_k$, one for each block of negative signs, where i_- is the first index belonging to a block k, and $l \geq 1$ is the block's length. The iteration over all rows and the conditional sign changes as above can be replaced by iteration over all such blocks. For each block, iterate sequentially in the range of indices i from i_- to i_-+l-1 , and change the signs unconditionally, thus eliminating the conditional branching based on the sign of \widetilde{J}_{ii} .

Such run-length-like encoding is employed in Phase 3, where it also accelerates the hyperbolic dot products in the case where the positive signs precede the negative ones on the diagonal of a sign matrix (i.e., at most one negative block exists) given by ZHEBPJ when forming the square factors for the inner Hari–Zimmermann method.

5. Phase 2 – optional (J, I) URV factorization. The one-sided Jacobi-type algorithms are fastest if they work on square matrices, since the column dot-products and updates are the shortest possible. If the square factors F, G, and the corresponding J of the matrix pair $(\widetilde{F}^*\widetilde{J}\widetilde{F}, \widetilde{G}^*\widetilde{G})$, can be found, instead of the rectangular factors $\widetilde{F}, \widetilde{G}$, and the corresponding \widetilde{J} , we expect that the overhead of such a shortening will be less than the computational time saved by avoiding the rectangular factors. To this end, matrix \widetilde{F} is shortened by using the hyperbolic QR factorization (also called the JQR factorization), according to the given \widetilde{J} :

$$(5.1) P_1 \widetilde{F} P_2 = Q_F F; \quad \widetilde{Q}_F^* \widetilde{J} \widetilde{Q}_F = J, \quad \widetilde{Q}_F := P_1^T Q_F,$$

where $F \in \mathbb{C}^{n \times n}$ is block upper triangular with diagonal blocks of order 1 or 2, $J = \operatorname{diag}(\pm 1) \in \mathbb{R}^{n \times n}$ is the shortened signature matrix, $P_2 \in \mathbb{R}^{n \times n}$ and $P_1 \in \mathbb{R}^{m \times m}$ are the column and the row permutation matrix, respectively. Our application does not use $Q_F \in \mathbb{C}^{m \times n}$, so it is not explicitly formed. From (5.1) it holds

$$P_2^T \widetilde{F}^* \widetilde{J} \widetilde{F} P_2 = F^* Q_F^* P_1 \widetilde{J} P_1^T Q_F F = F^* \widetilde{Q}_F^* \widetilde{J} \widetilde{Q}_F F = F^* J F.$$

Since the JQR requires both row and column pivoting (see [24]), matrix G, with its columns preparemented according to P_2 , the column pivoting of the JQR, will then be factored by the ordinary (tall-and-skinny) QR factorization (e.g., by the LAPACK routine ZGEQR). The latter QR factorization does not employ column pivoting, but in principle the row pivoting or presorting may be used:

$$P_3(\widetilde{G}P_2) = Q_GG; \quad \widetilde{Q}_G^*\widetilde{Q}_G = I_n, \quad \widetilde{Q}_G := P_3^TQ_G,$$

where $G \in \mathbb{C}^{n \times n}$ is upper triangular, and $Q_G \in \mathbb{C}^{m \times n}$, which is not needed in our application. We also do not depend on the special forms of F and G later on.

From (3.1) it follows that the \widetilde{J} -GHSVD of \widetilde{F} and \widetilde{G} , and that of $\widetilde{F}P_2$ and $\widetilde{G}P_2$, differ only in the column permutation of the right singular vectors, i.e., $\widetilde{X} = XP_2^T$, or, from (3.2), the row permutation of the eigenvectors, i.e., $\widetilde{Z} = P_2Z$, while Σ , and thus Λ , stay the same. Therefore, the square factors F, G, and J can be used in place of \widetilde{F} , \widetilde{G} , and \widetilde{J} throughout the rest of the computation, and then the results could be easily converted back to the ones of the original problem.

Whenever \widetilde{G} might be badly conditioned, Phase 2 could be skipped, or we should resort to a slower but more stable QR factorization of $\widetilde{G}P_2$ with the column pivoting,

$$P_4(\widetilde{G}P_2)P_5 = Q_{G'}G'; \quad \widetilde{Q}_{G'}^*\widetilde{Q}_{G'} = I_n, \quad \widetilde{Q}_{G'} := P_4^T Q_{G'},$$

where P_5 has to be applied back to F, and P_4 comes from an optional row pivoting. The column-pivoted QR factorization is provided by the LAPACK routine ZGEQP3.

Then, $F' := FP_5$ and G' (provided that it is not rank-deficient according to a user-defined tolerance) could be substituted for F and G in the rest of the computation. For X (or Z) thus obtained it holds $\widetilde{X}P' = X$ (or, $P'Z = \widetilde{Z}$), where $P' := P_2P_5$. Such an approach is not required for our datasets, and therefore it was not tested.

5.1. \widetilde{J} -dot products and norms. Since \widetilde{J} can in principle contain the positive and the negative signs in any order, and vectorization is strongly desired, a \widetilde{J} -dot product of two vectors, $f^*\widetilde{J}g$, is computed as \mathbf{v} piecewise sums Σ_j , $1 \leq j \leq \mathbf{v}$,

$$Re(\Sigma_j) = Re(\Sigma_j) + \widetilde{J}_i(Re(f_i) Re(g_i) + Im(f_i) Im(g_i)),$$

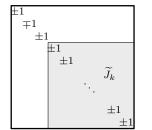
$$Im(\Sigma_j) = Im(\Sigma_j) + \widetilde{J}_i(Re(f_i) Im(g_i) - Im(f_i) Re(g_i)),$$

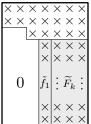
where i starts with a value of j and increments in steps of \mathbf{v} up to m. The components $\mathrm{Re}(\Sigma)$ and $\mathrm{Im}(\Sigma)$ of the resulting Σ are obtained by SUM-reducing $\mathrm{Re}(\Sigma_j)$ and $\mathrm{Im}(\Sigma_j)$, respectively. Similarly, the square of the \widetilde{J} -norm, $f^*\widetilde{J}f$ is computed by SUM-reducing Σ'_j , where

$$\Sigma'_j = \Sigma'_j + \widetilde{J}_i(\operatorname{Re}(f_i)^2 + \operatorname{Im}(f_i)^2).$$

The square of the \widetilde{J} -"norm" of a vector thus obtained can be positive or negative, with a possibility of cancellations inadvertently occurring in the summations. It is an open question how to compute the squares of the \widetilde{J} -"norms" both efficiently and accurately, though one possible speed improvement might be to encode \widetilde{J} as described in subsection 4.5.1 and simplify the above three piecewise summations accordingly.

- **5.2. Pivoting.** To achieve the maximal numerical stability, the JQR factorization is usually performed with complete pivoting. In the first step the pivot column(s) are chosen from the \widetilde{J} -Grammian matrix $H = \widetilde{F}^* \widetilde{J} \widetilde{F}$, and later on, in the kth step, from the \widetilde{J}_k -Grammian matrix $H_k = \widetilde{F}_k^* \widetilde{J}_k \widetilde{F}_k$, where \widetilde{F}_k the is a part of the matrix yet to be reduced, and \widetilde{J}_k is the matrix of signs that corresponds to the unreduced matrix \widetilde{F}_k (see Fig. 5.1). The complete pivoting in the first step needs formation of the whole H, i.e., $\mathcal{O}(mn^2)$ floating-point operations. Such an approach, consistently implemented throughout the algorithm, leads to $\mathcal{O}(m^2n^2)$ operations solely for the choice of pivots. Therefore, we relaxed the pivoting strategy to the diagonal pivoting supplemented with the partial pivoting [3, Algorithm C].
- **5.2.1. Diagonal and partial pivoting.** First, n-k+1 squares of the \widetilde{J}_k -norms $h_{ii}^{[k]} := \widetilde{f}_i^* \widetilde{J}_k \widetilde{f}_i$, where \widetilde{f}_i is the *i*th column of \widetilde{F}_k and $i \geq 1$, are computed in a parallel-do loop over i and stored in a work array. Since all columns have the length





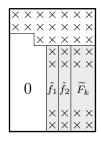


Fig. 5.1. Choosing a single pivot column or two pivot columns. Matrices \widetilde{J}_k and \widetilde{F}_k are shaded.

of m-k+1, each parallel loop iteration executes (sequentially) in approximately the same time and the work is therefore well balanced among the threads.

Let $j \geq 1$ be the smallest index such that $|h_{jj}^{[k]}| \geq |h_{ii}^{[k]}|$ for all i. If j > 1, the kth and the (j + k - 1)th columns of \widetilde{F} (and thus also the first and the jth columns of \widetilde{F}_k) are swapped. If k = n, the column pivoting is completed.

Otherwise, n-k \widetilde{J}_k -dot products $h_{1j}^{[k]}:=\widetilde{f}_1^*\widetilde{J}_k\widetilde{f}_j$, where \widetilde{f}_j is the jth column of \widetilde{F}_k and j>1, are computed in a parallel-do loop over j and stored in a complex workspace, while their magnitudes $|h_{1j}^{[k]}|$ are placed in a real workspace. Same as above, this work is well balanced among the threads.

Let i > 1 be the smallest index such that $|h_{1i}^{[k]}| \ge |h_{1j}^{[k]}|$ for all j > 1. As in [3], if $|h_{1i}^{[k]}| \ge \alpha |h_{1i}^{[k]}|$, with $\alpha := (1 + \sqrt{17})/8$, the column pivoting in the step k is completed.

Otherwise, n-k \widetilde{J}_k -dot products $h_{il}^{[k]} := \widetilde{f}_i^* \widetilde{J}_k \widetilde{f}_l$, where \widetilde{f}_l is the lth column of \widetilde{F}_k and $i \neq l \geq 1$, are computed in a parallel-do loop over l and their magnitudes $|h_{il}^{[k]}|$ are stored in a real workspace. This work is only slightly imbalanced among threads, since for l=i a thread assigned to the iteration sets $|h_{ll}^{[k]}|=0$, excluding the value and its index from the search for a maximum, unless all other values are also 0.

Let $j \geq 1$ be the smallest index such that $|h_{ij}^{[k]}| \geq |h_{il}^{[k]}|$ for all l. As in [3], if $|h_{11}^{[k]}||h_{ij}^{[k]}| \geq \alpha |h_{1i}^{[k]}|^2$, the column pivoting for the step k is completed; else, if $|h_{ii}^{[k]}| \geq \alpha |h_{ij}^{[k]}|$, the kth and the (i+k-1)th columns of \widetilde{F} (and thus also the first and the ith columns of \widetilde{F}_k) are swapped and the column pivoting in the step k is completed.

Otherwise, a 2×2 pivot is chosen, by taking the first column of \widetilde{F}_k and swapping the (k+1)th and (i+k-1)th columns of \widetilde{F} (and thus also the second and the *i*th columns of \widetilde{F}_k), if $i \neq 2$; else, the second pivot column is already in place.

The pivot column(s) have thus been brought to the front of the matrix \widetilde{F}_k by at most two column swaps. The ensuing row pivoting is explained further below.

5.3. Hyperbolic Householder reflectors. If a single pivot is chosen, the first column \tilde{f}_1 of \tilde{F}_k is reduced by a hyperbolic Householder reflector [25, Theorem 4.4] to a vector $f_1 = c_1 e_1$, where $c_1 \in \mathbb{C}$ and e_1 is the first vector of the canonical base. A variant of [25, Theorem 4.4] for the hyperbolic scalar product and a simple shape of f_1 follows.

THEOREM 5.1. Let \widetilde{J}_k be a hyperbolic scalar product matrix of order ℓ . Let $\widetilde{f}_1, f_1 \in \mathbb{C}^{\ell}$ be two distinct vectors. There exists a basic \widetilde{J}_k reflector H(w),

(5.2)
$$H(w) = I - 2w(w^*\widetilde{J}_k w)^+ w^*\widetilde{J}_k,$$

such that $H(w)\tilde{f}_1 = f_1$ if and only if \tilde{f}_1 and f_1 satisfy the \tilde{J}_k -isometry and \tilde{J}_k -symmetry property, respectively

$$(5.3) \widetilde{f}_1^* \widetilde{J}_k \widetilde{f}_1 = f_1^* \widetilde{J}_k f_1,$$

(5.4)
$$\tilde{f}_1^* \tilde{J}_k f_1 = f_1^* \tilde{J}_k \tilde{f}_1,$$

and $d = \tilde{f}_1 - f_1 \neq 0$ is nondegenerate, i.e., $d^* \tilde{J}_k d \neq 0$. Furthermore, whenever H(w) exists, it is unique. H(w) can be generated by any $w \in \mathbb{C}^\ell$ such that $w = \lambda d$, $\lambda \in \mathbb{C} \setminus \{0\}$. Finally, the same remains valid if we replace f_1 by $-f_1$, and d by $s = f_1 + \tilde{f}_1$.

If $\tilde{f}_1 = f_1$, there is nothing to do in this step, so we take H(w) = I. Otherwise, since we want to obtain f_1 in the form $f_1 = c_1 e_1$, equation (5.3) is equivalent to a requirement that the sign of \tilde{j}_{11} , the first diagonal element of \tilde{J}_k , is equal to the sign of $\tilde{f}_1^* \tilde{J}_k \tilde{f}_1 = |c_1|^2 \tilde{j}_{11} = f_1^* \tilde{J}_k f_1$. As $\tilde{f}_1^* \tilde{J}_k \tilde{f}_1$ has already been computed by the diagonal pivoting (see subsection 5.2.1), it is trivial to check if the requirement holds.

If not, it can be shown that there exists at least one element with the correct sign in \widetilde{J}_k . When there is more than one such element, our implementation sequentially finds the one that corresponds to the largest element in \widetilde{f}_1 by magnitude (say, lth). By permuting the diagonal of \widetilde{J}_k , this element can be brought to the first diagonal position. This implies a corresponding row permutation of \widetilde{F}_k that swaps the first and the lth rows of \widetilde{F}_k . For that, we employ the parallel ZSWAP routine, but should the rows be short enough, a sequential version of the routine could be considered instead.

Relation (5.3) implies that $|c_1| = |\tilde{f}_1^* \tilde{J}_k \tilde{f}_1|^{1/2}$. In general, c_1 is a complex number, $c_1 = re^{i\delta}$, and we have already determined $r = |c_1|$. It remains to find $\delta = \arg(c_1)$. From (5.4) it follows

(5.5)
$$\bar{\tilde{f}}_{11}\tilde{j}_{11}c_1 = \bar{c}_1\tilde{j}_{11}\tilde{f}_{11},$$

where $\tilde{f}_{11} := r_1 e^{\mathrm{i}\delta_1}$ is the first element in \tilde{f}_1 . Relation (5.5) can be divided by \tilde{j}_{11} and written as $r_1 r e^{\mathrm{i}(\delta - \delta_1)} = r_1 r e^{-\mathrm{i}(\delta - \delta_1)}$. Since $r_1, r \neq 0$, we may choose $\delta = \delta_1$ for $\arg(c_1)$. We only need to compute $e^{\mathrm{i}\delta}$, so $e^{\mathrm{i}\delta} = \tilde{f}_{11}/|\tilde{f}_{11}|$. Now we have satisfied the conditions (5.3)–(5.4) for construction of a reflector that maps \tilde{f}_1 to c_1e_1 , or to $-c_1e_1$.

We aim to compute d or s accurately. First, $w^*\widetilde{J}_k w$ is needed, where $w = \tilde{f}_1 \pm f_1$ (with the addition for s and the subtraction for d). Since \tilde{f}_1 and f_1 satisfy (5.3)–(5.4),

$$\begin{split} w^* \widetilde{J}_k w &= \left(\widetilde{f}_1 \pm f_1 \right)^* \widetilde{J}_k \left(\widetilde{f}_1 \pm f_1 \right) = 2 \left(\widetilde{f}_1^* \widetilde{J}_k \widetilde{f}_1 \pm f_1^* \widetilde{J}_k \widetilde{f}_1 \right) = 2 \left(\widetilde{f}_1^* \widetilde{J}_k \widetilde{f}_1 \pm \overline{c}_1 \widetilde{\jmath}_{11} \widetilde{f}_{11} \right) \\ &= 2 \left(\widetilde{f}_1^* \widetilde{J}_k \widetilde{f}_1 \pm |\widetilde{f}_1^* \widetilde{J}_k \widetilde{f}_1|^{1/2} \frac{\widetilde{f}_{11}}{|\widetilde{f}_{11}|} \widetilde{\jmath}_{11} \widetilde{f}_{11} \right) = 2 \left(\widetilde{f}_1^* \widetilde{J}_k \widetilde{f}_1 \pm |\widetilde{f}_1^* \widetilde{J}_k \widetilde{f}_1|^{1/2} |\widetilde{f}_{11}| \widetilde{\jmath}_{11} \right). \end{split}$$

Since (5.3) should be valid, we have $\operatorname{sgn}(\tilde{f}_1^* \widetilde{J}_k \tilde{f}_1) = \tilde{j}_{11}$, and both terms in the previous relation have the same sign. Therefore, to avoid unnecessary cancellation, our choice is w = s, and H(s) from (5.2) is then equal to

$$H(s) = I + \tau s s^* \widetilde{J}_k, \quad \tau = -1/\big(\widetilde{f}_1^* \widetilde{J}_k \widetilde{f}_1 + |\widetilde{f}_1^* \widetilde{J}_k \widetilde{f}_1|^{1/2} \, |\widetilde{f}_{11}| \widetilde{\jmath}_{11}\big).$$

The update of a column \tilde{f}_j , j > 1, is performed as $H(s)\tilde{f}_j = \tilde{f}_j + \tau s(s^*\tilde{J}_k\tilde{f}_j)$, which involves computing the \tilde{J}_k -dot product of s and \tilde{f}_j , scaling it by τ , and calling the ZAXPY BLAS 1 routine. We opted for a sequential ZAXPY version, but it can be argued that for the extremely long columns a parallel version would be better suited.

The column updates are mutually independent, and therefore can be performed in a parallel-do loop over j, with the work being well balanced among the threads.

In our application it is not required that the reflector generators $(s_k \text{ and } \tau_k \text{ in the step } k)$ are preserved, but that is nevertheless done with time overhead close to none in a separate complex matrix (the vector s_k is stored in its kth column, with the leading rows set to 0) and a real vector (the scalar τ_k is stored as its kth element).

The situation is more complicated if a pair of pivot columns is chosen. According to [25], the reduction can be performed by the block Householder matrix defined by these two columns. However, the computation of a block Householder reflector is a more difficult approach than the computation of a variant of the URV factorization with U hyperbolic, not unitary. The proposed factorization is similar, but not equal to the hyperbolic (sometimes also called signed) URV factorization presented in [30].

DEFINITION 5.2. Let $F \in \mathbb{C}^{m \times n}$ and $J \in \mathbb{Z}^{m \times m}$, $J = \operatorname{diag}(\pm 1)$, be given matrices. Let $J' = P^T J P$ for any permutation matrix P. A factorization F = U R V, where $U \in \mathbb{C}^{m \times m}$ is J'-unitary, $V \in \mathbb{C}^{n \times n}$ is unitary, and $R = \begin{bmatrix} R_0 \\ 0 \end{bmatrix} \in \mathbb{C}^{m \times n}$, with $R_0 \in \mathbb{C}^{n \times n}$ upper triangular, is called a (J, I) URV factorization according to J.

Note that the (J, I) URV factorization is not unique.

The Grammian matrix A_{12} of the pivot block $\widetilde{F}_{12} := [\widetilde{f}_1 \ \widetilde{f}_2]$, i.e., $A_{12} := \widetilde{F}_{12}^* \widetilde{J}_k \widetilde{F}_{12}$ (see Fig. 5.1) is nonsingular. Since the pivot strategy has chosen this block for the pivot block, the off-diagonal elements are larger in magnitude than the diagonal elements of A_{12} . The matrix A_{12} will be diagonalized by a single Jacobi rotation R_k ,

$$R_k^* A_{12} R_k = R_k^* \widetilde{F}_{12}^* \widetilde{J}_k \widetilde{F}_{12} R_k = \text{diag}(d_k, d_{k+1}).$$

This shows that the columns $[\hat{f}_1 \ \hat{f}_2] := \widetilde{F}_{12} R_k$ are mutually \widetilde{J}_k orthogonal – their scalar product is zero, and the squares of their \widetilde{J}_k norms are d_k and d_{k+1} , respectively. Since these columns are mutually \widetilde{J}_k orthogonal, the only way to handle them is applying two successive hyperbolic Householder transformations to reduce the pivot matrix \widetilde{F}_{12} to an upper triangular matrix \widehat{F}_{12} .

The first hyperbolic transformation, by $H(s_k)$, reduces the first column to a single element. The second hyperbolic transformation, by $H(s_{k+1})$, then reduces the elements, already transformed by $H(s_k)$, of the shortened (the first row of \widehat{F}_{12} is not changed anymore) second column. We should then multiply \widehat{F}_{12} by R_k^* to obtain the reduced matrix F_{12} with, generally full, topmost 2×2 block

$$F_{12} := \widehat{F}_{12} R_k^* = \begin{bmatrix} f_{11} & f_{21} & 0 & \cdots & 0 \\ f_{12} & f_{22} & 0 & \cdots & 0 \end{bmatrix}^T.$$

Finally, the step counter k is incremented by 2 (instead of 1 for a single pivot), and the process continues again with the column pivoting, as in subsection 5.2.

When present, the multiplications from right, first by the Jacobi rotation R_k , and then by R_k^* , cancel each other (since $R_k R_k^* = I_2$), so the right matrix V in this URV-like factorization is in fact identity.

The Jacobi rotations are applied by calling the parallel ZROT routine, but should the columns be short enough, a sequential version of the routine could be warranted.

After the final step of the reduction (for k = n), the new F and J are square matrices of order n. More precisely, J is the leading part of \widetilde{J} as left after all permutations due to the row pivoting, and F, unlike the standard compact representation of the LAPACK's QR factorizations, has zeros below the diagonal block set explicitly.

5.4. Testing. In Table 5.1 the wall execution times of both the JQR and the TSQR are shown. It is evident that there is still room for the future JQR's efficiency improvement, which might be achieved by blocking and delaying the columns updates.

Table 5.1

The wall execution time (wtime) in seconds of the Phase 2 steps: the hyperbolic QR (JQR), the tall-and-skinny QR (ZGEQR from LAPACK), and the prepermutation of the columns of \widetilde{G} by P_2 (max. wtime with 32, 64, and 48 threads). In the last column is the number of 2×2 pivots chosen.

ID	JO	QR wtime	[s]	TS	QR wtime	$\widetilde{G}P_2$	2×2	
110	32 thr.	64 thr.	48 thr.	32 thr.	$64 \mathrm{thr.}$	48 thr.	wtime [s]	pivots
A1	162.52	179.02	66.52	10.53	7.19	1.80	0.31	11
A2	449.13	482.63	184.46	14.06	8.26	4.07	0.87	9
A3	1041.92	1104.48	383.60	22.01	14.71	8.03	1.14	7
A4	2137.63	2249.65	753.71	42.82	27.38	16.54	2.37	6
B1	181.13	177.44	64.72	12.00	12.21	3.41	0.06	16
B2	520.15	495.70	186.48	15.63	10.29	4.28	0.69	16
В3	1263.68	1182.83	420.02	36.66	20.22	10.29	1.59	20
B4	2780.19	2439.50	841.83	49.04	33.12	19.35	3.41	16

5.4.1. Prepermuting of \widetilde{G} . Table 5.1 also contains the wall time of preparing $\widetilde{G}P_2$ in parallel. The fastest way to prepermute the columns of \widetilde{G} is to copy them to another matrix of the same size, with the column j going to $\pi_2(j)$, where π_2 denotes the permutation represented by P_2 . Such copying occurs in a parallel-do loop over j.

6. Phase 3 – generalized hyperbolic SVD. Vjeran Hari in his PhD thesis [9] developed a method for solving the generalized eigenproblem, when at least one of the two matrices is positive definite. The method is based on the ideas from the PhD thesis of Katharina Zimmermann [31], and has been revisited recently in [10, 11].

Based on the Hari-Zimmermann algorithm for the generalized eigenproblem, in [21] a one-sided method for computing the real generalized SVD has been derived. The main trick, how to obtain a one-sided method for the SVD from the two-sided method for the eigenproblem is always the same: think about the transformations in the two-sided fashion, and apply them from one (right or left) side on a matrix factor.

Since the elements of the pivot submatrices \widehat{H} of H and \widehat{S} of S are the scalar products of the columns of F and G, respectively, it is easier to write the transformations in terms of the elements of \hat{H} and \hat{S} ,

$$(6.1) \qquad \widehat{H} = \begin{bmatrix} h_{pp} & h_{pq} \\ \overline{h}_{pq} & h_{qq} \end{bmatrix} = \begin{bmatrix} f_p^* J f_p & f_p^* J f_q \\ \overline{f}_p^* J f_q & f_q^* J f_q \end{bmatrix}, \quad \widehat{S} = \begin{bmatrix} s_{pp} & s_{pq} \\ \overline{s}_{pq} & s_{qq} \end{bmatrix} = \begin{bmatrix} g_p^* g_p & g_p^* g_q \\ \overline{g}_p^* g_q & g_q^* g_q \end{bmatrix},$$

instead of in terms of the columns f_p , f_q , g_p , and g_q . The original method consists of 3 active transformations, and an auxiliary transformation that helps in coupling them all together.

6.1. Pointwise algorithm. The whole pointwise algorithm (with a pair of 2×2 pivot submatrices in each annihilation step) is taken from the PhD thesis of Vieran Hari [9]. However, we feel that the algorithm should be presented succinctly here, to aid its implementors, and also to incorporate some minor corrections.

6.1.1. Preprocessing. In the preprocessing step, H and S are scaled by a diagonal matrix D such that diag(DSD) = I, i.e.,

$$H_0 := DHD, \quad S_0 := DSD, \quad D = \operatorname{diag}\left(s_{11}^{-1/2}, s_{22}^{-1/2}, \dots, s_{nn}^{-1/2}\right).$$

Such preprocessing can be done only once, at the start of the algorithm, or it can be done before each annihilation step for the pivot column pair in question. The latter might seem redundant, but in a floating-point realization of the algorithm, after enough steps, diag(DSD) could veer off the identity enough to warrant such a rescaling. In that case, form a matrix D_0 that has, for the chosen pivot indices (p,q), as its pth and qth diagonal entries $(s_{pp})^{-1/2}$ and $(s_{qq})^{-1/2}$, respectively, while being equal to the identity elsewhere. For the approach with a single prescaling, let $D_0 = I$.

In both cases, let D_0 be a 2×2 restriction of D_0 to the pth and the qth rows and columns. We have implemented the pivot pair prescaling in each annihilation step.

6.1.2. Diagonalization of \widehat{S}_0 . In the first step the 2×2 pivot submatrix \widehat{S}_0 of S_0 (at the crossings of the pth and the qth rows and columns) is diagonalized by a complex Jacobi rotation \widehat{R}_1 , where \widehat{R}_k for $k \geq 1$ is

(6.2)
$$\widehat{R}_k = \begin{bmatrix} \cos \varphi_k & e^{i\alpha_k} \sin \varphi_k \\ -e^{-i\alpha_k} \sin \varphi_k & \cos \varphi_k \end{bmatrix}.$$

The same transformation is then applied to H_0 , to keep the new pair equivalent to the original one. After that, the new pair is $(H_1, S_1) := (R_1^* H_0 R_1, R_1^* S_0 R_1)$, where $R_1 = I$, except at the pivot positions, where $R_1 = \widehat{R}_1$. Since H and S have been preprocessed as in subsection 6.1.1, the diagonal elements of \widehat{S}_0 are the same, and we may choose $\varphi_1 = -\pi/4$ in (6.2). Now it is easy to determine that $\alpha_1 = \arg(s_{pq})$. If s_{pq} , written in the trigonometric form, was $s_{pq} = xe^{i\alpha_1}$, with $x = |s_{pq}|$, before the transformation, then after it we obtain

(6.3)
$$\widehat{S}_1 = \text{diag}(1+x, 1-x).$$

6.1.3. Rescaling of \widehat{S}_1 . The second step rescales the diagonal of S_1 to ones, and rescales H_1 with the same diagonal matrix. After the transformation, similar to the preprocessing step, we obtain $H_2 := D_2H_1D_2$, $S_2 := D_2S_1D_2$. From (6.3) we conclude that $\widehat{D}_2 = \operatorname{diag}((1+x)^{-1/2}, (1-x)^{-1/2})$. Elements of \widehat{H}_2 are

$$(6.4) \quad \widehat{H}_2 = \begin{bmatrix} h_{pp}^{(2)} & h_{pq}^{(2)} \\ \bar{h}_{pq}^{(2)} & h_{qq}^{(2)} \end{bmatrix} = \frac{1}{2} \begin{bmatrix} \frac{1}{1+x} (h_{pp} + h_{qq} + 2u) & \frac{e^{\mathrm{i}\alpha_1}}{\sqrt{1-x^2}} (h_{qq} - h_{pp} + 2\mathrm{i}v) \\ \frac{e^{-\mathrm{i}\alpha_1}}{\sqrt{1-x^2}} (h_{qq} - h_{pp} - 2\mathrm{i}v) & \frac{1}{1-x} (h_{pp} + h_{qq} - 2u) \end{bmatrix},$$

where

(6.5)
$$u + iv = e^{-i\alpha_1} h_{pq} = e^{-i\arg(s_{pq})} h_{pq}.$$

6.1.4. Diagonalization of \widehat{H}_2 . In the third step the pivot submatrix \widehat{H}_2 of H_2 is diagonalized by a complex Jacobi rotation \widehat{R}_3 of the form (6.2). The third transformation can be written as $H_3 = R_3^* H_2 R_3$, $S_3 = R_3^* S_2 R_3$, where $R_3 = I$, except at the pivot positions, where $R_3 = \widehat{R}_3$. Then φ_3 in (6.2) is written as $\varphi_3 = \vartheta + \pi/4$ to express the transformations in terms of ϑ . The relations for the angles of a rotation that annihilates the off-diagonal element of \widehat{H}_2 are given by

$$(6.6) \qquad \tan\left(2\vartheta + \frac{\pi}{2}\right) = \sigma \frac{2|h_{pq}^{(2)}|}{h_{qq}^{(2)} - h_{pp}^{(2)}}, \qquad \alpha_3 = \alpha_1 + \arg\left(\frac{h}{2} + \mathrm{i}v\right) + (1 - \sigma)\frac{\pi}{2},$$

where $h = h_{qq} - h_{pp}$. The requirement $\gamma := \alpha_3 - \alpha_1 \in (-\pi/2, \pi/2]$ yields $\sigma = \operatorname{sgn}(h)$. This choice of σ as the sign of h and the constraint $\vartheta \in (-\pi/4, \pi/4]$ ensures the convergence of the algorithm (see [9]).

Since $\tan(2\vartheta + \pi/2) = -\cot(2\vartheta) = -1/\tan(2\vartheta)$, from (6.4) and (6.6), we obtain

(6.7)
$$\tan(2\vartheta) = \sigma \frac{2u - (h_{pp} + h_{qq})x}{t\sqrt{h^2 + 4v^2}}, \qquad t := \sqrt{1 - x^2}.$$

After the first three steps, the pivot submatrix \widehat{S}_3 is still diagonal (in fact identity), i.e., $\widehat{S}_3 = \widehat{R}^* \widehat{S}_0 \widehat{R} = I$, where

$$\widehat{R} = \widehat{R}_1 \widehat{D}_2 \widehat{R}_3.$$

This constructively shows that \widehat{R} diagonalizes the matrix pair $(\widehat{H}_0, \widehat{S}_0)$.

6.1.5. Forming \widehat{R} and \widehat{Z}' . Hari in [9, Theorem 2.2] has proved the form of the general nonsingular matrix that diagonalizes a 2×2 Hermitian positive definite matrix. The intention of representing \widehat{R} as

$$\widehat{R} = \frac{1}{t} \begin{bmatrix} \cos \varphi & e^{\mathrm{i}\alpha} \sin \varphi \\ -e^{-\mathrm{i}\beta} \sin \psi & \cos \psi \end{bmatrix} \mathrm{diag}(e^{\mathrm{i}\sigma_p}, e^{\mathrm{i}\sigma_q})$$

is to simplify (6.8). Comparing the elements of (6.8) and (6.9) we obtain

$$(6.10) \begin{aligned} \cos \varphi &= \left(1/\sqrt{2}\right) \cdot \sqrt{1 + x \sin(2\vartheta) + t \cos \gamma \cos(2\vartheta)}, \quad 0 \le \varphi < \pi/2, \\ \cos \psi &= \left(1/\sqrt{2}\right) \cdot \sqrt{1 - x \sin(2\vartheta) + t \cos \gamma \cos(2\vartheta)}, \quad 0 \le \psi < \pi/2, \\ e^{\mathrm{i}\alpha} \sin \varphi &= e^{\mathrm{i}\alpha_1} \cdot \left(\left(\sin(2\vartheta) - x\right) + \mathrm{i}t \sin \gamma \cos(2\vartheta)\right) / \left(2\cos\psi\right), \\ e^{-\mathrm{i}\beta} \sin \psi &= e^{-\mathrm{i}\alpha_1} \cdot \left(\left(\sin(2\vartheta) + x\right) - \mathrm{i}t \sin \gamma \cos(2\vartheta)\right) / \left(2\cos\varphi\right). \end{aligned}$$

Since in (6.10) we need only $\sin \gamma$ and $\cos \gamma$, from (6.6) it then follows

$$(6.11) \tan \gamma = 2v/h.$$

The fourth step only deals with a formal simplification of \widehat{R} , by introducing a transformation Φ_4 such that $H_4 = \Phi_4^* H_3 \Phi_4$, $S_4 = \Phi_4^* S_3 \Phi_4$. The matrix Φ_4 is a diagonal matrix equal to identity, except at pivot positions, where $\widehat{\Phi}_4 = \operatorname{diag}(e^{-\mathrm{i}\sigma_p}, e^{-\mathrm{i}\sigma_q})$. Obviously, if \widehat{R} diagonalizes the pair $(\widehat{H}_0, \widehat{S}_0)$, then the transformation $\widehat{Z} = \widehat{R}\widehat{\Phi}_4$ will leave the final diagonal matrices intact. Then, let $\widehat{Z}' = \widehat{D}_0 \widehat{Z}$.

6.1.6. Exceptional cases. There can be a few exceptions in the computations of the elements of the matrix \hat{Z} which have to be accounted for in the algorithm.

If $h_{pq} = s_{pq} = 0$, we set $\widehat{Z} = I$, since both pivot submatrices are already diagonal. We could still apply the scaling by \widehat{D}_0 , but not count that as a transformation.

If $h_{pq} \neq 0$, but $s_{pq} = 0$ (i.e., x = 0), we set $\alpha_1 = 0$ and proceed as described above to determine \hat{Z} as an ordinary Jacobi rotation that diagonalizes \hat{H}_0 .

If h = v = 0 in (6.11), i.e., when $\arg s_{pq} = \arg h_{pq}$ and $h_{pp} = h_{qq}$, it can be shown that \widehat{R}_1 , the Jacobi rotation that diagonalizes \widehat{S}_0 , also diagonalizes \widehat{H}_0 , so $\widehat{Z} = \widehat{R}_1 \widehat{D}_2$.

6.1.7. Convergence criterion and finalization. As in [21], the convergence criterion in floating-point arithmetic has to take into account the relative magnitudes

of the off-diagonal elements, compared to the diagonal ones, in the pivot pairs. Therefore, a pivot pair undergoes the transformation if, for the machine precision ε ,

$$(6.12) |h_{pq}| \ge \left(\max\{|h_{pp}|, |h_{qq}|\} \cdot \left(\varepsilon\sqrt{n}\right)\right) \cdot \min\{|h_{pp}|, |h_{qq}|\} \text{or} |s_{pq}| \ge \varepsilon\sqrt{n}.$$

If \widehat{Z}' turns out to be identity, it is not applied. If $\cos \varphi = \cos \psi = 1$, such transformation is considered "small", and "big" otherwise. Near the end of the process, the transformations turn out to be small, and (in a blocking variant, the last level of) the algorithm is stopped when no big transformations are encountered in a sweep, to avoid perpetually applying the transformations that spring only from the accumulation of the rounding errors. With blocking, the inner level(s) of the algorithm count all transformations (big and small) in a sweep for stopping. For details, see [18, 21].

Outputs. The algorithm stops when the columns of the in-place transformed F' and G' are numerically mutually J-orthogonal, and orthogonal, respectively. Let Z' be the accumulated product of the applied transformations. Then,

$$\Sigma_F' = \operatorname{diag}(|f_1^* J f_1|^{1/2}, \dots, |f_n^* J f_n|^{1/2}), \quad \Sigma_G' = \operatorname{diag}((g_1^* g_1)^{1/2}, \dots, (g_n^* g_n)^{1/2}),$$

and $U=F'\Sigma_F'^{-1}$, $V=G'\Sigma_G'^{-1}$, $\Sigma_j=\left((\Sigma_F')_j^2+(\Sigma_G')_j^2\right)^{1/2}$, $\Sigma=\mathrm{diag}(\Sigma_1,\ldots,\Sigma_n)$, $\Sigma_F=\Sigma_F'\Sigma^{-1}$, $\Sigma_G=\Sigma_G'\Sigma^{-1}$, $Z=Z'\Sigma^{-1}$. For the inner levels of blocking, only the matrix Z (and therefore also Σ_F' , Σ_G' , and Σ , but not U and V) is required. Accumulating Z'^{-1} . Optionally, Z'^{-1} could be obtained by accumulating transformations of the second second

Accumulating Z'^{-1} . Optionally, Z'^{-1} could be obtained by accumulating transformations \widehat{Z}' from the right and their inverses \widehat{Z}'^{-1} from the left. Then, $X = Z^{-1} = \Sigma$, so Phase 4 would not be needed for the full G(H)SVD. This has not been implemented, since our main concern is a solution of the generalized eigenproblem.

6.1.8. The G(H)SVD algorithm. After obtaining the matrix \hat{Z}' , the pointwise implicit Hari–Zimmermann G(H)SVD algorithm can be written similarly to [21, Algorithm 3.1]. In Algorithm 6.1 we take into account the signature matrix J, since $H = F^*JF$, but with J = I it reduces to a generalized SVD method.

Z'=I; optional computing of D and prescaling $H_0=DHD$, $S_0=DSD$, Z'=D;

```
Algorithm 6.1 The pointwise implicit Hari–Zimmermann G(H)SVD algorithm.
```

```
it=1; {The sweep counter. Maximal number of sweeps, C_{\max}, is usually \approx 30.} repeat for all pairs (p,q), 1 \leq p < q \leq n do compute \widehat{H} and \widehat{S} from (6.1); compute the elements of \widehat{Z}' from (6.5), (6.7), (6.10)-(6.11); [f_p,f_q]=[f_p,f_q]\cdot\widehat{Z}'; \quad [g_p,g_q]=[g_p,g_q]\cdot\widehat{Z}'; \quad [z'_p,z'_q]=[z'_p,z'_q]\cdot\widehat{Z}'; end for it=it+1; until (no transformations in this sweep) or (it>C_{\max}) Output: Z and (optionally) U,V,\Sigma_F,\Sigma_G, and \Sigma;
```

If the prescaling as described in subsection 6.1.1 is performed only once, then in Algorithm 6.1 only $g_p^*g_q$ is computed, since the diagonal elements of \widehat{S} are assumed to be unity. Otherwise, by passing once through the columns g_p and g_q all three dot products can be formed. Similar holds for f_p , f_q , and the three J-dot products of \widehat{H} .

6.2. Vectorization. Many computational building blocks of Algorithm 6.1 provide both the challenges and the opportunities for the SIMD vectorization. The most

obvious such primitives are the J-dot products (and norms), which are computed combining the approaches presented in subsections 5.1 (for the unstructured patterns of signs in J) and 4.5.1 for a compact, partitioned representation of J with only a number n_+ of the leading positive (and therefore $n - n_+$ tailing negative) signs.

- **6.2.1.** ZVROTM. The column updates by \widehat{Z}' cannot be realized by a single call to a BLAS or a LAPACK routine (e.g., ZROT), since there are two angles involved in a transformation, and the columns are meant to be transformed in-place (overwritten). For that purpose, a ZVROTM routine has been implemented as a vectorized loop, that resembles a simplified version of the BLAS routine DROTM, but with the complex sines.
- **6.2.2. Transformations.** The greatest challenge lies in computing the transformations in the SIMD-parallel way, where each vector lane i computes \widehat{Z}'_i for its own pivot pair with indices (p_i,q_i) , as in [19] for the rotations in the Jacobi SVD method. Assume that \mathbf{v} pivot pairs $(\widehat{H}_i,\widehat{S}_i)$ have been obtained, with their corresponding p_i and q_i indices all different. If there are fewer than \mathbf{v} such pairs, let $\widehat{H}_i = \widehat{S}_i = I$ for the missing indices i. A vector (e.g., $\widehat{S}_{12}^{[i]} = g_{p_i}^* g_{q_i}$) has to be kept in an array of length \mathbf{v} , properly aligned in memory, in which the ith position holds data for the ith lane. To help the compiler, the complex arithmetic operations have been written in terms of the real and the imaginary parts of the complex numbers in the vectorizable regions.

For the start, check for which lanes their pivot pair has to be transformed by evaluating the criterion (6.12) in each lane. To aid the compiler, (6.12) can be rewritten as a branch-free arithmetic expression that has a non-zero value if and only if the criterion is fulfilled. If those values constitute a zero vector, no transformation is required for any i, and a fresh set of pivot pairs (if any remain) should be considered.

Then, compute \widehat{Z}'_i unconditionally (i.e., for all i). The idea is that the unneeded computation comes at no cost, while its results can be discarded afterwards. However, the exceptional cases from subsection 6.1.6 should be carefully dealt with, since a naïve branching might spoil the vectorization opportunities. The logical conditions are therefore arithmetized, while halting on the arithmetic exceptions is suppressed.

We assume that the second argument of the intrinsics MIN and MAX is returned when their first argument is a NaN. Such a behavior is not mandated by the Fortran standard (contrary to fmin and fmax in C), but is checked for at runtime by our code.

Let B(i) = (1-MAX(V(i)/V(i),0))*(1-MAX(H(i)/H(i),0)) and note that B(i) is 1 if, in (6.11), $v_i = h_i = 0$ (so $\tan \gamma_i = \text{NaN}$), and also if $h_i = \pm \infty$ (due to an overflow) with $v_i = 0$; otherwise, it is 0. The computation resumes regardless of the value of B(i). However, \widehat{Z}'_i in those lanes where B(i) = 1 is useless, so a new, correct \widehat{Z}'_i is taken according to the rules of subsection 6.1.6 at the end, sequentially for all such i. This is the only situation when \widehat{Z}'_i is not computed in one go for all i, but it should occur rarely in practice. A branch-free exception handling is also required when $\widehat{S}^{[i]}_{12} = 0$, since in calculating the polar form $|\widehat{S}^{[i]}_{12}|e^{i\phi_i}$ of $\widehat{S}^{[i]}_{12}$ the divisions of the real and the imaginary parts by the absolute value result in NaNs. If $\cos \phi_i$ is denoted by C(i), then selecting a default of $\cos \phi_i = 1$ can be done by presetting the variable to 1, and taking C(i) = MIN(x(i), C(i)), where x is obtained by the vector division, or the multiplication by the reciprocal. For $\sin \phi_i$, the default value can be also set to 1, and the correct value of 0 is obtained by multiplying the intermediate result by a variable, set in the process of checking (6.12), that is 0 if $\widehat{S}^{[i]}_{12} = 0$, and 1 otherwise.

6.3. Parallelization. There are N := n(n-1)/2 pairs of indices (p_i, q_i) such that they belong to a strictly upper triangle (i.e., $1 \le p_i < q_i \le n$) of the square

matrices of order n. Let $P := \{(p_i, q_i) : 1 \le i \le N\}$ be a set of those index pairs. At most $\lfloor n/2 \rfloor$ such pairs can be chosen from P so that all their indices are distinct.

Let S_j , for some $j \geq 1$, be a set of at most $\lfloor n/2 \rfloor$ index pairs, with all indices distinct. Then, the pivot pairs formed from the columns of F and G and indexed by the elements of S_j can be transformed concurrently. We call S_j the jth parallel Jacobi step, and a sequence of steps $S := (S_1, S_2, \ldots, S_{\bar{n}})$ a parallel (quasi-)cyclic Jacobi strategy if $\bigcup_j S_j = P$, and under assumption that the sequence is repeated forever in principle (or, in practice, until the convergence criteria are met). A strategy is called cyclic if its steps are mutually disjoint; else, it is called quasi-cyclic. In a cycle (also called a sweep) all pivot pairs are accessed (at least once, but maybe more under a quasi-cyclic strategy), in $\bar{n} \geq n-1$ steps. We aim for the steps as large as possible.

- **6.3.1.** Parallel strategies. Two classes of the parallel Jacobi strategies were under test: the modified modulus strategy (MM), described in, e.g., [21], and a generalization of the Mantharam–Eberlein strategy [17] (ME), described in [18]. The former is a quasi-cyclic strategy with $\bar{n}=n$, but easily generated on-the-fly as the computation progresses. The latter is cyclic, attains $\bar{n}=n-1$, has provided a faster execution than MM to the one-sided Jacobi SVD, with more accurate results, but is not readily available for all even n, and its convergence has not yet been proven.
- **6.3.2.** Bordering. For an odd n, no more than (n-1)/2 pairs fit into a step, so $\bar{n} \geq n$. We therefore consider the strategies for even n only, with the steps of size n/2, and when necessary border the matrices by appending a zero column and a zero row, except for the new element at position (n+1, n+1), which is set to unity.
- **6.3.3.** Vector-Parallel algorithm. Let a step S_j , with $k_1 := n/2$ index pairs, be given. Then, partition S_j into $k_v := \lceil k_1/v \rceil$ disjoint subsets V_k , where each subset has at most v pairs, i.e., $S_j = (V_1, V_2, \dots, V_{k_v})$. For each subset, the requirements for the vectorized computation of \widehat{Z}' as described in subsection 6.2 are satisfied.

Now, let $t \ge 1$ be a number of available (OpenMP) threads. Then, S_j is traversed with a parallel-do loop over the subsets, where a thread takes a chunk of subsets to be processed independently, while the subsets within a chunk are handled in sequence.

A sweep of such Vector-Parallel (VP) variant is shown in Algorithm 6.2. Note that each thread has to have a private set of vector variables, which is most easily done by reserving a $v \times t$ rank-2 array for each variable and making the lth thread access the lth column of such an array, where l is a thread's unique number, $1 \le l \le t$.

- **6.4. Blocking.** To better exploit the memory hierarchy by keeping data in the cache(s) longer, a multilevel blocking principle can be applied, with the Level 1 being the pointwise algorithm in its VP variant. In the next, second level of the algorithm the block columns of width $w \ge 1$ take place of the single columns, and the 2×2 pivot pairs are replaced by $(2w) \times (2w)$ block pivots. The same principle can be applied recursively further (e.g., see [18]), but we consider only the Level 2 algorithms here.
- **6.4.1. Block-column partitioning.** There are two ways a matrix can be partitioned into block columns. The first one is to prescribe w, and then split the matrix into at least $\lceil n/w \rceil$ block columns of width at most w, bearing in mind that the number of block columns has to be even, as explained in subsection 6.3.2, and reducing the maximal width accordingly. For $w \ge 2$, all block columns can be made to contain either w or w 1 (but no less) columns by redistributing their individual widths.

The second way, and the one we have chosen to implement, is to query at runtime a number ${\sf t}$ of threads to be used. A thread l is to be assigned one pair of block

Algorithm 6.2 A sweep of the Vector-Parallel implicit HZ algorithm.

```
for all steps S_j \in S, 1 \le j \le \bar{n} do {a sequential loop over the steps of S}
   for all subsets V_k \in S_i, 1 \le k \le k_v do {an OpenMP parallel do with t threads}
       for all (p_i, q_i) \in V_k, 1 \le i \le |V_k| \le v do {a sequential loop over V_k}

get \widehat{H}_i = \begin{bmatrix} f_{p_i}^* J f_{p_i} & f_{p_i}^* J f_{q_i} \\ f_{q_i}^* J f_{q_i} \end{bmatrix}; \widehat{S}_i = \begin{bmatrix} g_{p_i}^* g_{p_i} & g_{p_i}^* g_{q_i} \\ g_{q_i}^* g_{q_i} \end{bmatrix}; {vectorized a^*(J)b}
        for all (p_i, q_i) \in V_k, 1 \le i \le |V_k| \le v do {a SIMD parallel do over V_k}
           check the transformation criterion (6.12);
       end for
       if no pivot pairs have to be transformed, cycle; {a reduction}
       for all (p_i, q_i) \in V_k, 1 \le i \le |V_k| \le v do {a SIMD parallel do over V_k}
           compute the elements of \widehat{Z}'_i from (6.5), (6.7), (6.10)–(6.11);
        end for
        for all (p_i, q_i) \in V_k, 1 \le i \le |V_k| \le v do {a sequential loop over V_k}
           check if \widehat{Z}'_i has to be corrected and cycle if \widehat{Z}'_i = I;
           \begin{split} [f_{p_i},f_{q_i}] &= [f_{p_i},f_{q_i}] \cdot \widehat{Z}_i'; \quad [g_{p_i},g_{q_i}] = [g_{p_i},g_{q_i}] \cdot \widehat{Z}_i'; \quad [z_{p_i}',z_{q_i}'] = [z_{p_i}',z_{q_i}'] \cdot \widehat{Z}_i'; \\ \{\text{Two ZVROTMs of length } m \text{ and one of length } n.\} \end{split}
        end for
   end for
end for
```

columns with the block indices (p_{jl}, q_{jl}) in the jth block step, so the maximal width w is computed as $\lceil n/(2t) \rceil$, with exactly 2t block columns. The block column widths w_i are non-increasing across the whole partition and are equal to either w or w-1.

The blocking overhead can dominate the actual computation time for the matrices small enough, so we assume that n > 2t, and suggest a pointwise algorithm otherwise.

Each thread allocates a contiguous storage (in its own NUMA region, but visible to the other threads) for 2w columns (two block columns of the maximal width) of F, and similarly for G and Z. The same amount of memory, and of the same shape, is additionally allocated for the columns of the "shadow" matrices F, G, and G. Let, for the Gth thread in the Gth step of the Gth step of the sth block sweep, the contents of that storage be named $X_{jl}^{[s]} := [X_{pjl}^{[s]} X_{qjl}^{[s]}]$, with $X \in \{F, G, Z, F, G, Z\}$. If $W_{pjl} = W - 1$, the columns of X_{pjl} are placed from the second column of X_{jl} , while the columns of X_{qjl} are always placed from the column W + 1 of X_{jl} . That way all the columns of a block column pair are stored contiguously and can be viewed by the BLAS routines as a single matrix.

The storage for the block pivots and for the workspaces, along with a copy of J and the strategy tables, is also preallocated per thread, in MCDRAM if possible. There are two strategy tables; one for the outer, Level 2 Jacobi strategy $S^{[2]}$, and one for the inner, Level 1 strategy $S^{[1]}$, which do not have to belong to the same class. The tables are fully initialized, with all the (block) steps, before the start of the iterations.

tables are fully initialized, with all the (block) steps, before the start of the iterations. The initial data for $F_{1l}^{[1]}$ and $G_{1l}^{[1]}$ is loaded from the block columns $(\mathsf{p}_{1l},\mathsf{q}_{1l}) \in \mathsf{S}_1^{[2]},$ $Z_{1l}^{[1]}$ is initialized to a corresponding part of I_n , while $\mathsf{F}_{1l}^{[1]}$, $\mathsf{G}_{1l}^{[1]}$, and $\mathsf{Z}_{1l}^{[1]}$ are undefined.

6.4.2. Processing the block pivots. In a step j (of a block sweep s, which index we omit from the superscripts of the matrices for simplicity when it is implied by the context), the lth thread shortens its J to \widehat{J}_{jl} and its block column pairs from F_{jl} to \widehat{F}_{jl} such that $\widehat{H}_{jl} := \widehat{F}_{jl}^* \widehat{J}_{jl} \widehat{F}_{jl} = F_{jl}^* J F_{jl}$, and G_{jl} to \widehat{G}_{jl} such that $\widehat{S}_{jl} := \widehat{G}_{jl}^* \widehat{G}_{jl} = G_{jl}^* G_{jl}$, where \widehat{F}_{il} and \widehat{G}_{il} are both square, of order 2w or 2(w - 1) (the order in between the

two is necessarily odd, so the bordering from subsection 6.3.2 should then be applied).

There are two ways to shorten the block column pairs. The more efficient one is to form \widehat{H}_{jl} and \widehat{S}_{jl} explicitly. To do that, F_{jl} is copied to F_{jl} , and then the rows of F_{jl} are scaled by J. Here, it is beneficial to have J in the compact, run-length encoded form. Then, a single (parallel, or in our case, sequential) ZGEMM call computes $\widehat{H}_{jl} = F_{jl}^*(J\mathsf{F}_{jl})$ and stores it temporarily into \widehat{G}_{jl} . Then, \widehat{F}_{jl}^* and \widehat{J}_{jl} (partitioned to the positive and the negative sign blocks) are obtained by the Hermitian indefinite factorization with complete pivoting (see section 4), $\widehat{H}_{jl} = \widehat{F}_{jl}^* \widehat{J}_{jl} \widehat{F}_{jl}$. The factorization should reveal if \widehat{H}_{jl} is rank deficient, in which case the process stops (the computation could be retried with the (J)QR approach, as below, if the input data has been preserved). Finally, \widehat{F}_{il}^* is copied, with the transposition and the complex conjugation applied, to \widehat{F}_{il} .

A simpler procedure is used for \widehat{G}_{jl} . It suffices to compute $G_{jl}^*G_{jl}$ by a single (parallel, or in our case, sequential) ZHERK call, store \widehat{S}_{jl} temporarily to \widehat{Z}_{jl} , and perform either the diagonally-pivoted Cholesky factorization, or even the Hermitian indefinite factorization with complete pivoting, to obtain $\widehat{S}_{jl} = \widehat{G}_{jl}^*\widehat{G}_{jl}$. If \widehat{S}_{jl} is rank deficient or indefinite, the algorithm stops with an error message. A similar treatment has been implemented for \widehat{H}_{jl} , if J=I. Otherwise, \widehat{G}_{jl}^* is copied, with the transposition and the complex conjugation applied, to \widehat{G}_{il} . A ZLASET call finally initializes \widehat{Z}_{jl} to I.

The (J)QR approach, similar to Phase 2, should be more accurate and even necessary when the conditions of the matrices \widehat{H}_{jl} and/or \widehat{S}_{jl} are so large that factorizing them after forming them explicitly might fail. Both F_{jl} and G_{jl} should then be copied to F_{jl} and G_{jl} , respectively, and the JQR factorization on F_{jl} , followed by the column prepermutation and the column-pivoted QR factorization on G_{jl} , should be performed to obtain \widehat{F}_{jl} with \widehat{J}_{jl} , and \widehat{G}_{jl} , respectively. This has not been implemented, though.

The block pivots \widehat{F}_{jl} with \widehat{J}_{jl} and \widehat{G}_{jl} are handed over to a version of the Level 1 (VP) algorithm to be transformed. This, single-threaded VP version executes in the contexts of the already running threads. The Level 1 algorithm can either fully (implicitly) diagonalize \widehat{H}_{jl} and \widehat{S}_{jl} , or at least iterate until a reasonably high number of the inner sweeps has been attained ($C_{max} = 30$), in which case we talk about the Full Block (FB) variant; or pass over the block pivots only a prescribed number of times, e.g., once ($C_{max} = 1$), in the Block-Oriented (BO) variant (for more details on both, see [21]). The former variant corresponds to a full two-sided annihilation of the off-diagonal of \widehat{H}_{il} and \widehat{S}_{il} , while the latter implicitly reduces their off-diagonal norms.

Also, the BO variant exhibits similar execution times for every call of the Level 1 routine across all threads in a block step, while those can vary significantly, due to data, among the threads in the FB variant. On some platforms (e.g., the GPUs), that does not pose a huge problem [18], but on the CPUs it can cause delays on the synchronization primitives (OpenMP barriers) required between the block steps [21].

In any case, the transformations applied in the Level 1 are accumulated in \widehat{Z}_{jl} , while the counters of all and of "big" transformations are added atomically to the respective counters (shared among the threads) for the current block sweep. With distributed memory, such counters can be updated by MPI_Allreduce collective calls.

6.4.3. Updating and exchanging the block columns. If no transformations have been applied in the Level 1, F_{jl} , G_{jl} , and Z_{jl} are then copied to F_{jl} , G_{jl} , and Z_{jl} , respectively. Otherwise, $\mathsf{F}_{jl} = F_{jl}\widehat{Z}_{jl}$, $\mathsf{G}_{jl} = G_{jl}\widehat{Z}_{jl}$, and $\mathsf{Z}_{jl} = Z_{jl}\widehat{Z}_{jl}$ (three ZGEMM calls).

All computation have thus far been local to a thread, apart from the atomic operations. Now, by looking to $S_{i+1}^{[2]}$, it is easy to figure out which of the two block

columns should be retained by the thread, and to which thread the other block column has to be sent (in fact, the whole communication pattern has been precomputed). Also, the block columns can swap roles: the first one in the current step can become the second one, in the presently owning or in the receiving thread, and vice versa, in the next step. A thread might even send away both its block columns (each to a different recipient), and receive two new block columns (each from a different sender).

It might seem unnecessary to perform the physical exchanges of data on the shared memory, but the reason behind them is twofold. First, most modern machines have their memory paritioned according to the speed of access by, or "proximity" to, each CPU (NUMA), and it is beneficial to bring the data close to (a thread bound to) the CPU that processes it. Second, such a design makes the algorithm convertible to a distributed memory one (e.g., by using an MPI process for what a thread does now).

To perform the block column copies, a thread has to know the memory addresses of the storage of the threads it communicates to. All such addresses are kept available to all threads, and the block columns from the "shadow" storage are copied by their present owners to the regular storage of their future ones. Before a block column pair can be copied (from or to), it has to be updated first, so there is an OpenMP barrier between the three updates above and the three copying actions (by two ZLACPY calls each, one per a block column). Also, a thread cannot continue with the next block step until it has the new data ready and its shadow storage available, which is enforced for all threads simultaneously by placing another barrier after the copying actions.

It is a legitimate question if it is worthwhile (and in what circumstances) to try hiding the communication behind the computation, maybe by relying on some tasking mechanism. For example, $\mathsf{F}_{\mathsf{p}_{jl}}$ can be copied to $F_{\mathsf{q}_{j+1,l'}}$, and $\mathsf{F}_{\mathsf{q}_{jl}}$ to $F_{\mathsf{p}_{j+1,l''}}$, while G_{jl} , $G_{il'}$, and $G_{il''}$ are being updated. We leave those considerations for a future work.

Completing a sweep. At the end of a block sweep, the transformation counters are read and reset. If no "big" transformations have been applied in any of the threads in the block sweep, the process stops, and the outputs, now including Λ_F , Λ_G , and Λ , are generated, piecewise per thread, as described in subsection 6.1.7 and section 3.

- **6.5. Testing.** The aim of testing was to establish what algorithm variant to recommend for practice, as well as should Phase 2 be employed, and if so, when.
- **6.5.1.** Blocked vs. pointwise algorithms. In Table 6.1 it is shown that the Level 2 (BO) algorithm is several times faster than both the Level 1 (VP) and the Level 2 (FB) algorithms, and that blocking in general gives a significant advantage over the pointwise approach. Also, having more threads, and thus the smaller block pivot orders, benefits the FB algorithm, since the speedup with 64 versus 32 threads is there more than twofold. On the contrary, the speedup with twice more threads is less than twofold for the BO algorithm, since the formation of the block pivots in the BO variant takes a bigger portion of the overall time, compared to the Level 1 inner iterations. Nevertheless, the Level 2 (BO) algorithm is our choice for Phase 3.
- **6.5.2. Phase 2 benefits.** Phase 3 might have been run directly on the \widetilde{F} , \widetilde{J} , and \widetilde{G} . In Table 6.2 it is shown that the preprocessing of the tall-and-skinny inputs by Phase 2 into the square ones for Phase 3 is preferred, both time-wise and sweep-wise, but that advantage diminishes as the inputs approach a square-like form $(m \gtrsim n)$.
- **6.5.3.** Generalized eigenvalues. Please see the Figures S.2 and S.3 in the supplementary material, depicting the generalized eigenvalues $\Lambda(A1)-\Lambda(A4)$ and $\Lambda(B1)-\Lambda(B4)$, computed with 64 threads by the Level 2 (BO) Phase 3 GHSVD after shortening in Phase 2. The eigenvalues obtained with 32 threads differ a few ulps at most.

Table 6.1
The wall execution time (wtime), with 32, 64, and 48 threads, of the three algorithm candidates.

	Level 1 (V	P) wtime [s]	Level 2 (F	B) wtime [s]	Level 2	(BO) wti	me [s]
ID	32 thr.	64 thr.	32 thr.	64 thr.	32 thr.	64 thr.	48 thr.
A1	828.58	836.50	372.00	222.81	125.38	128.51	39.39
A2	4528.13	4533.41	2215.07	914.32	634.72	481.92	169.97
A3	19094.72	19357.67	9253.03	3539.32	2392.01	1474.61	569.85
A4	58910.73	54938.04	30925.99	11449.48	8003.41	4528.69	2128.77
B1	254.69	266.73	135.87	101.21	46.36	50.43	16.91
B2	1379.01	1461.47	619.19	328.24	188.76	173.86	53.28
В3	6377.36	6249.12	2870.85	1142.89	789.13	529.44	190.45
B4	22002.45	22421.95	9620.59	4015.04	2361.69	1506.33	579.38

Table 6.2

The wall execution time (wtime) with 32, 64, and 48 threads, the speedups, and the numbers of block sweeps of Phase 3 on the tall-and-skinny and the square inputs (obtained by Phase 2).

ID	BO tall-an 32 thr.	d-skinny (\$) 64 thr.	Phase 2 32 thr.	& BO squ 64 thr.	are (\star) 48 thr.	speedup $(\diamond)/$ 32 thr.	(*) & sweeps 64 thr.
A1	590.92	1035.61	298.71	315.04	107.73	1.978; 15 14	3.287; 15 14
$\overline{A2}$	1813.66	1600.17	1098.79	973.41	358.56	1.651; 17 16	1.644; 16 16
A3	5393.96	3457.28	3457.09	2594.89	961.61	1.560; 19 17	1.332; 17 16
A4	13293.75	7734.08	10185.98	6808.09	2899.21	1.305; 19 18	1.136; 18 18
B1	609.11	992.56	239.55	240.14	85.10	2.543;12 13	4.133; 12 12
B2	1665.46	1716.25	725.15	680.55	244.13	2.297; 14 13	2.522; 13 13
В3	3595.73	3126.35	2091.06	1733.91	620.94	1.720; 15 15	1.803; 14 14
B4	8031.08	6147.36	5194.34	3981.79	1440.84	1.546; 15 15	1.544; 15 15

6.5.4. Comparison with $\mathsf{ZHEGV}(\mathsf{D})$. In Table 6.3 the wall times for the explicit formation of H (by the \widetilde{J} -scaling of \widetilde{F} and the ZGEMM matrix multiplication) and S (by the ZHERK matrix multiplication), as well as for the LAPACK's generalized Hermitian eigensolvers ZHEGV and ZHEGVD on (H,S), left in MCDRAM when possible, are shown, alongside the speedups of this approach (i.e., of forming of H and S and then calling either ZHEGV or ZHEGVD , with the eigenvectors also computed) vs. the Level 2 (BO) Phase 3 GHSVD on the inputs shortened by Phase 2, with 64 and 48 threads. Please see Table S.2 in the supplementary material for the results with 32 and 24 threads.

The results suggest that both the CPU's clock and the size of its low-level caches play a crucial role in performance of our approach for the problems large enough.

The main advantage of our approach is its ability to compute the generalized eigenproblem accurately when the LAPACK-based one can fail, not least due to the squaring of the condition number $\kappa_2(\widetilde{G})$ in the forming of $S = \widetilde{G}^*\widetilde{G}$, i.e.,

$$\kappa_2(S) = \lambda_{\max}(S)/\lambda_{\min}(S) = \sigma_{\max}^2(\widetilde{G})/\sigma_{\min}^2(\widetilde{G}) = \kappa_2^2(\widetilde{G}),$$

that consequently leads to a failure in the Cholesky factorization, or to an unacceptable inaccuracy in the generalized eigenvalues. We believe that the demonstrated slowdown is a reasonable tradeoff for a reliable backup alternative in those difficult cases.

Table 6.3

The wall times for the explicit formation of H and S, combined with those for a LAPACK's generalized Hermitian eigensolver (ZHEGV or ZHEGVD), and the speedup vs. (\star) , with 64 and 48 threads.

	# of	(may of 2 r	runs for H, S	C) wall tim	a [e] for	total wall t	time [s] with	speedup	770 (+)
ID	$^{_{+\!\!\!/}}$ or $^{_{-\!\!\!\!/}}$ thrs.	$H = \widetilde{F}^* \widetilde{J} \widetilde{F}$	$S = \widetilde{G}^* \widetilde{G}$	ZHEGV	ZHEGVD	ZHEGV (⊲)	ZHEGVD (⊳)	(*)/(⊲)	(*)/(▷)
	64	1.490	1.519	11.896	5.573	14.905	8.317	21.137	37.881
A1	48	1.316	0.536	2.875	1.800	4.727	3.185	22.793	33.825
A2	64	3.959	3.323	38.197	18.890	45.479	25.913	21.403	37.565
A2	48	2.368	1.277	12.722	8.504	16.349	12.148	21.931	29.515
A3	64	11.642	7.696	100.320	56.839	119.448	75.966	21.724	34.158
АЗ	48	5.808	3.201	53.464	37.089	62.472	46.098	15.393	20.860
A4	64	24.399	15.369	284.750	168.956	324.281	208.513	20.994	32.651
Λ4	48	12.692	7.095	160.374	119.576	180.157	139.290	16.093	20.814
В1	64	1.348	1.909	5.399	2.573	8.656	5.575	27.742	43.076
D1	48	0.866	0.507	1.055	0.603	2.427	1.954	35.058	43.550
В2	64	4.401	3.786	17.015	8.750	24.987	16.648	27.236	40.879
	48	2.270	1.281	4.284	2.791	7.799	6.342	31.302	38.493
ВЗ	64	11.213	10.047	45.139	24.654	63.692	45.914	27.223	37.764
	48	5.508	3.104	15.564	12.399	24.162	20.904	25.699	29.704
В4	64	23.092	16.446	114.787	62.469	150.798	102.006	26.405	39.035
D4	48	12.003	6.790	51.817	41.265	70.570	59.836	20.417	24.080

7. Phase 4 – optional computation of the full G(H)SVD. In Phase 4 the right generalized singular vector matrix $X = Z^{-1}$ is computed by the LU factorization with complete pivoting, $Z = P^T L U Q^T$, from the LAPACK routine ZGETC2, followed by solving a linear system $ZX = I_n$ for X, with Z factored as above, by calling the sequential routine ZGESC2 in a parallel do loop. Each of n loop iterations solves for one column of X, and the iteration space is divided among the same number of OpenMP threads, also available to ZGETC2, as used for the previous phases.

Should the G(H)SVD of the tall-and-skinny factors be required, $\widetilde{Z} = P_2^T Z$ and $\widetilde{X} = \widetilde{Z}^{-1}$ could also be computed in Phase 4, together with $\widetilde{U} = \widetilde{Q}_F U$ and $\widetilde{V} = \widetilde{Q}_G V$.

7.1. Relative errors in the full GHSVD. To measure the accuracy of the Phase 3 algorithm, we can look at the Frobenius norm of the error in the obtained GHSVD decomposition, relative to the Frobenius norm of the original matrix, i.e.,

(7.1)
$$||F - U\Sigma_F X||_F / ||F||_F$$
 and $||G - V\Sigma_G X||_F / ||G||_F$.

For a detailed description of computing the relative errors, see the supplementary material, subsection S.3.1, and for the full accuracy results, see Table S.3 therein.

The relative errors in the decomposition of F range from $1.25 \cdot 10^{-13}$ to $3.52 \cdot 10^{-13}$, and from $1.12 \cdot 10^{-13}$ to $7.22 \cdot 10^{-13}$, for the datasets A and B, respectively. In the decomposition of G, the relative errors are slightly lower, ranging from $9.47 \cdot 10^{-14}$ to $8.23 \cdot 10^{-13}$, and from $7.98 \cdot 10^{-14}$ to $4.73 \cdot 10^{-13}$ for the datasets A and B, respectively. The relative errors are similar regardless the number of threads used in Phase 3.

From the range of errors it can be concluded both that our datasets are not highly ill-conditioned, and that in those cases the Phases 3 and 4 do not behave erratically. A rigorous stability analysis of the GHSVD method remains open for the future work.

7.2. Comparison with ZGGSVD3. The Phase 3 algorithm can also be used for the ordinary GSVD by setting J = I. It is therefore reasonable to compare it to the ZGGSVD3 LAPACK routine (from the parallel Intel MKL) serving the same purpose.

- **7.2.1. Dataset.** A dataset C, comprising five Hermitian matrix pairs, has been generated by a call to the ZLATMS LAPACK testing routine for each matrix of the full bandwidth, with its pseudorandom eigenvalues uniformly distributed in (0,1). The matrix orders are $k \cdot 1000$, $1 \le k \le 5$, and the full GSVD is required.
- **7.2.2.** Timing results. In Table 7.1 the speedup of Phase 3 followed by Phase 4, both with 64 threads, vs. ZGGSVD3 is shown. The speedup with 32 threads reaches a lower peak for n = 5000, 140 times, as shown in the supplementary material, Table S.4. From the tests it is evident that Phase 3 (with Phase 4 if needed) forms a

Table 7.1

The wall execution time (wtime) and the speedup of the Phases 3 (with J = I) and 4 versus the ZGGSVD3 LAPACK routine on the set C, with 64 threads and n denoting the order of the matrices.

n	ZGGSVD3 wtime [s] (•)	Phase 3 wtime [s] & sweeps	Phase 4 wtime [s]	Phases 3 & 4 wtime [s] (o)	speedup $(\bullet)/(\circ)$
1000	399.47	8.32;13	2.84	11.16	35.78
2000	5935.03	35.52; 14	20.28	55.80	106.37
3000	21880.11	100.71; 16	69.48	170.19	128.57
4000	54233.01	191.32; 16	180.70	372.02	145.79
5000	107424.26	332.29; 17	327.97	660.26	162.70

very competitive, highly parallel algorithm for the ordinary GSVD as well, while being adaptable to a wide variety of the modern high performance computing hardware.

8. Future work. A (multi-)GPU version of the implicit Hari–Zimmermann algorithm for the ordinary GSVD [19, 20] shows the promising results and complements the CPU implementation described herein. Therefore, a GPU implementation of the GHSVD and Phase 1 might be a practical companion to the present research.

Acknowledgments. The authors would like thank the anonymous referees for their suggestions on improving the contents and the presentation of this manuscript.

REFERENCES

- [1] E. Anderson, Z. Bai, C. Bischof, S. Blackford, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney, and D. Sorensen, LAPACK Users' Guide, SIAM, Philadelphia, PA, 3rd ed., 1999.
- A. W. Bojanczyk, An implicit Jacobi-like method for computing generalized hyperbolic SVD, Linear Algebra Appl., 358 (2003), pp. 293-307, https://doi.org/10.1016/S0024-3795(02) 00394-4.
- [3] J. R. Bunch and L. C. Kaufman, Some stable methods for calculating inertia and solving symmetric linear systems, Math. Comp., 31 (1977), pp. 163-179, https://doi.org/10.1090/ S0025-5718-1977-0428694-0.
- [4] J. R. Bunch and B. N. Parlett, Direct methods for solving symmetric indefinite systems of linear equations, SIAM J. Numer. Anal., 8 (1971), pp. 639-655, https://doi.org/10.1137/ 0708060.
- [5] A. Canning, W. Mannstadt, and A. J. Freeman, Parallelization of the FLAPW method, Comput. Phys. Commun., 130 (2000), pp. 233-243, https://doi.org/10.1016/ S0010-4655(00)00120-X.
- [6] D. DAVIDOVIĆ, D. FABREGAT-TRAVER, M. HÖHNERBACH, AND E. DI NAPOLI, Accelerating the computation of FLAPW methods on heterogeneous architectures, Concurrency Computat. Pract. Exper., 30 (2018), https://doi.org/10.1002/cpe.4905. article no. e4905.
- [7] I. S. DHILLON AND B. N. PARLETT, Multiple representations to compute orthogonal eigenvectors of symmetric tridiagonal matrices, Linear Algebra Appl., 387 (2004), pp. 1–28, https://doi.org/10.1016/j.laa.2003.12.028.

- [8] E. DI NAPOLI, E. PEISE, M. HRYWNIAK, AND P. BIENTINESI, High-performance generation of the Hamiltonian and Overlap matrices in FLAPW methods, Comp. Phys. Comm., 211 (2017), pp. 61–72, https://doi.org/10.1016/j.cpc.2016.10.003.
- [9] V. Hari, On Cyclic Jacobi Methods for the Positive Definite Generalized Eigenvalue Problem, PhD thesis, FernUniversität-Gesamthochschule, Hagen, 1984.
- [10] V. Hari, Globally convergent Jacobi methods for positive definite matrix pairs, Numer. Algorithms, 79 (2019), pp. 221–249, https://doi.org/10.1007/s11075-017-0435-5.
- [11] V. Hari, On the global convergence of the complex HZ method, SIAM J. Matrix Anal. Appl., 40 (2019), pp. 1291–1310, https://doi.org/10.1137/19M1265594.
- [12] P. Hohenberg, Inhomogeneous electron gas, Phys. Rev., 136 (1964), pp. B864–B871, https://doi.org/10.1103/PhysRev.136.B864.
- [13] H. J. F. Jansen and A. J. Freeman, Total-energy full-potential linearized augmented-planewave method for bulk solids: Electronic and structural properties of tungsten, Phys. Rev. B, 30 (1984), pp. 561–569, https://doi.org/10.1103/PhysRevB.30.561.
- [14] JÜLICH SUPERCOMPUTING CENTRE, JUWELS: Modular Tier-0/1 Supercomputer at the Jülich Supercomputing Centre, Journal of large-scale research facilities, 5, A135 (2019), https://doi.org/10.17815/jlsrf-5-171.
- [15] W. Kohn and L. J. Sham, Self-consistent equations including exchange and correlation effects, Phys. Rev., 140 (1965), pp. A1133–A1138, https://doi.org/10.1103/PhysRev.140. A1133.
- [16] P. Kurz, Non-Collinear Magnetism at Surfaces and in Ultrathin Films, PhD thesis, RWTH Aachen, February 2000, http://juser.fz-juelich.de/record/30593.
- [17] M. Mantharam and P. J. Eberlein, Block recursive algorithm to generate Jacobi-sets, Parallel Comput., 19 (1993), pp. 481–496, https://doi.org/10.1016/0167-8191(93)90001-2.
- [18] V. Novaković, A hierarchically blocked Jacobi SVD algorithm for single and multiple graphics processing units, SIAM J. Sci. Comput., 37 (2015), pp. C1–C30, https://doi.org/10.1137/ 140952429.
- [19] V. Novaković, Parallel Jacobi-type algorithms for the singular and the generalized singular value decomposition, PhD thesis, University of Zagreb, December 2017, https://urn.nsk. hr/urn:nbn:hr:217:515320.
- [20] V. NOVAKOVIĆ AND S. SINGER, Implicit Hari-Zimmermann algorithm for the generalized SVD on the GPUs, arXiv:1909.00101 [math.NA], (2019), https://arxiv.org/abs/1909.00101.
- [21] V. Novaković, S. Singer, and S. Singer, Blocking and parallelization of the Hari– Zimmermann variant of the Falk-Langemeyer algorithm for the generalized SVD, Parallel Comput., 49 (2015), pp. 136–152, https://doi.org/10.1016/j.parco.2015.06.004.
- [22] B. N. Parlett, The Symmetric Eigenvalue Problem, no. 20 in Classics in Applied Mathematics, SIAM, Philadelphia, PA, 1998.
- [23] C. Rostgaard, The projector augmented-wave method, arXiv:0910.1912 [cond-mat], (2009), https://arxiv.org/abs/0910.1921.
- [24] S. Singer, Indefinite QR factorization, BIT, 46 (2006), pp. 141–161, https://doi.org/10.1016/j.parco.2015.06.004.
- [25] S. Singer and S. Singer, Orthosymmetric block reflectors, Linear Algebra Appl., 429 (2008),
 pp. 1354–1385, https://doi.org/10.1016/j.laa.2008.04.008.
- [26] D. J. Singh and L. Nordström, eds., Planewaves, Pseudopotentials, and the LAPW Method, Springer, 2006.
- [27] I. Slapničar, Componentwise analysis of direct factorization of real symmetric and Hermitian matrices, Linear Algebra Appl., 272 (1998), pp. 227–275, https://doi.org/10.1016/S0024-3795(97)00334-0.
- [28] E. Wimmer, H. Krakauer, M. Weinert, and A. J. Freeman, Full-potential self-consistent linearized-augmented-plane-wave method for calculating the electronic structure of molecules and surfaces: O₂ molecule, Phys. Rev. B, 24 (1981), pp. 864–875, https://doi.org/10.1103/PhysRevB.24.864.
- [29] J. WINKELMANN, P. SPRINGER, AND E. DI NAPOLI, ChASE: Chebyshev accelerated subspace iteration eigensolver for sequences of Hermitian eigenvalue problems, ACM Trans. Math. Software, 45 (2019), pp. 21:1–21:34, https://doi.org/10.1145/3313828.
- [30] M. Zhou and A.-J. van der Veen, Stable subspace tracking algorithm based on a signed URV decomposition, IEEE Trans. Signal Process., 60 (2012), pp. 3036–3051, https://doi. org/10.1109/TSP.2012.2190732.
- [31] K. ZIMMERMANN, Zur Konvergenz eines Jacobiverfahren für gewönliche und verallgemeinerte Eigenwertprobleme, dissertation no. 4305, ETH, Zürich, 1969.