

# Deterministic complexity analysis of Hermitian eigenproblems

Aleksandros Sobczyk  
IBM Research and ETH Zurich  
Zurich, Switzerland

## Abstract

In this work we revisit the arithmetic and bit complexity of Hermitian eigenproblems. Recently, [BGVKS, FOCS 2020] proved that a (non-Hermitian) matrix  $A$  can be diagonalized with a randomized algorithm in  $O(n^\omega \log^2(\frac{n}{\epsilon}))$  arithmetic operations, where  $\omega \lesssim 2.371$  is the square matrix multiplication exponent, and [Shah, SODA 2025] significantly improved the bit complexity for the Hermitian case. Our main goal is to obtain similar deterministic complexity bounds for various Hermitian eigenproblems. In the Real RAM model, we show that a Hermitian matrix can be diagonalized deterministically in  $O(n^\omega \log(n) + n^2 \text{polylog}(\frac{n}{\epsilon}))$  arithmetic operations, improving the classic deterministic  $\tilde{O}(n^3)$  algorithms, and derandomizing the aforementioned state-of-the-art. The main technical step is a complete, detailed analysis of a well-known divide-and-conquer tridiagonal eigensolver of Gu and Eisenstat [GE95], when accelerated with the Fast Multipole Method, asserting that it can accurately diagonalize a symmetric tridiagonal matrix in nearly- $O(n^2)$  operations. In finite precision, we show that an algorithm by Schönhage [Sch72] to reduce a Hermitian matrix to tridiagonal form is stable in the floating point model, using  $O(\log(\frac{n}{\epsilon}))$  bits of precision. This leads to a deterministic algorithm to compute all the eigenvalues of a Hermitian matrix in  $O(n^\omega \mathcal{F}(\log(\frac{n}{\epsilon})) + n^2 \text{polylog}(\frac{n}{\epsilon}))$  bit operations, where  $\mathcal{F}(b) \in \tilde{O}(b)$  is the bit complexity of a single floating point operation on  $b$  bits. This improves the best known  $\tilde{O}(n^3)$  deterministic and  $O(n^\omega \log^2(\frac{n}{\epsilon}) \mathcal{F}(\log(\frac{n}{\epsilon})))$  randomized complexities. We conclude with some other useful subroutines such as computing spectral gaps, condition numbers, and spectral projectors, and with some open problems.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Models of computation and complexity . . . . .	3
1.2	Notation . . . . .	4
1.3	Methods and Contributions . . . . .	4
1.3.1	Real RAM . . . . .	4
1.3.2	Finite precision . . . . .	7
1.4	Outline . . . . .	8
<b>2</b>	<b>Diagonalization of symmetric tridiagonal and arrowhead matrices in Real RAM</b>	<b>9</b>
2.1	Symmetric arrowhead diagonalization . . . . .	9
2.2	Tridiagonal diagonalization . . . . .	10
2.3	Hermitian diagonalization . . . . .	11
<b>3</b>	<b>Stability of tridiagonal reduction</b>	<b>11</b>
3.1	Matrix nomenclature . . . . .	11
3.2	Rotations . . . . .	12
3.3	Recursive bandwidth halving . . . . .	12
3.4	Eigenvalues of Hermitian matrices . . . . .	13
<b>4</b>	<b>Further applications of stable tridiagonal reduction and Hermitian eigenvalue solver</b>	<b>13</b>
4.1	Singular values and condition number . . . . .	13
4.2	Definite pencil eigenvalues . . . . .	14
4.3	Spectral gaps . . . . .	14
4.4	Spectral projectors and invariant subspaces . . . . .	15
4.5	Stable inversion and factorizations . . . . .	15
<b>5</b>	<b>Conclusion</b>	<b>16</b>
<b>A</b>	<b>Preliminaries for symmetric arrowhead diagonalization</b>	<b>22</b>
A.1	Deflation . . . . .	22
A.2	Reconstruction from approximate eigenvalues . . . . .	23
<b>B</b>	<b>Fast Arrowhead diagonalization</b>	<b>23</b>
B.1	Fast Multipole Method . . . . .	24
B.2	Computing the eigenvalues with bisection . . . . .	25
B.3	Approximating the elements of the shaft . . . . .	25
B.4	Approximating inner products with the eigenvectors . . . . .	25
B.5	Proof of Theorem 2.1 . . . . .	26
<b>C</b>	<b>Tridiagonal diagonalization</b>	<b>26</b>
C.1	Omitted proofs . . . . .	26
C.2	Approximating only the eigenvalues of a tridiagonal matrix . . . . .	27
C.3	Application: Hermitian diagonalization . . . . .	27
C.4	Singular Value Decomposition . . . . .	28
<b>D</b>	<b>Floating point arithmetic</b>	<b>28</b>
<b>E</b>	<b>Reduction to tridiagonal form - omitted proofs and definitions</b>	<b>29</b>
E.1	Imported subroutines . . . . .	29
E.2	Rotations . . . . .	29
E.3	Bandwidth halving . . . . .	31

# 1 Introduction

Eigenproblems arise naturally in many applications. Given a matrix  $A$ , the goal is to compute (a subset of) the eigenvalues  $\lambda$  and/or the eigenvectors  $v$ , which satisfy

$$Av = \lambda v.$$

The properties of the given matrix, as well as the quantities that need to be computed can vary depending on the application, giving rise to different variants of the eigenproblem. These include (i) *eigenvalue problems*, such as the approximation of eigenvalues, singular values, spectral gaps, and condition numbers, (ii) *eigenspace problems*, which refer to the approximation of eigenvectors, spectral projectors, and invariant subspaces, and (iii) *diagonalization problems*, which involve the (approximate) computation of all the eigenvalues and eigenvectors of a matrix (or pencil), i.e., a full spectral factorization and/or the Singular Value Decomposition (SVD).

In this work we focus on algorithms for *Hermitian eigenproblems*, i.e., the special case where the input matrix is Hermitian. Our motivation to dedicate the analysis to this special class is twofold. First, they arise in many fundamental applications in Machine Learning, Spectral Graph Theory, and Scientific Computing. For example, the SVD, which is ubiquitous for many applications such as low rank approximations [73, 41, 37, 22], directly reduces to a Hermitian eigenproblem. Second, the spectral theorem states that a Hermitian matrix  $A$  can always be written in a factored form  $A = Q\Lambda Q^*$ , where  $Q$  is unitary and  $\Lambda$  is diagonal with real entries. This alleviates several difficulties of the non-Hermitian case, which can lead to efficient dedicated algorithms.

Algorithms for eigenproblems have been studied for decades, some of the earliest being attributed to Jacobi [55, 45]. We refer to standard textbooks for an overview of the rich literature [33, 46, 74, 12, 77]. Some landmark works include the power method [64], the Lanczos algorithm [59], and the paramount QR algorithm [39, 40, 58], which has been recognized as one of the “top ten algorithms of the twentieth century” [35], signifying the importance of the eigenproblem in science and engineering. Given a Hermitian tridiagonal matrix  $T$  with size  $n \times n$ , the algorithm can compute a set of approximate eigenvalues in  $O(n^2 \log(\frac{n}{\epsilon}))$  arithmetic operations, based on the seminal analyses of Wilkinson [88] and Dekker and Traub [27]. A set of approximate eigenvectors can be also returned in  $O(n^3 \log(\frac{n}{\epsilon}))$  operations. In conjunction with the classic unitary similarity transformation algorithm of Householder [54], the shifted-QR algorithm has heavily lifted the computational burden of solving eigenvalue problems for decades, both in theory and in practice. A detailed bit complexity analysis was provided recently in [8, 7, 9].

Despite the daunting literature, several details regarding the computational complexity of many algorithms remain unclear. It is well-known, for example, that the cubic arithmetic complexity to compute the eigenvalues of a dense matrix is not optimal: a classic work by Pan and Chen [72] showed that the eigenvalues can be approximated in  $O(n^\omega)$  arithmetic operations, albeit without detailing how to also compute eigenvectors, or to fully diagonalize a matrix. Here  $\omega \geq 2$  is the *matrix multiplication exponent*, i.e. the smallest number such that two  $n \times n$  matrices can be multiplied in  $O(n^{\omega+\eta})$  arithmetic operations, for any  $\eta > 0$ . The current best known upper bound is  $\omega < 2.371339$  [1], and we will write  $n^\omega$  instead of  $n^{\omega+\eta}$  hereafter for simplicity. Recently, Banks, Garza-Vargas, Kulkarni, and Srivastava [6] described a numerically stable randomized algorithm to compute a provably accurate diagonalization, in  $\tilde{O}(n^\omega)$  operations, improving the previous best-known bounds, specifically,  $O(n^3)$  (Hermitian) and  $O(n^{10})$  (non-Hermitian) [3]. [82] further improved the analysis for the Hermitian case, and several works have studied extensions and related applications [30, 31, 78, 83]. To date, we are not aware of any *deterministic* algorithm that achieves the same arithmetic (or bit) complexity with provable approximation guarantees, even for the Hermitian case. In this work, we proceed step-by-step and analyze several variants of the aforementioned eigenvalue, eigenspace, and diagonalization problems, in different *models of computation*, and report complexity upper bounds with provable approximation guarantees.

## 1.1 Models of computation and complexity

From the Abel-Ruffini theorem it is known that the eigenvalues and/or the eigenvectors of matrices with size larger than  $n = 5$  can not be computed exactly. Even in exact arithmetic, they can only be approximated. Before analyzing algorithms, we first need to clarify what are the quantities of interest, to define how accuracy is measured, and what is the underlying model of computation. The main two models that we use to analyze algorithms are the Real RAM and the Floating Point model, described below.

**Exact real arithmetic (Real RAM):** For the Real RAM model we follow the definition of [38]. The machine has a memory (RAM) and registers that store real numbers in infinite precision. Moving real numbers between the memory and the registers takes constant time. A processor can execute arithmetic operations  $\{+, -, \times, /, \sqrt{\cdot}, >\}$  on the real numbers stored in the registers exactly, without any errors, in constant time. Other functions such as  $\log(x)$ ,  $\exp(x)$ , and trigonometric functions, are not explicitly available, but they can be efficiently approximated up to some additive error, e.g. with a truncated polynomial expansion, using a polylogarithmic number of basic arithmetic operations. Bit-wise operations on real numbers are forbidden, since this can give the machine unreasonable computing power [50, 80, 38].

**Floating point:** In this model, a real number  $\alpha \in \mathbb{R}$  is rounded to a floating point number  $\text{fl}(\alpha) = s \times 2^{e-t} \times m$ , where  $s = \pm 1$ ,  $e$  is the exponent,  $t$  is the bias, and  $m$  is the significand. Floating point arithmetic operations also introduce rounding errors, i.e., for two floating point numbers  $\alpha$  and  $\beta$ , each operation  $\odot \in \{+, -, \times, /\}$  satisfies:

$$\text{fl}(\alpha \odot \beta) = (1 + \theta)(\alpha \odot \beta), \quad \text{and also} \quad \text{fl}(\sqrt{a}) = (1 + \theta)\sqrt{a}, \quad (1)$$

where  $|\theta| \leq \mathbf{u}$ , and  $\mathbf{u}$  is the machine precision. Assuming a total of  $b$  bits for each number, every floating point operation costs  $\mathcal{F}(b)$  bit operations, where typically  $\mathcal{F}(b) \in O(b^2)$ , or even  $\tilde{O}(b)$ , with more advanced algorithms [81, 42, 51]. More details can be found in Appendix D.

In Section 3.4 we will also use as a subroutine an algorithm from [13, 15], which was originally analyzed in Boolean RAM model. We describe in detail how to use it in the corresponding section.

**Arithmetic and boolean complexity** Given a model of computation, the *arithmetic complexity* of an algorithm is quantified in terms of the arithmetic operations executed. The *boolean* (or *bit*) *complexity*, accounts for the total number of boolean operations (i.e. boolean gates with maximum fan-in two).

## 1.2 Notation

Matrices and vectors are denoted with bold capital and small letters, respectively.  $\|\mathbf{A}\|$  is the spectral norm of  $\mathbf{A}$ ,  $\|\mathbf{A}\|_F$  its Frobenius norm,  $\kappa(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^\dagger\|$  is the condition number, and  $\Lambda(\mathbf{A})$  its spectrum. For the complexity analysis we denote the geometric series  $S_x(m) = \sum_{l=1}^m (2^{x-2})^l$ . As already mentioned, for simplicity, the complexity of multiplying two  $n \times n$  matrices will be denoted as  $O(n^\omega)$ , instead of  $O(n^{\omega+\eta})$ , for arbitrarily small  $\eta > 0$ . We also use the standard notation  $\omega(a, b, c)$  for the complexity of multiplying two rectangular matrices with sizes  $n^a \times n^b$  and  $n^b \times n^c$  in time  $O(n^{\omega(a,b,c)})$ , and therefore  $\omega := \omega(1, 1, 1)$ . For example, for  $a = c = 1$  and  $b = 2$ , the best known bound for  $\omega(1, 2, 1) \approx 3.250035$  [1], which is slightly better than naively performing  $n$  square multiplications in  $O(n \cdot n^\omega) \lesssim O(n^{3.371339})$ .  $\mathcal{F}(b)$  denotes the bit complexity of a single floating point operation on  $b$  bits.

## 1.3 Methods and Contributions

### 1.3.1 Real RAM

We start with the analysis in exact arithmetic, which is the simplest model to analyze numerical algorithms. The goal is to obtain end-to-end upper bounds for the arithmetic complexity of approximately diagonalizing symmetric tridiagonal matrices and, ultimately, dense Hermitian matrices, as well as to approximate the SVD. To measure the accuracy, we follow the notion of *backward-stability* (or *backward-approximation*) for Hermitian diagonalization from [67]. Formally, the following problems are considered.

**Problem 1.1.** *Backward-approximate diagonalization problems in exact arithmetic.*

- (i) **Symmetric arrowhead/tridiagonal diagonalization:** Given a symmetric arrowhead or tridiagonal matrix  $\mathbf{A}$  with size  $n \times n$ ,  $\|\mathbf{A}\| \leq 1$ , and accuracy  $\epsilon \in (0, 1)$ , compute a diagonal matrix  $\tilde{\Lambda}$  and a matrix  $\tilde{\mathbf{U}}$ , such that  $\tilde{\mathbf{U}} = \mathbf{U} + \mathbf{E}_U$ , where  $\mathbf{U}$  is orthogonal and  $\|\mathbf{E}_U\| \leq \epsilon$ , and  $\left\| \mathbf{A} - \mathbf{U} \tilde{\Lambda} \mathbf{U}^\top \right\| \leq \epsilon$ .

- (ii) **Hermitian diagonalization:** Given a Hermitian matrix  $\mathbf{A}$  with size  $n \times n$ ,  $\|\mathbf{A}\| \leq 1$ , and accuracy  $\epsilon \in (0, 1)$ , compute a diagonal matrix  $\tilde{\Lambda}$  and a matrix  $\tilde{\mathbf{U}}$ , such that  $\tilde{\mathbf{U}} = \mathbf{U} + \mathbf{E}_{\mathbf{U}}$ , where  $\mathbf{U}$  is unitary and  $\|\mathbf{E}_{\mathbf{U}}\| \leq \epsilon$ , and  $\|\mathbf{A} - \mathbf{U}\tilde{\Lambda}\mathbf{U}^*\| \leq \epsilon$ .
- (iii) **SVD:** Given a matrix  $\mathbf{A} \in \mathbb{C}^{m \times n}$ ,  $m = n^k$ ,  $k \geq 1$ ,  $\|\mathbf{A}\| \leq 1$ , and accuracy  $\epsilon \in (0, 1)$  compute a diagonal matrix  $\tilde{\Sigma}$ , an  $m \times n$  matrix  $\tilde{\mathbf{U}} = \mathbf{U} + \mathbf{E}_{\mathbf{U}}$ , and an  $n \times n$  matrix  $\tilde{\mathbf{V}} = \mathbf{V} + \mathbf{E}_{\mathbf{V}}$ , such that  $\mathbf{U}^*\mathbf{U} = \mathbf{V}^*\mathbf{V} = \mathbf{I}$ ,  $\|\mathbf{E}_{\{\mathbf{U}, \mathbf{V}\}}\| \leq \epsilon$ , and  $\|\mathbf{A} - \mathbf{U}\tilde{\Sigma}\mathbf{V}^*\| \leq \epsilon$ .

To obtain rigorous complexity guarantees for these problems, with respect to all the involved parameters, we start from Problem 1.1-(i), namely, diagonalization of symmetric tridiagonal matrices. To that end, we revisit the so-called fast tridiagonal eigensolvers, which aim to reduce the complexity from cubic to  $\tilde{O}(n^2)$  operations. Many such algorithms have been studied in the literature [36, 34, 14, 13, 43, 87, 70, 84, 10], most of which are based on the divide-and-conquer (DC) strategy of Cuppen [24]. The algorithms in all of the aforementioned works have been rigorously analyzed, however, explicit complexity bounds in terms of strictly solving Problem 1.1-(i) are not detailed. We resolve this by providing an end-to-end complexity analysis of the algorithm of Gu and Eisenstat [48]. In their original work, the authors outlined how to accelerate several parts of the algorithm with the Fast Multipole Method (FMM) [75], which could eventually lead to a final complexity of  $\tilde{O}(n^2)$ . However, the actual analysis of this approach and the FMM details were not provided. [65] further extended the analysis in floating point, but it also relies on a numerically stable FMM implementation, which is not detailed. In this work, we use the elegant FMM analysis of [47, 60, 21], which is particularly suited for the problems considered. It is detailed in the following Proposition 1.1.

**Proposition 1.1** (FMM). *There exists an algorithm, which we refer to as  $(\epsilon, n)$ -approximate FMM (or  $(\epsilon, n)$ -FMM, for short), which takes as input:*

- a kernel function  $k(x) \in \{\log(|x|), \frac{1}{x}, \frac{1}{x^2}\}$ ,
- $2n + m$  real numbers:  $\{x_1, \dots, x_m\} \cup \{c_1, \dots, c_n\} \cup \{y_1, \dots, y_n\}$ , and a constant  $C$ , such that  $m \leq n$  and for all  $i \in [m], j \in [n]$  it holds that

$$|x_i|, |c_j|, |y_j| < C \quad \text{and} \quad |x_i - y_j| \geq \Omega(\text{poly}(\frac{\epsilon}{n})).$$

It returns  $m$  values  $\tilde{f}(x_1), \dots, \tilde{f}(x_m)$  such that  $|\tilde{f}(x_i) - f(x_i)| \leq \epsilon$ , for all  $i \in [m]$ , where  $f(x) = \sum_{j=1}^n c_j k(x_i - y_j)$ , in a total of  $O\left(n \log^\xi\left(\frac{n}{\epsilon}\right)\right)$  arithmetic operations, where  $\xi \geq 1$  is a small constant that is independent of  $\epsilon, n$ .

*Proof.* **TOPROVE 0** □

By taking advantage of the FMM, the analysis from Section 2 and the supporting Appendices A and B.1 leads to the following Theorem 1.1, whose proof can be found in Section 2.2.

**Theorem 1.1.** *Given an unreduced symmetric tridiagonal matrix  $\mathbf{T}$  with size  $n \times n$ ,  $\|\mathbf{T}\| \leq 1$ , an accuracy  $\epsilon \in (0, 1/2)$ , and an  $(\epsilon, n)$ -FMM implementation as in Proposition 1.1, the recursive algorithm of [48] (Algorithm 1), returns an approximately orthogonal matrix  $\tilde{\mathbf{U}}$  and a diagonal matrix  $\tilde{\Lambda}$  such that*

$$\|\mathbf{T} - \tilde{\mathbf{U}}\tilde{\Lambda}\tilde{\mathbf{U}}^\top\| \leq \epsilon, \quad \|\tilde{\mathbf{U}}^\top\tilde{\mathbf{U}} - \mathbf{I}\| \leq \epsilon/n^2,$$

or, stated alternatively,

$$\tilde{\mathbf{U}} = \mathbf{U} + \mathbf{E}_{\mathbf{U}}, \quad \mathbf{U}^\top\mathbf{U} = \mathbf{I}, \quad \|\mathbf{E}_{\mathbf{U}}\| \leq \epsilon/n^2, \quad \|\mathbf{T} - \mathbf{U}\tilde{\Lambda}\mathbf{U}^\top\| \leq \epsilon,$$

using a total of  $O\left(n^2 \log^{\xi+1}\left(\frac{n}{\epsilon}\right)\right)$  arithmetic operations and comparisons, where  $\xi \geq 1$  is a small constant that depends on the specific FMM implementation and it is independent of  $\epsilon, n$ .

Table 1: Comparison of algorithms for the Problems 1.1 (i)-(iii) in the Real RAM model. Here  $\xi \geq 1$  is a small constant which depends on the FMM implementation (see Prop. 1.1). The algorithms marked with **(R)** are randomized and succeed with high probability (at least  $1 - 1/\text{poly}(n)$ ).

	Arithmetic Complexity	Comment
Prob. 1.1-(i)		
Arrowhead/Trid. diagonalization		
Refs. [24, 69, 48]	$O(n^3) + \tilde{O}(n^2)$	Conjectured $\tilde{O}(n^2)$
Theorem 1.1	$O\left(n^2 \log^{\xi+1}\left(\frac{n}{\epsilon}\right)\right)$	Req. FMM (Prop. 1.1)
Prob. 1.1-(ii)		
Hermitian diagonalization		
Refs. [6, 82] <b>(R)</b>	$O\left(n^\omega \log^2\left(\frac{n}{\epsilon}\right)\right)$	-
Ref. [11] <b>(R)</b>	$\tilde{O}\left(n^{\omega+1}\right)$	Conjectured $\tilde{O}(n^\omega)$
Refs. [27, 88, 53]	$O\left(n^3 \log\left(\frac{n}{\epsilon}\right)\right)$	Shifted-QR
Ref. [67]	$\tilde{O}\left(n^3\right)$	Req. separated spectrum
Ref. [72]	$O\left(n^\omega + n \text{polylog}\left(\frac{n}{\epsilon}\right)\right)$	Only eigenvalues
Corollary 2.1	$O\left(n^\omega \log(n) + n^2 \log^{\xi+1}\left(\frac{n}{\epsilon}\right)\right)$	Req. FMM (Prop. 1.1)
Prob. 1.1-(iii), SVD		
Shifted-QR on $A^*A$	$O\left(n^{\omega(1,k,1)} + n^3 \log\left(\frac{n}{\epsilon}\right)\right)$	Partial error analysis
Refs. [6, 82, 57] <b>(R)</b>	$O\left(n^{\omega(1,k,1)} + n^\omega \log^2\left(\frac{n}{\epsilon}\right)\right)$	Partial error analysis
Theorem 1.2	$O\left(n^{\omega(1,k,1)} + n^\omega \log(n) + n^2 \text{polylog}\left(\frac{n\kappa(A)}{\epsilon}\right)\right)$	Req. FMM (Prop. 1.1)

This result has a direct application to the dense Hermitian diagonalization problem. The best-known complexities of deterministic algorithms, with end-to-end analysis, are at least cubic. For example, the aforementioned shifted-QR algorithm requires  $O(n^3 \log(n/\epsilon))$  arithmetic operations, even for tridiagonal matrices. The cubic barrier originates from the accumulation of the elementary rotation matrices to form the final eigenvector matrix. The so-called *spectral divide-and-conquer* methods (see e.g. [33, 67]) have also at least cubic complexities in the deterministic case. The two main difficulties in their analysis are the basis computation of spectral projectors, for which the best-known deterministic complexity bounds are  $\Omega(n^3)$ , e.g. via Strong Rank-Revealing QR (RRQR) factorization [49], and the choice of suitable splitting points, which relies on the existence of spectral gaps.

Randomness can help to overcome certain difficulties in the analysis. [28] analyzed a numerically stable Randomized URV decomposition, which can be used to replace RRQR for basis computations in spectral DC algorithms. A highly parallelizable Hermitian diagonalization algorithm with end-to-end analysis was proposed in [11]. While the reported sequential arithmetic complexity is  $O(n^{\omega+1})$ , the authors conjectured that it can be further reduced to  $\tilde{O}(n^\omega)$ . The first end-to-end  $\tilde{O}(n^\omega)$  complexity upper bound was achieved in [6]. One of the main techniques was to use random perturbations to ensure that the pseudospectrum is well-separated, which helps to find splitting points in spectral DC. [82] further improved the analysis for Hermitian matrices.

Theorem 1.1 can be directly used to obtain a simple and deterministic solution for the Hermitian diagonalization problem. Specifically, Corollary 2.1 states that the problem can be solved in  $O(n^\omega \log(n) + n^2 \text{polylog}(n/\epsilon))$  arithmetic operations, which is both faster and fully deterministic. This is achieved by combining Theorem 1.1 with the (rather overlooked) algorithm of Schönhage [79], who proved that a Hermitian matrix can be reduced to tridiagonal form with unitary similarity transformations in  $O(n^\omega \log(n)c(\omega))$  arithmetic operations, where  $c(\omega) = O(1)$  if  $\omega$  is treated as a fixed constant larger than 2, while  $c(\omega) = O(\log(n))$  if it turns out that  $\omega = 2$ .

As a consequence, Theorem 1.2 reports similar deterministic complexity results for the SVD. The SVD is a fundamental computational kernel in many applications, such as Principal Component Analysis [56], and it is also widely used as a subroutine for many other advanced algorithms, e.g. [73, 41, 37, 18, 20, 23, 19, 22], to name a few. A straightforward algorithm to compute it is to first form the Gramian matrix  $A^*A$ , and then diagonalize it. Other

classic SVD algorithms, such as the widely adopted Golub-Kahan bidiagonalization and its variants [44], or polar decomposition-based methods [67, 66], avoid the computation of  $A^*A$  for numerical stability reasons, but they also rely on a diagonalization algorithm as a subroutine. [57] elaborated on a matrix multiplication time SVD algorithm, by using [6] to diagonalize  $A^\top A$ , albeit without fully completing the backward stability analysis. The following Theorem 1.2, which is proved in Appendix C.4, states our main result.

**Theorem 1.2 (SVD).** *Let  $A \in \mathbb{C}^{m \times n}$ ,  $m = n^k$ ,  $k \geq 1$ . Assume that  $1/n^c \leq \|A\| \leq 1$ , for some constant  $c$ . Given accuracy  $\epsilon \in (0, 1/2)$  and an  $(\epsilon, n)$ -FMM (see Prop. 1.1), we can compute three matrices  $\tilde{U} \in \mathbb{C}^{m \times n}$ ,  $\tilde{\Sigma} \in \mathbb{R}^{n \times n}$ ,  $\tilde{V} \in \mathbb{C}^{n \times n}$  such that  $\tilde{\Sigma}$  is diagonal with positive diagonal elements and*

$$A = \tilde{U}\tilde{\Sigma}\tilde{V}^*, \quad \|\tilde{U}^*\tilde{U} - I\| \leq \epsilon, \quad \|\tilde{V}^*\tilde{V} - I\| \leq \epsilon/(\kappa(A)n^2) \ll \epsilon,$$

or, stated alternatively,

$$\begin{aligned} \|A - U\tilde{\Sigma}V^*\| &\leq \epsilon, \quad \tilde{U} = U + E_U, \quad \|E_U\| \leq \epsilon, \quad U^*U = I, \\ \tilde{V} &= V + E_V, \quad \|E_V\| \leq \epsilon/n^2, \quad V^*V = I. \end{aligned}$$

The algorithm requires a total of at most  $O\left(n^{\omega(1,k,1)} + n^\omega \log(n) + n^2 \text{polylog}\left(\frac{n\kappa(A)}{\epsilon}\right)\right)$  arithmetic operations, where  $\kappa(A) = \|A\|\|A^\dagger\|$ .

To summarize this section, in Table 2 we list all the deterministic arithmetic complexities achieved in the Real RAM model, for all the aforementioned problems, and compare with some important existing algorithms.

### 1.3.2 Finite precision

Similar deterministic complexity upper bounds are obtained for several problems in finite precision. In particular, we study the following problems, for which we seek to bound the boolean complexity, i.e., the total number of bit operations.

**Problem 1.2.** *Main problems in finite precision.*

- (i) **Tridiagonal reduction:** *Given a Hermitian matrix  $A$  with floating point elements, reduce  $A$  to tridiagonal form using (approximate) unitary similarity transformations. In particular, return a tridiagonal matrix  $\tilde{T}$ , and (optionally) an approximately unitary matrix  $\tilde{Q}$ , such that*

$$\|\tilde{Q}\tilde{Q}^* - I\| \leq \epsilon, \quad \text{and} \quad \|A - \tilde{Q}\tilde{T}\tilde{Q}^*\| \leq \epsilon \|A\|,$$

- (ii) **Hermitian eigenvalues:** *Given a Hermitian matrix  $A$ ,  $\|A\| \leq 1$ , and accuracy  $\epsilon \in (0, 1)$ , compute a set of approximate eigenvalues  $\tilde{\lambda}_i$  such that  $|\tilde{\lambda}_i - \lambda_i(A)| \leq \epsilon$ .*

Regarding deterministic algorithms, with end-to-end-analysis, the standard approach is to first reduce the Hermitian matrix to tridiagonal form with Householder transformations [54], which can be done stably in  $O(n^3)$  arithmetic operations using  $O(\log(n/\epsilon))$  bits of precision; see e.g. [52]. Thereafter, there is a plethora of algorithms (e.g. the ones mentioned in the previous section) for the eigenvalues of the tridiagonal matrix, with varying complexities and stability properties. However, the total boolean complexity cannot be lower than  $\Omega(n^3 \mathcal{F}(\log(\frac{n}{\epsilon})))$  due to the Householder reduction step. Other well-known deterministic and numerically stable algorithms in the literature also require at least  $n^3$  arithmetic operations to compute all the eigenvalues [71, 32, 67], and at least  $\text{polylog}(n, 1/\epsilon)$  bits of precision in a floating point machine. The arithmetic complexity of the algorithm of [72] scales as  $O(n^\omega)$  with respect to the matrix size  $n$ , but the boolean complexity can increase up to  $O(n^{\omega+1})$  in rational arithmetic. [61] described a randomized algorithm to compute only the largest eigenvalue in nearly  $O(n^\omega)$  bit complexity. The fastest algorithm to compute all the eigenvalues of a Hermitian matrix is [82], which requires  $O(n^\omega \text{polylog}(n/\epsilon))$  boolean operations and succeeds with high probability.



Table 2: Boolean complexity for Problems 1.2, for matrix size  $n \times n$  and accuracy  $\epsilon \in (0, 1)$ .

	Boolean Complexity	Comment
Prob. 1.2-(i)		
Tridiag. Reduction		
Refs. [54, 52]	$O\left(n^3 \mathcal{F}\left(\log\left(\frac{n}{\epsilon}\right)\right)\right)$	Standard Householder reduction
Theorem 3.1	$O\left(n^\omega \log(n) \mathcal{F}\left(\log\left(\frac{n}{\epsilon}\right)\right)\right)$	[79] with stable fast QR [28]
Prob. 1.2-(ii)		
Herm. Eigenvalues		
Ref. [82]	$O\left(n^\omega \log^2\left(\frac{n}{\epsilon}\right) \mathcal{F}\left(\log\left(\frac{n}{\epsilon}\right)\right)\right)$	Randomized, $\Pr[\text{success}] \geq 1 - \frac{1}{n}$
Theorem 3.3	$O\left(n^\omega \mathcal{F}\left(\log\left(\frac{n}{\epsilon}\right)\right) + n^2 \text{polylog}\left(\frac{n}{\epsilon}\right)\right)$	Deterministic, Thm. 3.1 + [15]

Randomized eigenvalue algorithms have also been studied in the sketching/streaming setting [2, 68, 86]. The (optimal) algorithm of [86] has not been analyzed in finite-precision, but, due to its simplicity, it should be straightforward to achieve. The algorithm approximates all the eigenvalues of a Hermitian matrix  $A$  up to additive error  $\epsilon \|A\|_F$ . However, to reduce the error to a spectral-norm bound  $\epsilon \|A\|$ , the algorithm internally needs to diagonalize a matrix with size  $\Omega(n)$ , and therefore it does not provide any improvement against any other Hermitian eigenvalue solver. Nevertheless, our main results can be directly applied as the main eigenvalue subroutine of this algorithm, and to help analyze its bit complexity.

To improve the aforementioned Hermitian eigenvalue algorithms, we first prove in Theorem 3.1 it is proved that Schönhage’s algorithm is numerically stable in the floating point model of computation. Thereafter, we carefully combine it with the algorithms of [13, 14, 15] which provably and deterministically approximate all the eigenvalues of a symmetric tridiagonal matrix in  $\tilde{O}(n^2)$  boolean operations. The latter algorithm is analyzed in the Boolean RAM model, therefore, in order to use it we need to convert the floating point elements of the tridiagonal matrix that is returned by Theorem 3.1 to integers, which can be done efficiently under reasonable assumptions on the initial number of bits used to represent the floating point numbers. This is described in detail in the proof of our main Theorem 3.3. Our result derandomizes and slightly improves the final bit complexity of the algorithm of [82], which has the currently best known bit complexity for this problem. Table 2 summarizes this discussion.

As a direct consequence of Theorem 3.3, we also provide the analysis for several other useful subroutines related to eigenvalue/eigenvector computations, including:

- (i) **Singular values and condition number:** In Proposition 4.1 we describe how to obtain relative error approximations of singular values. In Corollary 4.1 we show how to compute the condition number.
- (ii) **Definite pencil eigenvalues:** In Corollary 4.2 we demonstrate how to extend Theorem 3.3 to compute the eigenvalues of Hermitian-definite pencils.
- (iii) **Spectral gaps:** In Corollary 4.4 we show how to compute the spectral gap and the midpoint between any two eigenvalues of a Hermitian-definite pencil. Our algorithm is deterministic and it requires significantly less bits of precision than the algorithm of [83], who described a randomized algorithm for this problem that is slightly faster than applying [6] as a black-box, but it only computes a single spectral gap.
- (iv) **Spectral projector:** Corollary 4.5 details how to compute spectral projectors on invariant subspaces of Hermitian-definite pencils, which are important for many applications.

## 1.4 Outline

The paper is organized as follows. In Section 2 we analyze the algorithm of [48] when implemented with the FMM and its applications (see also Appendices A, B.1, and C). In Section 3 it is proved that Schönhage’s algorithm is numerically stable in floating point, and it is used as a preprocessing step to compute the eigenvalues of Hermitian matrices. For the proof we use the technical lemmas that are proved in the supporting Appendix E. In Section 4



we mention some direct applications to compute singular values, pencil eigenvalues, spectral gaps, and spectral projectors. We finally conclude and state some open problems in Section 5.

## Acknowledgements

I am grateful to Efstratios Gallopoulos, Daniel Kressner, and David Woodruff for helpful discussions.

## 2 Diagonalization of symmetric tridiagonal and arrowhead matrices in Real RAM

Our main target is to compute the eigenvalues and the eigenvectors of tridiagonal symmetric matrices in nearly linear time. To derive the desired result, we provide an analysis the divide-and-conquer algorithm of Gu and Eisenstat [48], when implemented with the FMM from Proposition 1.1.

The algorithm of [48] first “divides” the problem by partitioning the (unreduced) tridiagonal matrix  $T$  as follows:

$$T = \begin{pmatrix} T_1 & \beta_{k+1}\mathbf{e}_k & \\ \beta_{k+1}\mathbf{e}_k^\top & \alpha_{k+1} & \beta_{k+2}\mathbf{e}_1^\top \\ & \beta_{k+2}\mathbf{e}_1 & T_2 \end{pmatrix}.$$

If one has access to the spectral decomposition of  $T_1$  and  $T_2$ , i.e.  $T_1 = Q_1 D_1 Q_1^\top$  and  $T_2 = Q_2 D_2 Q_2^\top$ , then  $T$  can be factorized as

$$\begin{pmatrix} & Q_1 \\ 1 & \\ & Q_2 \end{pmatrix} \begin{pmatrix} \alpha_{k+1} & \beta_{k+1}\mathbf{l}_1^\top & \beta_{k+2}\mathbf{f}_2^\top \\ \beta_{k+1}\mathbf{l}_1 & D_1 & \\ \beta_{k+2}\mathbf{f}_2 & & D_2 \end{pmatrix} \begin{pmatrix} Q_1^\top & 1 \\ & Q_2^\top \end{pmatrix} = QHQ^\top, \quad (2)$$

where  $\mathbf{l}_1^\top$  is the last row of  $Q_1$  and  $\mathbf{f}_2^\top$  is the first row of  $Q_2$ , and  $H$  has the so-called *arrowhead* structure. Thus, given this form, to compute the spectral decomposition of  $T$ , it suffices to diagonalize  $H$ . One can then apply recursively the algorithm to compute the spectral decompositions of  $T_1$  and  $T_2$ , and, finally, at the “conquering stage,” combine the solutions with the eigendecomposition of  $H$ .

For the individual steps to be computed efficiently, we will need to use the FMM. Specifically, we will need to use it to evaluate the functions that are listed in Equations (11)-(15), where the evaluation points will also satisfy certain criteria that are detailed in Lemma A.1. We ensure that these criteria are met by using a deflation pre-processing step (see also Appendix A.1), which allows us to take advantage of the FMM, specifically, Proposition 1.1.

### 2.1 Symmetric arrowhead diagonalization

The first step is to provide an end-to-end complexity analysis for the arrowhead diagonalization algorithm of [48], when implemented with the FMM. We start with an  $n \times n$  arrowhead matrix of the form

$$H = \begin{pmatrix} \alpha & \mathbf{z}^\top \\ \mathbf{z} & D \end{pmatrix}, \quad (3)$$

where  $D$  is a diagonal matrix,  $\mathbf{z}$  is a vector of size  $n - 1$  and  $\alpha$  is a scalar. Without loss of generality we assume that  $\|H\| \leq 1$ . The main result is stated in Theorem 2.1.

In order to prove Theorem 2.1, we expand in detail the following methodology which is outlined in [48], by proving several technical lemmas in Appendices A and B.1 that leverage the FMM.

1. Deflation: The matrix  $H$  is preprocessed to ensure that it satisfies the following:  $d_{i+1} - d_i \geq \tau$ , and  $|\mathbf{z}_i| \geq \tau$ , where  $\tau \in (0, 1)$ . This assumption implies several useful properties, described in Lemma A.1, and it allows us to efficiently utilize the FMM in the subsequent steps. The deflation procedure is described in Proposition A.1.

2. Eigenvalues: The eigenvalues of the deflated matrix can be conveniently approximated as the roots of the corresponding secular equation (Eq. (6)):  $f(\lambda) = \lambda - \alpha + \sum_{j=2}^n \frac{z_j^2}{d_j - \lambda}$ . An FMM-based root finder is detailed in Lemma B.1.
3. From Lemma A.2, the approximate eigenvalues returned by Lemma B.1 are the exact eigenvalues of another arrowhead matrix  $\widehat{\mathbf{H}} = \begin{pmatrix} \widehat{\alpha} & \widehat{\mathbf{z}}^\top \\ \widehat{\mathbf{z}} & \mathbf{D} \end{pmatrix}$ . There is an analytical expression for the elements of  $\widehat{\mathbf{H}}$  (see also [17, 48]). Lemma B.2 describes how to compute those elements with the FMM.
4. Given the exact eigenvalues and the approximate elements of the matrix  $\widehat{\mathbf{H}}$ , we can focus on its eigenvectors. In particular, in Lemma B.3 it is shown how to use the FMM to approximate the inner products between the eigenvectors of  $\widehat{\mathbf{H}}$  and some arbitrary unit vector  $\mathbf{b}$ . Computing such inner products with all the columns of the identity, we obtain the final approximate eigenvector matrix of  $\mathbf{H}$  and, ultimately, an approximate diagonalization of  $\mathbf{H}$ .

**Remark 2.1.** We note that, if we are only interested in a full diagonalization, Lemma B.2 is redundant, i.e., we can naively compute the elements exactly without the FMM with the same complexity. However, it is useful if we need to compute only a few matrix-vector products with the eigenvector matrix. Lemma B.3, details how to approximate such matrix-vector products efficiently with the FMM.

**Theorem 2.1.** Given a symmetric arrowhead matrix  $\mathbf{H} \in \mathbb{R}^{n \times n}$  as in Eq. (3), with  $\|\mathbf{H}\| \leq 1$ , an accuracy parameter  $\epsilon \in (0, 1)$ , a matrix  $\mathbf{B}$  with  $r$  columns  $\mathbf{B}_i, i \in [r]$ , where  $\|\mathbf{B}_i\| \leq 1$ , and an  $(\epsilon, n)$ -FMM implementation (see Prop. 1.1), we can compute a diagonal matrix  $\widetilde{\Lambda}$ , and a matrix  $\widetilde{\mathbf{Q}}_{\mathbf{B}}$ , such that

$$\left\| \mathbf{H} - \mathbf{Q} \widetilde{\Lambda} \mathbf{Q}^\top \right\| \leq \epsilon, \quad \left| \left( \mathbf{Q}^\top \mathbf{B} - \widetilde{\mathbf{Q}}_{\mathbf{B}} \right)_{i,j} \right| \leq \epsilon/n^2,$$

where  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  is (exactly) orthogonal, in  $O\left(nr \log^{\xi+1}\left(\frac{n}{\epsilon}\right)\right)$  arithmetic operations and comparisons, where  $\xi \geq 1$  is a small constant that depends on the specific FMM implementation and it is independent of  $\epsilon, n$ .

Alternatively, if only want to compute a set of approximate values  $\widetilde{\lambda}_1, \dots, \widetilde{\lambda}_n$ , such that  $|\lambda_i(\mathbf{H}) - \widetilde{\lambda}_i| \leq \epsilon$ , the complexity reduces to  $O\left(n \log\left(\frac{1}{\epsilon}\right) \log^\xi\left(\frac{n}{\epsilon}\right)\right)$  arithmetic operations.

*Proof.* **TOPROVE 1** □

## 2.2 Tridiagonal diagonalization

Given the analysis for arrowhead diagonalization, we can now proceed to tridiagonal matrices. The next lemma bounds the error of the reduction to arrowhead form when the spectral factorizations of the matrices  $\mathbf{T}_1$  and  $\mathbf{T}_2$  in Equation (2) are approximate rather than exact. This will be used as an inductive step for the final algorithm.

**Lemma 2.1.** Let  $\epsilon \in (0, 1/2)$  be a given accuracy parameter and  $\mathbf{T} = \begin{pmatrix} \mathbf{T}_1 & \beta_{k+1} \mathbf{e}_k & \\ \beta_{k+1} \mathbf{e}_k^\top & \alpha_{k+1} & \beta_{k+2} \mathbf{e}_1^\top \\ & \beta_{k+2} \mathbf{e}_1 & \mathbf{T}_2 \end{pmatrix}$  be a tridiagonal matrix with size  $n \geq 3$  and  $\|\mathbf{T}\| \leq 1$ , where  $\mathbf{T}_1 = \mathbf{U}_1 \mathbf{D}_1 \mathbf{U}_1^\top$  and  $\mathbf{T}_2 = \mathbf{U}_2 \mathbf{D}_2 \mathbf{U}_2^\top$  be the exact spectral factorizations of  $\mathbf{T}_1$  and  $\mathbf{T}_2$ . Let  $\widetilde{\mathbf{U}}_1, \widetilde{\mathbf{D}}_1, \widetilde{\mathbf{U}}_2, \widetilde{\mathbf{D}}_2$  be approximate spectral factorizations of  $\mathbf{T}_1, \mathbf{T}_2$ . If these factors satisfy

$$\left\| \mathbf{T}_{\{1,2\}} - \widetilde{\mathbf{U}}_{\{1,2\}} \widetilde{\mathbf{D}}_{\{1,2\}} \widetilde{\mathbf{U}}_{\{1,2\}}^\top \right\| \leq \epsilon_1, \quad \left\| \widetilde{\mathbf{U}}_{\{1,2\}} \widetilde{\mathbf{U}}_{\{1,2\}}^\top - \mathbf{I} \right\| \leq \epsilon_1/n,$$

for some  $\epsilon_1 \in (0, 1/2)$ , where  $\widetilde{\mathbf{D}}_{\{1,2\}}$  are both diagonal, then, assuming an  $(\epsilon, n)$ -FMM implementation as in Prop. 1.1, we can compute a diagonal matrix  $\widetilde{\mathbf{D}}$  and an approximately orthogonal matrix  $\widetilde{\mathbf{U}}$  such that

$$\left\| \widetilde{\mathbf{U}}^\top \widetilde{\mathbf{U}} - \mathbf{I} \right\| \leq 3(\epsilon_1 + \epsilon)/n, \quad \text{and} \quad \left\| \mathbf{T} - \widetilde{\mathbf{U}} \widetilde{\mathbf{D}} \widetilde{\mathbf{U}}^\top \right\| \leq 2\epsilon_1 + 7\epsilon,$$

in a total of  $O\left(n^2 \log^{\xi+1}\left(\frac{n}{\epsilon}\right)\right)$  arithmetic operations and comparisons, where  $\xi \geq 1$  is a small constant that depends on the specific FMM implementation and it is independent of  $\epsilon, n$ .

*Proof.* **TOPROVE 2** □

Lemma 2.1 gives rise to the following recursive algorithm. We can finally proceed with the proof of Theorem 1.1, which gives the complexity of Algorithm 1.

**Algorithm:**  $[\tilde{U}, \tilde{\Lambda}] \leftarrow \text{DIAGONALIZE}(\mathbf{T}, \epsilon)$

- 1: **if**  $n \leq 2$  :
- 2:   Compute  $\tilde{U}, \tilde{\Lambda}$  to be the exact diagonalization of  $\mathbf{T}$ .
- 3: **else:**
- 4:   Partition  $\mathbf{T} = \begin{pmatrix} \mathbf{T}_1 & \beta_{k+1}\mathbf{e}_k & \\ \beta_{k+1}\mathbf{e}_k^\top & \alpha_{k+1} & \beta_{k+2}\mathbf{e}_1^\top \\ & \beta_{k+2}\mathbf{e}_1 & \mathbf{T}_2 \end{pmatrix}$ .
- 5:    $[\tilde{U}_1, \tilde{D}_1] \leftarrow \text{DIAGONALIZE}(\mathbf{T}_1, \epsilon)$ .
- 6:    $[\tilde{U}_2, \tilde{D}_2] \leftarrow \text{DIAGONALIZE}(\mathbf{T}_2, \epsilon)$ .
- 7:   Assemble  $\tilde{U}, \tilde{\Lambda}$  from  $\mathbf{T}, \tilde{U}_1, \tilde{D}_1, \tilde{U}_2, \tilde{D}_2$  using Lemma 2.1 with parameter  $\epsilon$ .
- 8: **return**  $\tilde{U}, \tilde{\Lambda}$ .

**Algorithm 1:** Recursive algorithm based on [48] to diagonalize a symmetric tridiagonal matrix.

*Proof.* **TOPROVE 3** □

### 2.3 Hermitian diagonalization

Given an algorithm to diagonalize tridiagonal matrices, the following corollary is immediate.

**Corollary 2.1.** *Let  $\mathbf{A}$  be a Hermitian matrix of size  $n$  with  $\|\mathbf{A}\| \leq 1$ . Given accuracy  $\epsilon \in (0, 1/2)$ , and an  $(\epsilon, n)$ -FMM implementation of Prop. 1.1, we can compute a matrix  $\tilde{\mathbf{Q}}$  and a diagonal matrix  $\tilde{\mathbf{\Lambda}}$  such that*

$$\left\| \mathbf{A} - \tilde{\mathbf{Q}} \tilde{\mathbf{\Lambda}} \tilde{\mathbf{Q}}^* \right\| \leq \epsilon, \quad \left\| \tilde{\mathbf{Q}}^* \tilde{\mathbf{Q}} - \mathbf{I} \right\| \leq \epsilon/n^2.$$

The algorithm requires a total of  $O\left(n^\omega \log(n) + n^2 \log^{\xi+1}\left(\frac{n}{\epsilon}\right)\right)$  arithmetic operations and comparisons, where  $\xi \geq 1$  is a small constant that depends on the specific FMM implementation and it is independent of  $\epsilon, n$ .

*Proof.* **TOPROVE 4** □

## 3 Stability of tridiagonal reduction

In this section we analyze the numerical stability and the boolean complexity of Schönhage's algorithm the floating point model. For this we will use the following stable matrix multiplication and backward-stable QR factorization algorithms as subroutines from [28, 29]. The corresponding definitions and imported results for these subroutines are deferred to Appendix E.1.

### 3.1 Matrix nomenclature

Schönhage [79] used a block variant of Rutishauser's algorithm [76] to reduce a matrix to tridiagonal form, where elementary rotations are replaced with block factorizations; see also [16, 4, 5] for similar methodologies. We start with a  $n \times n$  block-pentadiagonal matrix  $\mathbf{A}^{(k,s,t)}$ , where  $k \in \{0, \dots, \log(n) - 2\}$  (we assume without loss of generality that  $n$  is a power of two). The matrix  $\mathbf{A}^{(k,s,t)}$  is partitioned in  $b_k \times b_k$  blocks of size  $n_k \times n_k$  each, where  $b_k = \frac{n}{n_k}$  and  $n_k = 2^k$ . The integer  $s \in 2, \dots, b_k$  denotes that all the blocks  $\mathbf{A}_{i,i-2}$  and  $\mathbf{A}_{i-2,i}$ , for all  $i = 2, \dots, s$  are equal to zero.

$s = 2$  is a special case to denote a full block pentadiagonal matrix. The integer  $t \in \{s + 2, \dots, b_k\}$  denotes that the matrix has two additional nonzero blocks in the third off-diagonals, specifically at positions  $\mathbf{A}_{t,t-3}$  and  $\mathbf{A}_{t-3,t}$ . These blocks are often called the “bulge” in the literature. When  $t = 0$ , there is no bulge. As a consequence, the matrix  $\mathbf{A}^{(k,2,0)}$  is full block-pentadiagonal, while the matrix  $\mathbf{A}^{(k,b_k,0)}$  is block-tridiagonal. An illustration of these definitions is shown in Equation (4). A box is placed around the bulge on the second matrix.

$$\begin{array}{c} \mathbf{A}^{(k,2,0)} \qquad \qquad \qquad \mathbf{A}^{(k,4,6)} \\ \left( \begin{array}{cccccccc} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} & \mathbf{A}_{1,3} & 0 & 0 & 0 & 0 & 0 \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} & \mathbf{A}_{2,3} & \mathbf{A}_{2,4} & 0 & 0 & 0 & 0 \\ \mathbf{A}_{3,1} & \mathbf{A}_{3,2} & \mathbf{A}_{3,3} & \mathbf{A}_{3,4} & \mathbf{A}_{3,5} & 0 & 0 & 0 \\ 0 & \mathbf{A}_{4,2} & \mathbf{A}_{4,3} & \mathbf{A}_{4,4} & \mathbf{A}_{4,5} & \mathbf{A}_{4,6} & 0 & 0 \\ 0 & 0 & \mathbf{A}_{5,3} & \mathbf{A}_{5,4} & \mathbf{A}_{5,5} & \mathbf{A}_{5,6} & \mathbf{A}_{5,7} & 0 \\ 0 & 0 & 0 & \mathbf{A}_{6,4} & \mathbf{A}_{6,5} & \mathbf{A}_{6,6} & \mathbf{A}_{6,7} & \mathbf{A}_{6,8} \\ 0 & 0 & 0 & 0 & \mathbf{A}_{7,5} & \mathbf{A}_{7,6} & \mathbf{A}_{7,7} & \mathbf{A}_{7,8} \\ 0 & 0 & 0 & 0 & 0 & \mathbf{A}_{8,6} & \mathbf{A}_{8,7} & \mathbf{A}_{8,8} \end{array} \right), \quad \left( \begin{array}{cccccccc} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} & 0 & 0 & 0 & 0 & 0 & 0 \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} & \mathbf{A}_{2,3} & 0 & 0 & 0 & 0 & 0 \\ 0 & \mathbf{A}_{3,2} & \mathbf{A}_{3,3} & \mathbf{A}_{3,4} & \mathbf{A}_{3,5} & \boxed{\mathbf{A}_{3,6}} & 0 & 0 \\ 0 & 0 & \mathbf{A}_{4,3} & \mathbf{A}_{4,4} & \mathbf{A}_{4,5} & \mathbf{A}_{4,6} & 0 & 0 \\ 0 & 0 & \mathbf{A}_{5,3} & \mathbf{A}_{5,4} & \mathbf{A}_{5,5} & \mathbf{A}_{5,6} & \mathbf{A}_{5,7} & 0 \\ 0 & 0 & \boxed{\mathbf{A}_{6,3}} & \mathbf{A}_{6,4} & \mathbf{A}_{6,5} & \mathbf{A}_{6,6} & \mathbf{A}_{6,7} & \mathbf{A}_{6,8} \\ 0 & 0 & 0 & 0 & \mathbf{A}_{7,5} & \mathbf{A}_{7,6} & \mathbf{A}_{7,7} & \mathbf{A}_{7,8} \\ 0 & 0 & 0 & 0 & 0 & \mathbf{A}_{8,6} & \mathbf{A}_{8,7} & \mathbf{A}_{8,8} \end{array} \right). \end{array} \quad (4)$$

### 3.2 Rotations

The algorithm defines two types of block rotations  $R_i$  and  $R'_j$ , which are unitary similarity transformations, with the following properties.

**Definition 3.1** (Rotations). *The algorithm of [79] uses the following two types of rotations:*

1.  $R_i(\mathbf{A}^{(k,i,0)})$ , for  $i = 2, \dots, b_k - 1$ , operates on a block-pentadiagonal matrix without a bulge. It transforms the matrix  $\mathbf{A}^{(k,i,0)}$  to a matrix  $\mathbf{A}^{(k,i+1,i+3)}$ . In particular, the block  $\mathbf{A}_{i,i-1}$  becomes upper triangular, the block  $\mathbf{A}_{i+1,i-1}$  becomes zero, and a new bulge block arises at  $\mathbf{A}_{i,i+3}$ . Due to symmetry,  $\mathbf{A}_{i-1,i}$  becomes lower triangular,  $\mathbf{A}_{i-1,i+1}$  is eliminated, and  $\mathbf{A}_{i+3,i}$  becomes non-zero.
2.  $R'_j(\mathbf{A}^{(k,s,j+1)})$ , for some  $j = s + 1, \dots, b_k - 1$ , operates on a block-pentadiagonal matrix with a bulge at positions  $\mathbf{A}_{j+1,j-2}$ ,  $\mathbf{A}_{j-2,j+1}$ . It transforms the matrix  $\mathbf{A}^{(k,s,j+1)}$  to a matrix  $\mathbf{A}^{(k,s,j+3)}$ , such that the bulge is moved two positions “down-and-right”, i.e. the blocks  $\mathbf{A}_{j-2,j+1}$  and  $\mathbf{A}_{j+1,j-2}$  become zero and the blocks  $\mathbf{A}_{j,i+3}$  and  $\mathbf{A}_{j,j+3}$  become the new bulge. In addition, the matrices  $\mathbf{A}_{j,j-2}$  and  $\mathbf{A}_{j+1,j-1}$  become upper triangular, and, by symmetry, the matrices  $\mathbf{A}_{j-2,j}$  and  $\mathbf{A}_{j-1,j+1}$  become lower triangular.

An example of the aforementioned rotations is illustrated in Equations (18) and (19) in the Appendix, and in Lemmas E.1 and E.2 it is proved that both types can be stably implemented in floating point using fast QR factorizations.

### 3.3 Recursive bandwidth halving

Using Lemmas E.1 and E.2, we can analyze the following Algorithm 2, which halves the bandwidth of a matrix. Its complexity and stability properties are stated in Lemma E.3 in the Appendix. Applying this algorithm recursively gives the main Theorem 3.1.

**Theorem 3.1.** *There exists a floating point implementation of the tridiagonal reduction algorithm of [79], which takes as input a Hermitian matrix  $\mathbf{A}$ , and returns a tridiagonal matrix  $\tilde{\mathbf{T}}$ , and (optionally) an approximately unitary matrix  $\tilde{\mathbf{Q}}$ . If the machine precision  $\mathbf{u}$  satisfies  $\mathbf{u} \leq \epsilon \frac{1}{cn^{\beta+4}}$ , where  $\epsilon \in (0, 1)$ ,  $c$  is a constant, and  $\beta$  is the same as in Corollary E.1, which translates to  $O(\log(n) + \log(1/\epsilon))$  bits of precision, then the following hold:*

$$\|\tilde{\mathbf{Q}}\tilde{\mathbf{Q}}^* - \mathbf{I}\| \leq \epsilon, \quad \text{and} \quad \|\mathbf{A} - \tilde{\mathbf{Q}}\tilde{\mathbf{T}}\tilde{\mathbf{Q}}^*\| \leq \epsilon \|\mathbf{A}\|.$$

The algorithm executes at most  $O(n^2 S_\omega(\log(n)))$  floating point operations to return only  $\tilde{\mathbf{T}}$ , where  $S_x(m) = \sum_{l=1}^m (2^{x-2})^l$ . If  $\mathbf{A}$  is banded with  $1 \leq d \leq n$  bands, the floating point operations reduce to  $O(n^2 S_\omega(\log(d)))$ . If  $\tilde{\mathbf{Q}}$  is also returned, the complexity increases to  $O(n^2 C_\omega(\log(n)))$ , where  $C_x(n) := \sum_{k=2}^{\log(n)-2} (S_x(\log(n) - 1) - S_x(k))$ . If  $\omega$  is treated as a constant  $\omega \approx 2.371$  the corresponding complexities are  $O(n^\omega)$ ,  $O(n^2 d^{\omega-2})$ , and  $O(n^\omega \log(n))$ , respectively.

*Proof.* **TOPROVE 5** □

**Algorithm:**  $[\widehat{\mathbf{Q}}^{(k)}, \mathbf{A}^{(k-1,2,0)}] \leftarrow \text{HALVE}(\mathbf{A}^{(k,2,0)}, k, n)$

- 1: Set  $n_k = 2^k$ ,  $b_k = n/n_k$ .
- 2: **for**  $i = 2, \dots, b_k$  :
- 3:     Compute  $\mathbf{Q}_{i,i}, \mathbf{A}^{(k,i+1,i+3)} \leftarrow R_i(\mathbf{A}^{(k,i,0)})$ .
- 4:     **for**  $j = i + 2, \dots, b_k$  with step 2 :
- 5:         Compute  $\mathbf{Q}_{i,j}, \mathbf{A}^{(k,i+1,j+3)} \leftarrow R'_j(\mathbf{A}^{(k,i+1,j+1)})$ .
- 6:     Stack together all the matrices  $\mathbf{Q}_{i,j}$  to form  $\mathbf{Q}_i$ .
- 7: Assemble the matrix  $\widehat{\mathbf{Q}}^{(k)}$  by multiplying the matrices  $\mathbf{Q}_i$ .
- 8: **return**  $\widehat{\mathbf{Q}}^{(k)}, \mathbf{A}^{(k-1,2,0)}$ .

**Algorithm 2:** Halves the bandwidth of a Hermitian matrix with unitary rotations.

### 3.4 Eigenvalues of Hermitian matrices

We now have all the prerequisites to compute the eigenvalues Hermitian matrices in nearly matrix multiplication time in finite precision. For this we can use the eigenvalue solver of [13], which has  $\widetilde{O}(n^2)$  boolean complexity, albeit in the Boolean RAM model. Specifically, the algorithm accepts as input symmetric tridiagonal matrices with bounded integer entries.

**Theorem 3.2** (Imported from [13, 15]). *Let  $\mathbf{T}$  be a symmetric tridiagonal matrix with integer elements bounded in magnitude by  $2^m$  for some  $m$ . Let  $\epsilon = 2^{-u} \in (0, 1)$  be a desired accuracy. Algorithm 4.1 of [13] computes a set of approximate eigenvalues  $\widetilde{\lambda}_i \in \mathbb{R}$  (which are returned as rationals) such that  $|\widetilde{\lambda}_i - \lambda_i(\mathbf{T})| < \epsilon$ . The algorithm requires  $O(n^2 b \log^2(n) \log(nb) (\log^2(b) + \log(n)) \log(\log(nb)))$  boolean operations, where  $b = m + u$ .*

**Theorem 3.3.** *Let  $\mathbf{A}$  be a (banded) Hermitian matrix, with  $\|\mathbf{A}\| \leq 1$ ,  $1 \leq d \leq n - 1$  off-diagonals, and let  $\epsilon \in (0, 1)$  be an accuracy parameter. Assume that the elements of  $\mathbf{A}$  are floating point numbers on a machine with precision  $\mathbf{u}$ ,  $t = \log(1/\mathbf{u})$  bits for the significand, and  $p = O(\log(\log(n)))$  bits for the exponent. There exists an algorithm that returns a set of  $n$  approximate eigenvalues  $\widetilde{\lambda}_1, \dots, \widetilde{\lambda}_n$  such that*

$$|\widetilde{\lambda}_i - \lambda_i(\mathbf{A})| \leq \epsilon$$

using at most

$$O(n^2 S_\omega(\log(d)) \cdot \mathcal{F}(\log(\frac{n}{\epsilon})) + n^2 \text{polylog}(\frac{n}{\epsilon}))$$

boolean operations, where  $\mathcal{F}(b)$  is the bit complexity of a floating point operation on  $b$  bits, and  $n^2 S_\omega(\log(d)) = O(n^2 d^{\omega-2})$  if  $\omega$  is treated as a constant greater than two.

*Proof.* **TOPROVE 6** □

## 4 Further applications of stable tridiagonal reduction and Hermitian eigenvalue solver

In this section we state some applications of the tridiagonal reduction algorithm to some eigenproblems.

### 4.1 Singular values and condition number

**Proposition 4.1.** *Given a matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$ , with  $\|\mathbf{A}\| \leq 1$ , an integer  $k \in [n]$ , and accuracy  $\epsilon \in (0, 1)$ , we can compute a value  $\widetilde{\sigma}_k \in (1 \pm \epsilon)\sigma_k(\mathbf{A})$  in*

$$O\left(\left[n^\omega \mathcal{F}(\log(\frac{n}{\epsilon \sigma_k})) + n^2 \text{polylog}(\frac{n}{\epsilon \sigma_k})\right] \log(\log(\epsilon \sigma_k))\right)$$

boolean operations, deterministically.

**Note:** If  $\|A\| > 1$ , we can scale with  $1/(n\|A\|_{\max})$  or  $1/\|A\|_F$ . The complexity is unaffected.

Proof. **TOPROVE 7** □

**Corollary 4.1.** Let  $A \in \mathbb{C}^{n \times n}$ ,  $\kappa = \kappa(A)$ , and  $\delta \in (0, 1/2)$ . We can compute  $\tilde{\kappa}$  which  $\kappa \leq \tilde{\kappa} \leq 3n\kappa$ , in

$$O\left(\left[n^\omega \mathcal{F}(\log(n\kappa)) + n^2 \text{polylog}(n\kappa)\right] \log(\log(n\kappa))\right)$$

boolean operations deterministically.

Proof. **TOPROVE 8** □

## 4.2 Definite pencil eigenvalues

Using the proposed algorithms we can compute the eigenvalues of a Hermitian definite pencil.

**Corollary 4.2.** Let  $H$  be Hermitian,  $S$  Hermitian positive-definite, both with size  $n$  and floating point elements, on a machine with precision  $u$ ,  $t = \log(1/u)$  bits for the significand, and  $p = O(\log(\log(n)))$  bits for the exponent. Assume that  $\|H\|, \|S^{-1}\| \leq 1$ ,  $\kappa(S) \in \text{poly}(n)$ , and that we have access to  $\tilde{\kappa} \in [\kappa(S), Z\kappa(S)]$ , where  $Z > 1$  might be a constant or a function of  $n$ . There exists an algorithm that returns a set of  $n$  approximate eigenvalues  $\tilde{\lambda}_1, \dots, \tilde{\lambda}_n$  such that

$$\left| \tilde{\lambda}_i - \lambda_i(H, S) \right| \leq \epsilon$$

using at most

$$O\left(n^\omega \mathcal{F}(\log(n) \log(Z\kappa(S)) + \log(\frac{1}{\epsilon})) + n^2 \text{polylog}(\frac{n}{\epsilon})\right)$$

boolean operations. If  $\tilde{\kappa}$  is not given, it can be computed up to a factor  $Z = 3n$  with Corollary 4.1.

Proof. **TOPROVE 9** □

**Remark 4.1.** In Theorem 4.2 we assumed that  $\|H\|, \|S^{-1}\| \leq 1$ . This is not a limitation since we can approximate  $\|S^{-1}\|$  with Proposition 4.1, and then scale accordingly. Formally, let  $\eta \gtrsim \|H\|$  and  $\sigma \gtrsim \|S^{-1}\|$ . Then we can rewrite the generalized eigenproblem

$$HC = SCA \Leftrightarrow \left(\frac{1}{\eta}H\right)C = (\sigma S)C\left(\Lambda \frac{1}{\eta\sigma}\right) \Leftrightarrow H'C = S'CA',$$

i.e. it is the same generalized eigenproblem only with scaled eigenvalues. Assuming that the matrices  $H$  and  $S$  are “well-conditioned,” i.e. their norms and condition numbers  $\in \text{poly}(n)$ , the eigenvalues are scaled by at most a  $1/\text{poly}(n)$  factor, and thus it suffices to scale  $\epsilon$  by  $1/\text{poly}(n)$  as well. We can thus safely make the unit-norms assumption.

## 4.3 Spectral gaps

In a similar way we can compute the spectral gap between a pair of eigenvalues of Hermitian matrices and Hermitian-definite pencils.

**Corollary 4.3.** Let  $A$  be a banded Hermitian matrix of size  $n$  with  $1 \leq d \leq n-1$  off-diagonals,  $\|A\| \leq 1$ , and its eigenvalues  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ . For some integer  $k \in [n-1]$ , let  $\mu_k = \frac{\lambda_k + \lambda_{k+1}}{2}$  and  $\Delta_k = \lambda_{k+1} - \lambda_k$ . Given an accuracy  $\epsilon \in (0, 1)$ , we can compute two values  $\tilde{\mu}_k$  and  $\tilde{\Delta}_k$  such that

$$\tilde{\mu}_k \in \mu_k \pm \epsilon \Delta_k, \quad \text{and} \quad \tilde{\Delta}_k \in (1 \pm \epsilon) \Delta_k,$$

in

$$O\left(n^2 \left[d^{\omega-2} \cdot \mathcal{F}(\log(\frac{n}{\epsilon \Delta_k})) + \text{polylog}(\frac{n}{\epsilon \Delta_k})\right] \log(\log(\frac{1}{\epsilon \Delta_k}))\right)$$

boolean operations.

*Proof.* TOPROVE 10 □

**Corollary 4.4.** *Let  $\mathbf{H}$  be Hermitian,  $\mathbf{S}$  Hermitian positive-definite, both with size  $n$  and  $\|\mathbf{H}\|, \|\mathbf{S}^{-1}\| \leq 1$ , which define a Hermitian-definite pencil  $(\mathbf{H}, \mathbf{S})$ . Given  $k \in [n-1]$ ,  $\tilde{\kappa} \in [\kappa(\mathbf{S}), Z\kappa(\mathbf{S})]$ , where  $Z > 1$  might be a constant or a function of  $n$ , and accuracy  $\epsilon \in (0, 1/2)$ , we can compute  $\tilde{\mu}_k = \mu_k \pm \epsilon \Delta_k$  and  $\tilde{\Delta}_k = (1 \pm \epsilon) \Delta_k$ , where  $\mu_k = \frac{\lambda_k + \lambda_{k+1}}{2}$  and  $\Delta_k = \lambda_k - \lambda_{k+1}$ . The algorithm requires*

$$O\left(\left\lceil n^\omega \mathcal{F}\left(\log(n) \log(Z\kappa(\mathbf{S})) + \log\left(\frac{1}{\epsilon \Delta_k}\right)\right) + n^2 \text{polylog}\left(\frac{n}{\epsilon \Delta_k}\right) \right\rceil \log\left(\log\left(\frac{1}{\epsilon \Delta_k}\right)\right)\right)$$

*boolean operations. If  $\tilde{\kappa}$  is not given, it can be computed up to a factor  $Z = 3n$  with Corollary 4.1.*

*Proof.* TOPROVE 11 □

#### 4.4 Spectral projectors and invariant subspaces

For a Hermitian matrix  $\mathbf{A}$ , the algorithm called GAP in [83] computes the spectral gap and the midpoint in the spirit of Corollary 4.3 using  $O\left(n^\omega \log\left(\frac{1}{\epsilon \Delta_k}\right) \log\left(\frac{1}{\delta \epsilon \Delta_k}\right) \cdot \mathcal{F}\left(\log^3\left(\frac{n}{\delta \epsilon \Delta_k}\right) \log\left(\frac{1}{\delta \epsilon \Delta_k}\right) \log(n)\right)\right)$  boolean operations and succeeds with probability  $1 - \delta$ . On the other hand, if  $\omega$  is treated as a constant larger than two, the algorithm of Corollary 4.3 requires

$$O\left(n^\omega \mathcal{F}(\log(n)) + n^2 \text{polylog}(n/\epsilon)\right)$$

boolean operations, which is significantly faster than GAP. Moreover, it is deterministic, and it has a lower complexity for banded matrices. A key difference between the two algorithms is that the GAP algorithm of [83] is fully analyzed in floating point, while the algorithm of 4.3 requires [15] which is analyzed in Boolean RAM.

Note that originally the GAP algorithm was used as a preliminary step to locate the gap and the midpoint in order to compute spectral projectors on invariant subspaces. Corollary 4.3 can serve as a direct replacement, providing an end-to-end deterministic algorithm for computing spectral projectors.

**Corollary 4.5.** *Let  $(\mathbf{H}, \mathbf{S})$  be a Hermitian definite pencil of size  $n$ , with  $\|\mathbf{H}\|, \|\mathbf{S}^{-1}\| \leq 1$ , and  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  its eigenvalues. Given as input  $\mathbf{H}, \mathbf{S}$ , an integer  $1 \leq k \leq n-1$ , an error parameter  $\epsilon \in (0, 1)$ , we can compute a matrix  $\tilde{\Pi}_k$  such that*

$$\left\|\tilde{\Pi}_k - \Pi_k\right\| \leq \epsilon,$$

*where  $\Pi_k$  is the true spectral projector on the invariant subspace that is associated with the  $k$  smallest eigenvalues, in*

$$O\left(n^\omega \mathcal{F}\left(\log(n) \log^3\left(\frac{1}{\Delta_k}\right) \log\left(\frac{n\kappa}{\epsilon \Delta_k}\right)\right) \left(\log\left(\frac{1}{\Delta_k}\right) + \log\left(\log\left(\frac{n\kappa}{\epsilon \Delta_k}\right)\right)\right) + n^2 \text{polylog}\left(\frac{n\kappa}{\Delta_k}\right)\right)$$

*boolean operations, where  $\kappa := \kappa(\mathbf{S}) = \|\mathbf{S}\| \|\mathbf{S}^{-1}\|$  and  $\Delta_k = \lambda_{k+1} - \lambda_k$ .*

*Proof.* TOPROVE 12 □

#### 4.5 Stable inversion and factorizations

Theorem 3.1 can be used to invert a matrix stably in floating point, improving the bit requirement of the logarithmically-stable algorithm of [28], however, at the cost of increasing the arithmetic complexity by a factor of  $\log(n)$ , which ultimately leads to a slower algorithm. To achieve this, we first form the matrix  $\begin{pmatrix} 0 & \mathbf{A}^* \\ \mathbf{A} & 0 \end{pmatrix}$ , and then reduce it to tridiagonal form. Afterwards we can solve  $O(n)$  linear systems stably with QR factorization of the tridiagonal matrix, where the right hand sides are the columns of  $\tilde{\mathbf{Q}}$ , in time  $O(n)$  each. We perform one final multiplication to obtain the inverse. We do not expand the analysis further since the boolean complexity of the algorithm is slower than the one of [28] (even if we use a slow algorithm for basic arithmetic operations). The stable inversion algorithm can be used to obtain stable factorization algorithms, such as LU [28] or Cholesky [83].



## 5 Conclusion

In this work we provided a deterministic complexity analysis for Hermitian eigenproblems. In the Real RAM model, we reported nearly-linear complexity upper bounds for the full diagonalization of arrowhead and tridiagonal matrices, and nearly matrix multiplication-type complexities for diagonalizing Hermitian matrices and for the SVD. This was achieved by analyzing the divide-and-conquer algorithm of [48], when implemented with the Fast Multipole Method. We also showed that the tridiagonal reduction algorithm of [79] is numerically stable in the floating point model. This paved the way to obtain improved deterministic boolean complexities for computing the eigenvalues, singular values, spectral gaps, and spectral projectors, of Hermitian matrices and Hermitian-definite pencils. Some interesting questions for future research are the following.

1. **Stability of arrowhead diagonalization:** The FMM-accelerated arrowhead diagonalization algorithm was only analyzed in the Real RAM model. Several works [48, 87, 70, 21] have provided stabilization techniques of related algorithms in floating point, albeit, without an end-to-end complexity analysis. Such an analysis will be insightful to better understand the boolean complexity of (Hermitian) diagonalization.
2. **Condition number in the SVD complexity:** The complexity of the SVD in Theorem 1.2 has a polylogarithmic dependence on the condition number. Frustratingly, we were not able to remove it at the time of this writing.
3. **Non-Hermitian diagonalization:** Schönhage also proved that a non-Hermitian matrix can be reduced to upper Hessenberg form in matrix multiplication time [79]. In this work we only provided the stability analysis for the Hermitian case. It would be interesting to investigate whether the Hessenberg reduction algorithm can be used to diagonalize non-Hermitian matrices in matrix multiplication time deterministically (e.g. in conjunction with [8, 7, 9]).
4. **Deterministic pseudospectral shattering:** One of the main techniques of the seminal work of [6] is called “pseudospectral shattering.” The main idea is to add a tiny random perturbation to the original matrix to ensure that the minimum eigenvalue gap between any pair of eigenvalues will be at least polynomial in  $1/n$ . We highlight that the deflation preprocessing step in Proposition A.1 has this precise effect: the pseudospectrum is shattered with respect to a deterministic grid. Can this be generalized to obtain *deterministic* pseudospectral shattering techniques, for Hermitian or even non-Hermitian matrices?

## References

- [1] Josh Alman, Ran Duan, Virginia Vassilevska Williams, Yinzhan Xu, Zixuan Xu, and Renfei Zhou. More asymmetry yields faster matrix multiplication. In *Proc. 2025 ACM-SIAM Symposium on Discrete Algorithms*, pages 2005–2039. SIAM, 2025.
- [2] Alexandr Andoni and Huy L Nguyen. Eigenvalues of a matrix in the streaming model. In *Proc. 24th ACM-SIAM Symposium on Discrete Algorithms*, pages 1729–1737. SIAM, 2013.
- [3] Diego Armentano, Carlos Beltrán, Peter Bürgisser, Felipe Cucker, and Michael Shub. A stable, polynomial-time algorithm for the eigenpair problem. *Journal of the European Mathematical Society*, 20(6):1375–1437, 2018.
- [4] Grey Ballard, James Demmel, and Nicholas Knight. Communication avoiding successive band reduction. *ACM SIGPLAN Notices*, 47(8):35–44, 2012.
- [5] Grey Ballard, James Demmel, and Nicholas Knight. Avoiding communication in successive band reduction. *ACM Transactions on Parallel Computing*, 1(2):1–37, 2015.
- [6] Jess Banks, Jorge Garza-Vargas, Archit Kulkarni, and Nikhil Srivastava. Pseudospectral Shattering, the Sign Function, and Diagonalization in Nearly Matrix Multiplication Time. *Foundations of Computational Mathematics*, pages 1–89, 2022.

- [7] Jess Banks, Jorge Garza-Vargas, and Nikhil Srivastava. Global Convergence of Hessenberg Shifted QR II: Numerical Stability. *arXiv*, 2022.
- [8] Jess Banks, Jorge Garza-Vargas, and Nikhil Srivastava. Global convergence of Hessenberg Shifted QR I: Exact Arithmetic. *Foundations of Computational Mathematics*, pages 1–34, 2024.
- [9] Jess Banks, Jorge Garza-Vargas, and Nikhil Srivastava. Global Convergence of Hessenberg Shifted QR III: Approximate Ritz Values via Shifted Inverse Iteration. *SIAM Journal on Matrix Analysis and Applications*, 46(2):1212–1246, 2025.
- [10] Jesse L Barlow. Error analysis of update methods for the symmetric eigenvalue problem. *SIAM Journal on Matrix Analysis and Applications*, 14(2):598–618, 1993.
- [11] Michael Ben-Or and Lior Eldar. A Quasi-Random Approach to Matrix Spectral Analysis. In *Proc. 9th Innovations in Theoretical Computer Science Conference*, pages 6:1–6:22. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2018.
- [12] Rajendra Bhatia. *Perturbation bounds for matrix eigenvalues*. SIAM, 2007.
- [13] Dario Bini and Victor Pan. Parallel complexity of tridiagonal symmetric eigenvalue problem. In *Proc. 2nd ACM-SIAM Symposium on Discrete Algorithms*, pages 384–393, 1991.
- [14] Dario Bini and Victor Pan. Practical improvement of the divide-and-conquer eigenvalue algorithms. *Computing*, 48(1):109–123, 1992.
- [15] Dario Bini and Victor Y Pan. Computing matrix eigenvalues and polynomial zeros where the output is real. *SIAM Journal on Computing*, 27(4):1099–1115, 1998.
- [16] Christian H Bischof, Bruno Lang, and Xiaobai Sun. A framework for symmetric band reduction. *ACM Transactions on Mathematical Software*, 26(4):581–601, 2000.
- [17] D Boley and Gene Howard Golub. Inverse eigenvalue problems for band matrices. In *Numerical Analysis: Proceedings of the Biennial Conference Held at Dundee, June 28–July 1, 1977*, pages 23–31. Springer, 1977.
- [18] Christos Boutsidis and David P Woodruff. Optimal CUR matrix decompositions. In *Proc. 46th ACM Symposium on Theory of Computing*, pages 353–362, 2014.
- [19] Christos Boutsidis, David P Woodruff, and Peilin Zhong. Optimal principal component analysis in distributed and streaming models. In *Proc. 48th ACM Symposium on Theory of Computing*, pages 236–249, 2016.
- [20] Christos Boutsidis, Anastasios Zouzias, Michael W Mahoney, and Petros Drineas. Randomized dimensionality reduction for  $k$ -means clustering. *IEEE Transactions on Information Theory*, 61(2):1045–1062, 2014.
- [21] Difeng Cai and Jianlin Xia. A stable matrix version of the fast multipole method: stabilization strategies and examples. *ETNA-Electronic Transactions on Numerical Analysis*, 54, 2020.
- [22] Kenneth L Clarkson and David P Woodruff. Low-rank approximation and regression in input sparsity time. *Journal of the ACM*, 63(6):1–45, 2017.
- [23] Michael B Cohen, Sam Elder, Cameron Musco, Christopher Musco, and Madalina Persu. Dimensionality reduction for  $k$ -means clustering and low rank approximation. In *Proc. 47th ACM Symposium on Theory of Computing*, pages 163–172, 2015.
- [24] Jan JM Cuppen. A divide and conquer method for the symmetric tridiagonal eigenproblem. *Numerische Mathematik*, 36:177–195, 1980.
- [25] Eric Darve. The fast multipole method i: Error analysis and asymptotic complexity. *SIAM Journal on Numerical Analysis*, 38(1):98–128, 2000.

- [26] Eric Darve. The fast multipole method: numerical implementation. *Journal of Computational Physics*, 160(1):195–240, 2000.
- [27] TJ Dekker and JF Traub. The shifted QR algorithm for Hermitian matrices. *Linear Algebra and its Applications*, 4(3):137–154, 1971.
- [28] James Demmel, Ioana Dumitriu, and Olga Holtz. Fast linear algebra is stable. *Numerische Mathematik*, 108(1):59–91, 2007.
- [29] James Demmel, Ioana Dumitriu, Olga Holtz, and Robert Kleinberg. Fast matrix multiplication is stable. *Numerische Mathematik*, 106(2):199–224, 2007.
- [30] James Demmel, Ioana Dumitriu, and Ryan Schneider. Fast and inverse-free algorithms for deflating subspaces. *arXiv*, 2024.
- [31] James Demmel, Ioana Dumitriu, and Ryan Schneider. Generalized Pseudospectral Shattering and Inverse-Free Matrix Pencil Diagonalization. *Foundations of Computational Mathematics*, pages 1–77, 2024.
- [32] James Demmel and Krešimir Veselić. Jacobi’s method is more accurate than QR. *SIAM Journal on Matrix Analysis and Applications*, 13(4):1204–1245, 1992.
- [33] James W Demmel. *Applied numerical linear algebra*. SIAM, 1997.
- [34] Inderjit Singh Dhillon. *A new  $O(N^2)$  algorithm for the symmetric tridiagonal eigenvalue/eigenvector problem*. University of California, Berkeley, 1997.
- [35] Jack Dongarra and Francis Sullivan. Guest Editors Introduction to the top 10 algorithms. *Computing in Science & Engineering*, 2(01):22–23, 2000.
- [36] Jack J Dongarra and Danny C Sorensen. A fully parallel algorithm for the symmetric eigenvalue problem. *SIAM Journal on Scientific and Statistical Computing*, 8(2):s139–s154, 1987.
- [37] Petros Drineas, Michael W Mahoney, and Shan Muthukrishnan. Relative-error CUR matrix decompositions. *SIAM Journal on Matrix Analysis and Applications*, 30(2):844–881, 2008.
- [38] Jeff Erickson, Ivor Van Der Hoog, and Tillmann Miltzow. Smoothing the gap between NP and ER. *SIAM Journal on Computing*, 0(0):FOCS20–102–FOCS20–138, 2022.
- [39] John GF Francis. The QR transformation a unitary analogue to the LR transformation—Part 1. *The Computer Journal*, 4(3):265–271, 1961.
- [40] John GF Francis. The QR transformation—Part 2. *The Computer Journal*, 4(4):332–345, 1962.
- [41] Alan Frieze, Ravi Kannan, and Santosh Vempala. Fast Monte-Carlo algorithms for finding low-rank approximations. *Journal of the ACM*, 51(6):1025–1041, 2004.
- [42] Martin Fürer. Faster integer multiplication. In *Proc. 39th ACM Symposium on Theory of Computing*, pages 57–66, 2007.
- [43] Doron Gill and Eitan Tadmor. An  $O(N^2)$  Method for Computing the Eigensystem of  $N \times N$  Symmetric Tridiagonal Matrices by the Divide and Conquer Approach. *SIAM Journal on Scientific and Statistical Computing*, 11(1):161–173, 1990.
- [44] Gene Golub and William Kahan. Calculating the singular values and pseudo-inverse of a matrix. *Journal of the Society for Industrial and Applied Mathematics, Series B: Numerical Analysis*, 2(2):205–224, 1965.
- [45] Gene H Golub and Henk A Van der Vorst. Eigenvalue computation in the 20th century. *Journal of Computational and Applied Mathematics*, 123(1-2):35–65, 2000.

- [46] Gene H Golub and Charles F Van Loan. *Matrix Computations*. Johns Hopkins University Press, 2013.
- [47] Ming Gu and Stanley C Eisenstat. A stable and fast algorithm for updating the singular value decomposition, 1993.
- [48] Ming Gu and Stanley C Eisenstat. A divide-and-conquer algorithm for the symmetric tridiagonal eigenproblem. *SIAM Journal on Matrix Analysis and Applications*, 16(1):172–191, 1995.
- [49] Ming Gu and Stanley C Eisenstat. Efficient algorithms for computing a strong rank-revealing qr factorization. *SIAM Journal on Scientific Computing*, 17(4):848–869, 1996.
- [50] Juris Hartmanis and Janos Simon. On the power of multiplication in random access machines. In *15th Annual Symposium on Switching and Automata Theory*, pages 13–23. IEEE, 1974.
- [51] David Harvey and Joris Van Der Hoeven. Integer multiplication in time  $O(n \log n)$ . *Annals of Mathematics*, 193(2):563–617, 2021.
- [52] Nicholas J Higham. *Accuracy and stability of numerical algorithms*. SIAM, 2002.
- [53] Walter Hoffmann and Beresford N Parlett. A new proof of global convergence for the tridiagonal QL algorithm. *SIAM Journal on Numerical Analysis*, 15(5):929–937, 1978.
- [54] Alston S Householder. Unitary triangularization of a nonsymmetric matrix. *Journal of the ACM*, 5(4):339–342, 1958.
- [55] C.G.J. Jacobi. Über ein leichtes Verfahren die in der Theorie der Säcularstörungen vorkommenden Gleichungen numerisch aufzulösen. *Journal für die reine und angewandte Mathematik*, 30:51–94, 1846.
- [56] Ian T Jolliffe. *Principal component analysis for special types of data*. Springer, 2002.
- [57] Praneeth Kacham and David P Woodruff. Faster Algorithms for Schatten-p Low Rank Approximation. In *Proc. Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 2024.
- [58] Vera N Kublanovskaya. On some algorithms for the solution of the complete eigenvalue problem. *USSR Computational Mathematics and Mathematical Physics*, 1(3):637–657, 1962.
- [59] Cornelius Lanczos. An iteration method for the solution of the eigenvalue problem of linear differential and integral operators. *Journal of Research of the National Bureau of Standards*, 45(4), 1950.
- [60] Oren E Livne and Achi Brandt.  $N$  roots of the secular equation in  $O(N)$  operations. *SIAM Journal on Matrix Analysis and Applications*, 24(2):439–453, 2002.
- [61] Anand Louis and Santosh S Vempala. Accelerated newton iteration for roots of black box polynomials. In *Proc. 57th IEEE Symposium on Foundations of Computer Science*, pages 732–740. IEEE, 2016.
- [62] Per-Gunnar Martinsson and Vladimir Rokhlin. An accelerated kernel-independent fast multipole method in one dimension. *SIAM Journal on Scientific Computing*, 29(3):1160–1178, 2007.
- [63] A Melman. Numerical solution of a secular equation. *Numerische Mathematik*, 69:483–493, 1995.
- [64] RV Mises and Hilda Pollaczek-Geiringer. Praktische Verfahren der Gleichungsauflösung. *ZAMM-Journal of Applied Mathematics and Mechanics/Zeitschrift für Angewandte Mathematik und Mechanik*, 9(1):58–77, 1929.
- [65] Cameron Musco, Christopher Musco, and Aaron Sidford. Stability of the Lanczos method for matrix function approximation. In *Proc. 29th ACM-SIAM Symposium on Discrete Algorithms*, pages 1605–1624. SIAM, 2018.
- [66] Yuji Nakatsukasa, Zhaojun Bai, and François Gygi. Optimizing Halley’s iteration for computing the matrix polar decomposition. *SIAM Journal on Matrix Analysis and Applications*, 31(5):2700–2720, 2010.

- [67] Yuji Nakatsukasa and Nicholas J Higham. Stable and efficient spectral divide and conquer algorithms for the symmetric eigenvalue decomposition and the SVD. *SIAM Journal on Scientific Computing*, 35(3):A1325–A1349, 2013.
- [68] Deanna Needell, William Swartworth, and David P Woodruff. Testing positive semidefiniteness using linear measurements. In *Proc. 2022 IEEE Symposium on Foundations of Computer Science*, pages 87–97. IEEE, 2022.
- [69] Dianne P O’Leary and Gilbert W Stewart. Computing the eigenvalues and eigenvectors of symmetric arrowhead matrices. *Journal of Computational Physics*, 90(2):497–505, 1990.
- [70] Xiaofeng Ou and Jianlin Xia. Superdc: superfast divide-and-conquer eigenvalue decomposition with improved stability for rank-structured matrices. *SIAM Journal on Scientific Computing*, 44(5):A3041–A3066, 2022.
- [71] Chris C Paige. Accuracy and effectiveness of the Lanczos algorithm for the symmetric eigenproblem. *Linear algebra and its Applications*, 34:235–258, 1980.
- [72] Victor Y Pan and Zhao Q Chen. The complexity of the matrix eigenproblem. In *Proc. 31st ACM Symposium on Theory of Computing*, pages 507–516, 1999.
- [73] Christos H Papadimitriou, Hisao Tamaki, Prabhakar Raghavan, and Santosh Vempala. Latent semantic indexing: A probabilistic analysis. In *Proc. 17th ACM Symposium on Principles of Database Systems*, pages 159–168, 1998.
- [74] Beresford N Parlett. *The Symmetric Eigenvalue Problem*. SIAM, 1998.
- [75] Vladimir Rokhlin. Rapid solution of integral equations of classical potential theory. *Journal of computational physics*, 60(2):187–207, 1985.
- [76] H Rutishauser. On Jacobi rotation patterns. In *Proceedings of Symposia in Applied Mathematics*, volume 15, pages 219–239, 1963.
- [77] Yousef Saad. *Numerical methods for large eigenvalue problems: revised edition*. SIAM, 2011.
- [78] Ryan Schneider. *Pseudospectral divide-and-conquer for the generalized eigenvalue problem*. University of California, San Diego, 2024.
- [79] Arnold Schönhage. Unitäre transformationen großer matrizen. *Numerische Mathematik*, 20:409–417, 1972.
- [80] Arnold Schönhage. On the power of random access machines. In *Proc. International Colloquium on Automata, Languages, and Programming*, pages 520–529. Springer, 1979.
- [81] Arnold Schönhage and Volker Strassen. Fast multiplication of large numbers. *Computing*, 7:281–292, 1971.
- [82] Rikhav Shah. Hermitian Diagonalization in Linear Precision. In *Proc. 2025 ACM-SIAM Symposium on Discrete Algorithms*, pages 5599–5615. SIAM, 2025.
- [83] Aleksandros Sobczyk, Marko Mladenovic, and Mathieu Luisier. Invariant subspaces and PCA in nearly matrix multiplication time. *Advances in Neural Information Processing Systems*, 37:19013–19086, 2024.
- [84] Nevena Jakovčević Stor, Ivan Slapničar, and Jesse L Barlow. Accurate eigenvalue decomposition of real symmetric arrowhead matrices and applications. *Linear Algebra and its Applications*, 464:62–89, 2015.
- [85] Xiaobai Sun and Nikos P Pitsianis. A matrix version of the fast multipole method. *SIAM Review*, 43(2):289–300, 2001.
- [86] William Swartworth and David P Woodruff. Optimal eigenvalue approximation via sketching. In *Proc. 55th ACM Symposium on Theory of Computing*, pages 145–155, 2023.

- [87] James Vogel, Jianlin Xia, Stephen Cauley, and Venkataramanan Balakrishnan. Superfast divide-and-conquer method and perturbation analysis for structured eigenvalue solutions. *SIAM Journal on Scientific Computing*, 38(3):A1358–A1382, 2016.
- [88] James Hardy Wilkinson. Global convergene of tridiagonal QR algorithm with origin shifts. *Linear Algebra and its Applications*, 1(3):409–420, 1968.

## A Preliminaries for symmetric arrowhead diagonalization

This section contains the required preliminaries to describe the algorithm for symmetric arrowhead diagonalization. Our analysis relies on the methodology of [48], but we refer also to [69, 84] for related techniques. We first state the following lemma which describes the desired properties of the input matrix.

**Lemma A.1.** *Let  $\mathbf{H} = \begin{pmatrix} \alpha & \mathbf{z}^\top \\ \mathbf{z} & \mathbf{D} \end{pmatrix}$ ,  $\|\mathbf{H}\| \leq 1$ , be an arrowhead matrix in the form of Eq. (3). Suppose that the elements of  $\mathbf{H}$  satisfy the following properties, for all  $i = 2, \dots, n$ :*

$$d_{i+1} - d_i \geq \tau, \quad \text{and} \quad |\mathbf{z}_i| \geq \tau, \quad (5)$$

for some  $\tau \in (0, 1)$ , where  $d_i$  is the  $i$ -th element of  $\mathbf{D}$ . The eigenvalues  $\lambda_i$  are the roots of the secular equation

$$f(\lambda) = \lambda - \alpha + \sum_{j=2}^n \frac{z_j^2}{d_j - \lambda}, \quad (6)$$

and they satisfy the following interlacing property:

$$\lambda_1 < d_2 < \lambda_2 < \dots < d_n < \lambda_n. \quad (7)$$

Moreover, it also holds that

$$\begin{aligned} \min \{ |d_i - \lambda_i|, |d_{i+1} - \lambda_i| \} &\geq \frac{\tau^3}{n+1}, \quad \text{for all } i = 2, \dots, n-1, \\ \min \{ |d_2 - \lambda_1|, |d_n - \lambda_n| \} &\geq \frac{\tau^3}{n}, \end{aligned} \quad (8)$$

i.e., there is a well defined gap between the eigenvalues and their boundaries.

*Proof.* **TOPROVE 13** □

### A.1 Deflation

In order to ensure that the given matrix satisfies the requirements of Equation (5) we will use deflation, specifically the methodology of section 4 in [48].

**Proposition A.1** (Arrowhead deflation). *Let  $\mathbf{H} = \begin{pmatrix} \alpha & \mathbf{z}^\top \\ \mathbf{z} & \mathbf{D} \end{pmatrix}$  be an arrowhead matrix in the form of Eq. (3). There exists an orthogonal matrix  $\mathbf{G}$  and a matrix  $\tilde{\mathbf{H}} = \begin{pmatrix} \tilde{\mathbf{H}}' & 0 \\ 0 & \tilde{\mathbf{D}} \end{pmatrix}$  such that  $\tilde{\mathbf{D}}$  is diagonal and  $\tilde{\mathbf{H}}' = \begin{pmatrix} \tilde{\alpha}' & \tilde{\mathbf{z}}'^\top \\ \tilde{\mathbf{z}}' & \tilde{\mathbf{D}}' \end{pmatrix}$  is an arrowhead matrix that satisfies the requirements of Equation (5), specifically:*

$$\tilde{d}'_{j+1} - \tilde{d}'_j \geq \tau, \quad |\tilde{\mathbf{z}}'_i| \geq \tau, \quad \text{and} \quad \|\mathbf{H} - \mathbf{G}\tilde{\mathbf{H}}\mathbf{G}^\top\| \leq n\tau.$$

$\tilde{\mathbf{H}}$  can be computed in  $O(n \log(n))$  arithmetic operations and comparisons, and the product  $\mathbf{G}^\top \mathbf{B}$  for some matrix  $\mathbf{B}$  with  $r$  columns can be computed in additional  $O(nr)$  arithmetic operations on-the-fly.

*Proof.* **TOPROVE 14** □



## A.2 Reconstruction from approximate eigenvalues

If we have access to a set of approximate eigenvalues  $\widehat{\lambda}_i$  of  $\mathbf{H}$  that also satisfy the same interlacing property, then we can construct a matrix  $\widehat{\mathbf{H}}$  that is close to  $\mathbf{H}$ , and  $\widehat{\lambda}_i$  are the eigenvalues of  $\widehat{\mathbf{H}}$ . This is achieved with the following lemma from [17]. We use its restatement from [48].

**Lemma A.2** ([17, 48]). *Given a set of  $n$  numbers  $\widehat{\lambda}_1, \dots, \widehat{\lambda}_n$  and a diagonal matrix  $\mathbf{D} = \text{diag}(d_2, \dots, d_n)$  such that*

$$\widehat{\lambda}_1 < d_2 < \widehat{\lambda}_2 < \dots < d_n < \widehat{\lambda}_n,$$

*there exists a symmetric arrowhead matrix  $\widehat{\mathbf{H}} = \begin{pmatrix} \widehat{\alpha} & \widehat{\mathbf{z}}^\top \\ \widehat{\mathbf{z}} & \mathbf{D} \end{pmatrix}$ , whose eigenvalues are  $\widehat{\lambda}_i$ . In this case*

$$\begin{aligned} |\widehat{\mathbf{z}}_i| &= \sqrt{(d_i - \widehat{\lambda}_1)(\widehat{\lambda}_n - d_i) \prod_{j=2}^{i-1} \frac{\widehat{\lambda}_j - d_i}{d_j - d_i} \prod_{j=i}^{n-1} \frac{\widehat{\lambda}_j - d_i}{d_{j+1} - d_i}}, \\ \widehat{\alpha} &= \widehat{\lambda}_1 + \sum_{j=2}^n (\widehat{\lambda}_j - d_j), \end{aligned} \tag{9}$$

where the sign of  $\widehat{\mathbf{z}}_i$  can be chosen arbitrarily.

We shall use the spectral decomposition of  $\widehat{\mathbf{H}}$  as a backward approximate spectral decomposition of  $\mathbf{H}$ . We write

$$\mathbf{H} = \begin{pmatrix} \alpha & \mathbf{z}^\top \\ \mathbf{z} & \mathbf{D} \end{pmatrix}, \quad \text{and} \quad \widehat{\mathbf{H}} = \begin{pmatrix} \widehat{\alpha} & \widehat{\mathbf{z}}^\top \\ \widehat{\mathbf{z}} & \mathbf{D} \end{pmatrix},$$

in which case  $\|\mathbf{H} - \widehat{\mathbf{H}}\| \leq |\alpha - \widehat{\alpha}| + \|\mathbf{z} - \widehat{\mathbf{z}}\|$ . It suffices to show that  $\widehat{\alpha} \approx \alpha$  and  $\widehat{\mathbf{z}} \approx \mathbf{z}$ .

**Lemma A.3.** *Let  $\mathbf{H} = \begin{pmatrix} \alpha & \mathbf{z}^\top \\ \mathbf{z} & \mathbf{D} \end{pmatrix}$  be a symmetric arrowhead matrix with  $\|\mathbf{H}\| \leq 1$ , satisfying the properties of Lemma A.1 with parameter  $\tau$ , and let  $\widehat{\lambda}_1, \widehat{\lambda}_2, \dots, \widehat{\lambda}_n$  be a set of approximate eigenvalues that satisfy*

$$\widehat{\lambda}_1 < d_2 < \widehat{\lambda}_2 < \dots < d_n < \widehat{\lambda}_n, \quad \text{and} \quad |\widehat{\lambda}_i - \lambda_i| \leq \epsilon \frac{\tau^3}{2(n+1)},$$

for some  $\epsilon \in (0, 1/n)$ . Then the quantities  $\widehat{\alpha}$ ,  $\widehat{\mathbf{z}}$ , and  $\widehat{\mathbf{H}}$  from Lemma A.2 satisfy:

$$|\alpha - \widehat{\alpha}| \leq \frac{\epsilon \tau^3}{2}, \quad \|\mathbf{z} - \widehat{\mathbf{z}}\| \leq \frac{n\epsilon}{1 - n\epsilon}, \quad \|\mathbf{H} - \widehat{\mathbf{H}}\| \leq \frac{\epsilon \tau^3}{2} + \frac{n\epsilon}{1 - n\epsilon}.$$

*Proof.* **TOPROVE 15** □

## B Fast Arrowhead diagonalization

In this section we provide the full analysis of the algorithm of [48] when accelerated with the FMM. The FMM was introduced in [75], to accelerate the evaluation of integrals between interacting bodies, and it has been comprehensively analyzed in the literature [26, 25, 21, 85, 62, 60, 47]. Consider a function of the form

$$f(x) = \sum_{j=1}^n c_j k(x - x_j), \tag{10}$$

where  $x_j, c_j$  are constants and  $k(x)$  is a suitably chosen kernel function, typically one of  $\{\log(x), \frac{1}{x}, \frac{1}{x^2}\}$ . The FMM can be used to approximately evaluate  $f(x)$  over  $m$  different points  $x_i$  in only  $\widetilde{O}(m+n)$  arithmetic operations (suppressing logarithmic terms in the accuracy), instead of the naive  $O(mn)$  evaluation.

## B.1 Fast Multipole Method

For our analysis, we will need to use the FMM to evaluate the following functions (11)-(15), each on  $O(n)$  points. For each function, there is a guarantee that the evaluation points will satisfy certain criteria. More precisely, the magnitude of the denominators and the logarithm arguments will have a well-defined lower bound. Here  $\tau \in (0, 1)$  is a parameter that is determined later.

Function:	Guarantees for FMM:	Used in:
$f(\lambda) = \sum_{j=1}^n \frac{\mathbf{z}_j^2}{d_j - \lambda},$	$ d_j - \lambda  \geq \Omega(\text{poly}(\frac{\tau}{n})),$	Lemma B.1, (11)
$f(d_i) = \sum_{j=2}^n \log( \widehat{\lambda}_j - d_i ),$	$ \widehat{\lambda}_j - d_i  \geq \Omega(\text{poly}(\frac{\tau}{n})),$	Lemma B.2, (12)
$f(d_i) = \sum_{j=2}^n \log( d_j - d_i ),$	$ d_j - d_i  \geq \tau,$	Lemma B.2, (13)
$f(\lambda) = \sum_{k=2}^n \frac{(1 + \epsilon_k) \widehat{\mathbf{z}}_k \mathbf{q}_k}{d_k - \lambda},$	$ d_k - \lambda  \geq \Omega(\text{poly}(\frac{\tau}{n})),$	Lemma B.3, (14)
$f(\lambda) = \sum_{k=2}^n \frac{(1 + \epsilon_k)^2 \widehat{\mathbf{z}}_k^2}{(d_k - \lambda)^2},$	$ d_k - \lambda  \geq \Omega(\text{poly}(\frac{\tau}{n})),$	Lemma B.3. (15)

The fact that the magnitudes of the denonimators and the logarithm arguments are bounded from below allows us to use the seminal FMM analysis of [47, 60, 21]. To the best of our knowledge, [47] is one of the first works to rigorously analyze the application of the FMM on such kernel functions with end-to-end approximation and complexity bounds. In Section 3.4 they achieved an error of  $O\left(\epsilon \sum_i \frac{|x_i|}{|\omega^2 - d_i^2|}\right)$  for the function  $\Phi(\omega) = \sum_i \frac{x_i}{d_i^2 - \omega^2}$ , assuming that  $d_i$  and  $\omega$  satisfy similar interlacing properties to ours, in  $O(n \log^2(1/\epsilon))$  arithmetic operations. This complexity translates to  $O(n \log^2(\frac{n}{\tau\epsilon}))$  arithmetic operations, if we rescale  $\epsilon$  appropriately to obtain an absolute error  $O(\epsilon)$ , i.e., it satisfies the guarantees of Proposition 1.1 with  $\xi = 2$ . [60] used a kernel-softening approach and report similar bounds in  $O(n \log(1/\epsilon))$  operations. More recently, [21] developed a general framework based on the matrix-version of the FMM [85], and provided a thorough analysis and explicit bounds similar to [47], which can be used to obtain guarantees for all the functions (11)-(15).

For completeness, we summarize the main ideas and results of the aforementioned works and provide a short proof for Proposition 1.1 below. We remind once more that we do not account for floating point errors, but rather focus only on the error FMM approximation errors. Obtaining full, end-to-end complexity analysis under floating point errors has its own merit, and it is left as future work.

**Proposition B.1** (FMM). *There exists an algorithm, which we refer to as  $(\epsilon, n)$ -approximate FMM (or  $(\epsilon, n)$ -FMM, for short), which takes as input*

- a kernel function  $k(x) = \{\log(|x|), \frac{1}{x}, \frac{1}{x^2}\}$ ,
- $2n + m$  real numbers:  $\{x_1, \dots, x_m\} \cup \{c_1, \dots, c_n\} \cup \{y_1, \dots, y_n\}$ , and a constant  $C$ , such that  $m \leq n$  and for all  $i \in [m], j \in [n]$  it holds that

$$|x_i|, |c_j|, |y_j| < C \quad \text{and} \quad |x_i - y_j| \geq \Omega(\text{poly}(\frac{\epsilon}{n})).$$

It returns  $m$  values  $\widetilde{f}(x_1), \dots, \widetilde{f}(x_m)$  such that  $|\widetilde{f}(x_i) - f(x_i)| \leq \epsilon$ , for all  $i \in [m]$ , where  $f(x) = \sum_{j=1}^n c_j k(x_i - y_j)$ , in a total of  $O\left(n \log^\xi(\frac{n}{\epsilon})\right)$  arithmetic operations, where  $\xi \geq 1$  is a small constant that is independent of  $\epsilon, n$ .

*Proof.* **TOPROVE 16**

□

## B.2 Computing the eigenvalues with bisection

Using the FMM as a black-box evaluation of the targeted kernel functions, we can compute all the eigenvalues of an arrowhead matrix which are given by the roots of the secular equation. We highlight that [60] provided a rigorous analysis of the Newton iteration, based on the results of [63], to compute all the roots of the secular equation. The reported complexity is  $O(n \log(1/\epsilon))$ , however, this does not include the number of Newton steps. Due to the well-known quadratic convergence properties, the number of Newton steps should be bounded by  $O(\log(\log(1/\epsilon)))$ . For the sake of simplicity and completeness, instead of repeating the analysis of [60], we will use the following standard bisection method to compute all the eigenvalues in a total of  $O(n \log^2(1/\epsilon))$  operations. It might be slightly slower than the Newton iteration (up to a log factor), but it is significantly simpler and it does not require the evaluation of derivatives. It should therefore be both easier to implement and to analyze in finite precision.

**Lemma B.1.** *Let  $\mathbf{H} = \begin{pmatrix} \alpha & \mathbf{z}^\top \\ \mathbf{z} & \mathbf{D} \end{pmatrix}$  be a symmetric arrowhead matrix that satisfies the properties of Lemma A.1 for some parameter  $\tau \in (0, 1)$ , with  $\|\mathbf{H}\| \leq 1$ , and assume that the diagonal elements of  $\mathbf{D}$  are sorted:  $d_2 < d_3 < \dots < d_n$ . Given  $\epsilon \in (0, 1)$ , and assuming an  $(\epsilon, n)$ -FMM implementation as in Proposition 1.1, we can compute approximate eigenvalues  $\hat{\lambda}_i$  such that*

$$\hat{\lambda}_1 < d_2 < \hat{\lambda}_2 < \dots < d_n < \hat{\lambda}_n, \quad \text{and} \quad |\hat{\lambda}_i - \lambda_i| \leq \epsilon,$$

in  $O\left(n \log\left(\frac{1}{\epsilon}\right) \log^\xi\left(\frac{n}{\tau\epsilon}\right)\right)$  arithmetic operations.

*Proof.* **TOPROVE 17** □

## B.3 Approximating the elements of the shaft

As a next step, we use the trick of [48] to approximate the elements of  $\hat{\mathbf{H}}$  from Lemma A.3.

**Lemma B.2.** *Let  $\mathbf{H} = \begin{pmatrix} \alpha & \mathbf{z}^\top \\ \mathbf{z} & \mathbf{D} \end{pmatrix}$  be a symmetric arrowhead matrix with  $\|\mathbf{H}\| \leq 1$ , that satisfies the requirements of Lemma A.1 with parameter  $\tau$ , and let  $d_2 < \dots < d_n$  be the diagonal elements of  $\mathbf{D}$ . From Lemma A.3, if we are given a set of approximate eigenvalues  $\hat{\lambda}_i$  that satisfy*

$$\hat{\lambda}_1 < d_2 < \hat{\lambda}_2 < \dots < d_n < \hat{\lambda}_n, \quad \text{and} \quad |\hat{\lambda}_i - \lambda_i| \leq \epsilon \frac{\tau^3}{2(n+1)},$$

for some  $\epsilon \in (0, 1/n)$ , then there exists a matrix  $\hat{\mathbf{H}} = \begin{pmatrix} \hat{\alpha} & \hat{\mathbf{z}}^\top \\ \hat{\mathbf{z}} & \mathbf{D} \end{pmatrix}$ , such that  $\hat{\lambda}_i$  are the exact eigenvalues of  $\hat{\mathbf{H}}$  and  $\|\mathbf{H} - \hat{\mathbf{H}}\| \leq \frac{\epsilon\tau^3}{2} + \frac{n\epsilon}{1-n\epsilon}$ . Assuming an  $(\epsilon, n)$ -FMM, for any  $\epsilon_z \in (0, 1/2)$ , we can compute an approximate vector  $\hat{\mathbf{z}}'$  such that  $|\hat{\mathbf{z}}_i - \hat{\mathbf{z}}'_i| \leq \epsilon_z |\hat{\mathbf{z}}_i|$  and  $\|\hat{\mathbf{z}}' - \hat{\mathbf{z}}\| \leq \epsilon_z (1 + \frac{n\epsilon}{1-n\epsilon})$  in  $O\left(n \log\left(\frac{1}{\epsilon_z}\right) \log^\xi\left(\frac{n}{\tau\epsilon_z}\right)\right)$  arithmetic operations. The matrix  $\hat{\mathbf{H}}' = \begin{pmatrix} \hat{\alpha} & \hat{\mathbf{z}}'^\top \\ \hat{\mathbf{z}}' & \mathbf{D} \end{pmatrix}$  satisfies

$$\|\mathbf{H} - \hat{\mathbf{H}}'\| \leq \frac{\epsilon\tau^3}{2} + \epsilon_z + (1 + \epsilon_z) \frac{n\epsilon}{1-n\epsilon}.$$

*Proof.* **TOPROVE 18** □

## B.4 Approximating inner products with the eigenvectors

In the last part of the analysis we provide bounds on the errors for inner products of the form  $\hat{\mathbf{u}}_i^\top \mathbf{q}$ , between the  $i$ -th eigenvector and some arbitrary vector  $\mathbf{q}$ . Following the procedure of Section 5 in [48], we write

$$\hat{\mathbf{u}}_i^\top \mathbf{q} = \frac{-\mathbf{q}_1 + \Phi(\hat{\lambda}_i)}{\sqrt{1 + \Psi(\hat{\lambda}_i)}},$$

where  $\Phi(\lambda) = \sum_{k=2}^n \frac{\widehat{z}_k \mathbf{q}_k}{d_k - \lambda}$  and  $\Psi(\lambda) = \sum_{k=2}^n \frac{\widehat{z}_k^2}{(d_k - \lambda)^2}$ . To compute all the products for all  $i$  with  $\mathbf{q}$ , we can use FMM to approximate  $\Phi(\lambda)$  and  $\Psi(\lambda)$  on  $n$  points.

**Lemma B.3.** Let  $\mathbf{H} = \begin{pmatrix} \alpha & \mathbf{z}^\top \\ \mathbf{z} & \mathbf{D} \end{pmatrix}$  be a symmetric arrowhead matrix with  $\|\mathbf{H}\| \leq 1$ , that satisfies the requirements of Lemma A.1 with parameter  $\tau$ , and let  $d_2 < d_3 < \dots < d_n$  be the diagonal elements of  $\mathbf{D}$ . Let  $\epsilon_z, \widehat{\lambda}_i, \widehat{\mathbf{z}}, \widehat{\mathbf{z}}', \widehat{\mathbf{H}}$ , and  $\widehat{\mathbf{H}}'$  be the same as in Lemma B.2. Let  $\mathbf{q}$  be a fixed vector with  $\|\mathbf{q}\| \leq 1$ , and  $\widehat{\mathbf{u}}_i, i \in [n]$ , be the eigenvectors of  $\widehat{\mathbf{H}}$ . Assuming an  $(\epsilon, n)$ -FMM, we can approximate all the inner products  $\widehat{\mathbf{u}}_i^\top \mathbf{q}$ , by some values  $x_i$  such that

$$|\widehat{\mathbf{u}}_i^\top \mathbf{q} - x_i| \leq 147\epsilon_z \frac{(n+1)^2}{\tau^6},$$

for all  $i \in [n]$  simultaneously, in  $O\left(n \log\left(\frac{1}{\epsilon_z}\right) \log^\xi\left(\frac{n}{\tau \epsilon_z}\right)\right)$  arithmetic operations.

Proof. [TOPROVE 19](#) □

## B.5 Proof of Theorem 2.1

We can finally combine all the results to prove Theorem 2.1, which we restate below for readability.

**Theorem B.1.** Given a symmetric arrowhead matrix  $\mathbf{H} \in \mathbb{R}^{n \times n}$  as in Eq. (3), with  $\|\mathbf{H}\| \leq 1$ , an accuracy parameter  $\epsilon \in (0, 1)$ , a matrix  $\mathbf{B}$  with  $r$  columns  $\mathbf{B}_i, i \in [r]$ , where  $\|\mathbf{B}_i\| \leq 1$ , and an  $(\epsilon, n)$ -FMM implementation (see Prop. 1.1), we can compute a diagonal matrix  $\widetilde{\Lambda}$ , and a matrix  $\widetilde{\mathbf{Q}}_{\mathbf{B}}$ , such that

$$\left\| \mathbf{H} - \mathbf{Q} \widetilde{\Lambda} \mathbf{Q}^\top \right\| \leq \epsilon, \quad \left| \left( \mathbf{Q}^\top \mathbf{B} - \widetilde{\mathbf{Q}}_{\mathbf{B}} \right)_{i,j} \right| \leq \epsilon/n^2,$$

where  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  is (exactly) orthogonal, in  $O\left(nr \log^{\xi+1}\left(\frac{n}{\epsilon}\right)\right)$  arithmetic operations and comparisons.

Alternatively, if only want to compute a set of approximate values  $\widetilde{\lambda}_1, \dots, \widetilde{\lambda}_n$ , such that  $|\lambda_i(\mathbf{H}) - \widetilde{\lambda}_i| \leq \epsilon$ , the complexity reduces to  $O\left(n \log\left(\frac{1}{\epsilon}\right) \log^\xi\left(\frac{n}{\epsilon}\right)\right)$  arithmetic operations.

Proof. [TOPROVE 20](#) □

## C Tridiagonal diagonalization

### C.1 Omitted proofs

The next lemma bounds the error of the reduction to arrowhead form when the spectral factorizations of the matrices  $\mathbf{T}_1$  and  $\mathbf{T}_2$  in Equation (2) are approximate rather than exact.

**Lemma C.1** (Restatement of Lemma 2.1). Let  $\epsilon \in (0, 1/2)$  be a given accuracy parameter and  $\mathbf{T} = \begin{pmatrix} \mathbf{T}_1 & \beta_{k+1} \mathbf{e}_k & \\ \beta_{k+1} \mathbf{e}_k^\top & \alpha_{k+1} & \beta_{k+2} \mathbf{e}_1^\top \\ & \beta_{k+2} \mathbf{e}_1 & \mathbf{T}_2 \end{pmatrix}$

be a tridiagonal matrix with size  $n \geq 3$  and  $\|\mathbf{T}\| \leq 1$ , where  $\mathbf{T}_1 = \mathbf{U}_1 \mathbf{D}_1 \mathbf{U}_1^\top$  and  $\mathbf{T}_2 = \mathbf{U}_2 \mathbf{D}_2 \mathbf{U}_2^\top$  be the exact spectral factorizations of  $\mathbf{T}_1$  and  $\mathbf{T}_2$ . Let  $\widetilde{\mathbf{U}}_1, \widetilde{\mathbf{D}}_1, \widetilde{\mathbf{U}}_2, \widetilde{\mathbf{D}}_2$  be approximate spectral factorizations of  $\mathbf{T}_1, \mathbf{T}_2$ . If these factors satisfy

$$\left\| \mathbf{T}_{\{1,2\}} - \widetilde{\mathbf{U}}_{\{1,2\}} \widetilde{\mathbf{D}}_{\{1,2\}} \widetilde{\mathbf{U}}_{\{1,2\}}^\top \right\| \leq \epsilon_1, \quad \left\| \widetilde{\mathbf{U}}_{\{1,2\}} \widetilde{\mathbf{U}}_{\{1,2\}}^\top - \mathbf{I} \right\| \leq \epsilon_1/n,$$

for some  $\epsilon_1 \in (0, 1/2)$ , where  $\widetilde{\mathbf{D}}_{\{1,2\}}$  are both diagonal, then, assuming an  $(\epsilon, n)$ -FMM implementation as in Prop. 1.1, we can compute a diagonal matrix  $\widetilde{\mathbf{D}}$  and an approximately orthogonal matrix  $\widetilde{\mathbf{U}}$  such that

$$\left\| \widetilde{\mathbf{U}}^\top \widetilde{\mathbf{U}} - \mathbf{I} \right\| \leq 3(\epsilon_1 + \epsilon)/n, \quad \text{and} \quad \left\| \mathbf{T} - \widetilde{\mathbf{U}} \widetilde{\mathbf{D}} \widetilde{\mathbf{U}}^\top \right\| \leq 2\epsilon_1 + 7\epsilon,$$

in a total of  $O\left(n^2 \log^{\xi+1}\left(\frac{n}{\epsilon}\right)\right)$  arithmetic operations and comparisons.

Proof. **TOPROVE 21** □

**Lemma C.2.** Let  $\epsilon \in (0, 1/2)$  be a given accuracy parameter and  $\mathbf{T} = \begin{pmatrix} \mathbf{T}_1 & \beta_{k+1}\mathbf{e}_k & \\ \beta_{k+1}\mathbf{e}_k^\top & \alpha_{k+1} & \beta_{k+2}\mathbf{e}_1^\top \\ & \beta_{k+2}\mathbf{e}_1 & \mathbf{T}_2 \end{pmatrix}$  be a tridiagonal matrix with size  $n \geq 2$ , where  $\mathbf{T}_1 = \mathbf{U}_1\mathbf{D}_1\mathbf{U}_1^\top$  and  $\mathbf{T}_2 = \mathbf{U}_2\mathbf{D}_2\mathbf{U}_2^\top$  be the exact spectral factorizations of  $\mathbf{T}_1$  and  $\mathbf{T}_2$ . Let  $\tilde{\mathbf{U}}_1, \tilde{\mathbf{D}}_1, \tilde{\mathbf{U}}_2, \tilde{\mathbf{D}}_2$  be approximate spectral factorizations of  $\mathbf{T}_1, \mathbf{T}_2$ . Assume that these factors satisfy

$$\left\| \mathbf{T}_{\{1,2\}} - \tilde{\mathbf{U}}_{\{1,2\}} \tilde{\mathbf{D}}_{\{1,2\}} \tilde{\mathbf{U}}_{\{1,2\}}^\top \right\| \leq \epsilon_1, \quad \left\| \tilde{\mathbf{U}}_{\{1,2\}} \tilde{\mathbf{U}}_{\{1,2\}}^\top - \mathbf{I} \right\| \leq \epsilon_1/n,$$

for some  $\epsilon_1 \in (0, 1/2)$ , where  $\tilde{\mathbf{D}}_{\{1,2\}}$  are both diagonal. Assume also that  $\tilde{\mathbf{D}}_1, \tilde{\mathbf{D}}_2$ , as well as the last row  $\tilde{\mathbf{l}}_1^\top$  of  $\tilde{\mathbf{U}}_1$ , and the first row  $\tilde{\mathbf{f}}_2^\top$  of  $\tilde{\mathbf{U}}_2$ , are explicitly available.

Then we can compute a diagonal matrix  $\tilde{\mathbf{D}}$ , as well as the first row  $\tilde{\mathbf{l}}$  and/or the last row  $\tilde{\mathbf{f}}$  of an approximately orthogonal matrix  $\tilde{\mathbf{U}}$ , which satisfy

$$\left\| \tilde{\mathbf{U}}^\top \tilde{\mathbf{U}} - \mathbf{I} \right\| \leq 3(\epsilon_1 + \epsilon)/n, \quad \text{and} \quad \left\| \mathbf{T} - \tilde{\mathbf{U}} \tilde{\mathbf{D}} \tilde{\mathbf{U}}^\top \right\| \leq 2\epsilon_1 + 7\epsilon,$$

in a total of  $O\left(n \log^{\xi+1}\left(\frac{n}{\epsilon}\right)\right)$  arithmetic operations.

Proof. **TOPROVE 22** □

## C.2 Approximating only the eigenvalues of a tridiagonal matrix

The following result gives the complexity of approximating only the eigenvalues of  $\mathbf{T}$ , instead of a full diagonalization. The main observation is that Lemma 2.1 can be simplified if we only need the eigenvalues. This simplification is listed in C.2 in the Appendix. Using this as an inductive step, we obtain the following Corollary.

**Corollary C.1.** Let  $\mathbf{T}$  be a symmetric unreduced tridiagonal matrix with  $\|\mathbf{T}\| \leq 1$  and  $\epsilon \in (0, 1/2)$  be a given accuracy parameter. Assuming access to an  $(\epsilon, n)$ -FMM, for  $\tau \in \Theta(\text{poly}(\frac{\epsilon}{n}))$ , we can compute approximate eigenvalues  $\tilde{\lambda}_1, \dots, \tilde{\lambda}_n$  such that  $|\tilde{\lambda}_i - \lambda_i(\mathbf{T})| \leq \epsilon$  in  $O\left(n \log^{\xi+2}\left(\frac{n}{\epsilon}\right)\right)$  arithmetic operations.

Proof. **TOPROVE 23** □

## C.3 Application: Hermitian diagonalization

Theorem 1.1 has a direct application to diagonalize Hermitian matrices. The following result is immediate.

**Corollary C.2.** Let  $\mathbf{A}$  be a Hermitian matrix of size  $n$  with  $\|\mathbf{A}\| \leq 1$ . Given accuracy  $\epsilon \in (0, 1/2)$ , and an  $(\epsilon, n)$ -FMM implementation of Prop. 1.1, we can compute a matrix  $\tilde{\mathbf{Q}}$  and a diagonal matrix  $\tilde{\mathbf{\Lambda}}$  such that

$$\left\| \mathbf{A} - \tilde{\mathbf{Q}} \tilde{\mathbf{\Lambda}} \tilde{\mathbf{Q}}^* \right\| \leq \epsilon, \quad \left\| \tilde{\mathbf{Q}}^* \tilde{\mathbf{Q}} - \mathbf{I} \right\| \leq \epsilon/n^2,$$

or, stated alternatively,

$$\tilde{\mathbf{Q}} = \mathbf{Q} + \mathbf{E}_Q, \quad \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}, \quad \|\mathbf{E}_Q\| \leq \epsilon/n^2, \quad \left\| \mathbf{A} - \mathbf{Q} \tilde{\mathbf{\Lambda}} \mathbf{Q}^\top \right\| \leq \epsilon.$$

The algorithm requires a total of  $O\left(n^\omega \log(n) + n^2 \log^{\xi+1}\left(\frac{n}{\epsilon}\right)\right)$  arithmetic operations and comparisons.

Proof. **TOPROVE 24** □

## C.4 Singular Value Decomposition

In this section we provide the proof of Theorem 1.2.

*Proof.* **TOPROVE 25** □

## D Floating point arithmetic

The sign  $s$  is + if the corresponding bit is one, and  $-$  if the bit is zero. The exponent  $e$  is stored as a binary number in the so-called *biased form*, and its range is  $e \in [-M, M]$ , where  $M = 2^{p-1}$ . The significand  $m$  is an integer that satisfies  $2^{t-1} \leq m \leq 2^t - 1$ , where the lower bound is enforced to ensure that the system is *normalized*, i.e. the first bit of  $m$  is always 1. We can therefore write  $\mathbf{fl}(\alpha)$  in a more intuitive representation

$$\mathbf{fl}(\alpha) = \pm 2^e \times \left( \frac{m_1}{2} + \frac{m_2}{2^2} + \dots + \frac{m_t}{2^t} \right),$$

where the first bit  $m_1$  of  $m$  is always equal to one for normalized numbers. The range of normalized numbers is therefore  $[2^{-M}, 2^M(2 - 2^{-t})]$ . Numbers that are smaller than  $2^{-M}$  are called *subnormal* and they will be ignored for simplicity, since we can either add more bits in the exponent. Similarly, numbers that are larger than  $2^M(2 - 2^{-t})$  are assumed to be numerically equal to infinity, denoted by INF.

From [52, Theorem 2.2], for all real numbers  $\alpha$  in the normalized range it holds that

$$\mathbf{fl}(\alpha) = (1 + \theta)\alpha,$$

where  $\theta \in \mathbb{R}$  satisfies  $|\theta| \leq 2^{-t} := \mathbf{u}$ , where  $\mathbf{u}$  is the *machine precision*. Clearly,  $t = O(\log(1/\mathbf{u}))$ , in which case we can always obtain a bound for the number of required bits of a numerical algorithm if we have an upper bound for the precision  $\mathbf{u}$ . We will write the same for complex numbers which are represented as a pair of normalized floating point numbers.

The floating point implementation of each arithmetic operation  $\odot \in \{+, -, \times, /\}$  also satisfies

$$\mathbf{fl}(\alpha \odot \beta) = (1 + \theta)(\alpha \odot \beta), \quad |\theta| \leq \mathbf{u}. \quad (16)$$

Divisions and multiplications with 1 and 2 do not introduce errors (for the latter we simply increase/decrease the exponent). We assume that we also have an implementation of  $\sqrt{\cdot}$  such that  $\mathbf{fl}(\sqrt{\alpha}) = (1 + \theta)\sqrt{\alpha}$  where  $|\theta| \leq \mathbf{u}$ . From [52, Lemma 3.1], we can bound products of errors as

$$\prod_{i=1}^n (1 + \theta_i)^{\rho_i} = 1 + \eta_n,$$

where  $\rho_i = \pm 1$  and  $|\eta_n| \leq \frac{n\mathbf{u}}{1 - n\mathbf{u}}$ .

The above can be extended also for complex arithmetic (see [52, Lemma 3.5]), where the bound becomes  $|\theta| \leq O(\mathbf{u})$ , but we will ignore the constant prefactor for simplicity.

Operations on matrices can be analyzed in a similar manner. Let  $\otimes$  denote the element-wise multiplication between two matrices and  $\oslash$  the element-wise division. The floating point representation of a matrix  $\mathbf{A}$  satisfies

$$\mathbf{fl}(\mathbf{A}) = \mathbf{A} + \Delta \otimes \mathbf{A}, \quad |\Delta_{i,j}| \leq \mathbf{u}.$$

It can be shown that  $\|\Delta\| \leq \mathbf{u}\sqrt{n}\|\mathbf{A}\|$ . For any operation  $\odot \in \{+, -, \otimes, \oslash\}$  and matrices  $\mathbf{A}$  and  $\mathbf{B}$  it holds that

$$\mathbf{fl}(\mathbf{A} \odot \mathbf{B}) = \mathbf{A} \odot \mathbf{B} + \Delta \otimes (\mathbf{A} \odot \mathbf{B}), \quad |\Delta_{i,j}| \leq \mathbf{u}, \quad \|\Delta \otimes (\mathbf{A} \odot \mathbf{B})\| \leq \mathbf{u}\sqrt{n}\|\mathbf{A} \odot \mathbf{B}\|. \quad (17)$$

## E Reduction to tridiagonal form - omitted proofs and definitions

### E.1 Imported subroutines

In this appendix we mention some preliminary results that we use in the analysis.

**Theorem E.1** (MM, stable fast matrix multiplication [29, 28]). *For every  $\eta > 0$ , there exists a fast matrix multiplication algorithm  $\text{MM}$  which takes as input two matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$  and returns  $\mathbf{C} \leftarrow \text{MM}(\mathbf{A}, \mathbf{B})$  such that*

$$\|\mathbf{C} - \mathbf{AB}\| \leq n^{c_\eta} \cdot \mathbf{u} \|\mathbf{A}\| \|\mathbf{B}\|,$$

*on floating point machine with precision  $\mathbf{u}$ , for some constant  $c_\eta$  independent of  $n$ . It requires  $O(n^{\omega+\eta})$  floating point operations.*

We can also assume that if the result  $\mathbf{AB}$  is Hermitian, then  $\text{MM}(\mathbf{A}, \mathbf{B})$  will be Hermitian as well (see e.g. [83]).

**Theorem E.2** (Cf. [28]). *Given a matrix  $\mathbf{A} \in \mathbb{C}^{m \times n}$ ,  $m \geq n$ , there exists an algorithm  $[\mathbf{Q}, \mathbf{R}] \leftarrow \text{QR}(\mathbf{A})$  which returns an upper triangular matrix  $\mathbf{R} \in \mathbb{C}^{m \times n}$  and a matrix  $\mathbf{Q} \in \mathbb{C}^{m \times m}$ , such that:*

$$(\mathbf{Q} + \mathbf{E}_\mathbf{Q})^* (\mathbf{A} + \mathbf{E}_\mathbf{A}) = \mathbf{R}, \quad \|\mathbf{E}_\mathbf{Q}\| \leq n^{c_{\text{QR}}} \mathbf{u}, \quad \|\mathbf{E}_\mathbf{A}\| \leq n^{c_{\text{QR}}} \mathbf{u} \|\mathbf{A}\|,$$

*for some constant  $c_{\text{QR}}$ , where the matrix  $\mathbf{Q} + \mathbf{E}_\mathbf{Q}$  is unitary. The algorithm executes  $O(mn^{\omega-1})$  floating point operations on a machine with precision  $\mathbf{u}$ .*

In our analysis it will be useful the following re-statement of the aforementioned results.

**Corollary E.1.** *There exists a global constant  $\beta \geq \max\{c_{\text{QR}}, c_\eta, 1\}$ , such that the matrices from Theorems E.1 and E.2 satisfy the following properties*

$$\begin{aligned} \|\text{MM}(\mathbf{A}, \mathbf{B}) - \mathbf{AB}\| &\leq n^\beta \mathbf{u} \|\mathbf{A}\| \|\mathbf{B}\|, \\ \|\mathbf{A} - \mathbf{QR}\| &\leq n^\beta \mathbf{u} \|\mathbf{A}\|, \\ \|\mathbf{R}\| &\leq (1 + n^\beta \mathbf{u}) \|\mathbf{A}\|, \\ \|\mathbf{Q}\| &\leq \sqrt{1 + n^\beta \mathbf{u}}, \\ \|\mathbf{Q}^{-1}\| &\leq \frac{1}{\sqrt{1 - n^\beta \mathbf{u}}}, \\ \max\{\|\mathbf{I} - \mathbf{QQ}^*\|, \|\mathbf{I} - \mathbf{Q}^* \mathbf{Q}\|\} &\leq n^\beta \mathbf{u}. \end{aligned}$$

*Proof.* **TOPROVE 26** □

Using square fast matrix multiplication we can obtain an algorithm to compute two banded matrices with the same bandwidth

**Corollary E.2.** *Let  $\mathbf{A}, \mathbf{B}$  in  $\mathbb{C}^{n \times n}$ , where, without loss of generality,  $n$  is a power of two. Assume that  $\mathbf{A}$  and  $\mathbf{B}$  are block-tridiagonal, with block size  $n_k = 2^k$  for some  $k \in \{0, 1, \dots, \log(n) - 2\}$ . We can compute a matrix  $\mathbf{C}'$  such that  $\|\mathbf{C}' - \mathbf{AB}\| \leq \mathbf{u} n_k^\beta \|\mathbf{A}\| \|\mathbf{B}\|$  in  $O(nn_k^{\omega-1})$  floating point operations.*

*Proof.* **TOPROVE 27** □

### E.2 Rotations

An example of the rotations is illustrated in Equations (18) and (19). In Equation (18), the matrices  $\mathbf{A}'_{1,2}$  and  $\mathbf{A}'_{2,1}$  are lower and upper triangular, respectively. In Equation (19), the matrices  $\mathbf{A}'_{4,2}$  and  $\mathbf{A}'_{5,3}$  are upper triangular, while  $\mathbf{A}'_{2,4}$



and  $\mathbf{A}'_{3,5}$  are lower triangular.

$$\begin{array}{c} \mathbf{A}^{(k,2,0)}, i=2 \\ \left( \begin{array}{cccccccc} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} & \boxed{\mathbf{A}_{1,3}} & 0 & 0 & 0 & 0 & 0 \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} & \mathbf{A}_{2,3} & \mathbf{A}_{2,4} & 0 & 0 & 0 & 0 \\ \boxed{\mathbf{A}_{3,1}} & \mathbf{A}_{3,2} & \mathbf{A}_{3,3} & \mathbf{A}_{3,4} & \mathbf{A}_{3,5} & 0 & 0 & 0 \\ 0 & \mathbf{A}_{4,2} & \mathbf{A}_{4,3} & \mathbf{A}_{4,4} & \mathbf{A}_{4,5} & \mathbf{A}_{4,6} & 0 & 0 \\ 0 & 0 & \mathbf{A}_{5,3} & \mathbf{A}_{5,4} & \mathbf{A}_{5,5} & \mathbf{A}_{5,6} & \mathbf{A}_{5,7} & 0 \\ 0 & 0 & 0 & \mathbf{A}_{6,4} & \mathbf{A}_{6,5} & \mathbf{A}_{6,6} & \mathbf{A}_{6,7} & \mathbf{A}_{6,8} \\ 0 & 0 & 0 & 0 & \mathbf{A}_{7,5} & \mathbf{A}_{7,6} & \mathbf{A}_{7,7} & \mathbf{A}_{7,8} \\ 0 & 0 & 0 & 0 & 0 & \mathbf{A}_{8,6} & \mathbf{A}_{8,7} & \mathbf{A}_{8,8} \end{array} \right) \xrightarrow{R_2} \left( \begin{array}{cccccccc} \mathbf{A}_{1,1} & \mathbf{A}'_{1,2} & 0 & 0 & 0 & 0 & 0 & 0 \\ \mathbf{A}'_{2,1} & \mathbf{A}_{2,2} & \mathbf{A}_{2,3} & \mathbf{A}_{2,4} & \boxed{\mathbf{A}_{2,5}} & 0 & 0 & 0 \\ 0 & \mathbf{A}_{3,2} & \mathbf{A}_{3,3} & \mathbf{A}_{3,4} & \mathbf{A}_{3,5} & 0 & 0 & 0 \\ 0 & \mathbf{A}_{4,2} & \mathbf{A}_{4,3} & \mathbf{A}_{4,4} & \mathbf{A}_{4,5} & \mathbf{A}_{4,6} & 0 & 0 \\ 0 & \boxed{\mathbf{A}_{5,2}} & \mathbf{A}_{5,3} & \mathbf{A}_{5,4} & \mathbf{A}_{5,5} & \mathbf{A}_{5,6} & \mathbf{A}_{5,7} & 0 \\ 0 & 0 & 0 & \mathbf{A}_{6,4} & \mathbf{A}_{6,5} & \mathbf{A}_{6,6} & \mathbf{A}_{6,7} & \mathbf{A}_{6,8} \\ 0 & 0 & 0 & 0 & \mathbf{A}_{7,5} & \mathbf{A}_{7,6} & \mathbf{A}_{7,7} & \mathbf{A}_{7,8} \\ 0 & 0 & 0 & 0 & 0 & \mathbf{A}_{8,6} & \mathbf{A}_{8,7} & \mathbf{A}_{8,8} \end{array} \right) \end{array} \quad (18)$$

$$\begin{array}{c} \mathbf{A}^{(k,3,5)} \\ \left( \begin{array}{cccccccc} \mathbf{A}_{1,1} & \mathbf{A}'_{1,2} & 0 & 0 & 0 & 0 & 0 & 0 \\ \mathbf{A}'_{2,1} & \mathbf{A}_{2,2} & \mathbf{A}_{2,3} & \mathbf{A}_{2,4} & \boxed{\mathbf{A}_{2,5}} & 0 & 0 & 0 \\ 0 & \mathbf{A}_{3,2} & \mathbf{A}_{3,3} & \mathbf{A}_{3,4} & \mathbf{A}_{3,5} & 0 & 0 & 0 \\ 0 & \mathbf{A}_{4,2} & \mathbf{A}_{4,3} & \mathbf{A}_{4,4} & \mathbf{A}_{4,5} & \mathbf{A}_{4,6} & 0 & 0 \\ 0 & \boxed{\mathbf{A}_{5,2}} & \mathbf{A}_{5,3} & \mathbf{A}_{5,4} & \mathbf{A}_{5,5} & \mathbf{A}_{5,6} & \mathbf{A}_{5,7} & 0 \\ 0 & 0 & 0 & \mathbf{A}_{6,4} & \mathbf{A}_{6,5} & \mathbf{A}_{6,6} & \mathbf{A}_{6,7} & \mathbf{A}_{6,8} \\ 0 & 0 & 0 & 0 & \mathbf{A}_{7,5} & \mathbf{A}_{7,6} & \mathbf{A}_{7,7} & \mathbf{A}_{7,8} \\ 0 & 0 & 0 & 0 & 0 & \mathbf{A}_{8,6} & \mathbf{A}_{8,7} & \mathbf{A}_{8,8} \end{array} \right) \xrightarrow{R'_4} \left( \begin{array}{cccccccc} \mathbf{A}_{1,1} & \mathbf{A}'_{1,2} & 0 & 0 & 0 & 0 & 0 & 0 \\ \mathbf{A}'_{2,1} & \mathbf{A}_{2,2} & \mathbf{A}_{2,3} & \mathbf{A}'_{2,4} & 0 & 0 & 0 & 0 \\ 0 & \mathbf{A}_{3,2} & \mathbf{A}_{3,3} & \mathbf{A}_{3,4} & \mathbf{A}'_{3,5} & 0 & 0 & 0 \\ 0 & \mathbf{A}'_{4,2} & \mathbf{A}_{4,3} & \mathbf{A}_{4,4} & \mathbf{A}_{4,5} & \mathbf{A}_{4,6} & \boxed{\mathbf{A}_{4,7}} & 0 \\ 0 & 0 & \mathbf{A}'_{5,3} & \mathbf{A}_{5,4} & \mathbf{A}_{5,5} & \mathbf{A}_{5,6} & \mathbf{A}_{5,7} & 0 \\ 0 & 0 & 0 & \mathbf{A}_{6,4} & \mathbf{A}_{6,5} & \mathbf{A}_{6,6} & \mathbf{A}_{6,7} & \mathbf{A}_{6,8} \\ 0 & 0 & 0 & \boxed{\mathbf{A}_{4,7}} & \mathbf{A}_{7,5} & \mathbf{A}_{7,6} & \mathbf{A}_{7,7} & \mathbf{A}_{7,8} \\ 0 & 0 & 0 & 0 & 0 & \mathbf{A}_{8,6} & \mathbf{A}_{8,7} & \mathbf{A}_{8,8} \end{array} \right) \end{array} \quad (19)$$

The rotations can be conveniently implemented in floating point using fast QR factorizations.

**Lemma E.1.** Let  $\mathbf{A}^{(k,i,0)}$  be a block-pentadiagonal matrix with no bulge of the form of Equation (4), consisting of  $b_k \times b_k$  blocks of size  $n_k \times n_k = 2^k \times 2^k$  each. The rotation  $[\mathbf{A}^{(k,i+1,i+3)}, \mathbf{Q}] \leftarrow R_i(\mathbf{A}^{(k,i,0)})$  can be implemented in  $O(n_k^\omega)$  floating point operations. The resulting matrix  $\mathbf{A}^{(k,i+1,i+3)}$  has a bulge at positions  $(i, i+3)$  and  $(i+3, i)$ , and  $\mathbf{Q}$  is approximately unitary. Moreover, if  $\mathbf{u} \leq 1/n^\beta$ , then

$$\left\| \mathbf{A}^{(k,i,0)} - \mathbf{Q} \mathbf{A}^{(k,i+1,i+3)} \mathbf{Q}^* \right\| \leq C \mathbf{u} n^\beta \|\mathbf{A}^{(k,i,0)}\|,$$

where  $C$  is a constant independent of  $n$ .

*Proof.* TOPROVE 28 □

**Lemma E.2.** Let  $\mathbf{A}^{(k,s,j+1)}$ ,  $j \geq s+1$ , be a block-pentadiagonal matrix of the form of Equation (4), with a bulge at position  $(j+1, j-2)$ , consisting of  $b_k \times b_k$  blocks of size  $n_k \times n_k = 2^k \times 2^k$  each. The rotation  $[\mathbf{A}^{(k,s,j+3)}, \mathbf{Q}] \leftarrow R'_j(\mathbf{A}^{(k,s,j+1)})$  can be implemented in  $O(n_k^\omega)$  floating point operations. The resulting matrix  $\mathbf{A}^{(k,s,j+3)}$  has a bulge at positions  $(j+3, j)$  and  $(j, j+3)$ , and  $\mathbf{Q}$  is approximately unitary. Moreover, if  $\mathbf{u} \leq 1/n^\beta$ , then

$$\left\| \mathbf{A}^{(k,s,j+1)} - \mathbf{Q} \mathbf{A}^{(k,s,j+3)} \mathbf{Q}^* \right\| \leq C' \mathbf{u} (2n_k)^\beta \|\mathbf{A}^{(k,s,j+1)}\|,$$

where  $C'$  is a constant independent of  $n$ .

*Proof.* TOPROVE 29 □

### E.3 Bandwidth halving

**Lemma E.3.** Let  $\mathbf{A}^{(k,2,0)}$  be a full block-pentadiagonal matrix of the form of Equation (4), with no bulge, consisting of  $b_k \times b_k$  blocks of size  $n_k \times n_k = 2^k \times 2^k$  each. For any  $k \in [\log(n) - 2]$ , if  $\mathbf{u} \leq \frac{1}{C_1 n^{\beta+3}}$ , where  $C_1$  is a constant, then Algorithm 2 returns a matrix  $\widehat{\mathbf{Q}}^{(k)}$  that is approximately orthogonal, and a matrix  $\mathbf{A}^{(k-1,2,0)}$  such that

$$\begin{aligned} \left\| \mathbf{A}^{(k,2,0)} - \widehat{\mathbf{Q}}^{(k)} \mathbf{A}^{(k-1,2,0)} \widehat{\mathbf{Q}}^{(k)*} \right\| &\leq C_2 \mathbf{u} n^{\beta+3} \|\mathbf{A}^{(k,2,0)}\|, \\ \left\| \widehat{\mathbf{Q}}^{(k)} \widehat{\mathbf{Q}}^{(k)*} - \mathbf{I} \right\| &\leq C_3 \mathbf{u} n^{\beta+3} \end{aligned}$$

where  $C_2, C_3$  are also constants independent of  $n$ . It requires  $O(n^2 (S_\omega(\log(n) - 1) - S_\omega(k)))$  floating point operations, where  $S_x(m) := \sum_{l=1}^m (2^{x-2})^l$ . If  $\widehat{\mathbf{Q}}^{(k)}$  is not required, the complexity reduces to  $O(n^2 n_k^{\omega-2}) = O(n^2 (2^{\omega-2})^k)$ .

*Proof.* **TOPROVE 30** □

**Theorem E.3** (Restatement of Theorem 3.1). There exists a floating point implementation of the tridiagonal reduction algorithm of [79], which takes as input a Hermitian matrix  $\mathbf{A}$ , and returns a tridiagonal matrix  $\widetilde{\mathbf{T}}$ , and (optionally) an approximately unitary matrix  $\widetilde{\mathbf{Q}}$ . If the machine precision  $\mathbf{u}$  satisfies  $\mathbf{u} \leq \epsilon \frac{1}{c n^{\beta+4}}$ , where  $\epsilon \in (0, 1)$ ,  $c$  is a constant, and  $\beta$  is the same as in Corollary E.1, which translates to  $O(\log(n) + \log(1/\epsilon))$  bits of precision, then the following hold:

$$\left\| \widetilde{\mathbf{Q}} \widetilde{\mathbf{Q}}^* - \mathbf{I} \right\| \leq \epsilon, \quad \text{and} \quad \left\| \mathbf{A} - \widetilde{\mathbf{Q}} \widetilde{\mathbf{T}} \widetilde{\mathbf{Q}}^* \right\| \leq \epsilon \|\mathbf{A}\|.$$

The algorithm executes at most  $O(n^2 S_\omega(\log(n)))$  floating point operations to return only  $\widetilde{\mathbf{T}}$ , where  $S_x(m) = \sum_{l=1}^m (2^{x-2})^l$ . If  $\mathbf{A}$  is banded with  $1 \leq d \leq n$  bands, the floating point operations reduce to  $O(n^2 S_\omega(\log(d)))$ . If  $\widetilde{\mathbf{Q}}$  is also returned, the complexity increases to  $O(n^2 C_\omega(\log(n)))$ , where  $C_x(n) := \sum_{k=2}^{\log(n)-2} (S_x(\log(n) - 1) - S_x(k))$ . If  $\omega$  is treated as a constant  $\omega \approx 2.371$  the corresponding complexities are  $O(n^\omega)$ ,  $O(n^2 d^{\omega-2})$ , and  $O(n^\omega \log(n))$ , respectively.

*Proof.* **TOPROVE 31** □

We have the following direct Corollary.

**Corollary E.3.** The eigenvalues of the matrix  $\widetilde{\mathbf{T}}$  returned by the algorithm of Theorem 3.1 satisfy

$$\left| \lambda_i(\widetilde{\mathbf{T}}) - \lambda_i(\mathbf{A}) \right| \leq \epsilon \|\mathbf{A}\|.$$

*Proof.* **TOPROVE 32** □