

A THEORY OF QUANTUM SUBSPACE DIAGONALIZATION *

ETHAN N. EPPERLY[†], LIN LIN[‡], AND YUJI NAKATSUKASA[§]

Abstract. Quantum subspace diagonalization methods are an exciting new class of algorithms for solving large-scale eigenvalue problems using quantum computers. Unfortunately, these methods require the solution of an ill-conditioned generalized eigenvalue problem, with a matrix pair corrupted by a non-negligible amount of noise that is far above the machine precision. Despite pessimistic predictions from classical worst-case perturbation theories, these methods can perform reliably well if the generalized eigenvalue problem is solved using a standard truncation strategy. By leveraging and advancing classical results in matrix perturbation theory, we provide a theoretical analysis of this surprising phenomenon, proving that under certain natural conditions, a quantum subspace diagonalization algorithm can accurately compute the smallest eigenvalue of a large Hermitian matrix. We give numerical experiments demonstrating the effectiveness of the theory and providing practical guidance for the choice of truncation level. Our new results can also be of independent interest to solving eigenvalue problems outside the context of quantum computation.

Key words. quantum subspace diagonalization, quantum linear algebra, generalized eigenvalue problem, matrix perturbation theory

AMS subject classifications. 68Q12, 65F15, 15A22, 15A45

1. Introduction. Quantum computing is a fundamentally new computational paradigm, which has the potential to have a transformative impact on certain areas of computational science [27, 28]. One particularly compelling use case for quantum computers is to solve eigenvalue problems related to quantum many-body systems, for which the dimension of the discretized matrix grows exponentially with respect to the number of particles.

Quantum subspace diagonalization (QSD) methods [6, 13, 15, 20, 21, 26, 31, 32], also known as quantum Krylov methods, are an exciting class of quantum algorithms for solving large-scale Hermitian eigenvalue problems. One common key step of these algorithms is to solve a nearly singular generalized eigenvalue problem, where each entry of the associated matrix pair can be corrupted by Monte Carlo errors many orders of magnitude larger than the round-off error typically seen in classical computation. For such noisy generalized eigenvalue problems, classical perturbation theory fails to explain why such problems could be solved accurately. Despite this, QSD methods appear to work in practice, at least on some examples and with some procedure to compensate for the measurement error. This article is addressed squarely at explaining why, in theory, QSD algorithms perform well and how, in practice, errors in the problem data can be effectively dealt with.

While our analysis is centered around the context of QSD, the underlying problem of solving an ill-conditioned generalized eigenvalue problem with noisy data is a fundamental linear algebra problem, which emerges in multiple application areas such as electronic structure theory [3, 14], control theory [9], and variational Monte Carlo

*This is a preprint version of *A Theory of Quantum Subspace Diagonalization* (<https://doi.org/10.1137/21M145954X>), which appeared in the SIAM Journal on Matrix Analysis and Applications on August 1, 2022.

[†]Department of Computing and Mathematical Sciences, California Institute of Technology, Pasadena, CA, USA (eepperly@caltech.edu).

[‡]Department of Mathematics, and Challenge Institute of Quantum Computation, University of California Berkeley, Berkeley, CA, USA and Computational Research Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA (linlin@math.berkeley.edu).

[§]Mathematical Institute, Oxford University, Oxford, UK (nakatsukasa@maths.ox.ac.uk).

optimization [30]. Our results may be of particular interest in filter diagonalization methods [18, 36] which were a classical antecedent to QSD methods [26] and also require the solution of an ill-conditioned generalized eigenvalue problem.

Notation. Vectors and matrices will be denoted by boldface lower- and upper-case letters respectively. We denote the conjugate transpose by $*$. All matrices and vectors are assumed to be over the complex numbers. The unembellished norm $\|\cdot\|$ shall refer to the Euclidean norm of a vector $\|\mathbf{x}\| = \sqrt{\mathbf{x}^* \mathbf{x}}$ or the spectral norm (largest singular value) of a matrix. At times, we shall also make use of the Frobenius norm $\|\mathbf{B}\|_F := \sqrt{\text{tr}(\mathbf{B}^* \mathbf{B})}$. The absolute values of the generalized eigenvalues of a pair (\mathbf{H}, \mathbf{S}) is denoted $|\Lambda(\mathbf{H}, \mathbf{S})|$. Relations \approx , \lesssim , \ll , etc. are informal, with no precise mathematical relation being implied.

1.1. Quantum subspace diagonalization and its numerical challenges.

We begin by describing a simple QSD algorithm developed in parallel by Parrish and McMahon [26] and Stair, Huang, and Evangelista [32] with a focus on its linear algebraic (Krylov) structure. For readers unfamiliar with quantum computation and why it might be advantageous over classical computing for certain problems, we recommend the classic textbook [24] as well as a recent tutorial aimed at a mathematical audience [22].

Suppose we are interested in the ground-state energy¹ (least eigenvalue) E_0 of a Hamiltonian operator (Hermitian matrix) $\widehat{\mathbf{H}}$. Our goal shall be to compute an approximation \tilde{E}_0 for which the (forward) error $|\tilde{E}_0 - E_0|$ is small. For the problem sizes for which QSD is appropriate, vectors with the dimension of the operator $\widehat{\mathbf{H}}$ are too large to store classically, so we would like to represent them as states in a quantum computer. The QSD method rests on two assumptions:

1. We can efficiently prepare a state φ_0 on the quantum device that has a nontrivial overlap $|\varphi_0^* \psi_0| \gg 0$ to the true ground-state eigenvector ψ_0 .
2. The time evolution $\varphi \mapsto e^{it\widehat{\mathbf{H}}} \varphi$ can be efficiently simulated on the quantum device.

To improve on our initial guess φ_0 , we enlarge to a subspace spanned by vectors

$$(1.1) \quad \varphi_j = e^{it_j \widehat{\mathbf{H}}} \varphi_0 \quad \text{for } j = 0, \dots, n-1$$

where $t_j = j\Delta t$ are a time sequence with step size $\Delta t > 0$. The subspace $\text{span}\{\varphi_j\}$ forms a “unitary Krylov space” and plays a role analogous to the Krylov subspace in the Lanczos method. Eigenvalue estimates for the operator $\widehat{\mathbf{H}}$ can be computed from this unitary Krylov subspace by applying the Rayleigh–Ritz method [25, §11].

In a classical Krylov method, one would usually orthogonalize the basis vectors (1.1). Unfortunately, this orthogonalization operation can be inherently difficult to perform on a quantum computer, so we instead work in the basis (1.1) as computed. In a non-orthogonal basis, the Rayleigh–Ritz eigenvalue estimates are obtained by solving a generalized eigenvalue problem

$$(1.2) \quad \mathbf{H}\mathbf{c} = E\mathbf{S}\mathbf{c}.$$

Here the *projected Hamiltonian* and *overlap* matrices \mathbf{H} and \mathbf{S} are Hermitian–Toeplitz matrices defined as

$$(1.3) \quad \mathbf{H}_{jk} = \varphi_j^* \widehat{\mathbf{H}} \varphi_k, \quad \mathbf{S}_{jk} = \varphi_j^* \varphi_k,$$

¹The ground-state energy, by itself, is a useful quantity for applications in electronic structure, as it determines the energy landscape for dynamical simulation.

\mathbf{c} is the reduced Ritz vector, and E is the Ritz value. Each matrix entry $\mathbf{H}_{jk}, \mathbf{S}_{jk}$ can be estimated via a Monte Carlo sampling procedure such as the Hadamard test on a quantum computer (see e.g., [24, Chapter 5], [6, App. D]). Once a sufficiently good estimate to \mathbf{H}, \mathbf{S} is obtained, the generalized eigenvalue problem (1.2) is solved on a classical computer.

Remark 1.1 (Other QSD methods). A number of QSD methods have been proposed which differ in how the basis states (1.1) are generated and how the eigenvalue estimates are obtained. Alternative methods for basis vector generation (1.1) from one or more initial guesses include: multiplication by creation and annihilation operators [5, 20] and imaginary-time evolution $\varphi_j := e^{i(it_j)\widehat{\mathbf{H}}}\varphi_0$ [21] (with $t_j \geq 0$). Methods also differ in whether the Rayleigh–Ritz procedure is applied to $\widehat{\mathbf{H}}$ itself [26, 32] or the time evolution operator $e^{i\Delta t\widehat{\mathbf{H}}}$ [6, 15]. For concreteness, we shall focus in this article on the QSD method as discussed, though our analysis should have insights for understanding the broader class of QSD algorithms.

There are two main questions in the analysis of the (forward) error of the QSD method:

- (A) How to analyze the Rayleigh–Ritz error due to the use of a finite-dimensional unitary Krylov subspace?
- (B) How to analyze the error of the generalized eigenvalue problem (1.2) in the presence of the Monte Carlo noise for estimating the matrix entries in \mathbf{H}, \mathbf{S} ?

Below we first discuss issues related to question (B), which are particularly challenging from the perspective of numerical linear algebra.

1.2. Numerical Issues with QSD. Numerical results indicate that the size of n needed to obtain desired accuracy can be very modest (e.g. $10 \leq n \leq 100$), so we are free to use any algorithm [11, 25] to solve the dense generalized eigenvalue problem (1.2). However, it is frequently observed that the states $\varphi_0, \dots, \varphi_{n-1}$ are very close to being linearly dependent, leading to the matrices \mathbf{H} and \mathbf{S} being nearly rank-deficient and the generalized eigenvalue problem (1.2) nearly singular. This ill-conditioning is an intrinsic feature of this method since the problem necessarily becomes ill-conditioned if the initial guess φ_0 possesses the desirable property of approximately lying in a low-dimensional invariant subspace.

The near-singularity of the problem (1.2) becomes particularly alarming when taken in conjunction with the fact that the matrix elements (1.3) will be corrupted by several types of error when measured from a quantum computer. Some forms of error, such as discretization error in evaluating the time evolution $e^{it_j\widehat{\mathbf{H}}}\varphi_0$ by a Trotter formula [4] and gate errors, can in principle be systematically controlled on a fault-tolerant quantum device. However, even on a flawless quantum device, the matrix elements (1.3) still need to be computed via sampling, which incurs Monte Carlo-type $\approx \delta^{-2}$ samples to compute each entry to δ -accuracy. We shall refer to all of these errors collectively as “noise”.

Classical perturbation theory [34, §VI.3] as well as modern improvements [19] do not apply when the perturbation is large enough to make the problem (1.2) singular, which is almost always the case for the QSD algorithm because of sampling error and ill-conditioning. Indeed, the following variant of the classical example [39, Eq. (4.10)] shows that a perturbation just a touch larger than the distance to singularity can

send the eigenvalues (originally, 1 and 2) to any pair of numbers $\alpha, \beta \in \mathbb{C}$:

$$(1.4) \quad \mathbf{H} = \begin{bmatrix} 2 & 0 \\ 0 & \epsilon \end{bmatrix}, \mathbf{S} = \begin{bmatrix} 1 & 0 \\ 0 & \epsilon \end{bmatrix} \xrightarrow{\mathcal{O}(\epsilon) \text{ error}} \widetilde{\mathbf{H}} = \begin{bmatrix} 2 & \alpha\epsilon \\ \beta\epsilon & 0 \end{bmatrix}, \widetilde{\mathbf{S}} = \begin{bmatrix} 1 & \epsilon \\ \epsilon & 0 \end{bmatrix}.$$

Even the well-conditioned eigenvalue 2 can be perturbed arbitrarily far if the noise is large enough to make the problem singular!

There is evidence that adversarially chosen perturbations such as (1.4) are pathologically unlikely to occur, with the well-conditioned eigenvalues of a pair (\mathbf{H}, \mathbf{S}) changing only modestly after perturbation. Indeed, Wilkinson showed that “most” $\mathcal{O}(\epsilon)$ perturbations to (1.4) have an eigenvalue near the well-conditioned eigenvalue of 2 [39], and recent analysis by Lotz and Noferini [16] show that some eigenvalues of genuinely singular generalized eigenvalue problems can, in effect, be locally well-conditioned with high probability.

Even if a good approximation to the ground-state energy is among the eigenvalues of $(\widetilde{\mathbf{H}}, \widetilde{\mathbf{S}})$, identifying it can be difficult. When the pair (\mathbf{H}, \mathbf{S}) is nearly singular, the perturbed problem is almost assured to possess spurious eigenvalues. To see why this is the case, observe that if (\mathbf{H}, \mathbf{S}) is nearly singular, there exists an eigenpair (\mathbf{c}, E) such that $\mathbf{H}\mathbf{c}, \mathbf{S}\mathbf{c} \approx \mathbf{0}$. Perturbations in \mathbf{H} and \mathbf{S} create large changes in both the numerator and the denominator of the Rayleigh quotient $E = \mathbf{c}^* \mathbf{H} \mathbf{c} / \mathbf{c}^* \mathbf{S} \mathbf{c}$; since eigenvalues of a pair with positive definite \mathbf{S} extremize the Rayleigh quotient, noise can easily introduce fake eigenvalues much smaller than the genuine least eigenvalue of (\mathbf{H}, \mathbf{S}) . Reliably distinguishing genuine eigenvalues from such fake eigenvalues is challenging; we tried many heuristics and all of them failed for a nontrivial fraction of random initializations of the measurement error (see section SM2).

To address these issues, we shall solve the eigenvalue problem (1.2) using the following truncation scheme: First, compute an eigendecomposition of the matrix \mathbf{S} and discard all eigenvalues smaller than or equal to a threshold $\epsilon > 0$. Then, letting $\mathbf{V}_{>\epsilon}$ denote a matrix whose columns are the non-discarded eigenvectors and $\mathbf{\Lambda}_{>\epsilon} := \mathbf{V}_{>\epsilon}^* \mathbf{S} \mathbf{V}_{>\epsilon}$, we solve the reduced generalized eigenvalue problem

$$\mathbf{V}_{>\epsilon}^* \mathbf{H} \mathbf{V}_{>\epsilon} \mathbf{c} = \widetilde{E} \mathbf{V}_{>\epsilon}^* \mathbf{S} \mathbf{V}_{>\epsilon} \mathbf{c},$$

or equivalently find the eigenvalues of $\mathbf{\Lambda}_{>\epsilon}^{-1/2} \mathbf{V}_{>\epsilon}^* \mathbf{S} \mathbf{V}_{>\epsilon} \mathbf{\Lambda}_{>\epsilon}^{-1/2}$. This procedure appears to have been first discovered in quantum physics by Löwdin in 1967 [17] (also rediscovered in [12]), where it is associated with the name *canonical orthogonalization*. We shall call this procedure *thresholding* and present it in Algorithm 1.1 for convenient reference in the rest of the document. A more careful variant of thresholding from the numerical analysis community was proposed by Fix and Heiberger in 1972 [10], though its authors expressly advise against using it in precisely the setting of the QSD algorithm where the pair (\mathbf{H}, \mathbf{S}) is nearly singular.

Despite appearing quite natural, there are examples where thresholding fails to work and is thus not appropriate for arbitrary Hermitian definite generalized eigenvalue problems. For instance, if one applies thresholding with parameter ϵ to the pair

$$(1.5) \quad \mathbf{H} = \begin{bmatrix} 1 & \epsilon \\ \epsilon & \epsilon^2 \end{bmatrix}, \quad \mathbf{S} = \begin{bmatrix} 1 & 0 \\ 0 & \epsilon^2 \end{bmatrix},$$

one recovers an eigenvalue of 1, which is far from the genuine eigenvalues 0 and 2 of the pair. Even with the existence of bad examples like (1.5), thresholding appears

Algorithm 1.1 Thresholding procedure for solving a noise-corrupted or nearly singular generalized eigenvalue problem $\mathbf{H}\mathbf{c} = E\mathbf{H}\mathbf{c}$.

```

procedure THRESHOLDING( $\mathbf{H}, \mathbf{S}, \epsilon$ )
   $(\mathbf{V}, \mathbf{D}) \leftarrow \text{eig}(\mathbf{S})$ 
   $I \leftarrow \{i : D_{ii} > \epsilon\}$ 
   $\mathbf{V} \leftarrow \mathbf{V}(:, I)$ 
  return smallest eigenvalue of  $(\mathbf{V}^* \mathbf{H} \mathbf{V}, \mathbf{V}^* \mathbf{S} \mathbf{V})$ 
end procedure

```

to be quite reliable at filtering out noise and dealing with the ill-conditioning of the overlap matrix for QSD-derived pairs (\mathbf{H}, \mathbf{S}) in our experiments, with similar results being observed for an SVD-based truncation strategy in [15, §II.F]. Despite truncation strategies such as canonical orthogonalization and Fix–Heiberger having a fifty year history, we are unaware of any general theory of why these methods work.

1.3. Overview and main results. We aim to elucidate why the QSD algorithm works when combined with the thresholding procedure, undeterred by the presence of negative examples such as (1.4) and (1.5). Our main three results are

- (i) In the absence of noise and with an appropriate choice of the time sequence, the QSD procedure with thresholding is accurate. This provides a positive answer to the question (A) in the error analysis of the QSD method due to the Krylov subspace approximation in the presence of thresholding.
- (ii) The thresholded problem is stable under noise in the sense that the thresholded problem we solve from the noise-corrupted pair $(\widetilde{\mathbf{H}}, \widetilde{\mathbf{S}})$ is close to the un-perturbed pair (\mathbf{H}, \mathbf{S}) .
- (iii) If the thresholded problem and its noisy perturbation are sufficiently close (as we establish with the previous result), the well-conditioned and well-separated eigenvalues of the thresholded problem are accurately computed in the presence of noise. Results (ii) and (iii) provide a positive answer to the question (B) in the error analysis of the QSD method due to the classical solution of the noisy generalized eigenvalue problem.

Together, these results paint a reasonably complete picture of why the QSD algorithm works in the presence of noise when thresholding is used. To our knowledge, this is the first work providing rigorous analysis of the theoretical efficacy of QSD-type algorithms, both in the noise-free and noisy settings.

Our first main result can be summarized informally as follows (a formal statement is presented as Theorem 2.7):

INFORMAL THEOREM 1.2. *Suppose the thresholding procedure (Algorithm 1.1) is applied to the perturbed pair $(\widetilde{\mathbf{H}}, \widetilde{\mathbf{S}}) = (\mathbf{H} + \Delta_{\mathbf{H}}, \mathbf{S} + \Delta_{\mathbf{S}})$, which are Hermitian matrices of size n . Consider the thresholded matrix pair $(\mathbf{A}, \mathbf{B}) := (\mathbf{V}_{>\epsilon}^* \mathbf{H} \mathbf{V}_{>\epsilon}, \mathbf{V}_{>\epsilon}^* \mathbf{S} \mathbf{V}_{>\epsilon})$ and let E_0 be its least eigenvalue. Assume E_0 is sufficiently well-separated from other eigenvalues of (\mathbf{A}, \mathbf{B}) , let d_0^{-1} denote the condition number of the eigenangle $\tan^{-1} E_0$, and suppose the perturbations $\Delta_{\mathbf{H}}$ and $\Delta_{\mathbf{S}}$ have spectral norm not exceeding η . There exists a constant $0 \leq \alpha \leq 1/2$ such that the recovered eigenvalue \widetilde{E}_0 from the noise-perturbed pair $(\widetilde{\mathbf{H}}, \widetilde{\mathbf{S}})$ using threshold parameter $\epsilon = \Theta\left(\eta^{\frac{1}{1+\alpha}}\right)$ satisfies the bound*

$$(1.6) \quad \left| \tan^{-1} \widetilde{E}_0^{\text{th}} - \tan^{-1} E_0^{\text{th}} \right| \leq \mathcal{O}\left(d_0^{-1} \eta^{\frac{1}{1+\alpha}}\right).$$

The implicit constant in the \mathcal{O} -notation depends on the eigenvalues of (\mathbf{H}, \mathbf{S}) , the spectrum of \mathbf{S} , and n .

Our result shows that, using the empirically observed value $\alpha = 1/4$ (see section SM6), we are able to recover the smallest eigenvalue of (\mathbf{H}, \mathbf{S}) (or more specifically its arctangent) with error proportional to $\eta^{4/5}$ times its condition number. Given examples (1.4) and (1.5) showing arbitrarily large errors for small perturbations of nearly singular generalized eigenvalue problems and thresholding, the fact that we are able to obtain any nontrivial error bounds for the QSD algorithm with thresholding may be regarded as surprising.

Our second main result provides an end-to-end bound for the QSD method.

INFORMAL THEOREM 1.3. *Let ΔE_j denote the difference between the j th smallest and the smallest eigenvalue E_0 of $\widehat{\mathbf{H}}$ and let γ_0 denote the inner product between the initial state φ_0 and the true ground-state eigenvector of $\widehat{\mathbf{H}}$. Let (\mathbf{H}, \mathbf{S}) denote the noise-free output of the QSD algorithm for a particular choice of time step, and instate the notation and assumptions of Informal Theorem 1.2. Then*

$$\left| \tan^{-1} \tilde{E}_0^{\text{th}} - \tan^{-1} E_0 \right| \leq \mathcal{O} \left(\frac{1 - |\gamma_0|^2}{|\gamma_0|^2} e^{-n \mathcal{O} \left(\frac{\Delta E_1}{\Delta E_{N-1}} \right)} + \left[\frac{\Delta E_{N-1}}{|\gamma_0|^2} + d_0^{-1} \right] \eta^{\frac{1}{1+\alpha}} \right).$$

The remainder of this paper is organized as follows. For expository reasons, we present our main results in the reverse order outlined here. Section 2 discusses perturbation analysis for the thresholded problem, leading to a formalization Theorem 2.7 of Informal Theorem 1.2 in section 2.3. Section 3 discusses Rayleigh–Ritz errors due to approximation by the finite-dimensional unitary Krylov subspace and thresholding procedure. We then present additional results in section 4 which are independent of the rest of the presentation. We draw particular attention to Theorem 4.2, which shows that thresholding applied to a general pair (\mathbf{H}, \mathbf{S}) recovers the least eigenvalue accurately if it is well-conditioned. This does not contradict the bad example (1.5) since both its eigenvalues are ill-conditioned with condition numbers $\Theta(\epsilon^{-1})$. We conclude with numerical experiments (section 5) and conclusions (section 6).

2. Perturbation Analysis for the Thresholded Problem. In this section, we analyze the effects of noise on the solution of the generalized eigenvalue problem (1.2) using thresholding Algorithm 1.1. The main result of this section is that well-separated, well-conditioned eigenvalues of the thresholded problem can be recovered accurately in the presence of noise. Together with section 3 which analyzes both the Rayleigh–Ritz and thresholding errors, this comprises a fairly complete explanation for the success of the QSD algorithm when implemented with thresholding.

Let \mathbf{H} and \mathbf{S} denote the exact outputs of the QSD algorithm (1.3) and \mathbf{V} the eigenvectors of \mathbf{S} with eigenvalues greater than ϵ . Dependence of $\mathbf{V} := \mathbf{V}_{>\epsilon}$ on ϵ , as in the introduction, has been suppressed for conciseness. The thresholded problem is described by the pair $(\mathbf{A}, \mathbf{B}) := (\mathbf{V}^* \mathbf{H} \mathbf{V}, \mathbf{V}^* \mathbf{S} \mathbf{V})$.² When implemented on a quantum computer, \mathbf{H} and \mathbf{S} are corrupted by noise as $\widetilde{\mathbf{H}} := \mathbf{H} + \Delta_{\mathbf{H}}$ and $\widetilde{\mathbf{S}} := \mathbf{S} + \Delta_{\mathbf{S}}$. As a simple measure of the size of the perturbation, we introduce

$$(2.1) \quad \eta := \sqrt{\eta_{\mathbf{H}}^2 + \eta_{\mathbf{S}}^2} := \sqrt{\|\Delta_{\mathbf{H}}\|^2 + \|\Delta_{\mathbf{S}}\|^2},$$

²To make the output of the thresholding procedure unambiguous, we assume eigenvectors are arranged left-to-right in decreasing order of the corresponding eigenvalues. The ordering convention does not effect the outputs of the thresholding procedure Algorithm 1.1.

which represents the noise level. In principle, one could undertake a careful analysis of the different sources of error (e.g., discretization, gate, and sampling) to obtain probabilistic bounds on η . (We provide such a bound for the sampling error alone in section 4.1.) For now, we shall just assume η or a good bound for it is known, and the threshold level ϵ is chosen to be (at least) larger than η . With the perturbations $\widetilde{\mathbf{H}}$ and $\widetilde{\mathbf{S}}$ in hand, the practitioner computes the large-eigenvalue eigenvectors $\widetilde{\mathbf{V}}$ of the perturbation $\widetilde{\mathbf{S}}$, and constructs the perturbed thresholded problem $(\widetilde{\mathbf{A}}, \widetilde{\mathbf{B}}) := (\widetilde{\mathbf{V}}^* \widetilde{\mathbf{H}} \widetilde{\mathbf{V}}, \widetilde{\mathbf{V}}^* \widetilde{\mathbf{S}} \widetilde{\mathbf{V}})$. We denote the dimension of $\widetilde{\mathbf{A}}$ and $\widetilde{\mathbf{B}}$ as q .

We hope to show that the smallest eigenvalue of the pair $(\widetilde{\mathbf{A}}, \widetilde{\mathbf{B}})$ —i.e., our computed approximation to the ground state energy—is close to the smallest eigenvalue of (\mathbf{A}, \mathbf{B}) . Unfortunately, there are a number of reasons to worry this might not be the case. First, even if $(\widetilde{\mathbf{H}}, \widetilde{\mathbf{S}})$ is close to (\mathbf{H}, \mathbf{S}) , it is still possible that $(\widetilde{\mathbf{A}}, \widetilde{\mathbf{B}})$ is not close to (\mathbf{A}, \mathbf{B}) . A small perturbation in just \mathbf{S} can lead to a large perturbation of \mathbf{A} : For a small parameter $\eta > 0$,

$$(2.2) \quad \begin{aligned} \mathbf{H} = \begin{bmatrix} 20 & 0 \\ 0 & 1 \end{bmatrix}, \mathbf{S} = \begin{bmatrix} 1 & 0 \\ 0 & 1 - \frac{\eta}{2} \end{bmatrix} &\xrightarrow{\eta \text{ error}} \widetilde{\mathbf{H}} = \begin{bmatrix} 20 & 0 \\ 0 & 1 \end{bmatrix}, \widetilde{\mathbf{S}} = \begin{bmatrix} 1 & 0 \\ 0 & 1 + \frac{\eta}{2} \end{bmatrix} \\ \mathbf{A} = \begin{bmatrix} 20 & 0 \\ 0 & 1 \end{bmatrix}, \mathbf{B} = \begin{bmatrix} 1 & 0 \\ 0 & 1 - \frac{\eta}{2} \end{bmatrix} &\longrightarrow \widetilde{\mathbf{A}} = \begin{bmatrix} 1 & 0 \\ 0 & 20 \end{bmatrix}, \widetilde{\mathbf{B}} = \begin{bmatrix} 1 + \frac{\eta}{2} & 0 \\ 0 & 1 \end{bmatrix}. \end{aligned}$$

Additionally, $\widetilde{\mathbf{A}}$ and \mathbf{A} can even be of different sizes if the perturbation causes the number of eigenvalues larger than ϵ to change. Fortunately, (2.2) suggests that the potential for small errors in \mathbf{S} to magnify into large errors in \mathbf{A} might have a benign source, with the error in this example caused simply by a reordering of the eigenvectors. The eigenvalues of the pair $(\widetilde{\mathbf{A}}, \widetilde{\mathbf{B}})$ are indifferent to a symmetric reordering of its rows and columns or, more generally, a $*$ -conjugation $(\widetilde{\mathbf{A}}, \widetilde{\mathbf{B}}) \mapsto (\mathbf{W}^* \widetilde{\mathbf{A}} \mathbf{W}, \mathbf{W}^* \widetilde{\mathbf{B}} \mathbf{W})$. Thus, it is sufficient for purposes of analysis to show that $(\widetilde{\mathbf{A}}, \widetilde{\mathbf{B}})$ are close to (\mathbf{A}, \mathbf{B}) after $(\widetilde{\mathbf{A}}, \widetilde{\mathbf{B}})$ is replaced by an appropriate $*$ -conjugation.

Assume this issue is addressed, and we obtain a $*$ -conjugation of $(\widetilde{\mathbf{A}}, \widetilde{\mathbf{B}})$ that is close to (\mathbf{A}, \mathbf{B}) . Classical worst-case perturbation theory still paints a grim portrait on the sensitivities of the eigenvalues of the thresholded pair (\mathbf{A}, \mathbf{B}) . After possibly replacing $(\widetilde{\mathbf{A}}, \widetilde{\mathbf{B}})$ by a $*$ -conjugation, let $(\widetilde{\mathbf{A}}, \widetilde{\mathbf{B}}) = (\mathbf{A} + \Delta_{\mathbf{A}}, \mathbf{B} + \Delta_{\mathbf{B}})$. A measure of the difference between (\mathbf{A}, \mathbf{B}) and $(\widetilde{\mathbf{A}}, \widetilde{\mathbf{B}})$ is given by

$$(2.3) \quad \chi := \sqrt{\|\Delta_{\mathbf{A}}\|^2 + \|\Delta_{\mathbf{B}}\|^2}.$$

Let E be the least eigenvalue of (\mathbf{A}, \mathbf{B}) , and \widetilde{E} be the least eigenvalue of $(\widetilde{\mathbf{A}}, \widetilde{\mathbf{B}})$. The classical perturbation theorem of Stewart [33, Thm. 3.2] can only show that

$$(2.4) \quad \left| \tan^{-1} \widetilde{E} - \tan^{-1} E \right| \leq \sin^{-1} \frac{\chi}{c(\mathbf{A}, \mathbf{B})} \leq \sin^{-1} \frac{\chi}{\epsilon}.$$

Here,

$$(2.5) \quad c(\mathbf{A}, \mathbf{B}) := \min_{\|\mathbf{x}\|=1} \sqrt{(\mathbf{x}^* \mathbf{A} \mathbf{x})^2 + (\mathbf{x}^* \mathbf{B} \mathbf{x})^2}$$

is the Crawford number, which is only guaranteed to be larger than ϵ under the standing assumptions. This leads to pessimistic error bounds on the order of χ/ϵ .

We would much prefer a bound which scales like χ times the condition number of $\tan^{-1} E$, without an explicit ϵ dependence.

In the rest of this section, we address these challenges. The analysis has two parts. In the first part, we use the Davis–Kahan $\sin \Theta$ theorem to show that, after replacing by a $*$ -conjugation, $(\tilde{\mathbf{A}}, \tilde{\mathbf{B}})$ is $\chi \leq \mathcal{O}(\eta/\epsilon^\alpha)$ -close to (\mathbf{A}, \mathbf{B}) , where $0 \leq \alpha \leq 1/2$ is a constant. In particular, α is assured to be no more than $1/2$, with $\alpha = 1/4$ appearing to be more representative in numerical experiments. For the second part, we use the perturbation theory of Mathias and Li [19] to improve on Stewart’s bound (2.4) in the case when the eigenvalue of interest is well-conditioned, obtaining the desired dependence on χ rather than χ/ϵ . We present these pieces in reverse order.

2.1. Eigenvalue Perturbation Bounds. Suppose that, potentially after a re-defining $(\tilde{\mathbf{A}}, \tilde{\mathbf{B}})$ to a $*$ -conjugation, $(\tilde{\mathbf{A}}, \tilde{\mathbf{B}})$ and (\mathbf{A}, \mathbf{B}) are separated by a small distance χ , as defined in (2.3). We shall show that the perturbation theory developed by Mathias and Li [19] allows us to significantly improve on the bound (2.4) furnished by Stewart’s theory. The beauty of the Mathias–Li theory is that the Crawford number in (2.4) can be replaced by a quantity related to the conditioning of E , provided a spectral gap condition is satisfied. We begin with a mildly specialized version of [19, Thm. 3.3]:

FACT 2.1. *Let (\mathbf{A}, \mathbf{B}) be a pair of $q \times q$ Hermitian matrices with all eigenvalues of \mathbf{B} larger than ϵ . Consider perturbations $(\tilde{\mathbf{A}}, \tilde{\mathbf{B}}) := (\mathbf{A} + \Delta_{\mathbf{A}}, \mathbf{B} + \Delta_{\mathbf{B}})$ where $\Delta_{\mathbf{A}}$ and $\Delta_{\mathbf{B}}$ are Hermitian matrices and χ is defined in (2.3) and satisfies $q\chi \leq \epsilon$. Let E_0, \dots, E_{q-1} be the eigenvalues of (\mathbf{A}, \mathbf{B}) with unit-norm eigenvectors $\mathbf{x}_0, \dots, \mathbf{x}_{q-1}$. Define*

$$(2.6) \quad \ell_j = \tan^{-1} E_j - \sin^{-1} \frac{q\chi}{d_j}, \quad u_j = \tan^{-1} E_j + \sin^{-1} \frac{q\chi}{d_j}.$$

where

$$(2.7) \quad d_j := |\mathbf{x}_j^*(\mathbf{A} + i\mathbf{B})\mathbf{x}_j|.$$

Let $\{\ell_j^\uparrow\}_{j=0}^{q-1}$ and $\{u_j^\uparrow\}_{j=0}^{q-1}$ denote the increasing rearrangements of the bounds $\{\ell_j\}_{j=0}^{q-1}$ and $\{u_j\}_{j=0}^{q-1}$. Then with $\tilde{E}_0 \leq \tilde{E}_1 \leq \dots \leq \tilde{E}_{q-1}$ the eigenvalues of $(\tilde{\mathbf{A}}, \tilde{\mathbf{B}})$, the following bound holds:

$$\ell_j^\uparrow \leq \tan^{-1} \tilde{E}_j \leq u_j^\uparrow.$$

We show a pictorial comparison of the proof of Stewart’s bound (2.4) and the Mathias–Li bound (2.9) in Figure 1. The perturbation theory of the definite generalized eigenvalues is naturally phrased in terms of the eigenangle $\theta_j := \tan^{-1} E_j$, which represents the angle of the ray through the complex number $\mathbf{x}_j^*(\mathbf{A} + i\mathbf{B})\mathbf{x}_j$ and the positive imaginary axis. In Stewart’s theory, one argues that, for a unit vector \mathbf{x} , the complex number $z := \mathbf{x}^*(\mathbf{A} + i\mathbf{B})\mathbf{x}$ must be a distance $c(\mathbf{A}, \mathbf{B})$ from the origin and the perturbed point $\tilde{z} := \mathbf{x}^*(\tilde{\mathbf{A}} + i\tilde{\mathbf{B}})\mathbf{x}$ is a distance at most χ from z . In view of a variational characterization Stewart proved for the eigenangles [33, Thm. 3.1], it follows that the eigenangles change by at most $\sin^{-1}(\chi/c(\mathbf{A}, \mathbf{B}))$. Mathias and Li instead consider the points $z_j := \mathbf{x}_j^*(\mathbf{A} + i\mathbf{B})\mathbf{x}_j$ for the unit-norm eigenvectors \mathbf{x}_j . A disk of radius χ centered at z_j is enclosed by rays with angles $\{\theta_j \pm \sin^{-1}(\chi/d_j)\} = \{\ell_j, u_j\}$. Using Stewart’s variational principle, Mathias and Li prove that, while the perturbed eigenangle $\theta_j := \tan^{-1} \tilde{E}_j$ need not lie within $[\ell_j, u_j]$, it must lie in $[\ell_j^\uparrow, u_j^\uparrow]$.

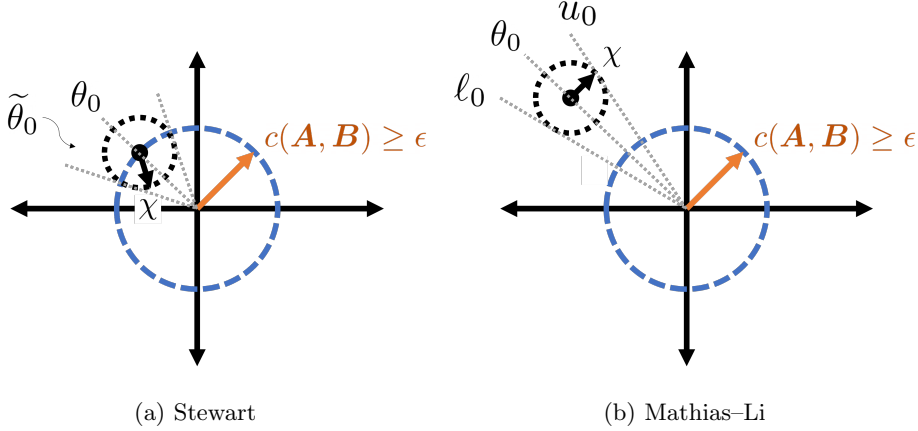


Fig. 1: Pictorial comparison of the proofs of Stewart's bound (2.4) and Mathias and Li's bound (2.9).

A consequence of Mathias and Li's analysis (see [19, Eq. (3.6)]) is that $\ell_j \leq \tilde{\theta}_j \leq u_j$ does hold if there is a large enough gap between θ_j and other eigenangles relative to the size of the perturbation. This bound $\ell_j \leq \tilde{\theta}_j \leq u_j$ is nearly as good a bound as one could hope since d_j^{-1} is the condition number of θ_j [34, Thm. VI.2.2].

COROLLARY 2.2. *Instate the notation and assumptions of Fact 2.1. Suppose that E_j satisfies the gap condition*

$$(2.8) \quad \min(\tan^{-1} E_j - \tan^{-1} E_{j-1}, \tan^{-1} E_{j+1} - \tan^{-1} E_j) \geq \sin^{-1} \frac{q\chi}{\epsilon} - \sin^{-1} \frac{q\chi}{d_j},$$

where the first or second term of the minimum can be ignored if $j = 0$ or $j = q - 1$, respectively. Then

$$(2.9) \quad \left| \tan^{-1} \tilde{E}_j - \tan^{-1} E_j \right| \leq \sin^{-1} \frac{q\chi}{d_j}.$$

A couple comments are in order before we proceed with the proof of Corollary 2.2. We have purged the suboptimal factor χ/ϵ from the eigenangle bound (2.9) and replaced it with the often much smaller quantity χ/d_j (tempered by a dimensional factor). However, χ/ϵ remains, just in the gap condition (2.8) (and the hypothesis $q\chi \leq \epsilon$). If the eigenangle gap on the left-hand side of (2.8) is reasonably large, then ϵ must merely be a modest multiple of χ for (2.8) to be satisfied.

Proof of Corollary 2.2. Denote $\tilde{\theta}_i := \tan^{-1} \tilde{E}_i$ and $\theta_i := \tan^{-1} E_i$ for every $i = 0, 1, \dots, q - 1$. We shall prove the upper bound $\tilde{\theta}_j - \theta_j \leq \sin^{-1}(q\chi/d_j)$ with the corresponding lower bound being proven in exactly the same way. We shall do this by showing $u_j^\uparrow \leq u_j$ under the gap assumption (2.8). To do this, we shall show that $u_i \leq \theta_j + \sin^{-1}(q\chi/d_j)$ for every $0 \leq i < j$. If $j = 0$, then no such i exists and the upper bound is automatically true. We thus continue in the case $j > 0$.

Fix $0 \leq i < j$. Since every eigenvalue of \mathbf{B} is at least ϵ , we have that $d_i = |\mathbf{x}_i^* (\mathbf{A} + i\mathbf{B}) \mathbf{x}_i| \geq \mathbf{x}_i^* \mathbf{B} \mathbf{x}_i > \epsilon$. Thus, we have $u_i \leq \theta_i + \sin^{-1}(q\chi/\epsilon)$. Since $E_i \leq E_{j-1}$,

we have $\theta_i = \tan^{-1} E_i \leq \tan^{-1} E_{j-1}$. Thus,

$$u_i \leq \tan^{-1} E_i + \sin^{-1} \frac{q\chi}{\epsilon} \leq \tan^{-1} E_{j-1} + \sin^{-1} \frac{q\chi}{\epsilon} \leq \tan^{-1} E_j + \sin^{-1} \frac{q\chi}{d_j}$$

by (2.8). From this we conclude $u_j^\uparrow \leq u_j$ so by Fact 2.1, $\tilde{\theta}_j \leq u_j^\uparrow \leq u_j = \theta_j + \sin^{-1}(q\chi/d_j)$. \square

2.2. How Do Perturbations Affect the Thresholded Problem? In this section, we seek to understand how perturbations of the pair (\mathbf{H}, \mathbf{S}) affect the thresholded problem (\mathbf{A}, \mathbf{B}) . As the example (2.2) shows, it need not be the case that $(\tilde{\mathbf{A}}, \tilde{\mathbf{B}})$ is close to (\mathbf{A}, \mathbf{B}) if $(\tilde{\mathbf{H}}, \tilde{\mathbf{S}})$ is close to (\mathbf{H}, \mathbf{S}) . Thus, in view of the fact that *-conjugations do not change the eigenvalues of a matrix pair [34, Thm. VI.1.8], our goal will be to show there exists a nonsingular matrix \mathbf{W} such that $(\mathbf{W}^* \tilde{\mathbf{A}} \mathbf{W}, \mathbf{W}^* \tilde{\mathbf{B}} \mathbf{W})$ is close to (\mathbf{A}, \mathbf{B}) . An instrumental tool in this goal will be the Davis-Kahan $\sin \Theta$ theorem [7], which we state a somewhat less general version here for reference.

FACT 2.3 (Davis-Kahan $\sin \Theta$ Theorem; see also [1, Thm. VII.3.1]). *Consider Hermitian matrices \mathbf{M} and $\tilde{\mathbf{M}}$ and let Π and $\tilde{\Pi}$ be the spectral projectors of \mathbf{M} and $\tilde{\mathbf{M}}$ associated with collections of eigenvalues of \mathbf{M} within an interval $[a, b]$ and of $\tilde{\mathbf{M}}$ outside the interval $[a - \delta, b + \delta]$ respectively. Then $\|\Pi\tilde{\Pi}\| \leq \|\tilde{\mathbf{M}} - \mathbf{M}\|/\delta$.*

We begin with some notation. Let $\mathbf{v}_1, \dots, \mathbf{v}_n$ denote the eigenvectors of \mathbf{S} with associated eigenvalues $\lambda_1, \dots, \lambda_n$ and Π denote the spectral projector associated with eigenvalues of \mathbf{S} which are larger than ϵ . We denote by m the critical index for which $\lambda_m > \epsilon \geq \lambda_{m+1}$. All quantities with tildes shall denote quantities defined as above for the perturbed problem.

We begin with the following simple bound:

PROPOSITION 2.4. *Suppose that*

$$(2.10) \quad \|\Pi\mathbf{H}\Pi - \tilde{\Pi}\mathbf{H}\tilde{\Pi}\| \leq \chi_{\mathbf{H}},$$

$$(2.11) \quad \|\Pi\mathbf{S}\Pi - \tilde{\Pi}\mathbf{S}\tilde{\Pi}\| \leq \chi_{\mathbf{S}},$$

and suppose that $\tilde{\mathbf{S}}$ and \mathbf{S} have the same number of eigenvalues larger than ϵ . Then for the nonsingular matrix $\mathbf{W} := \tilde{\mathbf{V}}^* \mathbf{V}$,

$$\sqrt{\|\mathbf{W}^* \tilde{\mathbf{A}} \mathbf{W} - \mathbf{A}\|^2 + \|\mathbf{W}^* \tilde{\mathbf{B}} \mathbf{W} - \mathbf{B}\|^2} \leq \sqrt{(\chi_{\mathbf{H}} + \eta_{\mathbf{H}})^2 + (\chi_{\mathbf{S}} + \eta_{\mathbf{S}})^2}.$$

This result gives a bound on χ as defined in (2.3) if one redefines $(\tilde{\mathbf{A}}, \tilde{\mathbf{B}})$ to its *-conjugation by \mathbf{W} .

Proof. First, note that \mathbf{W} is nonsingular because it is the product of a matrix with full row rank and full column rank. Note also that $\|\mathbf{V}\|, \|\tilde{\mathbf{V}}\|, \|\Pi\|, \|\tilde{\Pi}\| = 1$ since \mathbf{V} and $\tilde{\mathbf{V}}$ have orthonormal columns. The result then follows immediately by the bound

$$\begin{aligned} \|\mathbf{W}^* \tilde{\mathbf{A}} \mathbf{W} - \mathbf{A}\| &= \|\mathbf{V}^* \tilde{\mathbf{V}} \tilde{\mathbf{V}}^* \tilde{\mathbf{H}} \tilde{\mathbf{V}} \tilde{\mathbf{V}}^* \mathbf{V} - \mathbf{V}^* \mathbf{H} \mathbf{V}\| = \|\mathbf{V}^* (\tilde{\Pi} \tilde{\mathbf{H}} \tilde{\Pi} - \Pi \mathbf{H} \Pi) \mathbf{V}\| \\ &\leq \|\mathbf{V}^* (\tilde{\Pi} \tilde{\mathbf{H}} \tilde{\Pi} - \Pi \mathbf{H} \Pi) \mathbf{V}\| + \|\mathbf{V}^* \tilde{\Pi} (\tilde{\mathbf{H}} - \mathbf{H}) \tilde{\Pi} \mathbf{V}\| \\ &\leq \|\tilde{\Pi} \tilde{\mathbf{H}} \tilde{\Pi} - \Pi \mathbf{H} \Pi\| + \|\tilde{\mathbf{H}} - \mathbf{H}\| \leq \chi_{\mathbf{H}} + \eta_{\mathbf{H}} \end{aligned}$$

and similarly for $\|\mathbf{W}^* \tilde{\mathbf{B}} \mathbf{W} - \mathbf{B}\|$. \square

Note that a necessary condition for the hypotheses of our error bound Corollary 2.2 is for the distance χ (2.3) between $(\tilde{\mathbf{A}}, \tilde{\mathbf{B}})$ and (\mathbf{A}, \mathbf{B}) to be strictly smaller than ϵ . This forecloses our ability to use simple bounds for, e.g., (2.11) such as

$$\|\Pi \mathbf{S} \Pi - \tilde{\Pi} \tilde{\mathbf{S}} \tilde{\Pi}\| \leq \|\Pi \mathbf{S} \Pi - \mathbf{S}\| + \|\mathbf{S} - \tilde{\mathbf{S}}\| + \|\tilde{\mathbf{S}} - \tilde{\Pi} \tilde{\mathbf{S}} \tilde{\Pi}\| \leq 2(\epsilon + \eta).$$

We thus seek bounds of the form (2.10) and (2.11) without an additive $\mathcal{O}(\epsilon)$ term.

We begin with the more challenging of the two bounds, namely (2.10). Certainly, we should not expect a meaningful bound (2.10) if \mathbf{H} and \mathbf{S} have no relation to each other. For the sake of generality, we shall perform analysis under the assumption that (\mathbf{H}, \mathbf{S}) obey a weighted geometric mean inequality of the form

$$(2.12) \quad |\mathbf{v}_i^* \mathbf{H} \mathbf{v}_j| \leq \mu \min(\lambda_i, \lambda_j)^{1-\alpha} \max(\lambda_i, \lambda_j)^\alpha \quad \text{for all } 1 \leq i, j \leq n,$$

where $0 \leq \alpha \leq 1/2$ and $\mu > 0$ are constants. While this may appear strange, (2.12) necessarily holds with $\mu = \max |\Lambda(\mathbf{H}, \mathbf{S})|$ and $\alpha = 1/2$ by a direct application of the Courant–Fischer principle for generalized eigenvalue problems [34, Cor. VI.1.16]. Numerical experiments suggest that $\alpha = 1/4$ and $\mu \approx \max |\Lambda(\mathbf{H}, \mathbf{S})|$ appear to be more revealing of the empirically observed values of $|\mathbf{v}_i^* \mathbf{H} \mathbf{v}_j|$ for (\mathbf{H}, \mathbf{S}) computed by QSD for the physical models we tried; see section SM6.

Our most challenging technical result of this section will be a bound on the projection difference (2.10) under the assumption (2.12).

THEOREM 2.5. *Instate the prevailing notation and assume the bound (2.12) holds for $0 \leq \alpha \leq 1/2$. Assume the eigenvalue gap condition*

$$(2.13) \quad \lambda_{m+1} + \eta_{\mathbf{S}} \leq \epsilon < (1 + \rho)\epsilon \leq \lambda_m$$

for some $\rho > 0$. Suppose in addition that $\eta_{\mathbf{S}}$ is sufficiently small that $(1 + \rho^{-1})\eta_{\mathbf{S}}/\epsilon \leq 1$. Then the projection error (2.10) obeys the following bound

$$(2.14) \quad \|\Pi \mathbf{H} \Pi - \tilde{\Pi} \tilde{\mathbf{H}} \tilde{\Pi}\| \leq 3\mu n^3 (1 + \rho^{-1}) \left(\frac{\|\mathbf{S}\|}{\epsilon} \right)^\alpha \eta_{\mathbf{S}},$$

where $\eta_{\mathbf{S}}$ is defined in (2.1).

There are some unappealing features of this result, namely the cubic dependence on the problem size and the $\mathcal{O}(\eta_{\mathbf{S}}/\epsilon^\alpha)$ scaling. The first of these, the cubic dependence on n , we believe to likely be an artifact of our proof technique, applying Davis–Kahan “entry by entry”. Fortunately, numerical experiments do not suggest a dramatic dependence of the error on n . The second effect—the $\eta_{\mathbf{S}}/\epsilon^\alpha$ dependence rather than a more desirable $\eta_{\mathbf{S}}$ dependence—appears to be a genuine feature of this problem, at least without additional assumptions; see section SM5.³ Fortunately, we have yet to find an instance of a pair (\mathbf{H}, \mathbf{S}) generated by the QSD algorithm for which the ϵ^α factor appears to be necessary to understand the true error.

Finally, we note that the separation hypothesis (2.13) is relatively mild. The first inequality of (2.13) is necessary just to assure that \mathbf{A} and $\tilde{\mathbf{A}}$ have the same size. If we assume just a little more of a spectral gap around the thresholding level, quantified by the requirement that ρ is bounded away from zero, then we get a nice bound. A less careful application of Davis–Kahan would require that *all* the eigenvalues of \mathbf{S} are well-separated, so we consider a modest gap at the thresholding level to be a fairly mild requirement.

³Our evidence for this, presented in section SM5, is a synthetically generated pair (\mathbf{H}, \mathbf{S}) obeying the geometric mean condition (2.12); we did not obtain this pair from QSD.

Proof of Theorem 2.5. The proof shall be an enthusiastic exercise in applying the Davis–Kahan sin Θ theorem, Fact 2.3. We begin by bounding

$$(2.15) \quad \begin{aligned} \|\mathbf{\Pi H \Pi} - \tilde{\mathbf{\Pi H \Pi}}\|^2 &\leq \|\mathbf{\Pi H \Pi} - \tilde{\mathbf{\Pi H \Pi}}\|_{\mathbb{F}}^2 \\ &= \sum_{i,j=1}^n \left(\mathbf{v}_i^* (\mathbf{\Pi H \Pi} - \tilde{\mathbf{\Pi H \Pi}}) \mathbf{v}_j \right)^2 =: \sum_{i,j=1}^n I_{ij}^2. \end{aligned}$$

Our strategy will be to bound each of the terms I_{ij} .

For each i , we can expand $\tilde{\mathbf{\Pi}} \mathbf{v}_i = \sum_{k=1}^n c_{ik} \mathbf{v}_k$. Multiplying through by \mathbf{v}_i^* then gives that $c_{ik} = \mathbf{v}_i^* \tilde{\mathbf{\Pi}} \mathbf{v}_k$, which shows in particular that $c_{ik} = c_{ki}$. Our first goal will be to bound $|c_{ik}|$. We break into two cases, $k \leq m$ and $k > m$.

For case one, assume that $k \leq m$. By Weyl's inequality [34, Cor. IV.4.9], the $(m+1)$ st largest eigenvalue of $\tilde{\mathbf{S}}$ satisfies $\tilde{\lambda}_{m+1} \leq \lambda_{m+1} + \eta_S \leq \epsilon$. The Davis–Kahan theorem shows the difference $\boldsymbol{\delta}_k := \tilde{\mathbf{\Pi}} \mathbf{v}_k - \mathbf{v}_k$ satisfies

$$\|\boldsymbol{\delta}_k\| \leq \frac{\eta_S}{\lambda_k - \tilde{\lambda}_{m+1}} \leq \frac{\eta_S}{\lambda_k - \epsilon} \leq \frac{(1 + \rho^{-1})\eta_S}{\lambda_k}.$$

This gives a bound on the coefficients c_{ik} for $i \neq k$:

$$(2.16) \quad |c_{ik}| = \left| \mathbf{v}_i^* \tilde{\mathbf{\Pi}} \mathbf{v}_k \right| = |\mathbf{v}_i^* \mathbf{v}_k + \mathbf{v}_i^* \boldsymbol{\delta}_k| = |\mathbf{v}_i^* \boldsymbol{\delta}_k| \leq \|\boldsymbol{\delta}_k\| \leq \frac{(1 + \rho^{-1})\eta_S}{\lambda_k}.$$

For $i = k$, we have

$$(2.17) \quad |1 - c_{ii}| = |1 - \mathbf{v}_i^* (\mathbf{v}_i + \boldsymbol{\delta}_i)| \leq \|\boldsymbol{\delta}_i\| \leq \frac{(1 + \rho^{-1})\eta_S}{\lambda_i}.$$

Now consider case two where $k > m$. Since $\tilde{\lambda}_{m+1} \leq \epsilon$ as argued above, $\mathbf{\Pi}$ and $\tilde{\mathbf{\Pi}}$ are projections onto subspaces of the same dimensions so $\|\mathbf{\Pi} - \tilde{\mathbf{\Pi}}\| = \|\mathbf{\Pi}(\mathbf{I} - \tilde{\mathbf{\Pi}})\|$ by [34, Thm. I.5.5]. Applying Davis–Kahan then gives

$$\|\mathbf{\Pi} - \tilde{\mathbf{\Pi}}\| = \|\mathbf{\Pi}(\mathbf{I} - \tilde{\mathbf{\Pi}})\| \leq \frac{\eta_S}{\rho\epsilon}.$$

Then

$$(2.18) \quad |c_{ik}| = \left| \mathbf{v}_i^* \tilde{\mathbf{\Pi}} \mathbf{v}_k \right| \leq \left| \mathbf{v}_i^* (\tilde{\mathbf{\Pi}} - \mathbf{\Pi}) \mathbf{v}_k \right| + |\mathbf{v}_i^* \mathbf{\Pi} \mathbf{v}_k| \leq \|\mathbf{\Pi} - \tilde{\mathbf{\Pi}}\| \leq \frac{\eta_S}{\rho\epsilon}.$$

With these bounds in hand, we return to bounding I_{ij} as defined in (2.15). Let us introduce shorthand notation $a \wedge b$ and $a \vee b$ for the minimum and maximum of a and b respectively. Expanding $\tilde{\mathbf{\Pi}} \mathbf{v}_i$ and $\tilde{\mathbf{\Pi}} \mathbf{v}_j$ and using the bound (2.12), we obtain

$$(2.19) \quad \begin{aligned} I_{ij} &= \mathbf{v}_i^\dagger (\mathbf{\Pi H \Pi} - \tilde{\mathbf{\Pi H \Pi}}) \mathbf{v}_j = (1 - \overline{c_{ii}} c_{jj}) \cdot \mathbf{v}_i^\dagger \mathbf{H} \mathbf{v}_j - \sum_{\substack{k,\ell=1 \\ k \neq i \text{ or } \ell \neq j}}^n \overline{c_{ik}} c_{j\ell} \mathbf{v}_k^\dagger \mathbf{H} \mathbf{v}_\ell \\ &\leq \mu \left[|1 - \overline{c_{ii}} c_{jj}| (\lambda_i \wedge \lambda_j)^{1-\alpha} (\lambda_i \vee \lambda_j)^\alpha + \sum_{\substack{k,\ell=1 \\ k \neq i \text{ or } \ell \neq j}}^n |c_{ik} c_{j\ell}| (\lambda_k \vee \lambda_\ell)^{1-\alpha} (\lambda_k \wedge \lambda_\ell)^\alpha \right]. \end{aligned}$$

We thus turn our attention to bounding the summands in the second term of (2.19) for $k \neq i$ or $j \neq \ell$. First, assume that $k \neq i$ and suppose that $k \leq m$, we bound using (2.16):

$$\begin{aligned} |c_{ik}c_{j\ell}|(\lambda_k \wedge \lambda_\ell)^{1-\alpha}(\lambda_k \vee \lambda_\ell)^\alpha &\leq |c_{ik}|\lambda_k^{1-\alpha}\|\mathbf{S}\|^\alpha \leq \frac{(1+\rho^{-1})\eta_S\lambda_k^{1-\alpha}\|\mathbf{S}\|^\alpha}{\lambda_k} \\ &\leq (1+\rho^{-1})\eta_S(\|\mathbf{S}\|/\epsilon)^\alpha. \end{aligned}$$

Next assuming $k > m$, (2.18) yields

$$|c_{ik}c_{j\ell}|(\lambda_k \wedge \lambda_\ell)^{1-\alpha}(\lambda_k \vee \lambda_\ell)^\alpha \leq |c_{ik}|\epsilon^{1-\alpha}\|\mathbf{S}\|^\alpha \leq \eta_S\rho^{-1}(\|\mathbf{S}\|/\epsilon)^\alpha.$$

Turning our attention to the first term of the final bound in (2.19), (2.17) and the assumption $(1+\rho^{-1})\eta_S/\epsilon \leq 1$ give

$$\begin{aligned} |1 - \overline{c_{ii}c_{jj}}|(\lambda_i \wedge \lambda_j)^{1-\alpha}(\lambda_i \vee \lambda_j)^\alpha &\leq (|1 - c_{ii}| + |1 - c_{jj}| + |1 - c_{ii}| |1 - c_{jj}|)(\lambda_i \wedge \lambda_j)^{1-\alpha}\|\mathbf{S}\|^\alpha \\ &\leq 3\eta_S(1+\rho^{-1})(\|\mathbf{S}\|/\epsilon)^\alpha. \end{aligned}$$

Using the three previous displays to bound each of the n^2 summands in (2.19), we obtain

$$I_{ij} \leq 3\mu n^2(1+\rho^{-1})(\|\mathbf{S}\|/\epsilon)^\alpha\eta_S,$$

which then leads to the stated bound. \square

A bound for (2.11) is entirely analogous. The analysis is made significantly easier by the fact that the spectral projector $\mathbf{\Pi}$ is defined in terms of the matrix \mathbf{S} itself.

THEOREM 2.6. *Instate the prevailing notation. Assume that (2.13) holds for some $\rho > 0$. The projection error (2.11) satisfies the bound*

$$\|\mathbf{\Pi S \Pi} - \tilde{\mathbf{\Pi S \Pi}}\| \leq \|\mathbf{\Pi S \Pi} - \tilde{\mathbf{\Pi S \Pi}}\|_F \leq 2(1+\rho^{-1})\eta_S n + \epsilon^{-1}[(1+\rho^{-1})\eta_S n]^2.$$

In particular if $(1+\rho^{-1})\eta_S n/\epsilon \leq 1$, we have

$$\|\mathbf{\Pi S \Pi} - \tilde{\mathbf{\Pi S \Pi}}\| \leq \|\mathbf{\Pi S \Pi} - \tilde{\mathbf{\Pi S \Pi}}\|_F \leq 3(1+\rho^{-1})\eta_S n.$$

Proof. The proof is quite similar to Theorem 2.5 and we shall thus proceed more quickly. First, we bound

$$\|\mathbf{\Pi S \Pi} - \tilde{\mathbf{\Pi S \Pi}}\|^2 \leq \|\mathbf{\Pi S \Pi} - \tilde{\mathbf{\Pi S \Pi}}\|_F^2 = \sum_{i,j=1}^n \left(\mathbf{v}_i^* (\mathbf{\Pi S \Pi} - \tilde{\mathbf{\Pi S \Pi}}) \mathbf{v}_j \right)^2 =: \sum_{i,j=1}^n I_{ij}^2.$$

Consider the expansion $\tilde{\mathbf{\Pi}} \mathbf{v}_i = \sum_{k=1}^n c_{ik} \mathbf{v}_k$ as in the proof of Theorem 2.5. By the same arguments, the bounds (2.16), (2.17), and (2.18) hold under the respective hypotheses that $i \leq m$ and $i \neq k$, $i = k$, and $i > m$ respectively.

We now compute I_{ij} using the fact that $\mathbf{v}_i^* \mathbf{S} \mathbf{v}_j = \lambda_i \delta_{ij}$ where δ_{ij} denotes the Kronecker delta:

$$\begin{aligned} I_{ij} &= \mathbf{v}_i^* (\mathbf{\Pi S \Pi} - \tilde{\mathbf{\Pi S \Pi}}) \mathbf{v}_j = (1 - \overline{c_{ii}c_{jj}}) \cdot \mathbf{v}_i^* \mathbf{S} \mathbf{v}_j - \sum_{\substack{k,\ell=1 \\ k \neq i \text{ OR } \ell \neq j}}^n \overline{c_{ik}c_{j\ell}} \mathbf{v}_k^* \mathbf{H} \mathbf{v}_\ell \\ &= \lambda_i \delta_{ij} - \sum_{k=1}^n \overline{c_{ik}c_{jk}} \lambda_k. \end{aligned}$$

We now distinguish two cases. First suppose $i \neq j$. Then we bound using (2.16) and (2.18):

$$\begin{aligned} |I_{ij}| &\leq \sum_{k=1}^n |c_{ik}c_{jk}|\lambda_k \leq |c_{ij}|(\lambda_i + \lambda_j) + \sum_{\substack{k=1 \\ k \notin \{i,j\}}}^m |c_{ik}c_{jk}|\lambda_k + \sum_{m+1}^n |c_{ik}c_{jk}|\lambda_k \\ &\leq 2\eta_S(1 + \rho^{-1}) + \frac{n(1 + \rho^{-1})^2\eta_S^2}{\epsilon}. \end{aligned}$$

Next suppose $i = j$. Then applying (2.16), (2.17), and (2.18) gives the same bound

$$|I_{ij}| \leq 2\eta_S(1 + \rho^{-1}) + \frac{n(1 + \rho^{-1})^2\eta_S^2}{\epsilon}.$$

This entrywise bound immediately gives the desired result. \square

2.3. Main Result. We conclude this section by combining Corollary 2.2, Proposition 2.4, and Theorems 2.5 and 2.6 into our main result, which provides a formal statement of Informal Theorem 1.2 from the introduction.

THEOREM 2.7. *Let (\mathbf{H}, \mathbf{S}) be a pair of $n \times n$ Hermitian matrices perturbed to a pair $(\widetilde{\mathbf{H}}, \widetilde{\mathbf{S}})$ by perturbations $\Delta_{\mathbf{H}}$ and $\Delta_{\mathbf{S}}$ of spectral norms $\eta_{\mathbf{H}}$ and $\eta_{\mathbf{S}}$. Assume the following:*

- *The pair (\mathbf{H}, \mathbf{S}) satisfies the geometric mean bound (2.12) for some parameters $\mu > 0$ and $0 \leq \alpha \leq 1/2$.*
- *There exists an index m for which (2.13) holds for some $\rho > 0$.*
- *The noise $\eta_{\mathbf{S}}$ is sufficiently small so that $(1 + \rho^{-1})\eta_{\mathbf{S}}/\epsilon \leq 1$.*

Let $(\mathbf{A}, \mathbf{B}) = (\mathbf{V}^ \mathbf{H} \mathbf{V}, \mathbf{V}^* \mathbf{S} \mathbf{V})$ denote the thresholded matrix pair. The eigenvalues recovered by the thresholding procedure applied to the noise-perturbed pair $(\widetilde{\mathbf{H}}, \widetilde{\mathbf{S}})$ are the same as the eigenvalues of a pair $(\widetilde{\mathbf{A}}, \widetilde{\mathbf{B}})$ satisfying*

$$\sqrt{\|\widetilde{\mathbf{A}} - \mathbf{A}\|^2 + \|\widetilde{\mathbf{B}} - \mathbf{B}\|^2} \leq 3(2 + \mu)n^3(1 + \rho^{-1}) \left(\frac{\|\mathbf{S}\|}{\epsilon} \right)^\alpha \eta_{\mathbf{S}} + \eta_{\mathbf{H}} =: \chi.$$

Let E_0 and E_1 denote the least and second-to-least eigenvalues of (\mathbf{A}, \mathbf{B}) . Suppose further that

- *The error bound χ is sufficiently small: $n\chi \leq \epsilon$.*
- *The gap condition $\tan^{-1} E_1 - \tan^{-1} E_0 \geq \sin^{-1}(n\chi/\epsilon)$ holds.*

Then with d_0^{-1} the condition number of the eigenangle $\tan^{-1} E_0$ and \widetilde{E}_0 the eigenvalue recovered by thresholding applied to $(\widetilde{\mathbf{H}}, \widetilde{\mathbf{S}})$,

$$\left| \tan^{-1} E_0 - \tan^{-1} \widetilde{E}_0 \right| \leq \sin^{-1} \frac{n\chi}{d_0}.$$

In particular, for the theorem to hold, the threshold parameter must be chosen

$$\epsilon > n\chi = 3(2 + \mu)n^4(1 + \rho^{-1}) \left(\frac{\|\mathbf{S}\|}{\epsilon} \right)^\alpha \eta_{\mathbf{S}} + n\eta_{\mathbf{H}}.$$

This leads to the claimed value $\epsilon = \Theta\left(\eta^{\frac{1}{1+\alpha}}\right)$ and bound $\left| \tan^{-1} E_0 - \tan^{-1} \widetilde{E}_0 \right| \leq \mathcal{O}\left(d_0^{-1} \eta^{\frac{1}{1+\alpha}}\right)$ in Informal Theorem 1.2.

3. Analysis of unitary Krylov subspace approximation with thresholding. Having studied the errors due to noise in the previous section, we now turn to analyzing the Rayleigh–Ritz errors in computing the ground-state eigenvalue E_0 of $\widehat{\mathbf{H}}$ using the unitary Krylov space (1.1) and the thresholding procedure Algorithm 1.1. Specifically, we shall bound the difference between E_0 and the result \widetilde{E}_0 of applying thresholding to the pair (\mathbf{H}, \mathbf{S}) (1.3) computed on an error-free quantum computer.

If one considers our analysis with the threshold parameter ϵ set to zero, one obtains an analysis of the QSD method with no truncation. To our knowledge, this is the first quantitative error analysis of the QSD method even in the noise-free setting. This builds on two earlier explanations of the success of QSD. The first, by Stair, Huang, and Evangelista [32], argues based on Taylor series that the QSD subspaces approximately coincides with the classical polynomial Krylov space. This explanation has two drawbacks: (1) small timesteps make the ill-conditioning of \mathbf{H} and \mathbf{S} worse [32, §2.1] and (2) QSD often performs better with larger time steps [15, §II.B]. An alternate analysis based on filter diagonalization is provided by Klymko et al. [15], which provides an overcomplete set of *phase cancellation conditions* under which QSD computes the eigenvalues of interest with zero error. They then argue that these conditions hold approximately in the long-time limit for a randomly chosen timestep.

Our analysis is more direct than the two previous, emulating the classical analysis of the Lanczos method by Saad [29]. This leads to a quantitative error bound in terms of the distribution of the spectra which decreases exponentially in the number of timesteps used. The basic idea will be to use a linear combination of the QSD basis states φ_j which has the effect of applying a trigonometric polynomial to the eigenvalues E_0, \dots, E_{N-1} of $\widehat{\mathbf{H}}$. If φ_j approximately lies in the span of the first $M + 1$ eigenstates, we shall choose this trigonometric polynomial to be large at the eigenvalue of interest and exponentially small (in n) at eigenvalues E_1, \dots, E_M . The trigonometric polynomial will be bounded so it won't amplify any components of the eigenvector in the direction of any of the remaining eigenvectors. As an additional feature of this analysis, we are able to directly analyze thresholding “for free” in the noiseless setting, where thresholding has the effect of perturbing this trigonometric polynomial. The main theorem of this section is as follows:

THEOREM 3.1. *Let $\psi_0, \dots, \psi_{N-1}$ be the eigenvectors of a Hermitian operator $\widehat{\mathbf{H}}$ with eigenvalues E_0, \dots, E_{N-1} . Suppose the initial vector is expanded as*

$$(3.1) \quad \varphi_0 = \sum_{j=0}^{N-1} \gamma_j \psi_j.$$

Let $\Delta E_j := E_j - E_0$ and choose an index $0 \leq M \leq N - 1$. Suppose the QSD algorithm is implemented with time sequence $\{t_j\}_{j=-k}^k$ for $t_j = \pi j / \Delta E_M$. Suppose the generalized eigenvalue problem (1.2) is solved with thresholding parameter ϵ and let ϵ_{total} be the sum of the eigenvalues of \mathbf{S} discarded by thresholding. Then

$$(3.2) \quad 0 \leq \widetilde{E}_0 - E_0 \leq \frac{2 \left[\Delta E_{N-1} \epsilon_{\text{total}} + 4 \left(1 + \frac{\pi \Delta E_1}{\Delta E_M} \right)^{-2k} \sum_{i=1}^M \Delta E_i |\gamma_i|^2 + \sum_{i=M+1}^{N-1} \Delta E_i |\gamma_i|^2 \right]}{|\gamma_0|^2 - 2|\gamma_0| \sqrt{(2k+1)\epsilon}}.$$

Remark 3.2. The value $|\gamma_0|^2 = |\langle \varphi_0 | \psi_0 \rangle|^2$ is referred to as the initial overlap. In order for the QSD algorithm to succeed with a relatively small number of Krylov steps, $|\gamma_0|^2$ must be sufficiently large (for instance, ≥ 0.1). This is qualitatively different from

the assumption of classical Krylov subspace methods for solving eigenvalue problems, where the initial overlap $|\gamma_0|^2$ should also be nonzero but can be very small.

Let us note some salient features of this result. First, if we consider the noise-free case without thresholding, we obtain the bound

$$0 \leq \tilde{E}_0 - E_0 \leq 8 \left(1 + \frac{\pi \Delta E_1}{\Delta E_M}\right)^{-2k} \sum_{i=1}^M \Delta E_i \frac{|\gamma_i|^2}{|\gamma_0|^2} + 2 \sum_{i=M+1}^{N-1} \Delta E_i \frac{|\gamma_i|^2}{|\gamma_0|^2}.$$

The error bound has two terms, a term concerning the eigenvalues E_1, \dots, E_M which is damped exponentially fast with rate $\Delta E_1/\Delta E_M$ and an undamped (but also unamplified) term with the eigenvalues E_{M+1}, \dots, E_{N-1} . If the components of φ_0 in the directions of the $(M+1)$ st to $(N-1)$ st eigenvectors are small in the sense that $\Delta E_i |\gamma_i|^2 \ll 1$, this second term will be small. This bound thus has a tradeoff: If M is chosen larger (i.e., the timestep becomes smaller), more terms will be exponentially damped but at a slower rate as $\Delta E_1/\Delta E_M$ will decrease. One can obtain the simplest bound by choosing $M = N - 1$ and bounding $\Delta E_i \leq \Delta E_{N-1}$, leading to

$$(3.3) \quad 0 \leq \tilde{E}_0 - E_0 \leq 8 \Delta E_{N-1} \frac{1 - |\gamma_0|^2}{|\gamma_0|^2} \left(1 + \frac{\pi \Delta E_1}{\Delta E_{N-1}}\right)^{-2k}.$$

Unfortunately, this simplified bound often grossly overestimates the error.

Our analysis also indicates that thresholding has only a mild effect on the accuracy. We pick up a $1 + \mathcal{O}(\sqrt{k}\epsilon)$ prefactor and an additional term proportional to the sum of the discarded eigenvalues of \mathbf{S} , which in turn can be bounded $\epsilon_{\text{total}} \leq (2k+1)\epsilon$. In practice, we usually have $\epsilon_{\text{total}} \approx \epsilon$ due to rapid spectral decay.

Combining Theorem 3.1 with Informal Theorem 1.2 (formalized as Theorem 2.7) leads directly to Informal Theorem 1.3. More precise but also more complex error bounds can be obtained by using the full power (3.2) of Theorem 3.1 directly with Theorem 2.7.

3.1. Proof of Theorem 3.1. Our proof is based on the observation that thresholding is equivalent to applying the Rayleigh–Ritz procedure with a subspace spanned by the dominant left singular vectors of the Krylov matrix. Consider the Krylov matrix \mathbf{K} defined and factorized as

$$(3.4) \quad \begin{aligned} \mathbf{K} &:= [\varphi_{-k} \quad \cdots \quad \varphi_k] \\ &= \underbrace{[\psi_0 \quad \cdots \quad \psi_{N-1}]}_{:=\Psi} \underbrace{\begin{bmatrix} \gamma_0 & & \\ & \ddots & \\ & & \gamma_{N-1} \end{bmatrix}}_{:=\Gamma} \underbrace{\begin{bmatrix} e^{-ikE_0} & \cdots & e^{ikE_0} \\ \vdots & \ddots & \vdots \\ e^{-ikE_{N-1}} & \cdots & e^{ikE_{N-1}} \end{bmatrix}}_{:=\mathbf{F}}. \end{aligned}$$

Then we easily see that $\mathbf{H} = \mathbf{K}^* \widehat{\mathbf{H}} \mathbf{K}$ and $\mathbf{S} = \mathbf{K}^* \mathbf{K}$. Since the eigenvalues of $\mathbf{K}^* \mathbf{K}$ are the squares of singular values of \mathbf{K} with eigenvectors equal to the right singular vectors of \mathbf{K} , it follows that the thresholded problem is precisely the Rayleigh–Ritz procedure applied to the left singular subspace of \mathbf{K} with singular values larger than $\sqrt{\epsilon}$. From these left singular vectors, we are able to reconstruct the matrix \mathbf{K} up to a Frobenius norm error $\sqrt{\epsilon_{\text{total}}}$. We thus can analyze QSD with thresholding in much the same way as Saad’s analysis of the Lanczos method [29] with two twists. First, we have trigonometric polynomials in place of polynomials owing to the QSD algorithm’s

use of the unitary time-evolution operator. Second, our trigonometric basis functions are perturbed as a result of the truncated singular value decomposition.

With this roadmap in mind, we begin with a trigonometric version of a classic result in polynomial approximation theory.

LEMMA 3.3. *Let $0 < a < \pi$ and denote by \mathcal{T}_k the space of degree $\leq k$ trigonometric polynomials. The trigonometric polynomial minimax approximation problem*

$$\beta(a, k) = \min_{\substack{p \in \mathcal{T}_k \\ p(0)=1}} \max_{t \in (-\pi, \pi) \setminus (-a, a)} |p(t)|$$

is solved by

$$(3.5) \quad p^*(\theta) = \frac{T_k(1 + 2 \frac{\cos \theta - \cos a}{\cos a + 1})}{T_k(1 + 2 \frac{1 - \cos a}{\cos a + 1})}$$

where T_k denotes the k th Chebyshev polynomial. The optimum value is

$$(3.6) \quad \beta(a, k) = \left(T_k \left(1 + 2 \frac{1 - \cos a}{\cos a + 1} \right) \right)^{-1} \leq 2 \left(1 + 2 \sqrt{\frac{1 - \cos a}{\cos a + 1}} \right)^{-k} \leq 2(1 + a)^{-k}.$$

Proof. Once one convinces oneself that an optimal solution can be taken to be real and even, this result follows immediately from the analogous result for polynomial approximation [2, Thm. 4.1.11] together with the standard reparametrization $f(x) \mapsto f^\circ(\theta) := f(\cos \theta)$ which puts into bijection algebraic and even trigonometric polynomials. \square

We shall need a bound for the optimal trigonometric polynomial (3.5).

PROPOSITION 3.4. *The trigonometric polynomial p^* defined in (3.5) is bounded in absolute value by 1 and satisfies the L^2 bound*

$$(3.7) \quad \int_{-\pi}^{\pi} |p^*(\theta)|^2 d\theta \leq 2a + (2\pi - 2a)(\beta(a, k))^2 \leq 2\pi.$$

Proof. First, we show that p^* is monotone on $[0, a]$ by showing its derivative can't have a zero on $(0, a)$. Up to a scaling factor and a change of variables on the input, p^* coincides with T_k on $[a, \pi]$. Thus, p^* has $k - 1$ local extrema on (a, π) and symmetrically $k - 1$ on $(-\pi, -a)$. Since p^* is even and 2π -periodic, p^* has local extrema at 0 and π . Since $(p^*)'$ is a degree- k trigonometric polynomial, it has at most $2k$ zeros, all of which have already been accounted for. Thus, p^* is monotone on $[0, a]$.

Since $p^*(a) \leq \beta(a, k) < 1$, p^* is monotone decreasing on $[0, a]$ and thus achieves its maximum value on $[-a, a]$ of 1 at 0. On $(-\pi, \pi) \setminus [-a, a]$, $|p^*| \leq \beta(a, k)$, from which the bound (3.7) follows. \square

We shall be content with using the looser upper bound 2π in (3.7) in our subsequent analysis. With these approximation results in hand, we prove Theorem 3.1.

Proof of Theorem 3.1. Let φ_0 be expanded as (3.1) and consider the Krylov matrix and its factorization defined in (3.4). Let $\widetilde{\mathbf{K}}$ represent the truncation of \mathbf{K} by settings its singular values which are at most $\sqrt{\epsilon}$ to zero and factor it as $\widetilde{\mathbf{K}} = \Psi \Gamma \widetilde{\mathbf{F}}$ so that

$$(3.8) \quad \epsilon_{\text{total}} = \|\widetilde{\mathbf{K}} - \mathbf{K}\|_{\text{F}}^2 = \|\Gamma(\widetilde{\mathbf{F}} - \mathbf{F})\|_{\text{F}}^2 = \sum_{i=0}^{N-1} \sum_{j=-k}^k |\gamma_i|^2 \underbrace{\left| \widetilde{f}_{ij} - e^{ijE_i} \right|^2}_{:= \alpha_{ij}^2},$$

where \tilde{f}_{ij} denotes the ij entry of $\tilde{\mathbf{F}}$.

Let \mathcal{R}_ϵ denote the range of $\tilde{\mathbf{K}}$. By the Courant–Fischer theorem [34, Cor. IV.4.7],

$$(3.9) \quad \tilde{E}_0 - E_0 = \min_{|\xi\rangle \in \mathcal{R}_\epsilon \setminus \{0\}} \frac{\mathbf{x} \mathbf{i}^* (\widehat{\mathbf{H}} - E_0 \mathbf{I}) \xi}{\xi^* \xi},$$

where \mathbf{I} denotes the identity. Since $\widehat{\mathbf{H}} - E_0 \mathbf{I}$ is positive semidefinite, we have $\tilde{E}_0 - E_0 \geq 0$. The remainder of our effort will be dedicated to establishing an upper bound.

We shall use the minimax optimal trigonometric polynomial p^* (3.5) to construct an ansatz ξ_\star to plug into (3.9). Let $p^*(\theta - \pi E_0 / \Delta E_{M-1}) = \sum_{j=-k}^k c_j e^{ij\theta}$ be the Fourier series of the $\pi E_0 / \Delta E_{M-1}$ -translate of p^* . Choose as ansatz

$$\xi_\star := \sum_{i=0}^{N-1} \sum_{j=-k}^k \gamma_i \tilde{f}_{ij} c_j \psi_i \in \mathcal{R}_\epsilon.$$

Plugging this into (3.9), we obtain an upper bound

$$(3.10) \quad \tilde{E}_0 - E_0 \leq \frac{\xi_\star^* (\widehat{\mathbf{H}} - E_0 \mathbf{I}) \xi_\star}{\xi_\star^* \xi_\star} = \frac{\sum_{i=1}^{N-1} \Delta E_i |\gamma_i|^2 \left| \sum_{j=-k}^k \tilde{f}_{ij} c_j \right|^2}{\sum_{i=0}^{N-1} |\gamma_i|^2 \left| \sum_{j=-k}^k \tilde{f}_{ij} c_j \right|^2}.$$

First, focus on the numerator of the final bound in (3.10). We bound

$$\left| \sum_{j=-k}^k \tilde{f}_{ij} c_j \right|^2 \leq \left| \sum_{j=-k}^k (\tilde{f}_{ij} - e^{ijE_i}) c_j + \sum_{j=-k}^k e^{ijE_i} c_j \right|^2 \leq \left(\sum_{|j| \leq k} \alpha_{ij} c_j + p^*(E_i - E_0) \right)^2,$$

where $\alpha_{ij} \geq 0$ is defined in (3.8).

First suppose that $1 \leq i \leq M$. Then by the fact that $|p^*(E_i - E_0)| \leq \beta(a, k)$ as defined in (3.6) with $a := \pi \Delta E_1 / \Delta E_M$, the Parseval theorem (i.e., $2\pi \sum_{j=-k}^k |c_j|^2 = \int_{-\pi}^{\pi} |p^*(\theta)|^2 d\theta$), and the bound (3.7) we obtain

$$(3.11) \quad \left| \sum_{j=-k}^k \tilde{f}_{ij} c_j \right|^2 \leq 2 \left(\sum_{j=-k}^k |c_j|^2 \right) \left(\sum_{j=-k}^k \alpha_{ij}^2 \right) + 2(\beta(a, k))^2 \leq 2 \sum_{j=-k}^k \alpha_{ij}^2 + 2(\beta(a, k))^2.$$

For $M < i < N$, we get the same bound except with 1 in place of $\beta(a, k)$ since $|p^*| \leq 1$ by Proposition 3.4.

For the denominator of the final bound in (3.10), we bound

$$(3.12) \quad \left| \sum_{j=-k}^k \tilde{f}_{0j} c_j \right| \geq 1 - \sum_{j=-k}^k \alpha_{0j} \geq 1 - \sqrt{(2k+1) \sum_{j=-k}^k \alpha_{0j}^2} \geq 1 - \frac{1}{|\gamma_0|} \sqrt{(2k+1)\epsilon},$$

where we used a spectral norm bound similar to (3.8):

$$\epsilon \geq \|\tilde{\mathbf{K}} - \mathbf{K}\| = \|\mathbf{\Gamma}(\tilde{\mathbf{F}} - \mathbf{F})\| \geq \sqrt{|\gamma_0|^2 \sum_{j=-k}^k \left| \tilde{f}_{0j} - e^{ijE_0} \right|^2} = |\gamma_0| \sqrt{\sum_{j=-k}^k \alpha_{0j}^2}.$$

Plugging (3.11) and (3.12) into (3.10),

$$\begin{aligned} \tilde{E}_0 - E_0 &\leq \frac{2 \left[\sum_{i=1}^M \Delta E_i \left(\sum_{j=-k}^k \alpha_{ij}^2 + (\beta(a, k))^2 \right) + \sum_{i=M+1}^{N-1} \Delta E_i \left(\sum_{j=-k}^k \alpha_{ij}^2 + 1 \right) \right]}{\left(|\gamma_0| - \sqrt{(2k+1)\epsilon} \right)^2} \\ &\leq \frac{2 \left[\Delta E_{N-1} \epsilon_{\text{total}} + (\beta(a, k))^2 \sum_{i=1}^M \Delta E_i |\gamma_i|^2 + \sum_{i=M+1}^{N-1} \Delta E_i |\gamma_i|^2 \right]}{|\gamma_0|^2 - 2|\gamma_0| \sqrt{(2k+1)\epsilon}}. \end{aligned}$$

Using the bound (3.6) with $a = \pi \Delta E_1 / \Delta E_M$ leads precisely to (3.2). \square

4. Additional Results and Discussions. In this section, we include some additional results and discussions which follow from our analysis but are not directly germane to the main analysis of the QSD algorithm comprising Theorems 2.7 and 3.1.

4.1. On the Toeplitz Structure of \mathbf{H}, \mathbf{S} . With the choice of the basis vectors $\varphi_j = e^{it_j \widehat{H}} \varphi_0$ as in the Parrish–McMahon QSD procedure, the matrix elements of the projected matrices satisfy

$$\mathbf{H}_{jk} = \varphi_j^* \widehat{\mathbf{H}} \varphi_k = \varphi_0^* \widehat{\mathbf{H}} e^{i(t_k - t_j) \widehat{H}} \varphi_0 \quad \text{and} \quad \mathbf{S}_{jk} = \varphi_j^* \varphi_k = \varphi_0^* e^{i(t_k - t_j) \widehat{H}} \varphi_0.$$

Therefore, both \mathbf{H}, \mathbf{S} are Hermitian–Toeplitz matrices. Unfortunately, the Toeplitz structure relies on the assumption that the Hamiltonian simulation problem (i.e., $e^{it_j \widehat{H}} \varphi_0$) is computed exactly. In practice, the Hamiltonian simulation problem is often performed with approximate techniques (such as Trotter splitting), and the resulting projected matrices $\widetilde{\mathbf{H}}$ and $\widetilde{\mathbf{S}}$ may not have the Toeplitz structure. If this is the case, all n^2 entries of $\widetilde{\mathbf{H}}$ (and perhaps $\widetilde{\mathbf{S}}$ as well) need to be computed to apply the Rayleigh–Ritz procedure to the computed basis states $\varphi_0, \dots, \varphi_{n-1}$ in earnest. However, if one measures only the first row of \mathbf{H} and \mathbf{S} and imputes the remaining entries from the Hermitian–Toeplitz structure, then resulting recovered matrices $\widetilde{\mathbf{H}}$ and $\widetilde{\mathbf{S}}$ represent the true \mathbf{H} and \mathbf{S} corrupted by both Monte Carlo and discretization errors. Our main analysis makes no use of the Toeplitz structure.

An important question for the QSD procedure is how entrywise errors in the entries \mathbf{H} and \mathbf{S} correspond to errors in the spectral norm. For the standard QSD estimation procedure, the entries of $\mathbf{M} \in \{\mathbf{H}, \mathbf{S}\}$ are approximated by averaging m unbiased estimators each with maximum error B ($B = \mathcal{O}(1)$ for \mathbf{S} and $B = \mathcal{O}(\|\mathbf{H}\|)$ for \mathbf{H}). Consider the case where the Hamiltonian simulation problem is solved exactly and we compute estimates for $\mathbf{M} \in \{\mathbf{H}, \mathbf{S}\}$ by measuring the first row of \mathbf{M} and computing the remaining entries from the Hermitian–Toeplitz structure. Straightforward application of matrix concentration inequalities then shows that the approximation is $\mathcal{O}(B \sqrt{(n \log n)/m})$ -close to \mathbf{M} (see, e.g., [35, Thms. 3.6.1 and 4.6.1]).

4.2. Stability of Best Low-rank Approximation. Theorem 2.6 constitutes a stability result for the Eckart–Young best rank- m approximation $\llbracket \mathbf{S} \rrbracket_m = \mathbf{\Pi} \mathbf{S} \mathbf{\Pi}$ to \mathbf{S} . Since this result may be of independent interest, we state it here unburdened by the particularities of the QSD context.

THEOREM 4.1. *Let $\mathbf{A} \in \mathbb{C}^{n \times n}$ be a positive semidefinite matrix with eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$. Let $\llbracket \mathbf{A} \rrbracket_m$ denote the Eckart–Young best rank- m approximation to \mathbf{A} . Then, for any quadratic unitarily invariant norm $\|\cdot\|_{\text{QUI}}$ (such as the*

spectral or Frobenius norms; see [1, Def. IV.2.9]),

$$\|[\mathbf{A} + \mathbf{\Delta}]_m - [\mathbf{A}]_m\|_{\text{QUI}} \leq \|\mathbf{\Delta}\|_{\text{QUI}} + \frac{2n\lambda_m\|\mathbf{\Delta}\|}{\lambda_m - \lambda_{m+1} - \|\mathbf{\Delta}\|} \left(1 + \frac{0.5n\|\mathbf{\Delta}\|}{\lambda_m - \lambda_{m+1} - \|\mathbf{\Delta}\|}\right).$$

Proof. This follows immediately from Theorem 2.6 and the maximality of the Frobenius norm among quadratic unitarily invariant norms [1, Eq. (IV.38)]. \square

This bound is interesting because, in the setting where the gap $\lambda_m - \lambda_{m+1} - \|\mathbf{\Delta}\|$ is comparable in size to λ_m , the best rank- m approximation changes only by $\approx n\|\mathbf{\Delta}\|_{\text{QUI}}$ independent of the approximation error $\|\mathbf{A} - [\mathbf{A}]_m\|_{\text{QUI}}$. This can be a large improvement over the general-purpose bound (similar to [8, Cor. 2.4]) which holds for general unitarily invariant norms $\|\cdot\|_{\text{UI}}$:

$$\|[\mathbf{A} + \mathbf{\Delta}]_m - [\mathbf{A}]_m\|_{\text{UI}} \leq 2(\|\mathbf{A} - [\mathbf{A}]_m\|_{\text{UI}} + \|\mathbf{\Delta}\|_{\text{UI}}),$$

which does depend on the approximation error $\|\mathbf{A} - [\mathbf{A}]_m\|_{\text{UI}}$. We believe the dimensional factor n is likely quite pessimistic, and can hopefully be replaced by a small constant (at least) if the matrix enjoys a rapidly decaying spectrum.

4.3. General Analysis of Thresholding. We have presented an analysis in Theorem 3.1 of thresholding for the QSD problem, which simultaneously treats the Rayleigh–Ritz error from approximation from the titular quantum subspace and the thresholding procedure together. This is quite natural as the total error is what is most important to the practitioner, and one should in principle be able to obtain a more precise error bound by not decoupling these pieces. However, it remains a question of mathematical interest and of interest to the broader uses of thresholding beyond QSD to provide an analysis of thresholding for general matrices.

As the bad example (1.5) shows, thresholding may not work for general matrices, with thresholding with parameter ϵ able to introduce errors $\gg \epsilon$. This behavior is less surprising, however, when one notes that neither eigenvalue of (1.5) is well-conditioned, with both having a condition number $\Theta(\epsilon^{-1})$. In fact, this is the only obstruction to thresholding working (at least for recovering the smallest eigenvalue), as we shall show that the least eigenvalue is recovered accurately by thresholding if it is well-conditioned. (We recall [34, Thm. VI.2.2] that the condition number of the eigenangle $\tan^{-1} E_0$ associated with the least eigenvalue E_0 of (\mathbf{H}, \mathbf{S}) is $\|\mathbf{c}_0\|^2 \sqrt{1 + E_0^2}$ where \mathbf{c}_0 is the \mathbf{S} -normalized eigenvector, $\mathbf{c}_0^* \mathbf{S} \mathbf{c}_0 = 1$, associated with E_0 .)

THEOREM 4.2. *Let \mathbf{c}_0 be the \mathbf{S} -normalized eigenvector associated with the least eigenvalue E_0 of the generalized eigenvalue problem (1.2) for a Hermitian and Hermitian positive matrix \mathbf{H} and \mathbf{S} . Then, for \tilde{E}_0 the least eigenvalue recovered by thresholding with parameter ϵ and provided $2\sqrt{\epsilon}\|\mathbf{c}_0\| < 1$,*

$$0 \leq \tilde{E}_0 - E_0 \leq \frac{\Delta E \epsilon \|\mathbf{c}_0\|^2}{1 - 2\sqrt{\epsilon}\|\mathbf{c}_0\|},$$

where ΔE is the difference between the largest and smallest eigenvalue of (\mathbf{H}, \mathbf{S}) .

Proof. Let \mathcal{R}_ϵ be the span of the eigenvectors of \mathbf{S} with eigenvalue greater than ϵ . Using the same observation which motivated the proof of Theorem 3.1, we have

$$(4.1) \quad \tilde{E}_0 - E_0 = \min_{\mathbf{c} \in \mathcal{R}_\epsilon \setminus \{0\}} \frac{\mathbf{c}^* (\mathbf{H} - E_0 \mathbf{S}) \mathbf{c}}{\mathbf{c}^* \mathbf{S} \mathbf{c}}.$$

In particular, since $\mathbf{H} - E_0 \mathbf{S}$ is positive semidefinite, this implies $\tilde{E}_0 - E_0 \geq 0$ so we just need to concern ourselves with obtaining an upper bound. Letting $\tilde{\mathbf{S}}$ be the \mathbf{S} matrix with all its eigenvalues at most ϵ set to zero, we shall evaluate (4.1) at $\tilde{\mathbf{c}}_0 := \mathbf{S}^{-1/2} \tilde{\mathbf{S}}^{1/2} \mathbf{c}_0$, obtaining an upper bound. Defining the error $\boldsymbol{\delta} := \tilde{\mathbf{c}}_0 - \mathbf{c}_0$, we have that $\boldsymbol{\delta}$ satisfies the bound $\|\mathbf{S}^{1/2} \boldsymbol{\delta}\| \leq \|\mathbf{S}^{1/2} - \tilde{\mathbf{S}}^{1/2}\| \|\mathbf{c}_0\| \leq \sqrt{\epsilon} \|\mathbf{c}_0\|$. Thus

$$\tilde{\mathbf{c}}_0^* \mathbf{S} \tilde{\mathbf{c}}_0 \geq \mathbf{c}_0^* \mathbf{S} \mathbf{c}_0 - 2|\mathbf{c}_0^* \mathbf{S} \boldsymbol{\delta}| \geq 1 - 2(\mathbf{c}_0^* \mathbf{S} \mathbf{c}_0)^{1/2} (\boldsymbol{\delta}^* \mathbf{S} \boldsymbol{\delta})^{1/2} \geq 1 - 2\sqrt{\epsilon} \|\mathbf{c}_0\|$$

and

$$\tilde{\mathbf{c}}_0^* (\mathbf{H} - E_0 \mathbf{S}) \tilde{\mathbf{c}}_0 = \boldsymbol{\delta}^* (\mathbf{H} - E_0 \mathbf{S}) \boldsymbol{\delta} \leq \|\mathbf{S}^{-1/2} \mathbf{H} \mathbf{S}^{-1/2} - E_0 \mathbf{I}\| \|\mathbf{S}^{1/2} \boldsymbol{\delta}\|^2 \leq \Delta E \epsilon \|\mathbf{c}_0\|^2.$$

Plugging $\tilde{\mathbf{c}}_0$ into (4.1) and applying the previous two displays leads immediately to the stated result. \square

Theorem 4.2 yields an alternative analysis of the QSD algorithm with thresholding: Combine Theorem 3.1 with $\epsilon = 0$ (to measure the Rayleigh–Ritz error in isolation) together with Theorem 4.2 (to measure the thresholding error). Compared to using Theorem 3.1 alone, this alternative analysis has the advantage that the two types of error (Rayleigh–Ritz and thresholding) can be bounded independently of each other. However, to use Theorems 3.1 and 4.2 in this way, an additional piece of information is needed beyond what is required by Theorem 3.1 alone, namely the norm $\|\mathbf{c}_0\|$.⁴ See SM4 for a comparison of these two approaches on some numerical examples.

5. Numerical Experiments. In this section, we present some numerical experiments demonstrating the success of the theory in explaining the performance of QSD and other features of the numerical performance of the QSD method. We consider two sets of examples: the one-dimensional transverse field Ising model (TFIM) and the one-dimensional Hubbard model.

The Hamiltonian of the TFIM with L spins is

$$(5.1) \quad \widehat{\mathbf{H}} = - \sum_{i=1}^L \widehat{\mathbf{Z}}_i \widehat{\mathbf{Z}}_{i+1} - g \sum_{i=1}^L \widehat{\mathbf{X}}_i,$$

where $\widehat{\mathbf{X}}_i, \widehat{\mathbf{Z}}_i$ are the Pauli X, Z operators acting on the i th spin, respectively. We use the periodic boundary condition, and therefore $\widehat{\mathbf{Z}}_{L+1}$ is identified with $\widehat{\mathbf{Z}}_1$.

The Hamiltonian for the (spinful) Hubbard model of L sites is

$$(5.2) \quad \widehat{\mathbf{H}} = - \sum_{i=1}^L \sum_{\sigma \in \{\uparrow, \downarrow\}} \widehat{\mathbf{a}}_{i\sigma}^* \widehat{\mathbf{a}}_{i+1, \sigma} + U \sum_{i=1}^L \widehat{\mathbf{a}}_{i\uparrow}^* \widehat{\mathbf{a}}_{i\uparrow} \widehat{\mathbf{a}}_{i\downarrow}^* \widehat{\mathbf{a}}_{i\downarrow}.$$

Here $\widehat{\mathbf{a}}_{i\sigma}^*, \widehat{\mathbf{a}}_{i\sigma}$ are the fermionic creation and annihilation operators at site i with spin σ , which can be expressed in terms of spin operators following the Jordan–Wigner transformation (see, e.g., [23]). Similarly due to periodic boundary conditions, $\widehat{\mathbf{a}}_{L+1, \sigma}^*, \widehat{\mathbf{a}}_{L+1, \sigma}$ are identified with $\widehat{\mathbf{a}}_{1\sigma}^*, \widehat{\mathbf{a}}_{1\sigma}$, respectively. We are interested in finding the ground state energy of $\widehat{\mathbf{H}}$. The dimension of $\widehat{\mathbf{H}}$ of the TFIM is 2^L and that for the Hubbard model is 2^{2L} (due to spin degrees of freedom). Due to the high

⁴The norm of the smallest eigenvector of the perturbed thresholded pair $(\widetilde{\mathbf{H}}, \widetilde{\mathbf{S}})$ can provide a good *a posteriori* estimate of this quantity in the limit of a small perturbation (cf. [16, §6.2]).

	TFIM ($L = 10$)	Hubbard ($L = 10$)
Time step Δt	1.0	0.1
Initial overlap $ \gamma_0 ^2$	0.079	0.122
Exact E_0	-15.9799750	-3.31499673
Noiseless E_0 with $n = 40$	-15.9799748	-3.31499670

Table 1: Some parameters for the TFIM and the Hubbard models. Exact E_0 is obtained by diagonalizing $\widehat{\mathbf{H}}$. Noiseless E_0 is obtained by solving the projected generalized eigenvalue problem (of size $n \times n$) using the QSD algorithm in the noiseless setting.

dimensionality, these models are generally intractable to be solved directly when L is large.⁵ In all examples below, we set $L = 10$, $g = -\sqrt{2}$ for the TFIM model and $L = 10$, $U = 8$ for the Hubbard model.⁶ The Hubbard model also has an extra parameter called the total number of fermions denoted by N_e , which constrains $\widehat{\mathbf{H}}$ to a smaller diagonal block, and its value is set to the half filling with $N_e = L = 10$. We find that the numerical results do not depend sensitively to these parameters. Additional numerical results with other values of L , U , etc. can be found in the Supplementary Materials.

The number of time steps used shall be denoted by n , and the time grid is $t_j = j\Delta t$ where $j = 0, \dots, n-1$. The initial vector φ_0 of the TFIM model is taken to be a product state (an eigenstate with $g = 0$), and that of the Hubbard model is taken to be a Slater determinant state (an eigenstate with $U = 0$), respectively. We find that such a setup leads to a sufficiently large initial overlap $|\gamma_0|^2 = |\langle \varphi_0 | \psi_0 \rangle|^2$. This ensures that in the noiseless setting, the ground state energy of $\widehat{\mathbf{H}}$ can be estimated to high accuracy with a very modest value of n (see Table 1).

The Hamiltonian simulation and the computation of the projected matrix elements are performed using the QuSpin package [37, 38] in Python on a classical computer. All further experiments are performed in MATLAB, with the noise matrices $\Delta_{\mathbf{H}}$ and $\Delta_{\mathbf{S}}$ modeled as complex Gaussian Hermitian-Toeplitz matrices with the entries in the first rows independent with specified variances.

5.1. The Need for Thresholding. We first demonstrate why we advocate the use of thresholding by showing the potential pitfalls of some other strategies. Naturally, the first strategy one might attempt would be to do nothing at all: just solve the generalized eigenvalue problem

$$(5.3) \quad \widetilde{\mathbf{H}}\widetilde{\mathbf{c}} = \widetilde{E}\widetilde{\mathbf{S}}\widetilde{\mathbf{c}}$$

and return the least computed eigenvalue. The futility of this strategy is shown in Figure 2. Even for extremely low noise levels ($\sigma \approx 10^{-10}$), we see that the recovered least eigenvalue can deviate quite far from the genuine least eigenvalue with high probability. For nicer problem instances (see, e.g., Figure 2a), characterized by an only

⁵There exists special techniques that are particularly efficient for handling one-dimensional quantum systems. The main point of our numerical results is to demonstrate the performance of QSD algorithms, which does not rely on such special properties.

⁶In particular, $U = 8.0$ for the Hubbard model corresponds to a strongly correlated quantum system and is typically considered to be difficult.

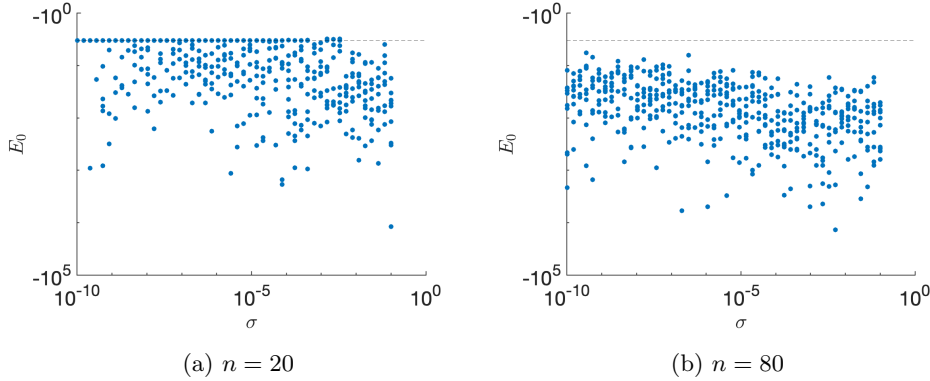


Fig. 2: Least eigenvalues computed from the perturbed pair $(\widetilde{\mathbf{H}}, \widetilde{\mathbf{S}})$ without any procedure to ameliorate the affects of noise. Shown are 10 random initializations of the noise for several random noise levels σ for the Hubbard example with $n = 20$ (left) and $n = 80$ (right). The true eigenvalue is shown as a horizontal dashed line.

modestly ill-conditioned (e.g., $\kappa(\mathbf{S}) = \|\mathbf{S}\| \|\mathbf{S}^{-1}\| \lesssim 10^{12}$) \mathbf{S} matrix, the eigenvalue of interest could perhaps be reliably recovered by taking a median over multiple trials (each of which requiring the quantum computation to be re-run). However, for problem instances with a more ill-conditioned \mathbf{S} matrix, such as is shown in Figure 2b, the probability of finding an eigenvalue close to the genuine smallest eigenvalue appears to occur with vanishingly small probability.

Alternately, one might try to apply “just a bit of thresholding” by setting the threshold parameter at a small constant value, independent of the noise level. This can be modestly effective for a well-conditioned \mathbf{S} matrix (particularly if combined with a median of multiple trials), but it falls down as soon as $\sigma \gtrsim \epsilon$ in general. See Figure SM6 for a demonstration of this. If one is to rely on thresholding alone to deal with the noise, then threshold parameter must be chosen large enough.

As another alternative to thresholding, one might attempt to solve the problem without explicitly filtering out the noise by thresholding (or only using a tiny threshold much smaller than the noise level) and attempting to systematically determine which eigenvalues are “real”. We investigate such strategies in section SM2 and ultimately conclude that these are less robust and less accurate than thresholding.

5.2. Choice of the Thresholding Parameter. Now that we have demonstrated why we prefer thresholding for solving the noise-perturbed QSD generalized eigenvalue problem (5.3), let us demonstrate the success of thresholding. In Figure 3, we demonstrate the error for the thresholding procedure with a threshold parameter ϵ proportional to the noise level.⁷ (See also Figure SM7 for more examples.) As the plots show, thresholding is robust on these examples in the sense that the maximum error over multiple trials is similar to the median, showing that thresholding is reliably able to filter out the noise over different random initializations.⁸ Since the norm of

⁷This is smaller than the theory (Theorem 2.5 and Corollary 2.2) predicts ϵ should be taken as, which suggests that one should choose $\epsilon \propto \sigma^{1/(1+\alpha)}$ ($\alpha \leq 1/2$ is the value for which (2.12) holds).

⁸This is seen to be not true for one example (Figure SM7d) in the Supplementary Material. This shows that thresholding is not infallible, but is still better than alternate strategies within our

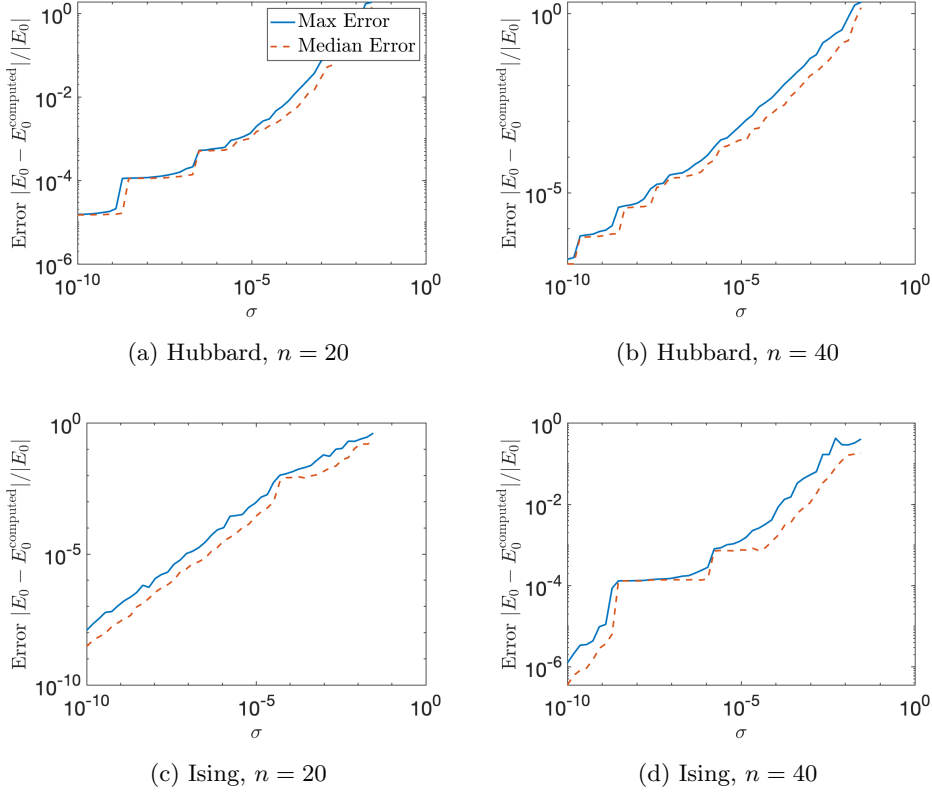


Fig. 3: Maximum (blue solid) and median (red dashed) error over 100 initializations for eigenvalues computed from the noise-perturbed pair $(\widetilde{\mathbf{H}}, \widetilde{\mathbf{S}})$ using thresholding with threshold parameter $25\sigma\|\widetilde{\mathbf{S}}\|$ for Hubbard model (top) and Ising model (bottom) with $n = 20$ (left) and 40 (right).

the noise matrix is a random quantity prone to occasionally being appreciably larger than its average value, the threshold parameter must be somewhat larger than the expected noise level to achieve good performance with high probability. These plots demonstrate that, for these examples at least, this multiple can be quite modest—just 25 is enough. In the error plots in Figure 3 we see two types of behavior: Sometimes the error increases relatively continuously with the thresholding parameter (e.g., most of Figure 3c) whereas other times it increases in a more stepwise fashion (e.g., low noise levels in Figure 3a). The first behavior is indicative of the error being dominated by the noise level (with the slope in log-log space being ≈ 1 demonstrating a linear dependence of the error on the noise level) with the second exemplifying the thresholding error being the dominant contribution. Behavior of the second type might suggest that the threshold parameter is being chosen conservatively and lower error could be achieved with a smaller threshold value.

The choice of the thresholding parameter is critical to the success of the method,

knowledge (section SM2).

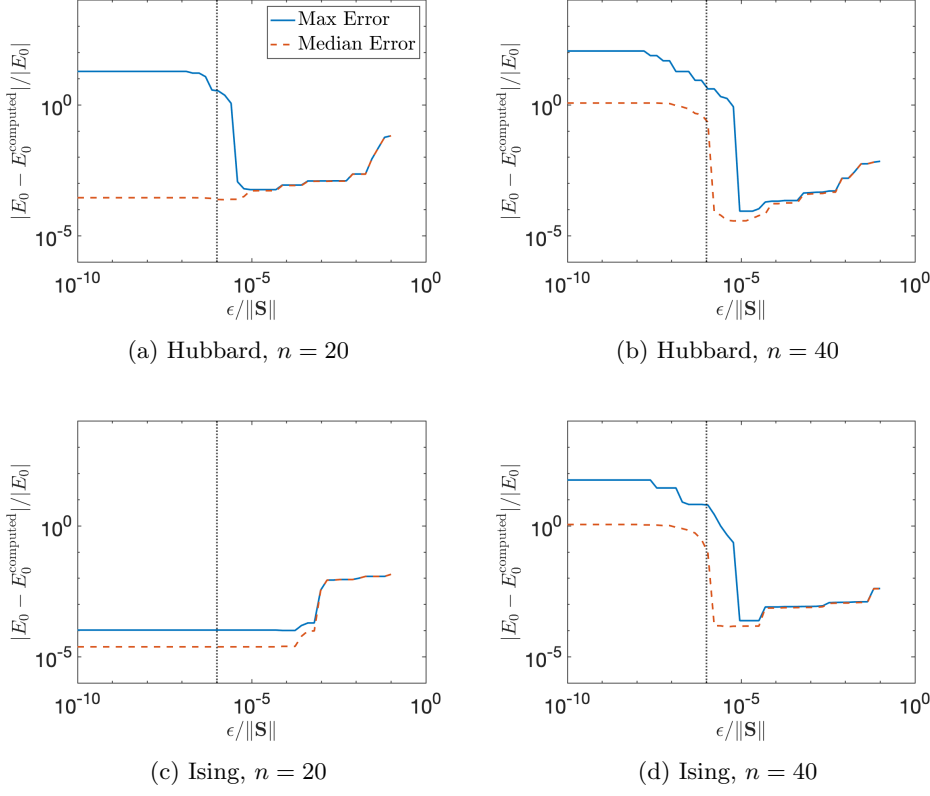


Fig. 4: Maximum (blue solid) and median (red dashed) error over 100 initializations for eigenvalues computed from the noise-perturbed pair $(\widetilde{\mathbf{H}}, \widetilde{\mathbf{S}})$ using thresholding for various values of the threshold ϵ for a fixed noise level $\sigma = 10^{-6}$ (dotted black line) for Hubbard model (top) and Ising model (bottom) with $n = 20$ (left) and 40 (right).

as is demonstrated in Figure 4. In these plots, we show the median and maximum error for the computed eigenvalue over 100 random initializations of the noise (fixed to $\sigma = 10^{-6}$) for different thresholding levels. In the normative case, the error decreases as the threshold parameter is decreased up until the threshold parameter reaches a modest multiple of the noise level, after which the error sharply rises. In the Supplementary Materials (section SM3), we discuss an automatically tuned variant of the thresholding procedure which can help in picking a good threshold level ϵ .

6. Conclusions. In this article, we have presented the first theoretical analysis of the accuracy of the quantum subspace diagonalization method with a thresholding procedure. Our explanation has two parts:

1. With an appropriate choice of time sequence $\{t_j\}$, QSD with thresholding (Algorithm 1.1) can compute an accurate approximation to the ground-state energy. (Theorem 3.1)
2. Under appropriate conditions, the thresholding procedure can robustly determine the smallest eigenvalue in the presence of perturbations to the pair

(\mathbf{H}, \mathbf{S}) . (Theorem 2.7)

These two pieces combine to give a bound on the total error of the QSD procedure (comprising Rayleigh–Ritz, thresholding, and perturbation errors) in Informal Theorem 1.3. The conditions of our theory are natural, and many of the parameters in our bounds can be estimated in the presence of noise, allowing our bounds to be able to give approximate bounds on the error *a posteriori*. Our numerical experiments (including additional experiments in the Supplementary Material) support the conclusion that QSD is accurate when implemented with thresholding (and not accurate when implemented without).

Our theoretical estimates can still be significantly improved, such as the bound $\chi_{\mathbf{H}} \leq \mathcal{O}(\eta_{\mathbf{S}}/\epsilon^\alpha + \eta_{\mathbf{H}})$ (see Theorem 2.5) for the discrepancy between the thresholded \mathbf{H} and $\widetilde{\mathbf{H}}$ matrices, where $0 \leq \alpha \leq 1/2$. This suggests that we need to take $\epsilon = \Omega(\eta_{\mathbf{S}}^{1/(1+\alpha)})$ for accurate recovery to be guaranteed by Corollary 2.2, which has hypothesis $\chi \leq \epsilon/q$. This is in contradiction to our numerical experiments, where ϵ can be chosen to be a small multiple of η . Synthetically generated worst-case examples (see section SM5) suggests the bound $\chi_{\mathbf{H}} \leq \mathcal{O}(\eta_{\mathbf{S}}/\epsilon^\alpha + \eta_{\mathbf{H}})$ is tight, but it remains possible a better bound can be derived for pairs (\mathbf{H}, \mathbf{S}) generated by QSD. An interesting open question is to give a convincing explanation for why we appear to have $\alpha = 1/4$ and $\mu \approx \max |\Lambda(\mathbf{H}, \mathbf{S})|$ in (2.12) for many QSD instances. Despite the modest size of n in practice, the polynomial dependence on n in Corollary 2.2 and Theorems 2.5, 2.6, and 4.1 can lead to significant overestimates of the error. Therefore, another natural question is whether these dimensional factors can be improved.

Natural extensions of this work are to generalized our analysis to excited states (interior eigenvalues) and to develop bounds on the accuracy of the computed eigenvectors. Although the QSD algorithm cannot be used to coherently prepare an eigenstate on a quantum computer, we may still compute other physical observables from the approximate eigenstate. Mathias and Li’s theory [19, §6] suggests eigenvectors might be more sensitive to the noise than the eigenvalues, and we plan to study the accuracy of the computed eigenvectors in future work.

Acknowledgments. This work is supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, Department of Energy Computational Science Graduate Fellowship under Award Number DE-SC0021110 (ENE), by the Department of Energy under Grant No. DE-SC0017867, and by the NSF Quantum Leap Challenge Institute (QLCI) program through grant number OMA-2016245 (LL). LL is a Simons Investigator. ENE thanks the Lawrence Berkeley Laboratory summer student program for providing a welcoming environment to perform this work. We thank Yulong Dong, Yu Tong, Norman Tubman, and Robert Webber for helpful discussions.

Disclaimer. This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor any agency thereof, nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof. The views and opinions of authors expressed herein do not necessarily state or reflect those of the

United States Government or any agency thereof.

REFERENCES

- [1] R. BHATIA, *Matrix Analysis*, vol. 169 of Graduate Texts in Mathematics, Springer-Verlag, New York, 1997, <https://doi.org/10.1007/978-1-4612-0653-8>. 10, 20, 29
- [2] Å. BJÖRCK, *Numerical Methods in Matrix Computations*, vol. 59 of Texts in Applied Mathematics, Springer International Publishing, 2015, <https://doi.org/10.1007/978-3-319-05089-8>. 17
- [3] Y. CAI, Z. BAI, J. E. PASK, AND N. SUKUMAR, *Hybrid preconditioning for iterative diagonalization of ill-conditioned generalized eigenvalue problems in electronic structure calculations*, J. Comput. Phys., 255 (2013), pp. 16–30. 1
- [4] A. M. CHILDS, Y. SU, M. C. TRAN, N. WIEBE, AND S. ZHU, *Theory of Trotter Error with Commutator Scaling*, Phys. Rev. X, 11 (2021), p. 011020, <https://doi.org/10.1103/PhysRevX.11.011020>. 3
- [5] J. I. COLLESS, V. V. RAMASESH, D. DAHLEN, M. S. BLOK, M. E. KIMCHI-SCHWARTZ, J. R. MCCLEAN, J. CARTER, W. A. DE JONG, AND I. SIDDIQI, *Computation of Molecular Spectra on a Quantum Processor with an Error-Resilient Algorithm*, Phys. Rev. X, 8 (2018), p. 011021, <https://doi.org/10.1103/PhysRevX.8.011021>. 3
- [6] C. L. CORTES AND S. K. GRAY, *Quantum Krylov subspace algorithms for ground and excited state energy estimation*, arXiv:2109.06868, (2021), <https://arxiv.org/abs/2109.06868>. 1, 3
- [7] C. DAVIS AND W. M. KAHAN, *The Rotation of Eigenvectors by a Perturbation. III*, SIAM J. Numer. Anal., 7 (1970), pp. 1–46, <https://doi.org/10.1137/0707001>. 10
- [8] P. DRINEAS AND I. C. F. IPSEN, *Low-Rank Matrix Approximations Do Not Need a Singular Value Gap*, SIAM J. Matrix Anal. Appl., 40 (2019), pp. 299–319, <https://doi.org/10.1137/18M1163658>. 20
- [9] A. EMAMI-NAEINI AND P. VAN DOOREN, *Computation of zeros of linear multivariable systems*, Automatica, 18 (1982), pp. 415–430, [https://doi.org/10.1016/0005-1098\(82\)90070-X](https://doi.org/10.1016/0005-1098(82)90070-X). 1
- [10] G. FIX AND R. HEIBERGER, *An Algorithm for the Ill-Conditioned Generalized Eigenvalue Problem*, SIAM J. Numer. Anal., 9 (1972), pp. 78–88, <https://doi.org/10.1137/0709009>. 4
- [11] G. H. GOLUB AND C. F. VAN LOAN, *Matrix computations*, Johns Hopkins Univ. Press, Baltimore, fourth ed., 2013. 3
- [12] L. A. GRIBOV AND B. K. NOVOSADOV, *Use of overcomplete basis sets in quantum-chemical calculations*, Journal of Molecular Structure: THEOCHEM, 136 (1986), pp. 387–389, [https://doi.org/10.1016/0166-1280\(86\)80152-X](https://doi.org/10.1016/0166-1280(86)80152-X). 4
- [13] W. J. HUGGINS, J. LEE, U. BAEK, B. O’GORMAN, AND K. B. WHALEY, *A non-orthogonal variational quantum eigensolver*, New Journal of Physics, 22 (2020), p. 073009, <https://doi.org/10.1088/1367-2630/ab867b>. 1
- [14] M. JUNGEN AND K. KAUFMANN, *The Fix-Heiberger procedure for solving the generalized ill-conditioned symmetric eigenvalue problem*, Int. J. Quantum Chem., 41 (1992), pp. 387–397. 1
- [15] K. KLYMKO, C. MEJUTO-ZAERA, S. J. COTTON, F. WUDARSKI, M. URBANEK, D. HAIT, M. HEAD-GORDON, K. B. WHALEY, J. MOUSSA, N. WIEBE, W. A. DE JONG, AND N. M. TUBMAN, *Real time evolution for ultracompact Hamiltonian eigenstates on quantum hardware*, arXiv:2103.08563, (2021), <https://arxiv.org/abs/2103.08563>. 1, 3, 5, 15
- [16] M. LOTZ AND V. NOFERINI, *Wilkinson’s Bus: Weak Condition Numbers, with an Application to Singular Polynomial Eigenproblems*, Found. Comput. Math., 20 (2020), pp. 1439–1473, <https://doi.org/10.1007/s10208-020-09455-y>. 4, 21
- [17] P.-O. LÖWDIN, *Group Algebra, Convolution Algebra, and Applications to Quantum Mechanics*, Reviews of Modern Physics, 39 (1967), pp. 259–287, <https://doi.org/10.1103/RevModPhys.39.259>. 4
- [18] V. A. MANDELSHTAM AND H. S. TAYLOR, *Harmonic inversion of time signals and its applications*, J. Chem. Phys., 107 (1997), pp. 6756–6769, <https://doi.org/10.1063/1.475324>. 2
- [19] R. MATHIAS AND C.-K. LI, *The definite generalized eigenvalue problem: A new perturbation theory*, T-NAREP No. 457, inst-MCCM, inst-MCCM:adr, Oct. 2004. 3, 8, 9, 26
- [20] J. R. MCCLEAN, M. E. KIMCHI-SCHWARTZ, J. CARTER, AND W. A. DE JONG, *Hybrid quantum-classical hierarchy for mitigation of decoherence and determination of excited states*, Phys. Rev. A, 95 (2017), p. 042308, <https://doi.org/10.1103/PhysRevA.95.042308>. 1, 3
- [21] M. MOTTA, C. SUN, A. T. K. TAN, M. J. O’ROURKE, E. YE, A. J. MINNICH, F. G. S. L. BRANDÃO, AND G. K.-L. CHAN, *Determining eigenstates and thermal states on a quantum computer using quantum imaginary time evolution*, Nature Physics, 16 (2020), pp. 205–

- 210, <https://doi.org/10.1038/s41567-019-0704-4>. 1, 3
- [22] G. NANNICINI, *An introduction to quantum computing, without the physics*, SIAM Rev., 62 (2020), pp. 936–981. 2
 - [23] J. W. NEGELE AND H. ORLAND, *Quantum many-particle systems*, Westview, 1988. 21
 - [24] M. A. NIELSEN AND I. CHUANG, *Quantum computation and quantum information*, 2000. 2, 3
 - [25] B. N. PARLETT, *The Symmetric Eigenvalue Problem*, Classics in Applied Mathematics, SIAM, 1998, <https://doi.org/10.1137/1.9781611971163>. 2, 3, 31
 - [26] R. M. PARRISH AND P. L. MCMAHON, *Quantum Filter Diagonalization: Quantum Eigendecomposition without Full Quantum Phase Estimation*, arXiv:1909.08925, (2019), <https://arxiv.org/abs/1909.08925>. 1, 2, 3
 - [27] J. PRESKILL, *Quantum computing in the NISQ era and beyond*, Quantum, 2 (2018), p. 79. 1
 - [28] J. PRESKILL, *Quantum computing 40 years later*, arXiv:2106.10522, (2021). 1
 - [29] Y. SAAD, *On the Rates of Convergence of the Lanczos and the Block-Lanczos Methods*, SIAM J. Numer. Anal., 17 (1980), pp. 687–706, <https://doi.org/10.1137/0717059>. 15, 16
 - [30] I. SABZEVARI, A. MAHAJAN, AND S. SHARMA, *An accelerated linear method for optimizing non-linear wavefunctions in variational Monte Carlo*, J. Chem. Phys., 152 (2020), p. 024111, <https://doi.org/10.1063/1.5125803>. 2
 - [31] K. SEKI AND S. YUNOKI, *Quantum power method by a superposition of time-evolved states*, Phys. Rev. X Quantum, 2 (2021), p. 010333. 1
 - [32] N. H. STAIR, R. HUANG, AND F. A. EVANGELISTA, *A Multireference Quantum Krylov Algorithm for Strongly Correlated Electrons*, J. Chem. Theory Comput., 16 (2020), pp. 2236–2245, <https://doi.org/10.1021/acs.jctc.9b01125>. 1, 2, 3, 15
 - [33] G. W. STEWART, *Perturbation bounds for the definite generalized eigenvalue problem*, Linear Algebra Appl., 23 (1979), pp. 69–85, [https://doi.org/10.1016/0024-3795\(79\)90094-6](https://doi.org/10.1016/0024-3795(79)90094-6). 7, 8
 - [34] G. W. STEWART AND J.-G. SUN, *Matrix Perturbation Theory*, Computer Science and Scientific Computing, Academic Press, 1st edition ed., 1990. 3, 9, 10, 11, 12, 18, 20, 29
 - [35] J. A. TROPP, *An introduction to matrix concentration inequalities*, Foundations and Trends in Machine Learning, 8 (2015), pp. 1–230. 19
 - [36] M. R. WALL AND D. NEUHAUSER, *Extraction, through filter-diagonalization, of general quantum eigenvalues or classical normal mode frequencies from a small number of residues or a short-time segment of a signal. I. Theory and application to a quantum-dynamics model*, The Journal of Chemical Physics, 102 (1995), pp. 8011–8022, <https://doi.org/10.1063/1.468999>. 2
 - [37] P. WEINBERG AND M. BUKOV, *QuSpin: a Python Package for Dynamics and Exact Diagonalisation of Quantum Many Body Systems part I: spin chains*, SciPost Phys., 2 (2017), p. 003. 22
 - [38] P. WEINBERG AND M. BUKOV, *QuSpin: a Python Package for Dynamics and Exact Diagonalisation of Quantum Many Body Systems. Part II: bosons, fermions and higher spins*, SciPost Phys., 7 (2019), p. 20. 22
 - [39] J. H. WILKINSON, *Kronecker’s canonical form and the QZ algorithm*, Linear Algebra Appl., 28 (1979), pp. 285–303, [https://doi.org/10.1016/0024-3795\(79\)90140-X](https://doi.org/10.1016/0024-3795(79)90140-X). 3, 4

SUPPLEMENTARY MATERIAL

SM1. Proof of Theorem 4.1. For reference, we provide a complete proof of Theorem 4.1.

Proof of Theorem 4.1. Let Π and $\tilde{\Pi}$ be the spectral projectors onto the dominant m -dimensional invariant subspaces of \mathbf{A} and $\mathbf{A} + \Delta$ respectively. First, we bound

$$\begin{aligned} \|\llbracket \mathbf{A} + \Delta \rrbracket_m - \llbracket \mathbf{A} \rrbracket_m\|_{\text{QUI}} &= \|\tilde{\Pi}(\mathbf{A} + \Delta)\tilde{\Pi} - \Pi\mathbf{A}\Pi\|_{\text{QUI}} \\ &\leq \|\tilde{\Pi}\Delta\tilde{\Pi}\|_{\text{QUI}} + \|\tilde{\Pi}\mathbf{A}\tilde{\Pi} - \Pi\mathbf{A}\Pi\|_{\text{QUI}}. \end{aligned}$$

For the first term, we bound $\|\tilde{\Pi}\Delta\tilde{\Pi}\|_{\text{QUI}} \leq \|\Delta\|_{\text{QUI}}$ using the fact $\|\cdot\|_{\text{QUI}}$ is *symmetric* in the sense that $\|\mathbf{B}_1\mathbf{B}_2\mathbf{B}_3\|_{\text{QUI}} \leq \|\mathbf{B}_1\| \|\mathbf{B}_2\|_{\text{QUI}} \|\mathbf{B}_3\|$ [34, Thm. 3.9]. Every quadratic unitarily invariant norm is bounded by the Frobenius norm, which follows from the definition of quadratic unitarily invariant norm and the fact that the nuclear norm bounds every unitarily invariant norm [1, Eq. (IV.38)]. Thus, the second term is bounded as

$$\|\tilde{\Pi}\mathbf{A}\tilde{\Pi} - \Pi\mathbf{A}\Pi\|_{\text{QUI}} \leq \|\tilde{\Pi}\mathbf{A}\tilde{\Pi} - \Pi\mathbf{A}\Pi\|_{\text{F}},$$

which is then bounded by Theorem 2.6 with $\epsilon = \lambda_{m+1} + \|\Delta\|$ and $\rho = (\lambda_m - \lambda_{m+1} - \|\Delta\|)/(\lambda_{m+1} - \|\Delta\|)$. \square

SM2. The Failure of Heuristics. In the main text, we show how accurate recovery for the QSD algorithm usually fails, if one solves the noisy generalized eigenvalue problem (5.3) with no special treatment. A forthcoming example (Figure SM6) shows that using a fixed threshold independent of the noise level may not perform much better. An alternate strategy, which we initially believed to be more promising than thresholding, is to compute the eigenvalues (either with no thresholding at all or with a small threshold independent of the noise level) and attempt to determine real from spurious eigenvalues by means of some property of the computed eigenvector. In this section, we shall consider a couple variants of such an approach and ultimately conclude the performance of these heuristics can still be unsatisfactory when compared to thresholding.

Two natural heuristics for whether \tilde{E} is a plausible candidate for the ground state energy suggest themselves. Let $\tilde{\mathbf{c}}$ be the unit-norm eigenvector associated with a computed eigenvalue \tilde{E} . The Ritz vector $\tilde{\psi}_0 := \sum_{j=0}^{n-1} \tilde{\mathbf{c}}_j \boldsymbol{\varphi}_j$ is supposed to be close to the true ground-state eigenvector ψ_0 of $\tilde{\mathbf{H}}$. Our heuristics are as follows:

1. **Require $h_1 := \tilde{\mathbf{c}}^* \tilde{\mathbf{S}} \tilde{\mathbf{c}}$ to be large.** The squared norm of the Ritz vector is precisely $\tilde{\mathbf{c}}^* \tilde{\mathbf{S}} \tilde{\mathbf{c}} \approx h_1$. If h_1 is small, then the norm of the Ritz vector is very small due to cancellations in the sum $\sum_{j=0}^{n-1} \tilde{\mathbf{c}}_j \boldsymbol{\varphi}_j$ and should thus be treated as suspect because of the noise. Thus, it is natural to insist on a large value of h_1 .⁹
2. **Require the estimated overlap $h_2 := |\mathbf{e}_0^* \tilde{\mathbf{S}} \tilde{\mathbf{c}}| \approx |\boldsymbol{\varphi}_0^* \tilde{\psi}_0|$ to be large.** It is important that the initial vector $\boldsymbol{\varphi}_0$ has a relatively large initial overlap $|\boldsymbol{\varphi}_0^* \psi_0|$ with the eigenvector of interest—indeed, our analysis suggests accurate recovery of the ground-state energy requires this (see Theorem 3.1). As such, it is natural to use the overlap (or its surrogate h_1 computable from the noise-corrupted $\tilde{\mathbf{S}}$ matrix) as a measure of whether an eigenvalue is a genuine candidate for the ground-state energy. Note that by unit-norm scaling

⁹ h_1 is also related to the conditioning of the eigenvalue [34, Eq. (VI.2.2)].

Algorithm SM3.1 Automatically tuned thresholding procedure for finding the least eigenvalue of a noise-corrupted generalized eigenvalue problem.

```

procedure AUTOTHRESHOLDING( $\mathbf{H}, \mathbf{S}, \epsilon_0, r$ )
   $E \leftarrow \text{THRESHOLDING}(\mathbf{H}, \mathbf{S}, \epsilon_0)$ 
   $\Lambda \leftarrow \{\lambda \in \text{eig}(\mathbf{S}) : \lambda < \epsilon_0\}$ 
  while  $\Lambda \neq \emptyset$  do
     $\epsilon \leftarrow \max \Lambda, \Lambda \leftarrow \Lambda \setminus \{\epsilon\}$ 
     $E' \leftarrow \text{THRESHOLDING}(\mathbf{H}, \mathbf{S}, \epsilon)$ 
    if  $|E - E'| / \min(|E|, |E'|) > r$  then
      break
    end if
     $E \leftarrow E'$ 
  end while
  return  $E$ 
end procedure

```

$\tilde{\mathbf{c}}$ (rather than adopting the normalization $\tilde{\mathbf{c}}^* \tilde{\mathbf{S}} \tilde{\mathbf{c}} = 1$), we are implicitly also incorporating the condition for $\tilde{\psi}_0$ to be a stable linear combination of the basis states which motivated our interest in h_1 .

There are several ways of using a heuristic $h \in \{h_1, h_2\}$ as an algorithm for computing the ground-state eigenvalue: (a) pick E with the highest h , (b) pick the smallest E of the eigenvalues with the top k values of h , and (c) pick the smallest E with h above some thresholding h_0 (or simply the largest h if none exceeds h_0).

Unfortunately, unlike thresholding where there is a natural choice of the threshold parameter (related to the noise level η (2.1) which can usually be reliably estimated), we are unaware of any good systematic ways to pick the parameters k and h_0 for strategies (b) and (c). These heuristics thus usually require some tuning to make them accurate for a given problem instance, with the parameters needing to be readjusted when a new problem is encountered. This reduces the reliability of these heuristics, when the ground truth is unavailable to compare against. The robustness of heuristics such as (a), (b), and (c) can be improved by medians of repeated trials or by comparing the results of different heuristics against each other. However, even with such improvements, without rigorous guarantees, the validity of these heuristics remains conjectural when the genuine ground-state energy is unavailable to be validated against.

Figure SM1 shows the suggested heuristics (a), (b), and (c) with the figures of merit h_1 and h_2 with $k = 5$ and $h_0 = 10^{-2} \|\tilde{\mathbf{S}}\|$. The first subfigure, Figure SM1a, shows a relatively optimistic case for the heuristics. For low levels of noise, the eigenvalue is generally recovered with low error with the exception of a few outliers which could be ameliorated by the median trick. Figure SM1b shows the potential danger of applying these heuristics; despite working well for the Hubbard example with $n = 20$ in Figure SM1a, the heuristics fail with the same parameter choices for the Ising example with $n = 40$. For this problem, the eigenvalues are observed to be recovered accurately only with very small probability. Improvements to both plots are likely possible by more careful choice of the heuristic parameters or more complicated heuristics, but this is a point against such heuristics rather than for them: we ideally want a method which works well without tuning problem-dependent parameters.

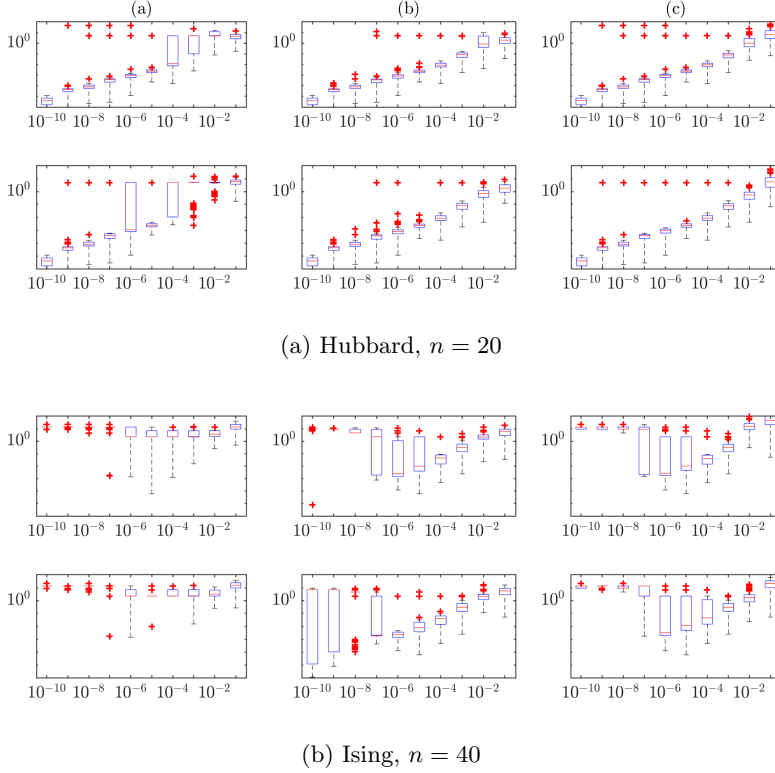


Fig. SM1: Errors (vertical axis) for eigenvalues computed from the perturbed pair $(\widetilde{\mathbf{H}}, \widetilde{\mathbf{S}})$ with heuristics (a), (b), and (c) described in the text for quality metrics h_1 (first row) and h_2 (second row). The noisy generalized eigenvalue problem (5.3) is solved using thresholding with a fixed threshold parameter $\epsilon = 10^{-12} \|\widetilde{\mathbf{S}}\|$. Shown are 100 random initializations of the noise for several random noise levels σ (horizontal axis) for the Hubbard example with $n = 20$ (top) and Ising example with $n = 40$ (bottom).

SM3. Automatic Thresholding. As we saw in the main text, choosing a good maximum thresholding level ϵ is critical to the success of the thresholding procedure Algorithm 1.1. A useful “sanity check” is thus to solve the problem using a handful of plausible thresholding parameters to make sure the computed eigenvalues are close to each other. (A variant of this strategy is proposed by Parlett for the Fix–Heiberger procedure [25, §15.5].) A more ambitious strategy is to solve with a range of thresholding parameters beginning with a conservative (but not comically large) threshold parameter ϵ_0 and then tuning it down until the eigenvalue “jumps” to a presumably spurious value. The best approximation to the ground-state energy suggested by this procedure is the last value before this jump. If one wishes to automate this procedure, one needs to have a mechanistic way of deciding whether a jump has occurred: For this purpose, we shall test whether the relative difference exceeds a cutoff r . This procedure is demonstrated in Algorithm SM3.1. The success of this procedure relies

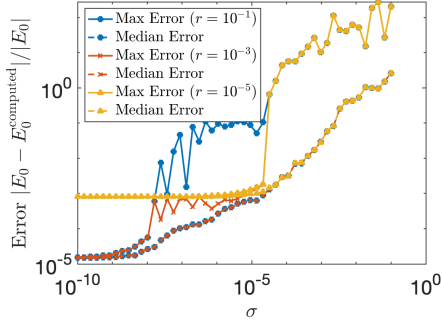
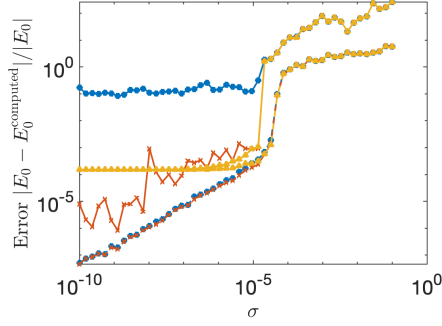
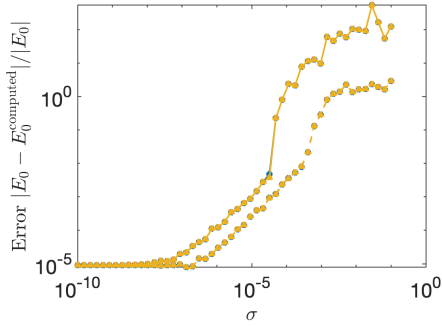
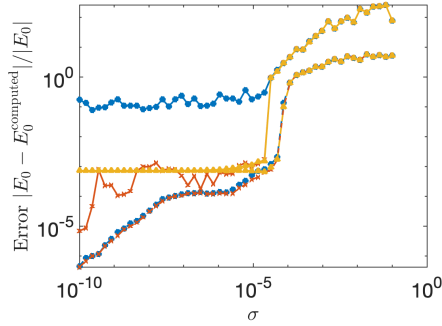
(a) Hubbard, $n = 20$ (b) Hubbard, $n = 40$ (c) Ising, $n = 20$ (d) Ising, $n = 40$

Fig. SM2: Maximum and median error over 100 initializations for eigenvalues computed from the noise-perturbed pair $(\hat{\mathbf{H}}, \hat{\mathbf{S}})$ using automatically tuned thresholding (Algorithm SM3.1) for three cutoffs $r \in \{10^{-1}, 10^{-3}, 10^{-5}\}$ for various values of the noise level σ for Hubbard model (top) and Ising model (bottom) with $n = 20$ (left) and 40 (right).

on the choice of ϵ_0 not being too large as the method uses the eigenvalue recovered with parameter ϵ_0 as a baseline, large deviations from which are characterized as erroneous. Usually, one will have some good estimate of the amount of noise so picking a sensible ϵ_0 should be possible.

The performance of the automatic thresholding procedure Algorithm SM3.1 with three choices of the parameter r are shown in Figure SM2. These plots represent the worst-case situation where the noise level is completely unknown and the choice one has available for ϵ_0 is a constant multiple of $\|\hat{\mathbf{S}}\|$. The best case scenario is shown in the $r = 10^{-3}$ lines in Figures SM2b and SM2c; in these cases, the error decays nicely as the noise does with the procedure being relatively robust (as shown by the error over a maximum over 100 trials being similar to the median). This automatic thresholding procedure can still be somewhat delicate, with the maximum error over 100 runs being near the cutoff r for the $r = 10^{-1}$ in Figures SM2a, SM2b, and SM2c; this suggests, in the worst case, one must be willing to accept an error level on the order r due to overly aggressive automatic tuning of the thresholding

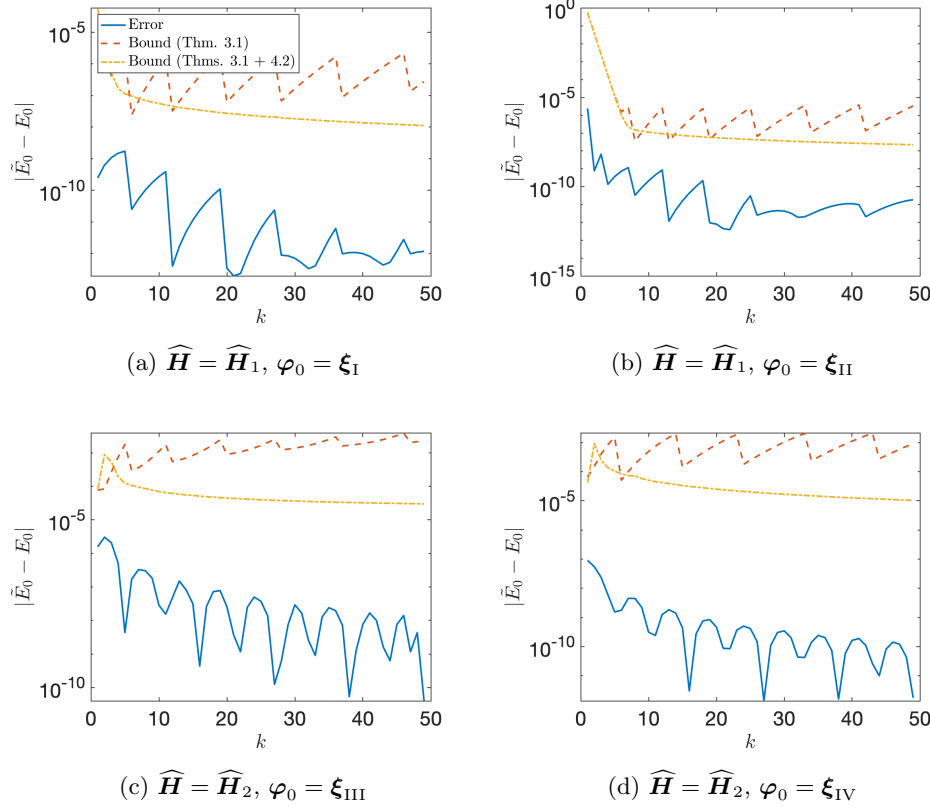


Fig. SM3: Error for QSD with the time sequence from the hypotheses of Theorem 3.1 with threshold parameter $\epsilon = 10^{-6}$ for various values k and four different sets of input data.

parameter. However, these same plots illustrate the importance of not being too cautious either, with the maximum error being $\approx 10^{-3}$ with $r = 10^{-5}$ due to overly conservative automatic tuning of the thresholding parameter. In totality, Figure SM2 shows that the automatic thresholding procedure cannot determine a near-optimal choice for the thresholding parameter in all cases, but it can be useful in “upgrading” an overly cautious threshold parameter ϵ_0 to a better choice for ϵ , obtaining a couple more decimal digits of accuracy in the best case. As a final comment, observe that thresholding is an inherently discrete process since each eigenvalue must either be discarded or not. Therefore, even with automatically tuned thresholding, there can be plateaus in the noise-level vs accuracy curve, owing to the importance of a single eigenvalue to the overall error; this is shown in Figure SM2c.

SM4. Validation of Theorems 3.1 and 4.2. Despite its desirable theoretical implications, Theorem 3.1 may still significantly overestimate the error incurred by thresholding in practice. Consider the following examples:

- (I) $\widehat{\mathbf{H}} = \widehat{\mathbf{H}}_1$ and $\varphi_0 = \xi_{\text{I}}$,
- (II) $\widehat{\mathbf{H}} = \widehat{\mathbf{H}}_1$ and $\varphi_0 = \xi_{\text{II}}$,

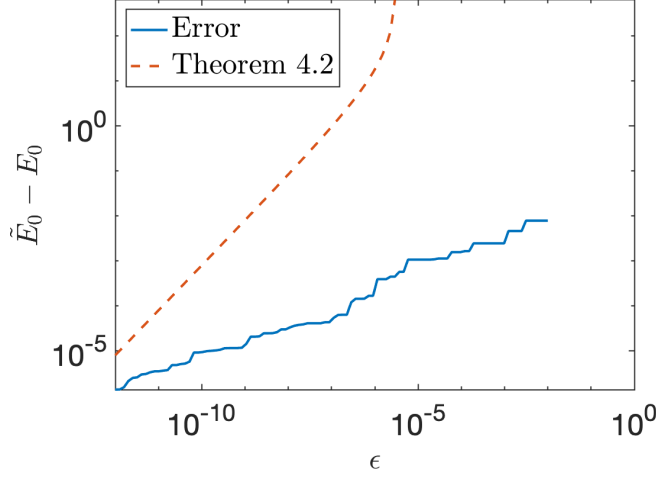


Fig. SM4: Error due to thresholding and error bound from Theorem 4.2 for least eigenvalue computed from a synthetically generated pair (\mathbf{H}, \mathbf{S}) for various threshold parameters ϵ .

$$\text{(III)} \quad \widehat{\mathbf{H}} = \widehat{\mathbf{H}}_2 \text{ and } \varphi_0 = \xi_{\text{III}},$$

$$\text{(IV)} \quad \widehat{\mathbf{H}} = \widehat{\mathbf{H}}_2 \text{ and } \varphi_0 = \xi_{\text{IV}}.$$

where

$$\xi_{\text{I}} = \left(\sqrt{1 - 10^{-4}}, \sqrt{\frac{10^{-4}}{998}}, \sqrt{\frac{10^{-4}}{998}}, \dots, \sqrt{\frac{10^{-4}}{998}} \right) \in \mathbb{C}^{999},$$

$$\xi_{\text{II}} = \left(\sqrt{0.5}, \sqrt{\frac{0.5}{998}}, \sqrt{\frac{0.5}{998}}, \dots, \sqrt{\frac{0.5}{998}} \right) \in \mathbb{C}^{999},$$

$$\xi_{\text{III}} = \left(\sqrt{1 - 10^{-4} - 10^{-8}}, \sqrt{\frac{10^{-4}}{998}}, \sqrt{\frac{10^{-4}}{998}}, \dots, \sqrt{\frac{10^{-4}}{998}}, 10^{-4} \right) \in \mathbb{C}^{1000}$$

$$\xi_{\text{IV}} = \left(1, \frac{0.01}{2}, \frac{0.01}{3}, \dots, \frac{0.01}{1000} \right) / \left\| \left(1, \frac{0.01}{2}, \frac{0.01}{3}, \dots, \frac{0.01}{1000} \right) \right\| \in \mathbb{C}^{1000},$$

and

$$\widehat{\mathbf{H}}_1 = \text{diag} \left(1, 2 + 0 \cdot \frac{0.1}{997}, 2 + 1 \cdot \frac{0.1}{997}, \dots, 2 + 997 \cdot \frac{0.1}{997} \right) \in \mathbb{C}^{999 \times 999},$$

$$\widehat{\mathbf{H}}_2 = \text{diag} \left(1, 2 + 0 \cdot \frac{0.1}{997}, 2 + 1 \cdot \frac{0.1}{997}, \dots, 2 + 997 \cdot \frac{0.1}{997}, 1000 \right) \in \mathbb{C}^{1000 \times 1000}.$$

These examples are artificial: They are engineered to have a large spectral gap ΔE_1 but a small spectral range ΔE_{999} . Even with these artificial examples, Theorem 3.1 (as well as Theorems 3.1 and 4.2 together) still overestimates the error by several orders of magnitude. See Figure SM3.

We consider the bound Theorem 4.2 by itself in Figure SM4. Shown is the error $\tilde{E}_0 - E_0$ between the recovered least eigenvalue \tilde{E}_0 and the true least eigenvalue E_0 for different choices of threshold parameter ϵ . For this example, we used $\mathbf{H} = \mathbf{K}^* \text{diag}(1, 2, \dots, 100) \mathbf{K}$ and $\mathbf{S} = \mathbf{K}^* \mathbf{K}$, where \mathbf{K} is a product of diagonal matrix with (j, j) th entry j^{-2} and a “randsvd” matrix from MATLAB’s gallery with approximate condition number 10^3 . This was chosen to give a fairly ill-conditioned “ \mathbf{S} ” matrix ($\kappa(\mathbf{S}) \approx 10^{12}$) in which the \mathbf{S} -normalized ground state eigenvector of fairly small norm ($\|\mathbf{c}_0\| \approx 10^2$). We find the bound from Theorem 4.2 becomes increasingly conservative as ϵ increases, diverging to $+\infty$ at $\epsilon \approx 10^{-6}$ which the true error $\tilde{E}_0 - E_0$ remains bounded $\leq 10^{-2}$ for $\epsilon \leq 10^{-2}$.

SM5. Evidence for Tightness of Theorem 2.5. First, we present a synthetically generated numerical example which suggests that the η/ϵ^α behavior in Theorem 2.5 is necessary, at least without further assumptions. As our example, we set $\mathbf{A} = (\mathbf{G} + \mathbf{G}^*)/2$ to be the Hermitian part of a 5×5 real standard Gaussian matrix \mathbf{G} and pick $\mathbf{S} = \text{diag}(1, 0.1, 3 \times 10^{-10}, 2 \times 10^{-10}, 10^{-10})$ and $\mathbf{H} = \mathbf{S}^{1/2} \mathbf{A} \mathbf{S}^{1/2}$. By construction, this example obeys the geometric mean bound (2.12) with $\alpha = 1/2$ and $\mu = 0.5 \lambda_{\max}(\mathbf{G} + \mathbf{G}^*)$ which is $\lesssim 10$ with high probability. We choose a threshold level of $\epsilon = 1.5 \times 10^{-10}$, so that the thresholded problem has dimension four. As perturbation, we take $\Delta_{\mathbf{S}} = 10^{-12} \cdot (\mathbf{\Gamma} + \mathbf{\Gamma}^*)/2$ (for a 5×5 real standard Gaussian matrix $\mathbf{\Gamma}$).

Let Π and $\tilde{\Pi}$ denote the spectral projectors onto the eigenvectors $> \epsilon$ for \mathbf{S} and $\tilde{\mathbf{S}} = \mathbf{S} + \Delta_{\mathbf{S}}$ respectively. For one random initialization of the Gaussian test matrices (which we find is broadly representative of repeat trials), we computed

$$(SM5.1) \quad \|\tilde{\Pi} \mathbf{H} \tilde{\Pi} - \Pi \mathbf{H} \Pi\| = 6.3 \times 10^{-8} \approx 10^{-7} \approx \|\Delta_{\mathbf{S}}\|/\epsilon^{1/2}.$$

Were the $\epsilon^{-\alpha}$ dependence in Theorem 2.5 unnecessary, we would expect that the projection error $\|\tilde{\Pi} \mathbf{H} \tilde{\Pi} - \Pi \mathbf{H} \Pi\|$ would be bounded by $\mu(1 + \rho^{-1})5^3 \|\Delta_{\mathbf{S}}\| \approx 10^{-9}$. We take this as evidence that the $\epsilon^{-\alpha}$ factor in the bound in Theorem 2.5 is necessary, at least without additional assumptions.

SM6. The Value of α in Eq. (2.12). As we argued in the main text, any pair (\mathbf{H}, \mathbf{S}) obeys the geometric mean bound (2.12) with $\alpha = 1/2$ and $\gamma = \max |\Lambda(\mathbf{H}, \mathbf{S})|$. In this section, we present numerical evidence that (2.12) often holds with $\alpha = 1/4$ and $\gamma \approx \max |\Lambda(\mathbf{H}, \mathbf{S})|$ for QSD problem instances, a substantial improvement on the provable bound. This $\alpha = 1/4$ behavior remains somewhat mysterious to us, and we have yet to discover a convincing explanation for why this behavior emerges.

The numerical validity of (2.12) with $\alpha = 1/4$ and $\gamma \approx \max |\Lambda(\mathbf{H}, \mathbf{S})|$ is demonstrated in Figure SM5. In these plots, we plot $y = |\mathbf{v}_j^* \mathbf{H} \mathbf{v}_k| / \max(\lambda_j, \lambda_k)$ against $x = \min(\lambda_j, \lambda_k) / \max(\lambda_j, \lambda_k)$ over all indices j and k for several different QSD instances, where we use the notation from section 2.2 that $(\lambda_j, \mathbf{v}_j)$ represents the j th largest eigenpair of \mathbf{S} . Since the accurately computable eigenvalues span a range roughly on the order of the inverse machine precision ($\approx 10^{16}$ in double precision), we only plot pairs (x, y) corresponding to indices j and k for which $\min(\lambda_j, \lambda_k) \geq 10^{-16} \lambda_1$. The bound (2.12) holds only if all points (x, y) (as well as those numerically incomputable) lie below a power law curve $y \leq \gamma x^{1-\alpha}$. The curve $\max |\Lambda(\mathbf{H}, \mathbf{S})| \cdot x^{0.75}$ is shown on each of the subplots in Figure SM5, and it lies above almost all of the pairs (x, y) . We consider this convincing evidence of the validity of (2.12) with $\alpha = 1/4$ and $\gamma \approx \max |\Lambda(\mathbf{H}, \mathbf{S})|$ for the QSD examples we tested. We suspect this relation will continue to hold for “reasonable” QSD instances, though we lack a precise definition of “reasonable” and a formal argument justifying this suspicion.

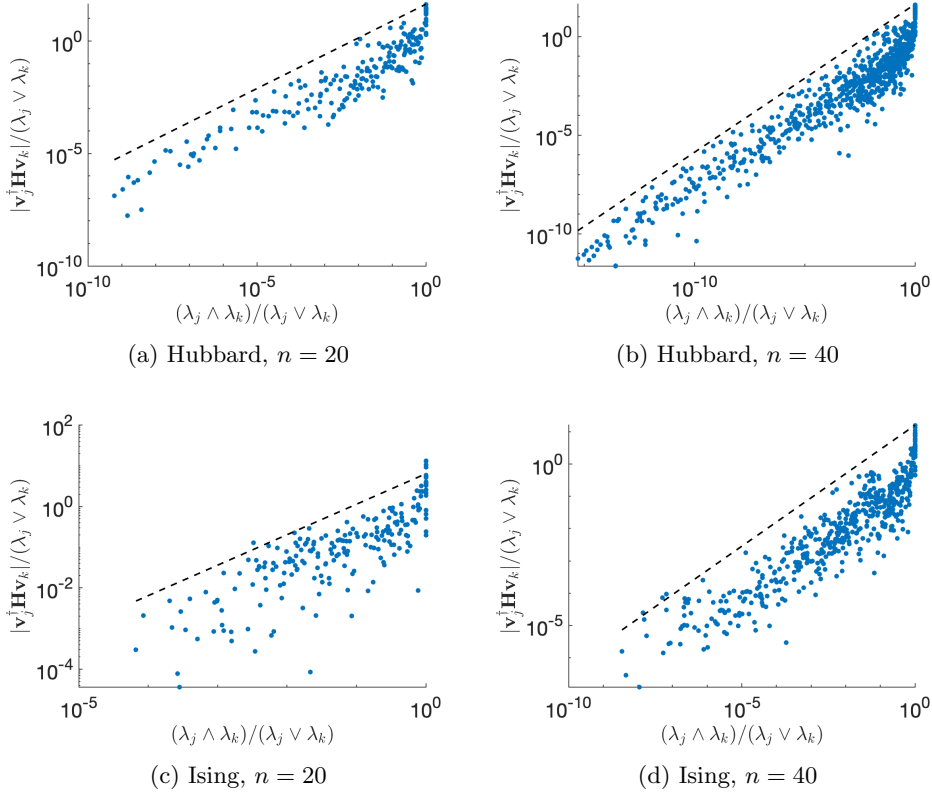


Fig. SM5: Scatter plot of values $y = |v_j^* H v_k| / (\lambda_j \vee \lambda_k)$ versus $x = \min(\lambda_j, \lambda_k) / \max(\lambda_j, \lambda_k)$ over all indices $j, k = 1, \dots, n$ for which $\min(\lambda_j, \lambda_k) \geq 10^{-16} \lambda_1$ for Hubbard model (top) and Ising model (bottom) with $n = 20$ (left) and 40 (right). Shown as a dashed black line is $\max |\Lambda(H, S)| \cdot x^{3/4}$, demonstrating that (2.12) holds numerically with $\alpha = 1/4$ and $\gamma \approx \max |\Lambda(H, S)|$.

SM7. Extra Figures. Finally, we conclude with some additional figures concerning additional numerical experiments. Figure SM6 shows the smallest eigenvalue computed when a fixed threshold parameter is used, independent of the noise level. Figures SM7, SM8, and SM9 provide more parameter settings for the Hubbard and Ising models for the Figures 3, 4, and SM2.

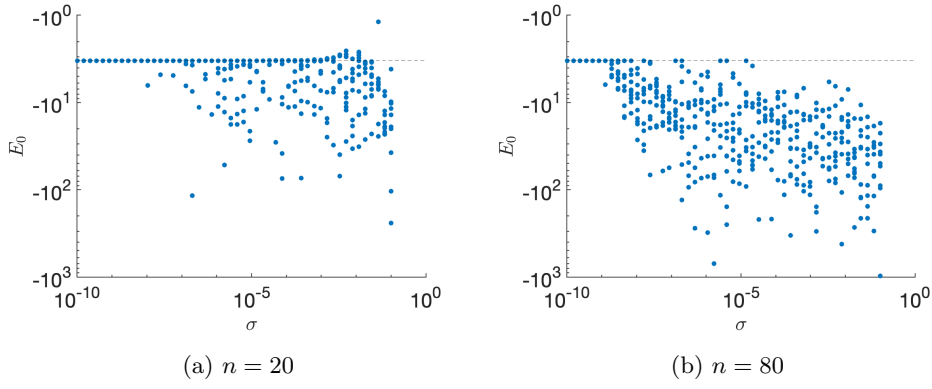


Fig. SM6: Least eigenvalues computed from the perturbed pair $(\widetilde{\mathbf{H}}, \widetilde{\mathbf{S}})$ with a fixed threshold $10^{-8}\|\mathbf{S}\|$. Shown are 10 random initializations of the noise for several random noise levels σ for the Hubbard example with $n = 20$ (left) and $n = 80$ (right). The true eigenvalue is shown for reference as a horizontal dashed line.

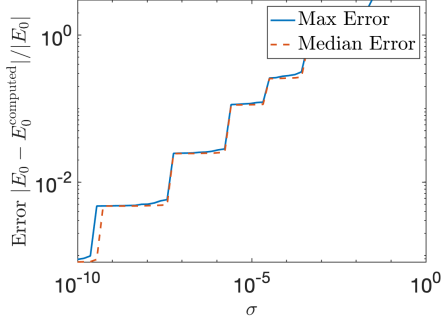
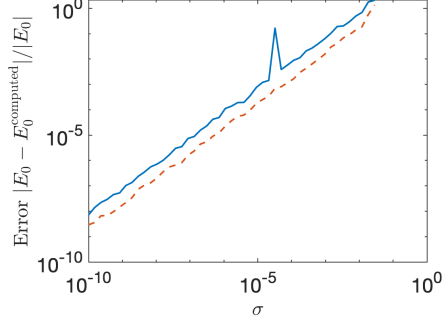
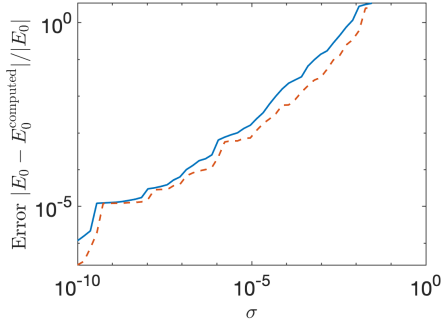
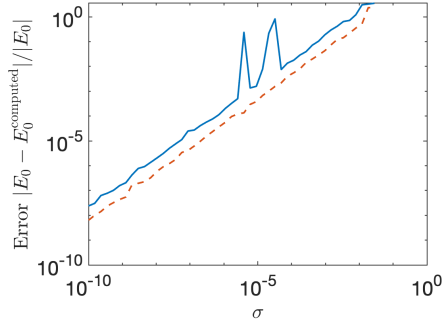
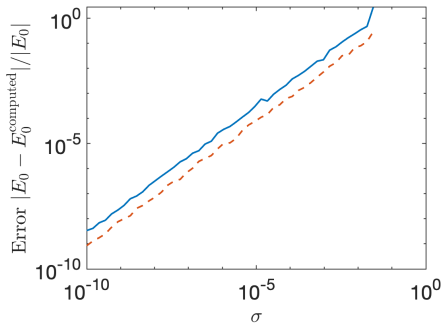
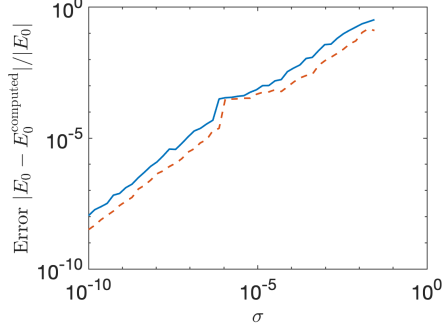
(a) Hubbard, $L = 10$, $U = 8$, $n = 10$ (b) Hubbard, $L = 10$, $U = 8$, $n = 80$ (c) Hubbard, $L = 10$, $U = 10$, $n = 30$ (d) Hubbard, $L = 10$, $U = 10$, $n = 80$ (e) Hubbard, $L = 6$, $U = 8$, $n = 40$ (f) Ising, $L = 8$, $n = 20$

Fig. SM7: Maximum (blue solid) and median (red dashed) error over 100 initializations for eigenvalues computed from the noise-perturbed pair $(\widehat{\mathbf{H}}, \widehat{\mathbf{S}})$ using thresholding with threshold parameter $25\sigma\|\widehat{\mathbf{S}}\|$ for Hubbard and Ising models for various parameters not considered in Figure 3.

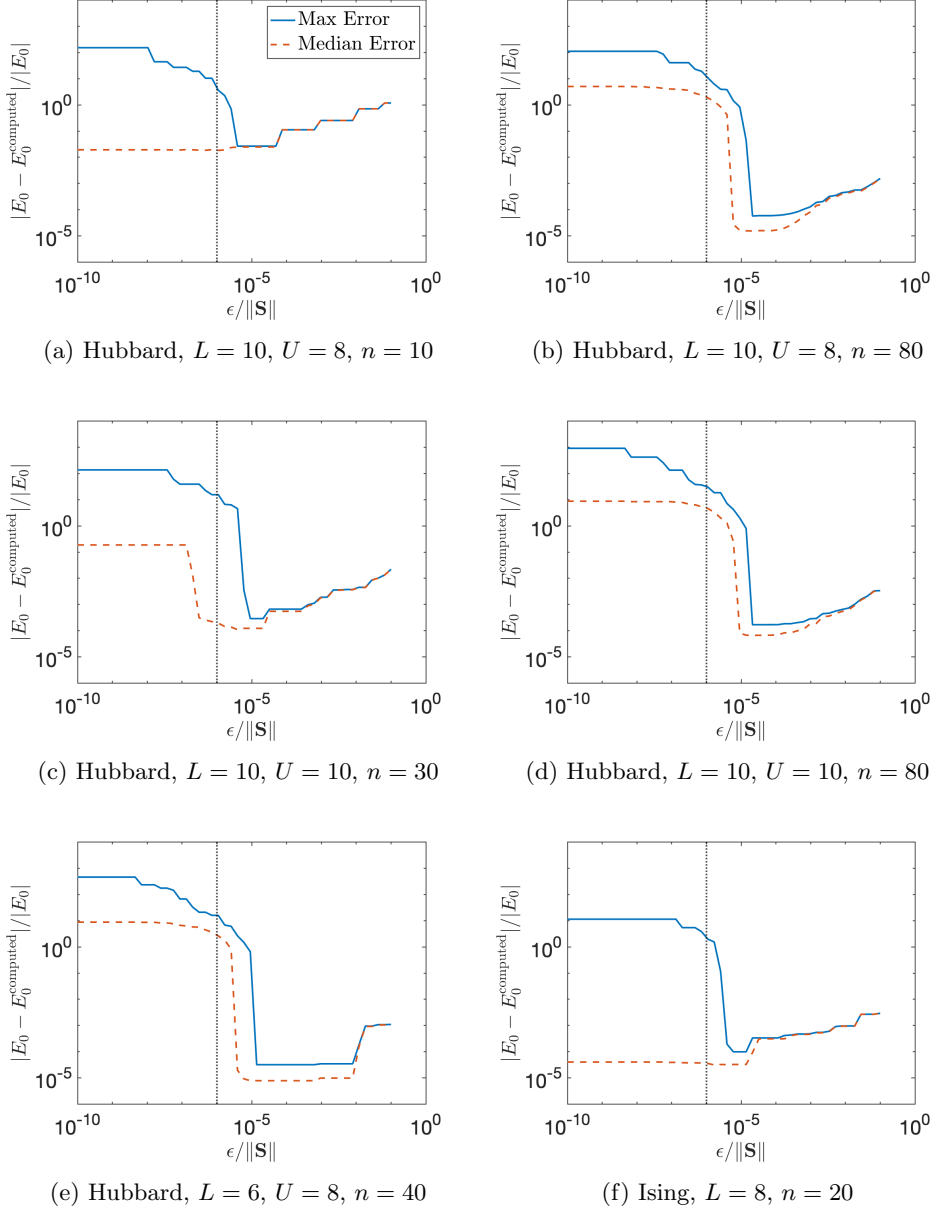


Fig. SM8: Maximum (blue solid) and median (red dashed) error over 100 initializations for eigenvalues computed from the noise-perturbed pair $(\tilde{\mathbf{H}}, \tilde{\mathbf{S}})$ using thresholding for various values of the threshold ϵ for a fixed noise level $\sigma = 10^{-6}$ (dotted black line) for Hubbard and Ising models for various parameters not considered in Figure 4.

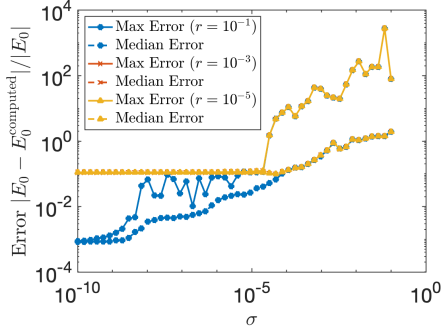
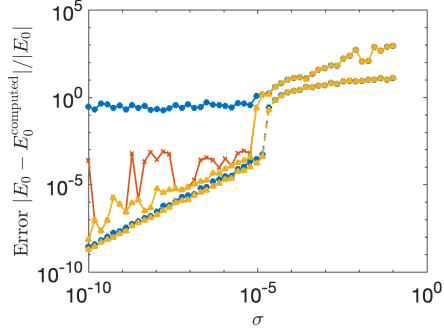
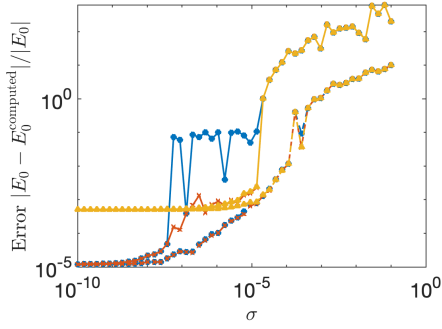
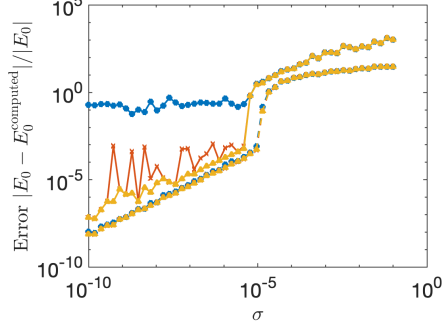
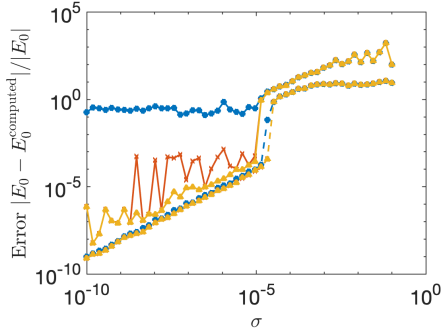
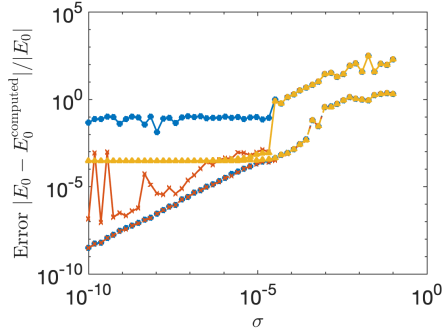
(a) Hubbard, $L = 10$, $U = 8$, $n = 10$ (b) Hubbard, $L = 10$, $U = 8$, $n = 80$ (c) Hubbard, $L = 10$, $U = 10$, $n = 30$ (d) Hubbard, $L = 10$, $U = 10$, $n = 80$ (e) Hubbard, $L = 6$, $U = 8$, $n = 40$ (f) Ising, $L = 8$, $n = 20$

Fig. SM9: Maximum and median error over 100 initializations for eigenvalues computed from the noise-perturbed pair $(\hat{\mathbf{H}}, \hat{\mathbf{S}})$ using automatically tuned thresholding (Algorithm SM3.1) for three cutoffs $r \in \{10^{-1}, 10^{-3}, 10^{-5}\}$ for various values of the noise level σ for Hubbard and Ising models for various parameters not considered in Figure SM2.

SUPPLEMENTARY MATERIALS: A THEORY OF QUANTUM SUBSPACE DIAGONALIZATION *

ETHAN N. EPPERLY[†], LIN LIN[‡], AND YUJI NAKATSUKASA[§]

In this supplement, we provide additional proofs and numerical experiments to support the claims made in the main text.

SM1. Proof of Theorem ??. For reference, we provide a complete proof of Theorem ??.

Proof of Theorem ??. Let Π and $\tilde{\Pi}$ be the spectral projectors onto the dominant m -dimensional invariant subspaces of \mathbf{A} and $\mathbf{A} + \mathbf{\Delta}$ respectively. First, we bound

$$\begin{aligned} \|\llbracket \mathbf{A} + \mathbf{\Delta} \rrbracket_m - \llbracket \mathbf{A} \rrbracket_m\|_{\text{QUI}} &= \|\tilde{\Pi}(\mathbf{A} + \mathbf{\Delta})\tilde{\Pi} - \Pi\mathbf{A}\Pi\|_{\text{QUI}} \\ &\leq \|\tilde{\Pi}\mathbf{\Delta}\tilde{\Pi}\|_{\text{QUI}} + \|\tilde{\Pi}\mathbf{A}\tilde{\Pi} - \Pi\mathbf{A}\Pi\|_{\text{QUI}}. \end{aligned}$$

For the first term, we bound $\|\tilde{\Pi}\mathbf{\Delta}\tilde{\Pi}\|_{\text{QUI}} \leq \|\mathbf{\Delta}\|_{\text{QUI}}$ using the fact $\|\cdot\|_{\text{QUI}}$ is *symmetric* in the sense that $\|\mathbf{B}_1\mathbf{B}_2\mathbf{B}_3\|_{\text{QUI}} \leq \|\mathbf{B}_1\| \|\mathbf{B}_2\|_{\text{QUI}} \|\mathbf{B}_3\|$ [?, Thm. 3.9]. Every quadratic unitarily invariant norm is bounded by the Frobenius norm, which follows from the definition of quadratic unitarily invariant norm and the fact that the nuclear norm bounds every unitarily invariant norm [?, Eq. (IV.38)]. Thus, the second term is bounded as

$$\|\tilde{\Pi}\mathbf{A}\tilde{\Pi} - \Pi\mathbf{A}\Pi\|_{\text{QUI}} \leq \|\tilde{\Pi}\mathbf{A}\tilde{\Pi} - \Pi\mathbf{A}\Pi\|_{\text{F}},$$

which is then bounded by Theorem ?? with $\epsilon = \lambda_{m+1} + \|\mathbf{\Delta}\|$ and $\rho = (\lambda_m - \lambda_{m+1} - \|\mathbf{\Delta}\|)/(\lambda_{m+1} - \|\mathbf{\Delta}\|)$. \square

SM2. The Failure of Heuristics. In the main text, we show how accurate recovery for the QSD algorithm usually fails, if one solves the noisy generalized eigenvalue problem (??) with no special treatment. A forthcoming example (Figure SM6) shows that using a fixed threshold independent of the noise level may not perform much better. An alternate strategy, which we initially believed to be more promising than thresholding, is to compute the eigenvalues (either with no thresholding at all or with a small threshold independent of the noise level) and attempt to determine real from spurious eigenvalues by means of some property of the computed eigenvector. In this section, we shall consider a couple variants of such an approach and ultimately conclude the performance of these heuristics can still be unsatisfactory when compared to thresholding.

Two natural heuristics for whether \tilde{E} is a plausible candidate for the ground state energy suggest themselves. Let $\tilde{\mathbf{c}}$ be the unit-norm eigenvector associated with a computed eigenvalue \tilde{E} . The Ritz vector $\tilde{\psi}_0 := \sum_{j=0}^{n-1} \tilde{\mathbf{c}}_j \boldsymbol{\varphi}_j$ is supposed to be close to the true ground-state eigenvector ψ_0 of \hat{H} . Our heuristics are as follows:

*This is a preprint version of *A Theory of Quantum Subspace Diagonalization* (<https://doi.org/10.1137/21M145954X>), which appeared in the SIAM Journal on Matrix Analysis and Applications on August 1, 2022.

[†]Department of Computing and Mathematical Sciences, California Institute of Technology, Pasadena, CA, USA (eepperly@caltech.edu).

[‡]Department of Mathematics, and Challenge Institute of Quantum Computation, University of California Berkeley, Berkeley, CA, USA and Computational Research Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA (linlin@math.berkeley.edu).

[§]Mathematical Institute, Oxford University, Oxford, UK (nakatsukasa@maths.ox.ac.uk).

1. **Require $h_1 := \tilde{\mathbf{c}}^* \tilde{\mathbf{S}} \tilde{\mathbf{c}}$ to be large.** The squared norm of the Ritz vector is precisely $\tilde{\mathbf{c}}^* \tilde{\mathbf{S}} \tilde{\mathbf{c}} \approx h_1$. If h_1 is small, then the norm of the Ritz vector is very small due to cancellations in the sum $\sum_{j=0}^{n-1} \tilde{\mathbf{c}}_j \boldsymbol{\varphi}_j$ and should thus be treated as suspect because of the noise. Thus, it is natural to insist on a large value of h_1 .¹
2. **Require the estimated overlap $h_2 := |\mathbf{e}_0^* \tilde{\mathbf{S}} \tilde{\mathbf{c}}| \approx |\boldsymbol{\varphi}_0^* \tilde{\boldsymbol{\psi}}_0|$ to be large.** It is important that the initial vector $\boldsymbol{\varphi}_0$ has a relatively large initial overlap $|\boldsymbol{\varphi}_0^* \tilde{\boldsymbol{\psi}}_0|$ with the eigenvector of interest—indeed, our analysis suggests accurate recovery of the ground-state energy requires this (see Theorem ??). As such, it is natural to use the overlap (or its surrogate h_1 computable from the noise-corrupted $\tilde{\mathbf{S}}$ matrix) as a measure of whether an eigenvalue is a genuine candidate for the ground-state energy. Note that by unit-norm scaling $\tilde{\mathbf{c}}$ (rather than adopting the normalization $\tilde{\mathbf{c}}^* \tilde{\mathbf{S}} \tilde{\mathbf{c}} = 1$), we are implicitly also incorporating the condition for $\tilde{\boldsymbol{\psi}}_0$ to be a stable linear combination of the basis states which motivated our interest in h_1 .

There are several ways of using a heuristic $h \in \{h_1, h_2\}$ as an algorithm for computing the ground-state eigenvalue: (a) pick E with the highest h , (b) pick the smallest E of the eigenvalues with the top k values of h , and (c) pick the smallest E with h above some thresholding h_0 (or simply the largest h if none exceeds h_0).

Unfortunately, unlike thresholding where there is a natural choice of the threshold parameter (related to the noise level η (??) which can usually be reliably estimated), we are unaware of any good systematic ways to pick the parameters k and h_0 for strategies (b) and (c). These heuristics thus usually require some tuning to make them accurate for a given problem instance, with the parameters needing to be readjusted when a new problem is encountered. This reduces the reliability of these heuristics, when the ground truth is unavailable to compare against. The robustness of heuristics such as (a), (b), and (c) can be improved by medians of repeated trials or by comparing the results of different heuristics against each other. However, even with such improvements, without rigorous guarantees, the validity of these heuristics remains conjectural when the genuine ground-state energy is unavailable to be validated against.

Figure SM1 shows the suggested heuristics (a), (b), and (c) with the figures of merit h_1 and h_2 with $k = 5$ and $h_0 = 10^{-2} \|\tilde{\mathbf{S}}\|$. The first subfigure, Figure SM1a, shows a relatively optimistic case for the heuristics. For low levels of noise, the eigenvalue is generally recovered with low error with the exception of a few outliers which could be ameliorated by the median trick. Figure SM1b shows the potential danger of applying these heuristics; despite working well for the Hubbard example with $n = 20$ in Figure SM1a, the heuristics fail with the same parameter choices for the Ising example with $n = 40$. For this problem, the eigenvalues are observed to be recovered accurately only with very small probability. Improvements to both plots are likely possible by more careful choice of the heuristic parameters or more complicated heuristics, but this is a point against such heuristics rather than for them: we ideally want a method which works well without tuning problem-dependent parameters.

SM3. Automatic Thresholding. As we saw in the main text, choosing a good maximum thresholding level ϵ is critical to the success of the thresholding procedure Algorithm ?. A useful “sanity check” is thus to solve the problem using a handful of plausible thresholding parameters to make sure the computed eigenvalues are close to

¹ h_1 is also related to the conditioning of the eigenvalue [?, Eq. (VI.2.2)].

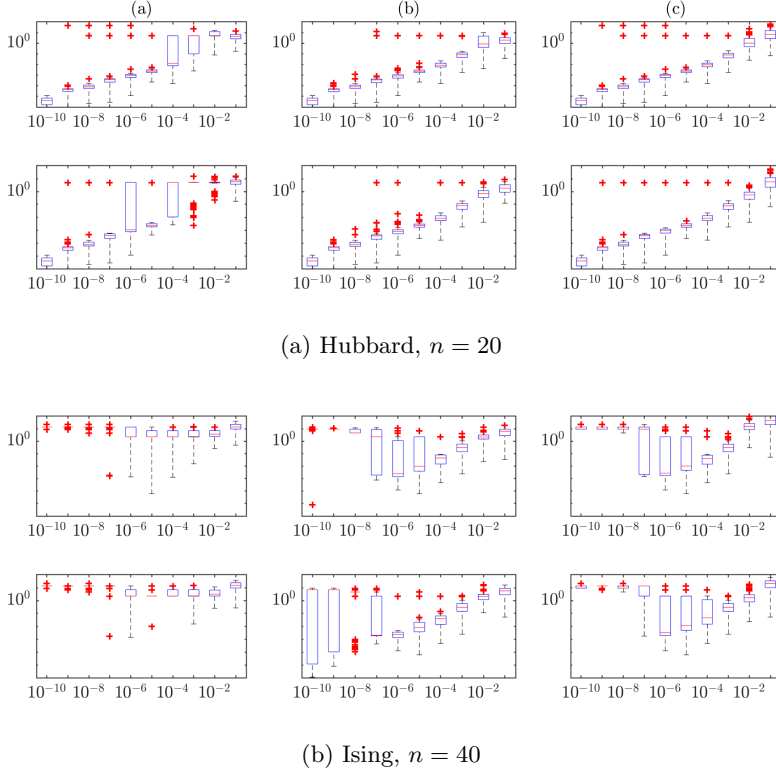


Fig. SM1: Errors (vertical axis) for eigenvalues computed from the perturbed pair $(\mathbf{H}, \tilde{\mathbf{S}})$ with heuristics (a), (b), and (c) described in the text for quality metrics h_1 (first row) and h_2 (second row). The noisy generalized eigenvalue problem (??) is solved using thresholding with a fixed threshold parameter $\epsilon = 10^{-12} \|\tilde{\mathbf{S}}\|$. Shown are 100 random initializations of the noise for several random noise levels σ (horizontal axis) for the Hubbard example with $n = 20$ (top) and Ising example with $n = 40$ (bottom).

each other. (A variant of this strategy is proposed by Parlett for the Fix–Heiberger procedure [?, §15.5].) A more ambitious strategy is to solve with a range of thresholding parameters beginning with a conservative (but not comically large) threshold parameter ϵ_0 and then tuning it down until the eigenvalue “jumps” to a presumably spurious value. The best approximation to the ground-state energy suggested by this procedure is the last value before this jump. If one wishes to automate this procedure, one needs to have a mechanistic way of deciding whether a jump has occurred: For this purpose, we shall test whether the relative difference exceeds a cutoff r . This procedure is demonstrated in Algorithm SM3.1. The success of this procedure relies on the choice of ϵ_0 not being too large as the method uses the eigenvalue recovered with parameter ϵ_0 as a baseline, large deviations from which are characterized as erroneous. Usually, one will have some good estimate of the amount of noise so picking a sensible ϵ_0 should be possible.

Algorithm SM3.1 Automatically tuned thresholding procedure for finding the least eigenvalue of a noise-corrupted generalized eigenvalue problem.

```

procedure AUTOTHRESHOLDING( $\mathbf{H}$ ,  $\mathbf{S}$ ,  $\epsilon_0$ ,  $r$ )
   $E \leftarrow \text{THRESHOLDING}(\mathbf{H}, \mathbf{S}, \epsilon_0)$ 
   $\Lambda \leftarrow \{\lambda \in \text{eig}(\mathbf{S}) : \lambda < \epsilon_0\}$ 
  while  $\Lambda \neq \emptyset$  do
     $\epsilon \leftarrow \max \Lambda$ ,  $\Lambda \leftarrow \Lambda \setminus \{\epsilon\}$ 
     $E' \leftarrow \text{THRESHOLDING}(\mathbf{H}, \mathbf{S}, \epsilon)$ 
    if  $|E - E'| / \min(|E|, |E'|) > r$  then
      break
    end if
     $E \leftarrow E'$ 
  end while
  return  $E$ 
end procedure

```

The performance of the automatic thresholding procedure Algorithm SM3.1 with three choices of the parameter r are shown in Figure SM2. These plots represent the worst-case situation where the noise level is completely unknown and the choice one has available for ϵ_0 is a constant multiple of $\|\tilde{\mathbf{S}}\|$. The best case scenario is shown in the $r = 10^{-3}$ lines in Figures SM2b and SM2c; in these cases, the error decays nicely as the noise does with the procedure being relatively robust (as shown by the error over a maximum over 100 trials being similar to the median). This automatic thresholding procedure can still be somewhat delicate, with the maximum error over 100 runs being near the cutoff r for the $r = 10^{-1}$ in Figures SM2a, SM2b, and SM2c; this suggests, in the worst case, one must be willing to accept an error level on the order r due to overly aggressive automatic tuning of the thresholding parameter. However, these same plots illustrate the importance of not being too cautious either, with the maximum error being $\approx 10^{-3}$ with $r = 10^{-5}$ due to overly conservative automatic tuning of the thresholding parameter. In totality, Figure SM2 shows that the automatic thresholding procedure cannot determine a near-optimal choice for the thresholding parameter in all cases, but it can be useful in “upgrading” an overly cautious threshold parameter ϵ_0 to a better choice for ϵ , obtaining a couple more decimal digits of accuracy in the best case. As a final comment, observe that thresholding is an inherently discrete process since each eigenvalue must either be discarded or not. Therefore, even with automatically tuned thresholding, there can be plateaus in the noise-level vs accuracy curve, owing to the importance of a single eigenvalue to the overall error; this is shown in Figure SM2c.

SM4. Validation of Theorems ?? and ??. Despite its desirable theoretical implications, Theorem ?? may still significantly overestimate the error incurred by thresholding in practice. Consider the following examples:

- (I) $\widehat{\mathbf{H}} = \widehat{\mathbf{H}}_1$ and $\varphi_0 = \xi_I$,
- (II) $\widehat{\mathbf{H}} = \widehat{\mathbf{H}}_1$ and $\varphi_0 = \xi_{II}$,
- (III) $\widehat{\mathbf{H}} = \widehat{\mathbf{H}}_2$ and $\varphi_0 = \xi_{III}$,
- (IV) $\widehat{\mathbf{H}} = \widehat{\mathbf{H}}_2$ and $\varphi_0 = \xi_{IV}$.

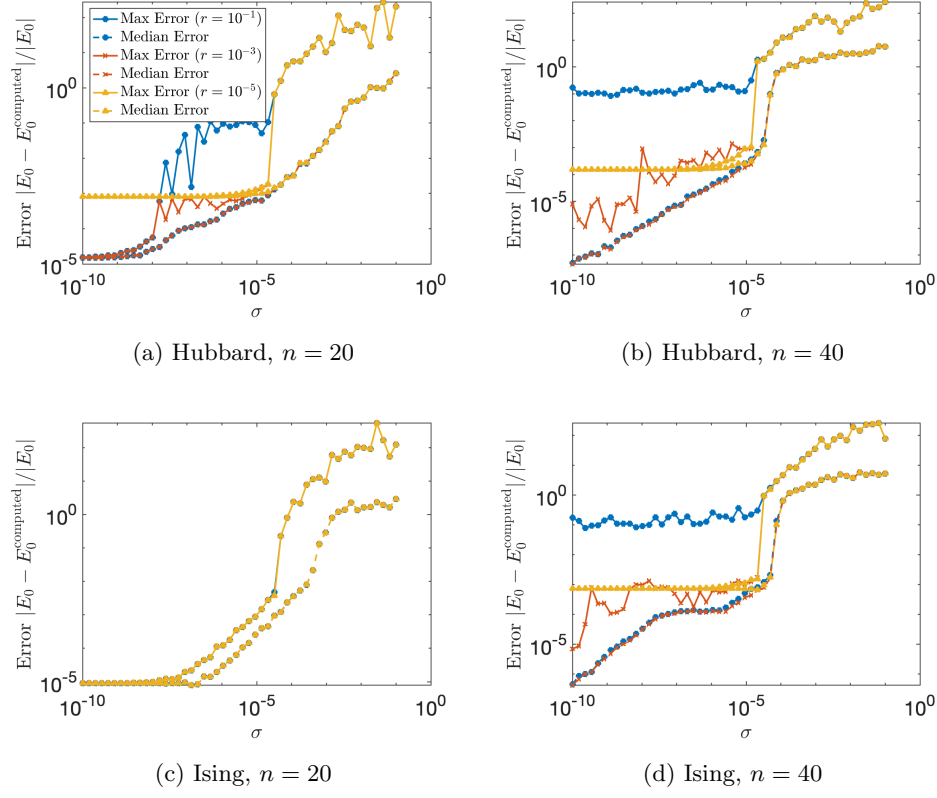


Fig. SM2: Maximum and median error over 100 initializations for eigenvalues computed from the noise-perturbed pair $(\hat{\mathbf{H}}, \hat{\mathbf{S}})$ using automatically tuned thresholding (Algorithm SM3.1) for three cutoffs $r \in \{10^{-1}, 10^{-3}, 10^{-5}\}$ for various values of the noise level σ for Hubbard model (top) and Ising model (bottom) with $n = 20$ (left) and 40 (right).

125 where

$$\begin{aligned}
 126 \quad \xi_{\text{I}} &= \left(\sqrt{1 - 10^{-4}}, \sqrt{\frac{10^{-4}}{998}}, \sqrt{\frac{10^{-4}}{998}}, \dots, \sqrt{\frac{10^{-4}}{998}} \right) \in \mathbb{C}^{999}, \\
 127 \quad \xi_{\text{II}} &= \left(\sqrt{0.5}, \sqrt{\frac{0.5}{998}}, \sqrt{\frac{0.5}{998}}, \dots, \sqrt{\frac{0.5}{998}} \right) \in \mathbb{C}^{999}, \\
 128 \quad \xi_{\text{III}} &= \left(\sqrt{1 - 10^{-4} - 10^{-8}}, \sqrt{\frac{10^{-4}}{998}}, \sqrt{\frac{10^{-4}}{998}}, \dots, \sqrt{\frac{10^{-4}}{998}}, 10^{-4} \right) \in \mathbb{C}^{1000} \\
 129 \quad \xi_{\text{IV}} &= \left(1, \frac{0.01}{2}, \frac{0.01}{3}, \dots, \frac{0.01}{1000} \right) / \left\| \left(1, \frac{0.01}{2}, \frac{0.01}{3}, \dots, \frac{0.01}{1000} \right) \right\| \in \mathbb{C}^{1000}, \\
 130
 \end{aligned}$$

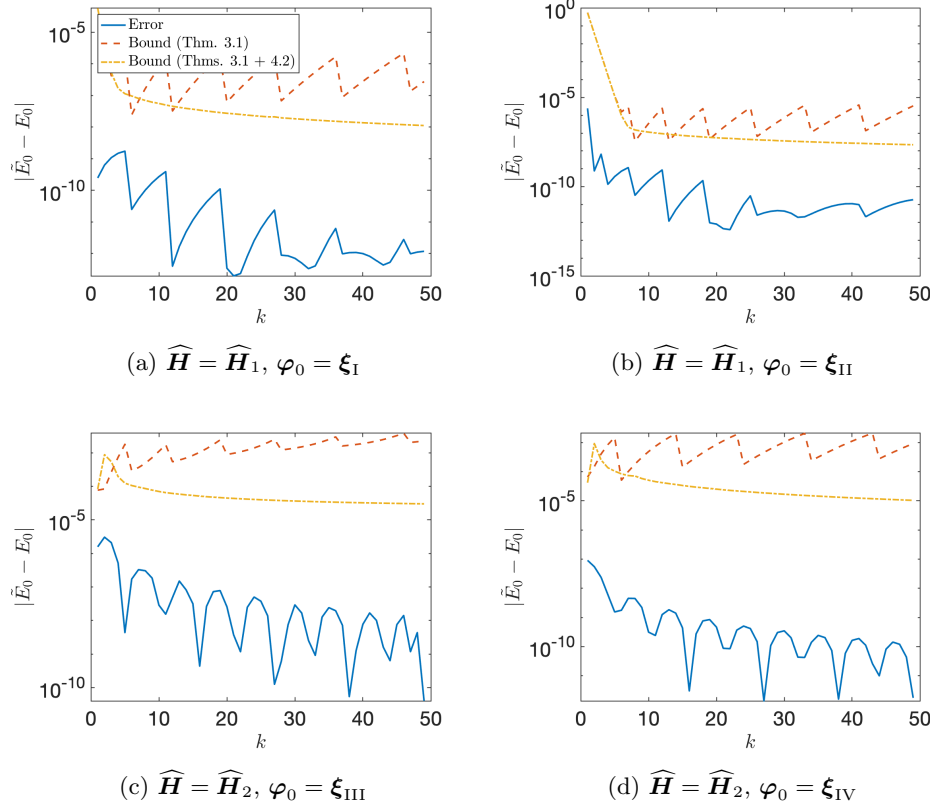


Fig. SM3: Error for QSD with the time sequence from the hypotheses of Theorem ?? with threshold parameter $\epsilon = 10^{-6}$ for various values k and four different sets of input data.

131 and

$$132 \quad \widehat{\mathbf{H}}_1 = \text{diag}\left(1, 2 + 0 \cdot \frac{0.1}{997}, 2 + 1 \cdot \frac{0.1}{997}, \dots, 2 + 997 \cdot \frac{0.1}{997}\right) \in \mathbb{C}^{999 \times 999},$$

$$133 \quad \widehat{\mathbf{H}}_2 = \text{diag}\left(1, 2 + 0 \cdot \frac{0.1}{997}, 2 + 1 \cdot \frac{0.1}{997}, \dots, 2 + 997 \cdot \frac{0.1}{997}, 1000\right) \in \mathbb{C}^{1000 \times 1000}.$$

134 These examples are artificial: They are engineered to have a large spectral gap ΔE_1
 135 but a small spectral range ΔE_{999} . Even with these artificial examples, Theorem ?? (as
 136 well as Theorems ?? and ?? together) still overestimates the error by several orders
 137 of magnitude. See Figure SM3.

138 We consider the bound Theorem ?? by itself in Figure SM4. Shown is the er-
 139 ror $\bar{E}_0 - E_0$ between the recovered least eigenvalue \bar{E}_0 and the true least eigen-
 140 value E_0 for different choices of threshold parameter ϵ . For this example, we used
 141 $\mathbf{H} = \mathbf{K}^* \text{diag}(1, 2, \dots, 100) \mathbf{K}$ and $\mathbf{S} = \mathbf{K}^* \mathbf{K}$, where \mathbf{K} is a product of diagonal matri-
 142 x with (j, j) th entry j^{-2} and a “randsvd” matrix from MATLAB’s gallery with
 143 approximate condition number 10^3 . This was chosen to give a fairly ill-conditioned
 144 “ \mathbf{S} ” matrix ($\kappa(\mathbf{S}) \approx 10^{12}$) in which the \mathbf{S} -normalized ground state eigenvector of

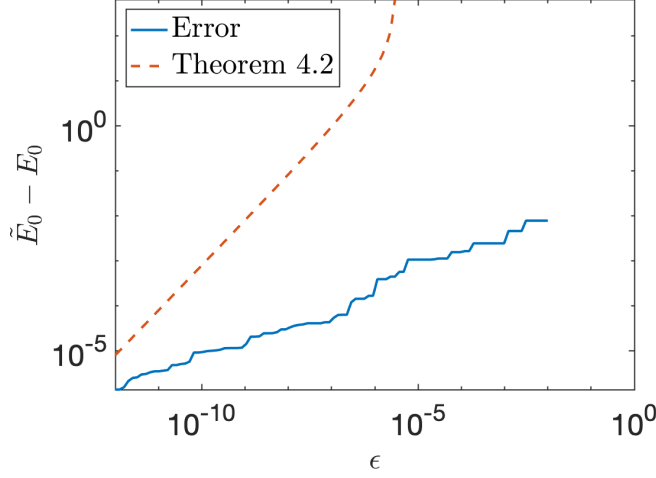


Fig. SM4: Error due to thresholding and error bound from Theorem ?? for least eigenvalue computed from a synthetically generated pair (\mathbf{H}, \mathbf{S}) for various threshold parameters ϵ .

fairly small norm ($\|\mathbf{c}_0\| \approx 10^2$). We find the bound from Theorem ?? becomes increasingly conservative as ϵ increases, diverging to $+\infty$ at $\epsilon \approx 10^{-6}$ which the true error $\tilde{E}_0 - E_0$ remains bounded $\leq 10^{-2}$ for $\epsilon \leq 10^{-2}$.

SM5. Evidence for Tightness of Theorem ??. First, we present a synthetically generated numerical example which suggests that the η/ϵ^α behavior in Theorem ?? is necessary, at least without further assumptions. As our example, we set $\mathbf{A} = (\mathbf{G} + \mathbf{G}^*)/2$ to be the Hermitian part of a 5×5 real standard Gaussian matrix \mathbf{G} and pick $\mathbf{S} = \text{diag}(1, 0.1, 3 \times 10^{-10}, 2 \times 10^{-10}, 10^{-10})$ and $\mathbf{H} = \mathbf{S}^{1/2} \mathbf{A} \mathbf{S}^{1/2}$. By construction, this example obeys the geometric mean bound (??) with $\alpha = 1/2$ and $\mu = 0.5 \lambda_{\max}(\mathbf{G} + \mathbf{G}^*)$ which is $\lesssim 10$ with high probability. We choose a threshold level of $\epsilon = 1.5 \times 10^{-10}$, so that the thresholded problem has dimension four. As perturbation, we take $\Delta_{\mathbf{S}} = 10^{-12} \cdot (\mathbf{\Gamma} + \mathbf{\Gamma}^*)/2$ (for a 5×5 real standard Gaussian matrix $\mathbf{\Gamma}$).

Let $\mathbf{\Pi}$ and $\tilde{\mathbf{\Pi}}$ denote the spectral projectors onto the eigenvectors $> \epsilon$ for \mathbf{S} and $\tilde{\mathbf{S}} = \mathbf{S} + \Delta_{\mathbf{S}}$ respectively. For one random initialization of the Gaussian test matrices (which we find is broadly representative of repeat trials), we computed

$$(SM5.1) \quad \|\tilde{\mathbf{\Pi}} \mathbf{H} \tilde{\mathbf{\Pi}} - \mathbf{\Pi} \mathbf{H} \mathbf{\Pi}\| = 6.3 \times 10^{-8} \approx 10^{-7} \approx \|\Delta_{\mathbf{S}}\|/\epsilon^{1/2}.$$

Were the $\epsilon^{-\alpha}$ dependence in Theorem ?? unnecessary, we would expect that the projection error $\|\tilde{\mathbf{\Pi}} \mathbf{H} \tilde{\mathbf{\Pi}} - \mathbf{\Pi} \mathbf{H} \mathbf{\Pi}\|$ would be bounded by $\mu(1 + \rho^{-1})5^3 \|\Delta_{\mathbf{S}}\| \approx 10^{-9}$. We take this as evidence that the $\epsilon^{-\alpha}$ factor in the bound in Theorem ?? is necessary, at least without additional assumptions.

SM6. The Value of α in Eq. (??). As we argued in the main text, any pair (\mathbf{H}, \mathbf{S}) obeys the geometric mean bound (??) with $\alpha = 1/2$ and $\gamma = \max |\Lambda(\mathbf{H}, \mathbf{S})|$. In this section, we present numerical evidence that (??) often holds with $\alpha = 1/4$ and $\gamma \approx \max |\Lambda(\mathbf{H}, \mathbf{S})|$ for QSD problem instances, a substantial improvement on

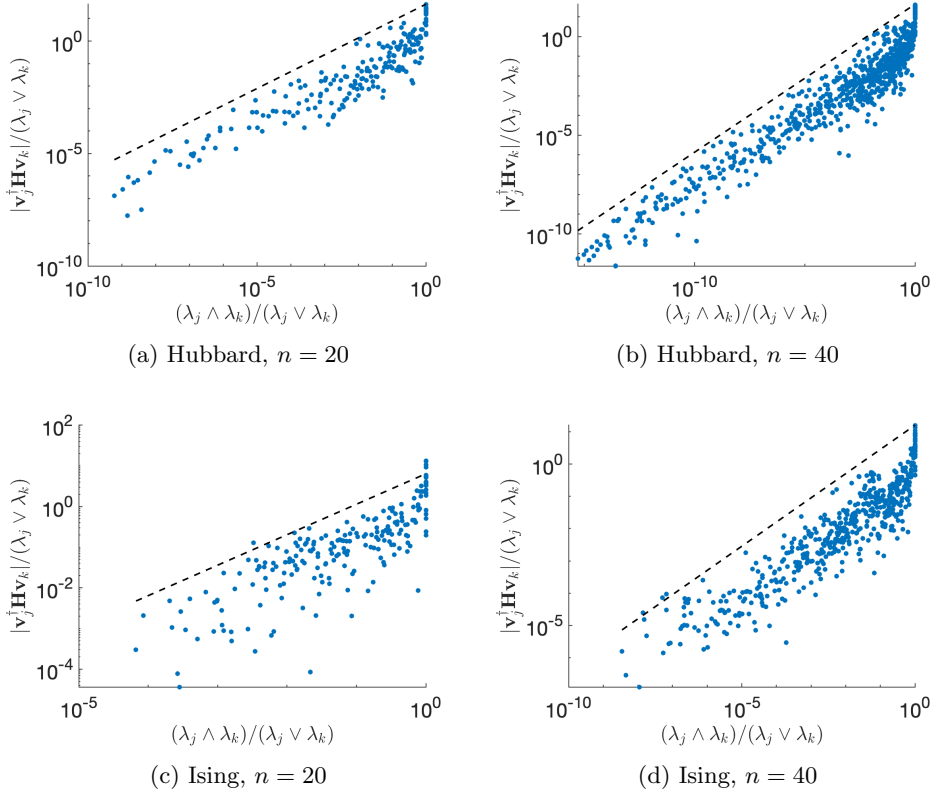


Fig. SM5: Scatter plot of values $y = |\mathbf{v}_j^* \mathbf{H} \mathbf{v}_k|/(\lambda_j \vee \lambda_k)$ versus $x = \min(\lambda_j, \lambda_k)/\max(\lambda_j, \lambda_k)$ over all indices $j, k = 1, \dots, n$ for which $\min(\lambda_j, \lambda_k) \geq 10^{-16} \lambda_1$ for Hubbard model (top) and Ising model (bottom) with $n = 20$ (left) and 40 (right). Shown as a dashed black line is $\max |\Lambda(\mathbf{H}, \mathbf{S})| \cdot x^{3/4}$, demonstrating that (??) holds numerically with $\alpha = 1/4$ and $\gamma \approx \max |\Lambda(\mathbf{H}, \mathbf{S})|$.

the provable bound. This $\alpha = 1/4$ behavior remains somewhat mysterious to us, and we have yet to discover a convincing explanation for why this behavior emerges.

The numerical validity of (??) with $\alpha = 1/4$ and $\gamma \approx \max |\Lambda(\mathbf{H}, \mathbf{S})|$ is demonstrated in Figure SM5. In these plots, we plot $y = |\mathbf{v}_j^* \mathbf{H} \mathbf{v}_k|/\max(\lambda_j, \lambda_k)$ against $x = \min(\lambda_j, \lambda_k)/\max(\lambda_j, \lambda_k)$ over all indices j and k for several different QSD instances, where we use the notation from section ?? that $(\lambda_j, \mathbf{v}_j)$ represents the j th largest eigenpair of \mathbf{S} . Since the accurately computable eigenvalues span a range roughly on the order of the inverse machine precision ($\approx 10^{16}$ in double precision), we only plot pairs (x, y) corresponding to indices j and k for which $\min(\lambda_j, \lambda_k) \geq 10^{-16} \lambda_1$. The bound (??) holds only if all points (x, y) (as well as those numerically incomputable) lie below a power law curve $y \leq \gamma x^{1-\alpha}$. The curve $\max |\Lambda(\mathbf{H}, \mathbf{S})| \cdot x^{0.75}$ is shown on each of the subplots in Figure SM5, and it lies above almost all of the pairs (x, y) . We consider this convincing evidence of the validity of (??) with $\alpha = 1/4$ and $\gamma \approx \max |\Lambda(\mathbf{H}, \mathbf{S})|$ for the QSD examples we tested. We suspect this relation will continue to hold for “reasonable” QSD instances, though we lack a precise definition

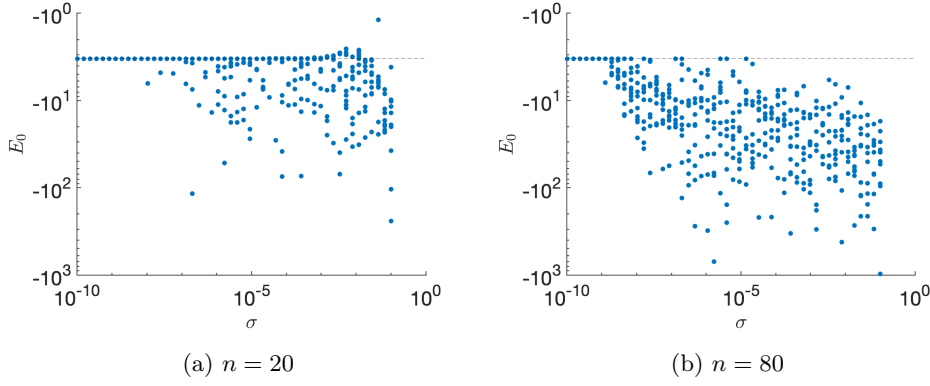


Fig. SM6: Least eigenvalues computed from the perturbed pair $(\widetilde{\mathbf{H}}, \widetilde{\mathbf{S}})$ with a fixed threshold $10^{-8}\|\mathbf{S}\|$. Shown are 10 random initializations of the noise for several random noise levels σ for the Hubbard example with $n = 20$ (left) and $n = 80$ (right). The true eigenvalue is shown for reference as a horizontal dashed line.

185 of “reasonable” and a formal argument justifying this suspicion.

186 **SM7. Extra Figures.** Finally, we conclude with some additional figures con-
 187 cerning additional numerical experiments. Figure SM6 shows the smallest eigenvalue
 188 computed when a fixed threshold parameter is used, independent of the noise level.
 189 Figures SM7, SM8, and SM9 provide more parameter settings for the Hubbard and
 190 Ising models for the Figures ??, ??, and SM2.

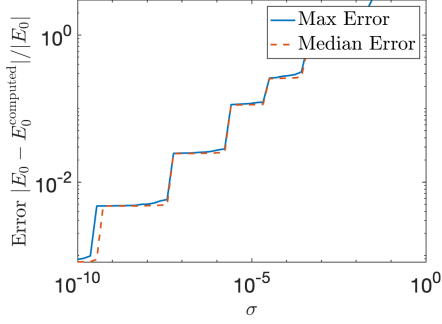
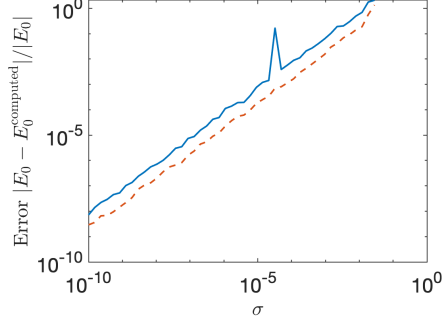
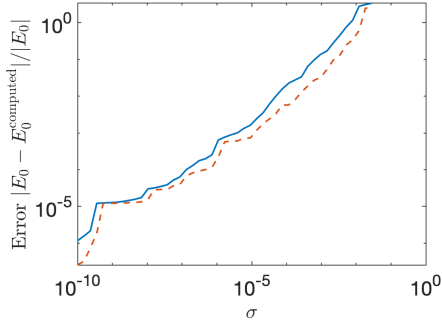
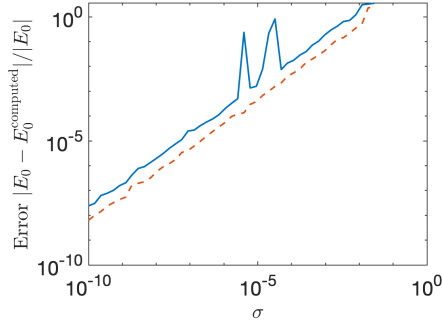
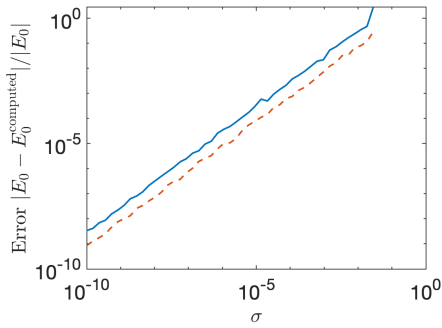
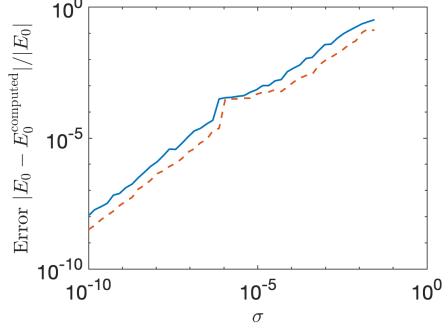
(a) Hubbard, $L = 10$, $U = 8$, $n = 10$ (b) Hubbard, $L = 10$, $U = 8$, $n = 80$ (c) Hubbard, $L = 10$, $U = 10$, $n = 30$ (d) Hubbard, $L = 10$, $U = 10$, $n = 80$ (e) Hubbard, $L = 6$, $U = 8$, $n = 40$ (f) Ising, $L = 8$, $n = 20$

Fig. SM7: Maximum (blue solid) and median (red dashed) error over 100 initializations for eigenvalues computed from the noise-perturbed pair $(\widehat{\mathbf{H}}, \widehat{\mathbf{S}})$ using thresholding with threshold parameter $25\sigma\|\widehat{\mathbf{S}}\|$ for Hubbard and Ising models for various parameters not considered in Figure ??.

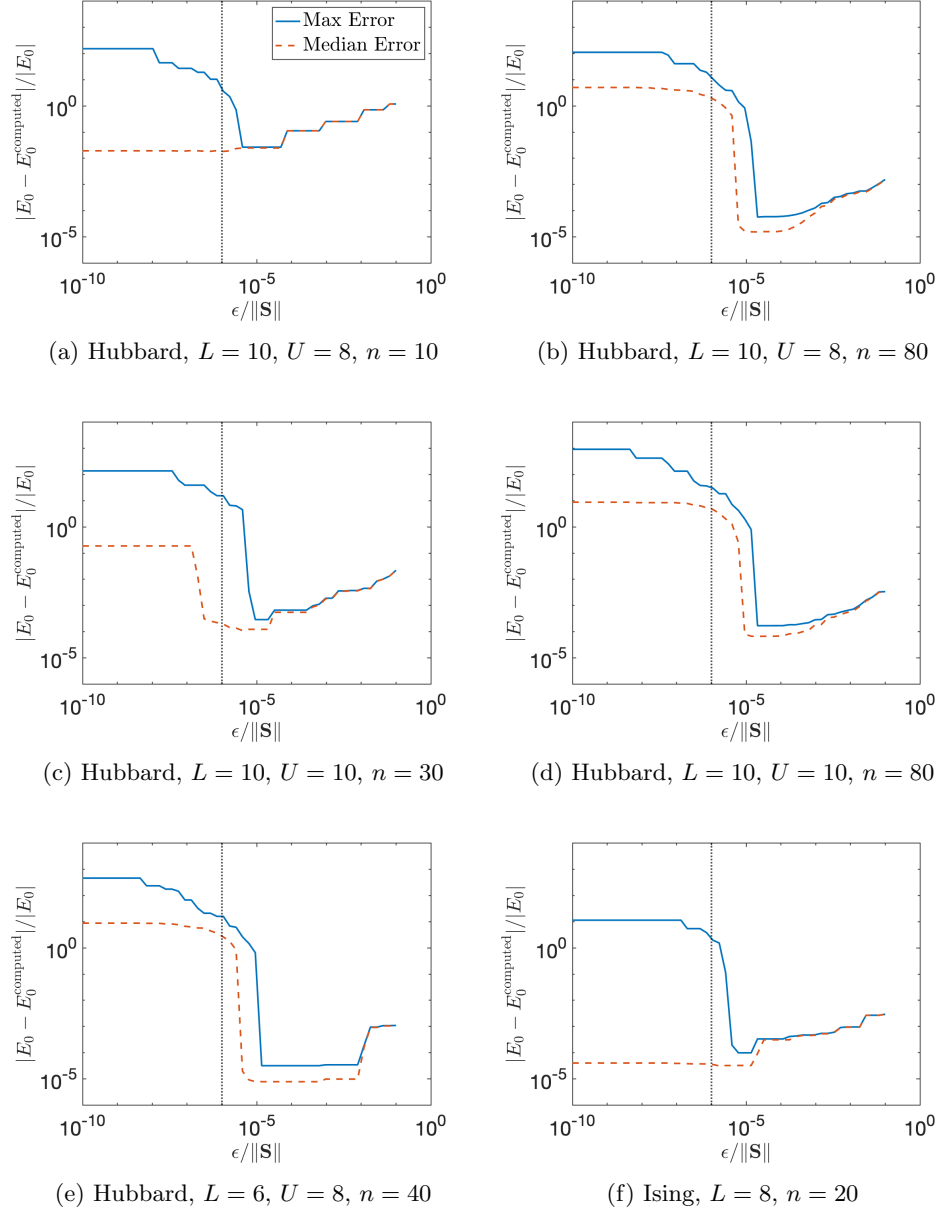


Fig. SM8: Maximum (blue solid) and median (red dashed) error over 100 initializations for eigenvalues computed from the noise-perturbed pair $(\tilde{\mathbf{H}}, \tilde{\mathbf{S}})$ using thresholding for various values of the threshold ϵ for a fixed noise level $\sigma = 10^{-6}$ (dotted black line) for Hubbard and Ising models for various parameters not considered in Figure ??.

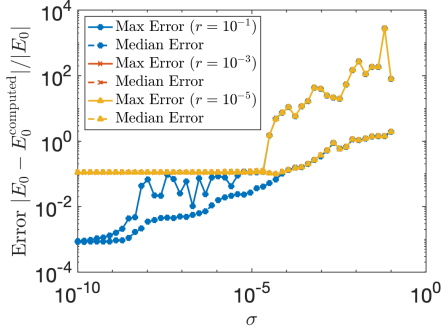
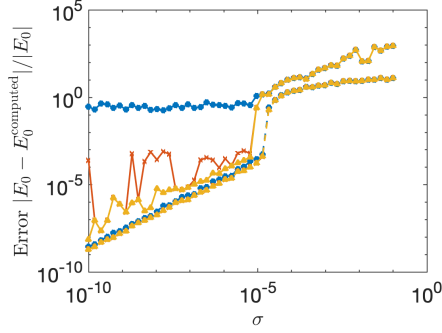
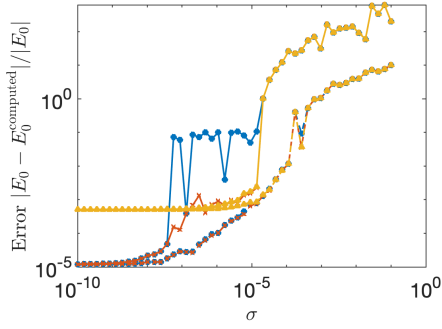
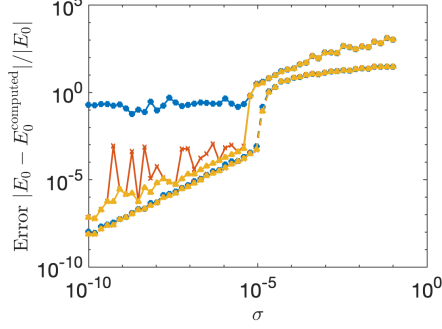
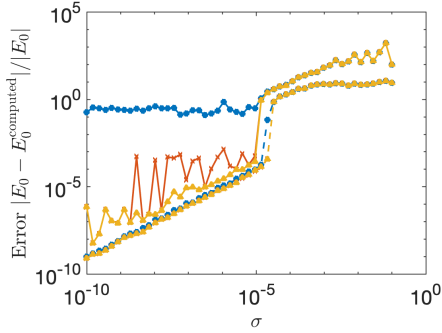
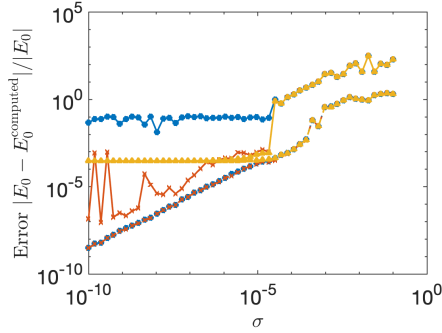
(a) Hubbard, $L = 10, U = 8, n = 10$ (b) Hubbard, $L = 10, U = 8, n = 80$ (c) Hubbard, $L = 10, U = 10, n = 30$ (d) Hubbard, $L = 10, U = 10, n = 80$ (e) Hubbard, $L = 6, U = 8, n = 40$ (f) Ising, $L = 8, n = 20$

Fig. SM9: Maximum and median error over 100 initializations for eigenvalues computed from the noise-perturbed pair $(\hat{\mathbf{H}}, \hat{\mathbf{S}})$ using automatically tuned thresholding (Algorithm SM3.1) for three cutoffs $r \in \{10^{-1}, 10^{-3}, 10^{-5}\}$ for various values of the noise level σ for Hubbard and Ising models for various parameters not considered in Figure SM2.