# UC Berkeley
## UC Berkeley Previously Published Works

**Title**
A Structure Preserving Lanczos Algorithm for Computing the Optical Absorption Spectrum

**Permalink**
https://escholarship.org/uc/item/1n29k4ww

**Journal**
SIAM Journal on Matrix Analysis and Applications, 39(2)

**ISSN**
0895-4798

**Authors**

Shao, Meiyue
da Jornada, Felipe H
Lin, Lin
et al.

**Publication Date**
2018

**DOI**
10.1137/16m1102641

Peer reviewed

# A structure preserving Lanczos algorithm for computing the optical absorption spectrum

Meiyue Shao[1], Felipe H. da Jornada[2,3], Lin Lin[1,4], Chao Yang[1], Jack Deslippe[5], and Steven G. Louie[2,3]

[1]*Computational Research Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720*
[2]*Department of Physics, University of California, Berkeley, CA 94720*
[3]*Materials Sciences Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720*
[4]*Department of Mathematics, University of California, Berkeley, CA 94720*
[5]*NERSC, Lawrence Berkeley National Laboratory, Berkeley, CA 94720*

### Abstract

We present a new structure preserving Lanczos algorithm for approximating the optical absorption spectrum in the context of solving full Bethe–Salpeter equation without Tamm–Dancoff approximation. The new algorithm is based on a structure preserving Lanczos procedure, which exploits the special block structure of Bethe–Salpeter Hamiltonian matrices. A recently developed technique of generalized averaged Gauss quadrature is incorporated to accelerate the convergence. We also establish the connection between our structure preserving Lanczos procedure with several existing Lanczos procedures developed in different contexts. Numerical examples are presented to demonstrate the effectiveness of our Lanczos algorithm.

**Keywords:** Bethe–Salpeter equation, Tamm–Dancoff approximation, optical absorption spectrum, Lanczos procedure, structure preserving algorithm, matrix functional, Gauss quadrature.

**MSC2010:** 65F15, 65F60

## 1 Introduction

Optical absorption and emission processes provide invaluable information to characterize the electronic properties of solids and molecules. At the same time, an accurate microscopic theory is also highly valuable to predict optical behavior of materials and help design more efficient photovoltaic and light-emitting devices. Physically, the optical spectra of

materials can be understood in terms of correlated electron–hole pairs known as excitons. When a photon gets absorbed by a molecule or solid, an electron can be promoted form an occupied to an unoccupied state [25, 30] in a process that creates both a negatively charge particle (known as quasielectron, or simply electron), and a positively charged particle (known as quasihole, or hole). The excitation energy required to produce such an electron–hole pair, or exciton, is directly related to the optical absorption and emission spectrum of the material. A two-particle collective excitation can be described by a two-particle Green's function, of which the real part of the poles give excitation energies. Since the two-particle Green's function satisfies the so called *Bethe–Salpeter equation* (BSE) [26, 30], the excitation energies can be obtained by solving the Bethe–Salpeter equation.

Under an appropriate discretization scheme, the *Bethe–Salpeter Hamiltonian* (BSH) matrix, which is the discrete representation of the Bethe–Salpeter Hamiltonian operator, has the block structure

$$H = \begin{bmatrix} A & B \\ -\overline{B} & -\overline{A} \end{bmatrix} \in \mathbb{C}^{2n \times 2n}, \tag{1}$$

where

$$A^{\mathsf{H}} = A, \qquad \overline{B}^{\mathsf{H}} = B. \tag{2}$$

We can rewrite $H$ as $H = C_n \Omega$ where

$$C_n = \begin{bmatrix} I_n & 0 \\ 0 & -I_n \end{bmatrix}, \qquad \Omega = \begin{bmatrix} A & B \\ \overline{B} & \overline{A} \end{bmatrix}. \tag{3}$$

For most physical systems, $\Omega$ is Hermitian positive definite (see, e.g., [35]), which we will denote by

$$\Omega \succ 0. \tag{4}$$

We define a BSH matrix $H$ that satisfies the condition (4) a *definite* Bethe–Salpeter Hamiltonian matrix. Throughout this paper, we assume that the BSH matrix $H$ is definite, that is, (4) is always assumed.

The matrices $A$ and $B$ are of size $n = n_v n_c n_k$, where $n_v$, $n_c$, and $n_k$ are the numbers of valence states, conduction states, and $k$-points, respectively. Both $n_v$ and $n_c$ are proportional to the number of electrons $n_e$ in the system. Therefore, the dimension $n = O(n_e^2 n_k)$ can be very large for systems of practical interest.

The optical absorption spectrum, which can be measured in optical absorption experiments, provides a global picture of all excitation states. Mathematically, the optical absorption spectrum is a matrix functional of the form $d_r^{\mathsf{H}} f(H; \omega) d_l$, where $f(H; \omega)$ is a function of $H$ depending on a parameter $\omega$, and $d_l$, $d_r \in \mathbb{C}^{2n}$. The peaks in the optical absorption spectrum correspond to the excitation energies.

Accurate computation of the optical absorption spectrum can be obtained by fully diagonalizing the BSH matrix [27]. However, when the problem size $n$ grows, diagonalizing the BSH matrix, whose complexity is $O(n^3)$, becomes increasingly expensive and eventually

2

unaffordable. In this paper we shall discuss how to quickly estimate the absorption without diagonalizing $H$.

In practice there is actually no need to accurately locate all peaks and determine the corresponding heights in the optical absorption spectrum—a reasonable approximation is often sufficient. Therefore Krylov subspace-based methods, such as the Lanczos algorithm, becomes very attractive for this purpose. In the context of *Tamm–Dancoff approximation* (TDA) [6, 31], which sets the off-diagonal blocks $B$ in $H$ to zero, Haydock's recursive algorithm [12, 13, 14, 15], which is essentially based on the symmetric Lanczos algorithm, can be used to solve this problem. For full BSE calculations, the non-symmetric Lanczos algorithm [21] can be applied, while the structure of $H$ is not taken into account. Recently a structure preserving Lanczos algorithm that is applicable to full BSE has been proposed in [11]. In this paper we shall develop a new structure preserving Lanczos algorithm to quickly estimate the optical absorption spectrum. Our new algorithm incorporates a recently developed technique of generalized averaged Gauss quadrature [19, 29] and largely improves the algorithm in [11] with negligible additional cost.

The rest of the paper is organized as follows. In Section 2, we review some basic properties of the definite Bethe–Salpeter Hamiltonian and the optical absorption spectrum. In Section 3, we describe how the standard Lanczos algorithm can be used to estimate the absorption spectrum in the context of TDA. Then in Section 4 we discuss how the Lanczos algorithm can be modified when it is applied to a full BSH. We compare several variants of the Lanczos algorithm. Finally, computational examples are presented in Section 5 to demonstrate the effectiveness and efficiency of the Lanczos algorithm.

## 2 Preliminaries

### 2.1 Properties of definite Bethe–Salpeter Hamiltonian matrices

We first briefly review some basic spectral properties of definite BSH matrices. Detailed discussion on these properties can be found in [3, 27, 28].

Although a definite BSH matrix $H$ defined in (1) is in general non-Hermitian, it is diagonalizable and has real spectrum. Moreover the special structure of the BSH leads to a structured spectral decomposition as stated in Theorem 1 below.

**Theorem 1** ([27, Theorem 3]). *Let $H$ be a definite Bethe–Salpeter Hamiltonian matrix. Then the spectral decomposition of $H$ is of the form $H = Z \operatorname{diag} \{\Lambda_+, \Lambda_-\} Z^{-1}$ where*

$$Z = \begin{bmatrix} X & \overline{Y} \\ Y & \overline{X} \end{bmatrix}, \qquad Z^{-1} = C_n Z^{\mathsf{H}} C_n = \begin{bmatrix} X & -\overline{Y} \\ -Y & \overline{X} \end{bmatrix}^{\mathsf{H}}, \tag{5}$$

$\Lambda_+ = \operatorname{diag} \{\lambda_1, \lambda_2, \ldots, \lambda_n\}$, *and* $\Lambda_- = \operatorname{diag} \{\lambda_{n+1}, \lambda_{n+2}, \ldots, \lambda_{2n}\}$ *with*

$$\lambda_1 = -\lambda_{n+1} \geq \lambda_2 = -\lambda_{n+2} \geq \cdots \geq \lambda_n = -\lambda_{2n} > 0.$$

3

Since the eigenvalues of $H$ appear in positive and negative pairs $\pm\lambda_j$, we use $\lambda_j^+ \equiv \lambda_j$ and $\lambda_j^- \equiv -\lambda_j$, for $1 \leq j \leq n$, in the following to emphasize on the signs of these eigenvalues. Let $X = [x_1, \ldots, x_n]$, $Y = [y_1, \ldots, y_n] \in \mathbb{C}^{n \times n}$ be the submatrices in (5). Theorem 1 suggests that the right and left eigenvectors associated with the positive eigenvalue $\lambda_j^+$ are $z_j = [x_j^{\mathsf{H}}, y_j^{\mathsf{H}}]^{\mathsf{H}}$ and $C_n z_j = [x_j^{\mathsf{H}}, -y_j^{\mathsf{H}}]^{\mathsf{H}}$ respectively, and the right and left eigenvectors associated with $\lambda_j^-$ are $z_{n+j} = [\bar{y}_j^{\mathsf{H}}, \bar{x}_j^{\mathsf{H}}]^{\mathsf{H}}$ and $-C_n z_{n+j} = [-\bar{y}_j^{\mathsf{H}}, \bar{x}_j^{\mathsf{H}}]^{\mathsf{H}}$ respectively. The normalization condition $(C_n Z C_n)^{\mathsf{H}} Z = I_{2n}$ implies that

$$x_j^{\mathsf{H}} x_j - y_j^{\mathsf{H}} y_j = 1$$

for $j = 1, \ldots, n$. As long as the right eigenvectors associated with the positive eigenvalues are properly normalized, other eigenvectors can be easily recovered.

From (5), we observe that the right eigenvectors of $H$ are orthonormal with respect to the $C$-inner product, $\langle u, v \rangle_C = v^{\mathsf{H}} C_n u$, which is an *indefinite inner product*. Another observation is

$$Z^{\mathsf{H}} \Omega Z = Z^{\mathsf{H}} C_n Z \operatorname{diag}\{\Lambda_+, -\Lambda_+\} = C_n \operatorname{diag}\{\Lambda_+, -\Lambda_+\} = \operatorname{diag}\{\Lambda_+, \Lambda_+\}, \qquad (6)$$

indicating that the right eigenvectors of $H$ are also orthogonal with respect to the $\Omega$-inner product $\langle u, v \rangle_\Omega = v^{\mathsf{H}} \Omega u$. These orthogonalities are crucial for developing structure preserving Lanczos procedures. By structure preserving, we mean that the positive and negative pairing of the eigenvalues is preserved in the approximations to the eigenvalues of BSH.

## 2.2 Optical absorption spectra

Let $(z_r)_j$ and $(z_l)_j$ be the right and left eigenvectors of $H$, respectively, associated with the eigenvalue $\lambda_j$, $(1 \leq j \leq 2n)$. We denote by $\varepsilon_2(\omega)$ the imaginary part of the macroscopic dielectric function; $\varepsilon_2(\omega)$ is also proportional to the optical absorption spectrum of a material, and can be computed in a straightforward way from the eigenvalues and eigenvectors of the BSH as

$$\varepsilon_2(\omega) = \frac{8\pi^2 e^2}{V_{\text{xtal}}} \epsilon(\omega),$$

$$\epsilon(\omega) := d_r^{\mathsf{H}} \delta(\omega I_{2n} - H) d_l = \sum_{j=1}^{2n} \frac{(d_r^{\mathsf{H}}(z_r)_j)((z_l)_j^{\mathsf{H}} d_l)}{(z_l)_j^{\mathsf{H}}(z_r)_j} \delta(\omega - \lambda_j), \qquad (7)$$

where $V_{\text{xtal}}$ is the crystal volume, $e$ is the elementary charge, and

$$d_r = \begin{bmatrix} d \\ -\bar{d} \end{bmatrix} \quad \text{and} \quad d_l = \begin{bmatrix} d \\ \bar{d} \end{bmatrix}$$
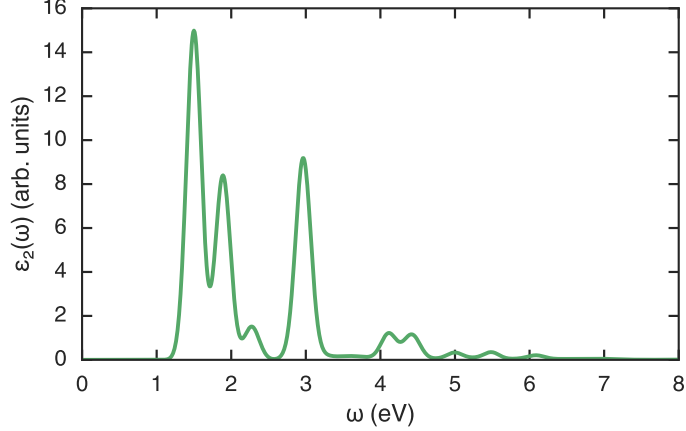
Figure 1: A typical curve for the imaginary part of the dielectric function. This curve is obtained from a single-wall $(8,0)$ carbon nanotube.

are the right and left *optical transition vectors*, respectively. Because the $d_r$ and $d_l$ depend solely on $d$, we will simply refer to $d$ as the optical transition vector. The coefficient, $(d_r^{\mathsf{H}}(z_r)_j)((z_l)_j^{\mathsf{H}} d_l)/((z_l)_j^{\mathsf{H}}(z_r)_j)$, of the Dirac delta function $\delta(\omega - \lambda_j)$ in (7) is known as the *oscillator strength* associated to the excitonic state $j$. Figure 1 shows a typical curve for the imaginary part of the dielectric function. In order to produce this plot, each Dirac delta function was broadened by a Gaussian function, as we will discuss below. The height of each peak in the spectrum is determined by the oscillator strength associated to each eigenvalue $\lambda_j$ and the number of eigenvalues clustered around an energy. Since the optical absorption spectrum is proportional to $\varepsilon_2(\omega)$, which is in turn proportional to $\epsilon(\omega)$, in this work we will broadly refer to both $\epsilon(\omega)$ and $\varepsilon_2(\omega)$ as the optical absorption spectrum of a material.

If $H$ can be fully diagonalized, we can compute $\epsilon(\omega)$ using the eigenpairs of $H$. However, diagonalizing $H$ is often costly, especially when the dimension of $H$ becomes large. For instance, for many low dimensional systems such as monolayer $MoS_2$, $n$ is on the order of 360,000 [18]. Similarly, large $n$'s are required to fully converge calculations on semiconducting carbon nanotubes and to obtain the correct order of the excited excitonic states in bulk semiconductors, for example. Therefore, it is natural to seek alternative approaches.

Using the structure of the eigenvectors of the BSH matrix $H$, we can simplify the expression of the absorption spectrum. For positive eigenpairs, we can choose $(z_r)_j = z_j$ and $(z_l)_j = C_n z_j$ so that $(z_l)_j^{\mathsf{H}}(z_r)_j = 1$. Then we have

$$(z_l)_j^{\mathsf{H}} d_l = (C_n z_j)^{\mathsf{H}}(C_n d_r) = z_j^{\mathsf{H}} d_r = \overline{d_r^{\mathsf{H}} z_j}.$$

5

It follows that

$$\epsilon^+(\omega) := \sum_{j=1}^{n} \frac{(d_r^{\mathsf{H}}(z_r)_j)((z_l)_j^{\mathsf{H}}d_l)}{(z_l)_j^{\mathsf{H}}(z_r)_j} \delta(\omega - \lambda_j)$$

$$= \sum_{j=1}^{n} |d_r^{\mathsf{H}} z_j|^2 \delta(\omega - \lambda_j^+)$$

$$= \sum_{j=1}^{n} |d^{\mathsf{H}} x_j - \overline{d}^{\mathsf{H}} y_j|^2 \delta(\omega - \lambda_j^+).$$

We remark that the oscillator strength $|d_r^{\mathsf{H}} z_j|^2$ is nonnegative. Similarly, for negative eigenpairs, we have

$$(z_l)_j^{\mathsf{H}} d_l = (-C_n z_j)^{\mathsf{H}} (C_n d_r) = -z_j^{\mathsf{H}} d_r = -\overline{d_r^{\mathsf{H}} z_j}$$

and

$$\epsilon^-(\omega) := -\sum_{j=n+1}^{2n} |d_r^{\mathsf{H}} z_j|^2 \delta(\omega - \lambda_j) = -\sum_{j=1}^{n} |d^{\mathsf{H}} x_j - \overline{d}^{\mathsf{H}} y_j|^2 \delta(\omega + \lambda_j^+) = -\epsilon^+(-\omega).$$

Therefore, the absorption spectrum

$$\epsilon(\omega) = \epsilon^+(\omega) + \epsilon^-(\omega) = \epsilon^+(\omega) - \epsilon^+(-\omega)$$

can be viewed as an odd function of the frequency $\omega$ in the distribution sense.

In practice, it is not desirable to plot the imaginary part of the polarizability as a sum of Dirac delta functions. A broadened peaked function, such as the Lorentzian function

$$L_\sigma(\omega) := \frac{1}{\pi} \cdot \frac{\sigma}{\omega^2 + \sigma^2} = \frac{1}{\pi} \mathrm{Im} \frac{1}{\omega - \mathrm{i}\sigma}$$

or the Gaussian function

$$G_\sigma(\omega) := \frac{1}{\sqrt{2\pi}\,\sigma} \mathrm{e}^{-\omega^2/(2\sigma^2)},$$

is used to replace the Dirac delta function, where the broadening factor $\sigma > 0$ is a small number. The first reason for doing so is because there is physically a lifetime associated to each excitonic state. The second reason is due to discretization procedures in performed in the calculations, such as employing a finite number of $k$-points in calculations on extended systems. If a calculation could be carried with infinitely many $k$-points, the optical absorption spectrum would consist of a few isolated low-energy sharp peaks, but the delta functions merge at higher energy and form a continuum spectrum. On the other hand, a calculation performed with a finite number of $k$-points only samples a finite number of transitions in this continuum region.

Therefore, we wish to plot the imaginary part of the dielectric function using a generic peaked function $g_\sigma(\omega)$ (either $L_\sigma(\omega)$ or $G_\sigma(\omega)$) characterized by a typical small width $\sigma$ instead of $\delta(\omega)$. The imaginary part of the dielectric matrix can be expressed now in terms of

$$\epsilon_\sigma(\omega) = d_r^{\mathsf{H}} g_\sigma(\omega I_{2n} - H) d_l = \sum_{j=1}^{n} \left| d^{\mathsf{H}} x_j - \overline{d}^{\mathsf{H}} y_j \right|^2 \left[ g_\sigma(\omega - \lambda_j^+) - g_\sigma(\omega + \lambda_j^+) \right], \qquad (8)$$

which is an odd function in $\omega$, that is, $\epsilon_\sigma(-\omega) = -\epsilon_\sigma(\omega)$. Thus it suffices to compute the function value for $\omega > 0$.

Note that (8), which is a scalar function of $\omega$, can be viewed as an expected value of a matrix function. We are interested in the positions and heights of the major peaks of this function, which are given by the eigenvalues and eigenvectors of the BSH. However, the precise position and height of each peak is seldom required, especially since the underlying theories employed to obtain these spectra are already themselves approximate. Therefore, efficient methods that can provide estimates of (8) without computing each individual eigenpair of $H$ are of great interest. In Sections 3 and 4, we discuss how to use Lanczos algorithms to estimate $\epsilon_\sigma(\omega)$ efficiently.

# 3   Tamm–Dancoff approximation

Tamm–Dancoff approximation (TDA) [6, 25, 31] is a technique often used to in practice to reduce the computational cost of the absorption spectrum calculation. For many systems, especially on bulk semiconductors and metals, the TDA incurs a very small error in the optical absorption spectrum, and for that reason it has been a widely used approximation in condensed-matter physics. In this section we discuss how to estimate the absorption spectrum with a Lanczos procedure within the TDA.

However, we remark that for many systems, including systems with reduced dimensionality optically excited with light polarized along a confined direction, the TDA may incur in large errors for the optical absorption spectrum. We shall discuss full BSE solvers in Section 4.

## 3.1   Lanczos algorithm

In TDA, the off-diagonal block of $H$, $B$, is set to zero. We denote the resulting block diagonal BSH by $H_{\mathrm{TDA}} = \mathrm{diag}\left\{ A, -\overline{A} \right\}$, which is a Hermitian matrix. It follows that the absorption spectrum associated with $H_{\mathrm{TDA}}$ becomes

$$\epsilon(\omega) = d_r^{\mathsf{H}} \delta(\omega I_{2n} - H_{\mathrm{TDA}}) d_l = d^{\mathsf{H}} \delta(\omega I_n - A) d - \overline{d^{\mathsf{H}} \delta(\omega I_n + A) d}.$$

As $d^{\mathsf{H}} \delta(\omega I_n \pm A) d$ is real and nonnegative, we can omit the complex conjugation in the second term. In practice, we compute

$$\epsilon_\sigma(\omega) = d^{\mathsf{H}} g_\sigma(\omega I_n - A) d - d^{\mathsf{H}} g_\sigma(\omega I_n + A) d =: d^{\mathsf{H}} f(A; \omega) d \qquad (9)$$

for $\omega > 0$, where $f(t; \omega) = g_\sigma(\omega - t) - g_\sigma(\omega + t)$.

Since $A$ is Hermitian and positive definite, the matrix functionals in (9) can be estimated using the Lanczos algorithm. Starting with $u_1 = d/\|d\|_2$, a $k$-step Lanczos procedure produces

$$AU_k = U_k T_k + \beta_k u_{k+1} e_k^{\mathsf{H}}, \qquad U_{k+1}^{\mathsf{H}} U_{k+1} = I_{k+1}, \tag{10}$$

where

$$T_k = \operatorname{tridiag} \left\{ \begin{matrix} & \beta_1 & \cdots & & \beta_{k-1} & \\ \alpha_1 & & \cdots & & \cdots & \alpha_k \\ & \beta_1 & \cdots & & \beta_{k-1} & \end{matrix} \right\} \tag{11}$$

is a real symmetric, tridiagonal, positive definite, and componentwise nonnegative matrix. Here we use the convention $U_j = [u_1, \ldots, u_j]$ (for $1 \leq j \leq k+1$) to represent the Lanczos vectors, and $e_j$ is the $j$th column of the identity matrix. Then the absorption spectrum can be estimated by

$$d^{\mathsf{H}} f(A; \omega) d = \|d\|_2^2 u_1^{\mathsf{H}} f(A; \omega) U_k e_1 \approx \|d\|_2^2 u_1^{\mathsf{H}} U_k f(T_k; \omega) e_1 = \|d\|_2^2 e_1^{\mathsf{H}} f(T_k; \omega) e_1. \tag{12}$$

As long as $k \ll n$, the matrix function of the projected matrix $T_k$, $f(T_k; \omega)$, can be easily evaluated by diagonalizing $T_k$. Moreover, there is no need to explicitly store the whole history of the Lanczos vectors because eventually only $T_k$ is used in (12). However, it is important to ensure columns of the generated $U_{k+1}$ matrix are orthonormal. A desired feature here is that the estimated absorption spectrum in this approach is nonnegative for $\omega > 0$. Clearly, the Lanczos algorithm possesses this desired feature.

Finally, we remark that when the Gaussian functions are replaced by Lorentzian functions Haydock's recursive algorithm [12, 13, 14, 15] is mathematically equivalent to the Lanczos algorithm. As the Lanczos algorithm is more general—it can handle any approximation to the Dirac delta function, it is a simple and flexible replacement of Haydock's recursive method in this context. Another advantage of the Lanczos algorithm will be discussed in the next subsection.

## 3.2 Generalized averaged Gauss quadrature

It is well known that the Lanczos algorithm for estimating matrix functionals can be interpreted as Gauss quadrature [8, 9]. In [19], a recently developed generalized averaged Gauss quadrature rule [29] has been adopted to improve the accuracy of the Lanczos algorithm with little extra effort. In the following we briefly describe the procedure of this approach.

After the $k$-step Lanczos procedure is performed, we can construct a $(2k-1) \times (2k-1)$ symmetric tridiagonal matrix $\widehat{T}_k$ as

$$\widehat{T}_k = \operatorname{tridiag} \left\{ \begin{matrix} & \beta_1 & \cdots & \beta_{k-1} & \beta_k & \beta_{k-2} & \cdots & \beta_1 & \\ \alpha_1 & & \cdots & \cdots & \alpha_k & \alpha_{k-1} & \cdots & \cdots & \alpha_1 \\ & \beta_1 & \cdots & \beta_{k-1} & \beta_k & \beta_{k-2} & \cdots & \beta_1 & \end{matrix} \right\}. \tag{13}$$

Then we replace $e_1^{\mathsf{H}} f(T_k; \omega) e_1$ in (12) by $e_1^{\mathsf{H}} f(\widehat{T}_k; \omega) e_1$,[1] that is,

$$d^{\mathsf{H}} f(A; \omega) d \approx \|d\|_2^2 e_1^{\mathsf{H}} f(\widehat{T}_k; \omega) e_1. \qquad (14)$$

When $k$ is not very large, the cost of computing $f(T_k; \omega)$ or $f(\widehat{T}_k; \omega)$ is negligible compared to that of forming $T_k$. As the spectrum of $\widehat{T}_k$ is a superset of that of $T_{k-1}$, and $\Lambda(\widehat{T}_k) \backslash \Lambda(T_{k-1})$ interlaces with $\Lambda(T_{k-1})$, (14) should be a better approximation compared to (12) with negligible computational overhead. We refer the readers to [19, 29] for detailed discussions.

If the Lanczos procedure breaks down at $k$th step, that is, $\beta_k = 0$, then (12) holds exactly instead of approximately. In this lucky breakdown, (14) also holds exactly because $\widehat{T}_k$ decouples into two tridiagonal submatrices. We remark that an extra benefit of using the generalized averaged Gauss quadrature is that, for the same number of quadrature points, the generalized averaged Gauss quadrature requires fewer Lanczos steps, hence the risk of loss of orthogonality among the Lanczos vectors is reduced.

Certainly the generalized averaged Gauss quadrature can be adopted here for the estimation of absorption spectrum. Similar to the Lanczos algorithm with standard Gauss quadrature, the generalized averaged Gauss quadrature also produces nonnegative oscillator strengths. Thus the estimated absorption spectrum is also nonnegative for $\omega > 0$ when $\widetilde{T}_k$ is positive definite. However, $\widehat{T}_k$ as defined in (13) can sometimes has one nonpositive eigenvalue. (The second smallest eigenvalue of $\widehat{T}_k$ is always positive since $\Lambda(\widehat{T}_k) \backslash \Lambda(T_{k-1})$ interlaces with $\Lambda(T_{k-1})$.) Such a nonpositive eigenvalue may violate the property $\epsilon_\sigma(\omega) \geq 0$ for $\omega > 0$. A simple remedy is to redefine $f(t; \omega)$ as

$$f(t; \omega) = \begin{cases} g_\sigma(\omega - t) - g_\sigma(\omega + t) & \text{if } t > 0, \\ 0 & \text{if } t \leq 0. \end{cases}$$

Then in the resulting generalized averaged Gauss quadrature (14) we can simply discard the term involving the nonpositive eigenvalue of $\widehat{T}_k$, if there is any. In fact, dropping the nonpositive eigenvalue does not affect the accuracy, because the eigenvalues of $T_{k-1}$, as the common Gauss quadrature nodes for both (12) and (14) (assuming in (12) we use the approximation from a $(k-1)$-step Lanczos procedure instead of a $k$-step one), have the same weights (up to scaling) in both quadrature rules [29]. We summarize the Lanczos algorithm with generalized averaged Gauss quadrature in Algorithm 1. The utility of generalized averaged Gauss quadrature provides another advantage of the Lanczos algorithm over Haydock's recursive algorithm.

## 4   Absorption spectrum for full BSE

In this section we investigate how to estimate the absorption spectrum without using the Tamm–Dancoff approximation. Like the Lanczos algorithm in the TDA setting, the

---

[1]The vector $e_1$ is of length $k$ in $e_1^{\mathsf{H}} f(T_k; \omega) e_1$, and is of length $2k - 1$ in $e_1^{\mathsf{H}} f(\widehat{T}_k; \omega) e_1$.

**Algorithm 1** Lanczos algorithm for estimating the absorption spectrum under TDA.

**Input:** A Hermitian positive definite matrix $A \in \mathbb{C}^{n \times n}$, an optical transition vector $d \in \mathbb{C}^n$, a broadening factor $\sigma > 0$, the number of Lanczos steps $k$, and a set of frequencies $\{\omega_i\}_{i=1}^m$.

**Output:** The estimated absorption spectrum $\epsilon_\sigma(\omega)$ sampled at $\omega_i$ (for $1 \leq i \leq m$).

1: Perform $k$ Lanczos steps using $d$ as the starting vector.
2: Formulate $\widehat{T}_k$ as defined in (13).
3: Compute the spectral decomposition $\widehat{T}_k = \widehat{S}_k \, \text{diag}\left\{\widehat{\theta}_1, \ldots, \widehat{\theta}_{2k-1}\right\} \widehat{S}_k^{\mathsf{H}}$, where $\widehat{S}_k^{\mathsf{H}} \widehat{S}_k = I_{2k-1}$.
4: Evaluate

$$\epsilon_\sigma(\omega_i) = \|d\|_2^2 \sum_{\substack{j=1 \\ \widehat{\theta}_j > 0}}^{2k-1} |\widehat{S}_k(1,j)|^2 \big[g_\sigma(\omega_i - \widehat{\theta}_j) - g_\sigma(\omega_i + \widehat{\theta}_j)\big]$$

for $i = 1, \ldots, m$.

following features are desired.

1. Any breakdown in the Lanczos procedure is a lucky breakdown.

2. The computed absorption spectrum is real and nonnegative for $\omega > 0$.

3. The full history of Lanczos vectors is not required.

4. The technique of generalized averaged Gauss quadrature can be applied.

We shall demonstrate that all these features are feasible for full BSE calculations.

## 4.1 Lanczos algorithm for real BSE

We first examine a simpler case in which both $A$ and $B$ are real symmetric matrices, and

$$H = \begin{bmatrix} A & B \\ -B & -A \end{bmatrix} \in \mathbb{R}^{2n \times 2n} \tag{15}$$

is real also. Such an $H$ results from systems with real-space inversion symmetry. It is not difficult to verify that the condition (4) is equivalent to the following conditions:

$$M := A + B \succ 0, \qquad K := A - B \succ 0. \tag{16}$$

We also assume that the optical transition vector $d$ is real. A Lanczos algorithm that can be used to estimate $\epsilon(\omega)$ for BSH matrices of this type has been studied in [5], in the context of linear response time-dependent density functional theory (TDDFT) based calculations. In the following, we briefly summarize this algorithm.

Using the spectral decomposition of $H$ as shown in Theorem 1, we can verify that

$$M = (X - Y)\Lambda_+(X - Y)^\mathsf{T}, \qquad K = (X + Y)\Lambda_+(X + Y)^\mathsf{T},$$

and

$$(X - Y)^\mathsf{T}(X + Y) = I_n.$$

Then we have

$$\epsilon(\omega) = \sum_{j=1}^{n} \left[d^\mathsf{T}(x_j - y_j)\right]^2 \left[\delta(\omega - \lambda_j^+) - \delta(\omega + \lambda_j^+)\right]$$

$$= 2\operatorname{sign}(\omega) \sum_{j=1}^{n} \lambda_j^+ \left[d^\mathsf{T}(x_j - y_j)\right]^2 \delta\left(\omega^2 - (\lambda_j^+)^2\right)$$

$$= 2\operatorname{sign}(\omega) d^\mathsf{T}(X - Y)\Lambda_+(X - Y)^\mathsf{T}(X + Y)\delta(\omega^2 I_n - \Lambda_+^2)(X - Y)^\mathsf{T}d$$

$$= 2\operatorname{sign}(\omega) d^\mathsf{T} M \delta(\omega^2 I_n - KM)d.$$

Therefore, we reduce this problem size from $2n \times 2n$ to $n \times n$. Although $KM$ is nonsymmetric in general, it is symmetric and positive definite with respect to the $M$-inner product because

$$\langle x, KMy \rangle_M = y^\mathsf{T} MKMx = \langle KMx, y \rangle_M.$$

A Lanczos procedure in which standard Euclidean inner product is replaced with an $M$-inner product reads

$$KMU_k = U_k T_k + \beta_k u_{k+1} e_k^\mathsf{T}, \tag{17}$$

with $u_1 = d/\|d\|_M$ and $U_{k+1}^\mathsf{T} MU_{k+1} = I_{k+1}$. Algorithm 2 outlines the computational procedure of calculating (17). We remark that in [5] full orthogonalization is used to retain numerical stability of the Lanczos procedure. In contrast, Algorithm 2 uses a careful formulation of short recurrence. The numerical stability is observed to be comparable with full orthogonalization if $k$ is reasonably small.

It follows from (17) and the identity

$$\delta(\omega - |\lambda|) - \delta(\omega + |\lambda|) = 2|\lambda|\operatorname{sign}(\omega)\delta(\omega^2 - \lambda^2)$$

that $\epsilon(\omega)$ can be approximated through

$$\epsilon(\omega) = 2\operatorname{sign}(\omega) d^\mathsf{T} M \delta(\omega^2 I_n - KM)d$$

$$\approx 2\operatorname{sign}(\omega) d^\mathsf{T} MU_k \delta\left(\omega^2 I_k - T_k\right)U_k^\mathsf{T} Md$$

$$= \|d\|_M^2 e_1^\mathsf{T} \left[\delta\left(\omega I_k - T_k^{1/2}\right) - \delta\left(\omega I_k + T_k^{1/2}\right)\right]T_k^{-1/2}e_1$$

$$\approx \|d\|_M^2 e_1^\mathsf{T} \left[g_\sigma\left(\omega I_k - T_k^{1/2}\right) - g_\sigma\left(\omega I_k + T_k^{1/2}\right)\right]T_k^{-1/2}e_1. \tag{18}$$

11

**Algorithm 2** Lanczos procedure in $M$-inner product for real full BSE.

---

**Input:** A definite Bethe–Salpeter Hamiltonian matrix $H \in \mathbb{R}^{2n \times 2n}$; a starting vector $u_1 \in \mathbb{R}^n$ satisfying $u_1^\mathsf{T}(A + B)u_1 = 1$; the number of Lanczos steps, $k$.

**Output:** $\alpha_1, \ldots, \alpha_k, \beta_1, \ldots, \beta_k \in \mathbb{R}$, and $u_1, \ldots, u_{k+1} \in \mathbb{R}^n$ satisfying (17) and $U_{k+1}^\mathsf{T}(A + B)U_{k+1} = I_{k+1}$.

1: $\beta_0 \leftarrow 0, \quad u_0 \leftarrow 0, \quad v_0 \leftarrow 0$.
2: $v_1 = (A + B)u_1$.
3: **for** $j = 1, \ldots, k$ **do**
4:      $x \leftarrow (A - B)v_j - \beta_{j-1}u_{j-1}$.
5:      $\alpha_j \leftarrow v_j^\mathsf{T} x$.
6:      $x \leftarrow x - \alpha_j u_j$.
7:      $y \leftarrow (A + B)x$.
8:      $\beta_j \leftarrow \sqrt{x^\mathsf{T} y}$.
9:      $u_{j+1} \leftarrow x/\beta_j, \quad v_{j+1} \leftarrow y/\beta_j$.
10: **end for**

---

Here $T_k$ is a real symmetric tridiagonal matrix as in (11). Similar to the Lanczos algorithm in the TDA setting, the approximate $\epsilon_\sigma(\omega)$ is nonnegative for $\omega > 0$, which is a desired property. There is also no need to keep the whole history of Lanczos vectors.

We have already seen in Section 3.2 that the generalized averaged Gauss quadrature can be incorporated in the Lanczos algorithm. This is also the case for (18). Let

$$f(t; \omega) = t^{-1/2}\big[g(\omega - t^{1/2}) - g(\omega + t^{1/2})\big].$$

Then the generalized averaged Gauss quadrature replaces $e_1^\mathsf{T} f(T_k; \omega)e_1$ in (18) by $e_1^\mathsf{T} f(\widehat{T}_k; \omega)e_1$, that is,

$$\epsilon_\sigma(\omega) \approx \|d_1\|_M^2 e_1^\mathsf{T} f(\widehat{T}_k; \omega)e_1, \tag{19}$$

where $\widehat{T}_k \in \mathbb{R}^{(2k-1)\times(2k-1)}$ is as defined in (13). It is expected that (19) in general provides a better approximation to $\epsilon_\sigma(\omega)$ compared to (18). Similar to the discussions in Section 3.2, $\widehat{T}_k$ can sometimes has one nonpositive eigenvalue. But in (18) $\widehat{T}_k$ needs to be positive definite so that $\widehat{T}_k^{1/2}$ is also positive definite. The remedy is to extend the definition of $f(t; \omega)$ as

$$f(t; \omega) = \begin{cases} t^{-1/2}\big[g_\sigma(\omega - t^{1/2}) - g_\sigma(\omega + t^{1/2})\big] & \text{if } t > 0, \\ 0 & \text{if } t \leq 0, \end{cases}$$

and discard the term involving the nonpositive eigenvalue of $\widehat{T}_k$, if there is any. Algorithm 3 summarizes the Lanczos algorithm for real full BSE incorporated with the generalized averaged Gauss quadrature.

We should point out that $\epsilon(\omega)$ can be obtained by computing the eigenpairs of $KM$ or $H$ directly. If one is only interested in the low energy region of the absorption spectrum,

**Algorithm 3** The Lanczos algorithm for estimating the optical absorption spectrum for real full BSE.

---

**Input:** Real symmetric positive definite matrices $M$, $K \in \mathbb{R}^{n \times n}$, an optical transition vector $d \in \mathbb{R}^n$, a broadening factor $\sigma > 0$, the number of Lanczos steps $k$, and a set of frequencies $\{\omega_i\}_{i=1}^m$.

**Output:** The estimated absorption spectrum $\epsilon_\sigma(\omega)$ sampled at $\omega_i$ (for $1 \leq i \leq m$).

1: Perform $k$ Lanczos steps in $M$-inner product using $d$ as the starting vector.
2: Formulate $\widehat{T}_k$ as defined in (13).
3: Compute the spectral decomposition $\widehat{T}_k = \widehat{S}_k \operatorname{diag}\left\{\widehat{\theta}_1^2, \ldots, \widehat{\theta}_{2k-1}^2\right\} \widehat{S}_k^{\mathsf{H}}$, where $\widehat{S}_k^{\mathsf{H}} \widehat{S}_k = I_{2k-1}$ and $\widehat{\theta}_{2k-1} \geq \cdots \geq \widehat{\theta}_2 > 0$.
4: Evaluate

$$\epsilon_\sigma(\omega_i) = d^{\mathsf{T}} M d \sum_{\substack{j=1 \\ \widehat{\theta}_j > 0}}^{2k-1} |\widehat{S}_k(1,j)|^2 \frac{g_\sigma(\omega_i - \widehat{\theta}_j) - g_\sigma(\omega_i + \widehat{\theta}_j)}{\widehat{\theta}_j}$$

for $i = 1, \ldots, m$.

---

iterative methods such as the ones proposed in [1, 20, 33] can be used to compute the first few eigenpairs. However, these methods can become costly when the absorption spectrum window becomes large, and more eigenpairs need to be computed.

## 4.2 Structure preserving Lanczos procedure for complex BSE

In this subsection we discuss how to develop a structure preserving Lanczos procedure for complex BSE. Just like the real case, we will try to reformulate the problem so that only $n$-dimensional matrices and vectors are involved. To this end, let us define

$$\mathcal{U}_\phi = \left\{ \begin{bmatrix} u \\ e^{i\phi}\overline{u} \end{bmatrix} : u \in \mathbb{C}^n \right\}, \qquad (\phi \in \mathbb{R}).$$

It can be easily verified that $H\mathcal{U}_\phi = \mathcal{U}_{\phi+\pi}$ and $H^2\mathcal{U}_\phi = \mathcal{U}_\phi$. However, we remark that $\mathcal{U}_\phi$ is *not* an invariant subspaces of $H^2$ as it is not a subspace of $\mathbb{C}^{2n}$ over $\mathbb{C}$; it can only be regarded as a linear space over $\mathbb{R}$. To approximate $d_r^{\mathsf{H}} \delta(\omega I_{2n} - H) d_l$ using a Lanczos procedure, it is natural to use $d_l$ as the starting vector. Note that $d_l$ and $d_r$ are structured because $d_l \in \mathcal{U}_0$ and $d_r \in \mathcal{U}_\pi$. In the following, we discuss how to preserve this type of structure in a Lanczos procedure.

It was observed in [10] that $H = C_n \Omega$ is self-adjoint with respect to the inner product defined by $\Omega$ in (3), because

$$\langle x, Hy \rangle_\Omega = y^{\mathsf{H}} \Omega C_n \Omega x = \langle Hx, y \rangle_\Omega.$$

We make another observation that $H^2 = (C_n \Omega C_n)\Omega$ is Hermitian and positive definite with respect to the $\Omega$-inner product. Thus there exists a Lanczos procedure associated with $H^2$ that is defined in terms of the $\Omega$-inner product. If we start with the vector $q_1 \in \mathcal{U}_0$, the recurrence relationship among the Lanczos vectors is characterized by the following theorem.

**Theorem 2.** *Let $H = C_n \Omega$ be a definite Bethe–Salpeter Hamiltonian matrix. Suppose that $u_1 \in \mathbb{C}^n$ satisfies $\mathrm{Re}(u_1^\mathsf{H} A u_1 + u_1^\mathsf{H} B \overline{u}_1) = 1$. Then for $k < n$, applying a $k$-step Lanczos procedure to $H^2$ in the $\Omega$-inner product with the starting vector $[u_1^\mathsf{H}, \overline{u}_1^\mathsf{H}]^\mathsf{H}$ produces*

$$H^2 \begin{bmatrix} U_k \\ \overline{U}_k \end{bmatrix} = \begin{bmatrix} U_k \\ \overline{U}_k \end{bmatrix} T_k + \beta_k \begin{bmatrix} u_{k+1} \\ \overline{u}_{k+1} \end{bmatrix} e_k^\mathsf{H}, \tag{20}$$

*where $U_k = [u_1, \ldots, u_k] \in \mathbb{C}^{n \times k}$, $T_k \in \mathbb{R}^{k \times k}$ is as defined in (11). The tridiagonal matrix $T_k$ is positive definite and componentwise nonnegative, and $\beta_k > 0$, if the Lanczos procedure does not break down. The Lanczos vectors satisfy the orthogonality condition*

$$\begin{bmatrix} u_i \\ \overline{u}_i \end{bmatrix}^\mathsf{H} \Omega \begin{bmatrix} u_j \\ \overline{u}_j \end{bmatrix} = 2\delta_{ij}, \qquad (1 \leq i, j \leq k+1), \tag{21}$$

*where $\delta_{ij}$ is the Kronecker delta notation.*

*Proof.* In the generic case (i.e., assuming no breakdown occurs), the Arnoldi procedure using the orthogonality condition (21) with starting vector $q_1 = [u_1^\mathsf{H}, \overline{u}_1^\mathsf{H}]$ reads

$$H^2 Q_k = Q_k T_k + \beta_k q_{k+1} e_k^\mathsf{H},$$

where $T_k$ is an upper Hessenberg matrix with positive subdiagonal entries, and $\beta_k > 0$. Multiplying from the left by $Q_k^\mathsf{H} \Omega$, we obtain that

$$2T_k = Q_k^\mathsf{H} \Omega H^2 Q_k = (C_n \Omega Q_k)^\mathsf{H} \Omega (C_n \Omega Q_k)$$

is Hermitian positive definite. Consequently the diagonal entries of $T_k$ are real and positive. Hence we conclude that $T_k$ is real symmetric, tridiagonal, positive definite, and componentwise nonnegative. The Arnoldi procedure is in fact a Lanczos procedure.

Let us denote by $\alpha_i$ and $\beta_i$, respectively, the $i$th diagonal and subdiagonal entries of $T_k$, i.e., $T_k$ is of the form (11). Notice that $q_1 \in \mathcal{U}_0$ implies $H^2 q_1 \in \mathcal{U}_0$. From the Lanczos procedure we have

$$q_2 = \frac{1}{\beta_1}(H^2 q_1 - \alpha_1 q_1) \in \mathcal{U}_0,$$

because both $\alpha_1$ and $\beta_1$ are real. By induction, we have

$$q_{i+1} = \frac{1}{\beta_i}(H^2 q_i - \alpha_i q_i - \beta_{i-1} q_{i-1}) \in \mathcal{U}_0$$

for $i = 2, \ldots, k$, as the linear combination on the vectors from $\mathcal{U}_0$ involves only real coefficients. This completes the proof. $\qquad\square$

14

In Section 4.4 we show that this Lanczos procedure reduces to the one given in Section 4.1 for real BSE. The additional factor of two in (21) is introduced to make the two Lanczos procedures identical.

It may appear that the Lanczos procedure associated with $H^2$ only provides one of the two sets of vectors required to construct approximations to the oscillator strength. The following observation shows that the other set of vectors can be easily recovered. Let

$$\begin{bmatrix} V_k \\ \overline{V}_k \end{bmatrix} = \Omega \begin{bmatrix} U_k \\ \overline{U}_k \end{bmatrix}, \tag{22}$$

or, equivalently,

$$\begin{bmatrix} V_k \\ -\overline{V}_k \end{bmatrix} = H \begin{bmatrix} U_k \\ \overline{U}_k \end{bmatrix}.$$

The $U_k$ and $V_k$ matrices can also be generated together from the following recurrence

$$H \begin{bmatrix} U_k & V_k \\ \overline{U}_k & -\overline{V}_k \end{bmatrix} = \begin{bmatrix} U_k & V_k \\ \overline{U}_k & -\overline{V}_k \end{bmatrix} \begin{bmatrix} 0 & T_k \\ I_k & 0 \end{bmatrix} + \beta_k \begin{bmatrix} u_{k+1} \\ \overline{u}_{k+1} \end{bmatrix} e_{2k}^{\mathsf{H}}. \tag{23}$$

The orthogonality condition (21) becomes

$$\begin{bmatrix} U_k \\ \overline{U}_k \end{bmatrix}^{\mathsf{H}} \begin{bmatrix} V_k \\ \overline{V}_k \end{bmatrix} = 2I_k. \tag{24}$$

However, this condition is not sufficient for constructing the (oblique) projector associated with the subspace

$$\mathrm{span} \begin{bmatrix} U_k & V_k \\ \overline{U}_k & -\overline{V}_k \end{bmatrix}.$$

We show a stronger result in the following theorem.

**Theorem 3.** *Under the same assumption given in Theorem 2. Let $U_k$ and $V_k$ be as defined in (22) and (23). Then*

$$\begin{bmatrix} V_k & U_k \\ \overline{V}_k & -\overline{U}_k \end{bmatrix}^{\mathsf{H}} \begin{bmatrix} U_k & V_k \\ \overline{U}_k & -\overline{V}_k \end{bmatrix} = 2I_{2k}. \tag{25}$$

*Proof.* Since

$$\begin{bmatrix} V_k & U_k \\ \overline{V}_k & -\overline{U}_k \end{bmatrix}^{\mathsf{H}} \begin{bmatrix} U_k & V_k \\ \overline{U}_k & -\overline{V}_k \end{bmatrix} = \begin{bmatrix} 2I_k & V_k^{\mathsf{H}}V_k - \overline{V_k^{\mathsf{H}}V_k} \\ U_k^{\mathsf{H}}U_k - \overline{U_k^{\mathsf{H}}U_k} & 2I_k \end{bmatrix},$$

it suffices to show that $u_i^{\mathsf{H}}u_j$ and $v_i^{\mathsf{H}}v_j$ are both real for all $i$ and $j$. The proof is based on the fact that

$$\begin{bmatrix} u \\ \pm\overline{u} \end{bmatrix}^{\mathsf{H}} (H^{\mathsf{H}})^{\ell_1} C_n H^{\ell_2} \begin{bmatrix} u \\ \pm\overline{u} \end{bmatrix} = \begin{bmatrix} u \\ \mp\overline{u} \end{bmatrix}^{\mathsf{H}} (C_n\Omega)^{(\ell_1+\ell_2)/2} C_n (\Omega C_n)^{(\ell_1+\ell_2)/2} \begin{bmatrix} u \\ \mp\overline{u} \end{bmatrix} = 0$$

15

holds for any $u \in \mathbb{C}^n$ and any nonnegative integers $\ell_1, \ell_2$ as long as $\ell_1 + \ell_2$ is even.

From (20) it can be verified that $[u_j^{\mathsf{H}}, \bar{u}_j^{\mathsf{H}}]^{\mathsf{H}}$ can be expressed as

$$\begin{bmatrix} u_j \\ \bar{u}_j \end{bmatrix} = p_j(H^2) \begin{bmatrix} u_1 \\ \bar{u}_1 \end{bmatrix},$$

where $p_j(\cdot)$ is a polynomial of degree $j$ with real coefficients. Then we obtain

$$2\mathrm{i} \cdot \mathrm{Im}(u_i^{\mathsf{H}} u_j) = \begin{bmatrix} u_i \\ \bar{u}_i \end{bmatrix}^{\mathsf{H}} C_n \begin{bmatrix} u_j \\ \bar{u}_j \end{bmatrix} = \begin{bmatrix} u_1 \\ \bar{u}_1 \end{bmatrix}^{\mathsf{H}} p_i(H^2)^{\mathsf{H}} C_n p_j(H^2) \begin{bmatrix} u_1 \\ \bar{u}_1 \end{bmatrix} = 0$$

by expanding $p_i(H^2)^{\mathsf{H}} C_n p_j(H^2)$ as the sum of monomials. Similarly, $[v_j^{\mathsf{H}}, -\bar{v}_j^{\mathsf{H}}]^{\mathsf{H}}$ can be expressed as

$$\begin{bmatrix} v_j \\ -\bar{v}_j \end{bmatrix} = H \begin{bmatrix} u_j \\ \bar{u}_j \end{bmatrix} = H p_j(H^2) \begin{bmatrix} u_1 \\ \bar{u}_1 \end{bmatrix} = p_j(H^2) H \begin{bmatrix} u_1 \\ \bar{u}_1 \end{bmatrix} = p_j(H^2) \begin{bmatrix} v_1 \\ -\bar{v}_1 \end{bmatrix},$$

and then

$$2\mathrm{i} \cdot \mathrm{Im}(v_i^{\mathsf{H}} v_j) = \begin{bmatrix} v_i \\ -\bar{v}_i \end{bmatrix}^{\mathsf{H}} C_n \begin{bmatrix} v_j \\ -\bar{v}_j \end{bmatrix} = \begin{bmatrix} v_1 \\ -\bar{v}_1 \end{bmatrix}^{\mathsf{H}} p_i(H^2)^{\mathsf{H}} C_n p_j(H^2) \begin{bmatrix} v_1 \\ -\bar{v}_1 \end{bmatrix} = 0. \qquad \square$$

From Theorem 3, we conclude that

$$\frac{1}{2} \begin{bmatrix} U_k & V_k \\ \bar{U}_k & -\bar{V}_k \end{bmatrix} \begin{bmatrix} V_k & U_k \\ \bar{V}_k & -\bar{U}_k \end{bmatrix}^{\mathsf{H}}$$

is the projector we seek, and

$$\begin{bmatrix} 0 & T_k \\ I_k & 0 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} V_k & U_k \\ \bar{V}_k & -\bar{U}_k \end{bmatrix}^{\mathsf{H}} H \begin{bmatrix} U_k & V_k \\ \bar{U}_k & -\bar{V}_k \end{bmatrix}.$$

is indeed a projected form of $H$.

The recurrence given by (23) is more desirable than that given by (20) because it removes the ambiguity introduced by squaring the eigenvalues of the projected matrix

$$\begin{bmatrix} 0 & T_k \\ I_k & 0 \end{bmatrix},$$

which appear in pairs $\pm\theta_i$, where $\theta_i^2$ is the eigenvalue of $T_k$. We regard (23) as a structure preserving Lanczos procedure as the spectrum of the projected matrix is real and symmetric with respect to the origin. Algorithm 4 outlines the structure preserving Lanczos procedure for complex BSE. Similar to Algorithm 2, a careful formulation of short recurrence instead of full orthogonalization is used to largely retain numerical stability.

16

---

**Algorithm 4** Lanczos procedure in $\Omega$-inner product for complex full BSE.

---

**Input:** A definite Bethe–Salpeter Hamiltonian matrix $H \in \mathbb{C}^{2n \times 2n}$; a starting vector $u_1 \in \mathbb{C}^n$ satisfying $\left[u_1^{\mathsf{H}}, \overline{u}_1^{\mathsf{H}}\right]^{\mathsf{H}} \Omega \left[u_1^{\mathsf{H}}, \overline{u}_1^{\mathsf{H}}\right]^{\mathsf{H}} = 2$; the number of Lanczos steps, $k$.

**Output:** $\alpha_1, \ldots, \alpha_k, \beta_1, \ldots, \beta_k \in \mathbb{R}$, and $u_1, \ldots, u_{k+1}, v_1, \ldots, v_{k+1} \in \mathbb{C}^n$ satisfying (22)–(24).

1: $\beta_0 \leftarrow 0, \quad u_0 \leftarrow 0, \quad v_0 \leftarrow 0.$
2: $v_1 = Au_1 + B\overline{u}_1.$
3: **for** $j = 1, \ldots, k$ **do**
4: $\quad x \leftarrow Av_j - B\overline{v}_j - \beta_{j-1}u_{j-1}.$
5: $\quad \alpha_j \leftarrow \mathrm{Re}(v_j^{\mathsf{H}} x).$
6: $\quad x \leftarrow x - \alpha_j u_j.$
7: $\quad y \leftarrow Ax + B\overline{x}.$
8: $\quad \beta_j \leftarrow \sqrt{\mathrm{Re}(x^{\mathsf{H}} y)}.$
9: $\quad u_{j+1} \leftarrow x/\beta_j, \quad v_{j+1} \leftarrow y/\beta_j.$
10: **end for**

---

Finally, we remark that this Lanczos procedure can be extended to have a starting vector from $\mathcal{U}_\phi$. Let $D_\phi = \mathrm{diag}\left\{I_n, e^{i\phi} I_n\right\}$. Notice that $D_\phi^{\mathsf{H}} H D_\phi = C_n(D_\phi^{\mathsf{H}} \Omega D_\phi)$ is also a definite BSH matrix. Thus, the Lanczos procedure of $D_\phi^{\mathsf{H}} H D_\phi$,

$$(D_\phi^{\mathsf{H}} H D_\phi) \begin{bmatrix} U_k & V_k \\ \overline{U}_k & -\overline{V}_k \end{bmatrix} = \begin{bmatrix} U_k & V_k \\ \overline{U}_k & -\overline{V}_k \end{bmatrix} \begin{bmatrix} 0 & T_k \\ I_k & 0 \end{bmatrix} + \beta_k \begin{bmatrix} u_{k+1} \\ \overline{u}_{k+1} \end{bmatrix} e_{2k}^{\mathsf{H}},$$

is equivalent to

$$H \begin{bmatrix} U_k & V_k \\ e^{i\phi}\overline{U}_k & -e^{i\phi}\overline{V}_k \end{bmatrix} = \begin{bmatrix} U_k & V_k \\ e^{i\phi}\overline{U}_k & -e^{i\phi}\overline{V}_k \end{bmatrix} \begin{bmatrix} 0 & T_k \\ I_k & 0 \end{bmatrix} + \beta_k \begin{bmatrix} u_{k+1} \\ e^{i\phi}\overline{u}_{k+1} \end{bmatrix} e_{2k}^{\mathsf{H}}.$$

## 4.3 Estimation of the absorption spectrum

In the following we describe how to use the Lanczos procedure defined by (23) to estimate the absorption spectrum. It follows from (23) and the orthogonality condition (25) that

$$\epsilon_\sigma(\omega) = d_r^{\mathsf{H}} g_\sigma(\omega I_{2n} - H) d_l$$

$$= \frac{1}{2} \|d_l\|_\Omega^2 \begin{bmatrix} u_1 \\ -\overline{u}_1 \end{bmatrix}^{\mathsf{H}} g_\sigma(\omega I_{2n} - H) \begin{bmatrix} u_1 \\ \overline{u}_1 \end{bmatrix}$$

$$\approx \frac{1}{4} \|d_l\|_\Omega^2 \begin{bmatrix} u_1 \\ -\overline{u}_1 \end{bmatrix}^{\mathsf{H}} \begin{bmatrix} U_k & V_k \\ \overline{U}_k & -\overline{V}_k \end{bmatrix} g_\sigma \left( \omega I_{2n} - \begin{bmatrix} 0 & T_k \\ I_k & 0 \end{bmatrix} \right) \begin{bmatrix} V_k & U_k \\ \overline{V}_k & -\overline{U}_k \end{bmatrix}^{\mathsf{H}} \begin{bmatrix} u_1 \\ \overline{u}_1 \end{bmatrix}. \quad (26)$$

In the proof of Theorem 3, we showed that $U_k^{\mathsf{H}} u_1$ is real. As a result, we obtain that

$$\begin{bmatrix} u_1 \\ -\overline{u}_1 \end{bmatrix}^{\mathsf{H}} \begin{bmatrix} U_k & V_k \\ \overline{U}_k & -\overline{V}_k \end{bmatrix} = 2 \begin{bmatrix} 0 \\ e_1 \end{bmatrix}^{\mathsf{H}}, \qquad \begin{bmatrix} V_k & U_k \\ \overline{V}_k & -\overline{U}_k \end{bmatrix}^{\mathsf{H}} \begin{bmatrix} u_1 \\ \overline{u}_1 \end{bmatrix} = 2 \begin{bmatrix} e_1 \\ 0 \end{bmatrix}.$$

17

These equations allow us to further simplify the expression given in (26). The simplification removes $U_k$ and $V_k$ in the approximation of $\epsilon_\sigma(\omega)$. Hence these vectors do not need to be explicitly stored. Let $T_k = S_k \Theta_k^2 S_k^{\mathsf{H}}$ be the spectral decomposition of $T_k$, where $\Theta_k = \mathrm{diag}\{\theta_1, \ldots, \theta_k\} \succ 0$. By simple calculation, we obtain

$$
f\left(\begin{bmatrix} 0 & T_k \\ I_k & 0 \end{bmatrix}\right) = \begin{bmatrix} S_k & 0 \\ 0 & S_k \end{bmatrix} f\left(\begin{bmatrix} 0 & \Theta_k^2 \\ I_k & 0 \end{bmatrix}\right) \begin{bmatrix} S_k & 0 \\ 0 & S_k \end{bmatrix}^{\mathsf{H}}
$$

$$
= \frac{1}{2} \begin{bmatrix} S_k & 0 \\ 0 & S_k \end{bmatrix} \begin{bmatrix} \Theta_k & -\Theta_k \\ I_k & I_k \end{bmatrix} \begin{bmatrix} f(\Theta_k) & 0 \\ 0 & f(-\Theta_k) \end{bmatrix} \begin{bmatrix} \Theta_k^{-1} & I_k \\ -\Theta_k^{-1} & I_k \end{bmatrix} \begin{bmatrix} S_k & 0 \\ 0 & S_k \end{bmatrix}^{\mathsf{H}}
$$

and

$$
\begin{bmatrix} 0 \\ e_1 \end{bmatrix}^{\mathsf{H}} f\left(\begin{bmatrix} 0 & T_k \\ I_k & 0 \end{bmatrix}\right) \begin{bmatrix} e_1 \\ 0 \end{bmatrix} = \frac{1}{2} e_1^{\mathsf{H}} S_k \left[ f(\Theta_k) - f(-\Theta_k) \right] \Theta_k^{-1} S_k^{\mathsf{H}} e_1 \tag{27}
$$

for any smooth function $f(t)$. Substituting $f(t) = f(t; \omega) = g_\sigma(\omega - t)$, we finally arrive at

$$
\epsilon_\sigma(\omega) \approx \frac{1}{2} \|d_l\|_\Omega^2 e_1^{\mathsf{H}} S_k \left[ g_\sigma(\omega I_k - \Theta_k) - g_\sigma(\omega I_k + \Theta_k) \right] \Theta_k^{-1} S_k^{\mathsf{H}} e_1
$$

$$
= \mathrm{Re}\left( d^{\mathsf{H}} A d + d^{\mathsf{H}} B \bar{d} \right) \sum_{j=1}^{k} |S_k(1,j)|^2 \frac{g_\sigma(\omega - \theta_j) - g_\sigma(\omega + \theta_j)}{\theta_j}. \tag{28}
$$

Again we have the desired property that $\epsilon_\sigma(\omega) \geq 0$ always holds for $\omega > 0$.

The technique of generalized averaged Gauss quadrature can also be adopted here. Notice that (27) can be interpreted as

$$
\begin{bmatrix} 0 \\ e_1 \end{bmatrix}^{\mathsf{H}} f\left(\begin{bmatrix} 0 & T_k \\ I_k & 0 \end{bmatrix}; \omega\right) \begin{bmatrix} e_1 \\ 0 \end{bmatrix} = e_1^{\mathsf{H}} h(T_k; \omega) e_1, \tag{29}
$$

where $h(t; \omega) = t^{-1/2} \left[ f(t^{1/2}; \omega) - f(-t^{1/2}; \omega) \right]$. We expect to obtain a better approximation by replacing $T_k$ in (29) with $\widehat{T}_k$ defined in (13). Certainly, the identity matrix $I_k$ needs to be replaced by $I_{2k-1}$ accordingly. Let $\widehat{T}_k = \widehat{S}_k \widehat{\Theta}_k^2 \widehat{S}_k^{\mathsf{H}}$ be the spectral decomposition of $\widehat{T}_k$, where $\widehat{\Theta}_k = \mathrm{diag}\left\{ \widehat{\theta}_1, \ldots, \widehat{\theta}_{2k-1} \right\}$ has at most one nonpositive eigenvalue. The generalized averaged Gauss quadrature produces

$$
\epsilon_\sigma(\omega) \approx \mathrm{Re}\left( d^{\mathsf{H}} A d + d^{\mathsf{H}} B \bar{d} \right) \sum_{\substack{j=1 \\ \widehat{\theta}_j > 0}}^{2k-1} |\widehat{S}_k(1,j)|^2 \frac{g_\sigma(\omega - \widehat{\theta}_j) - g_\sigma(\omega + \widehat{\theta}_j)}{\widehat{\theta}_j}, \tag{30}
$$

which is expected to be better than (28) in general. Algorithm 5 summarizes the Lanczos algorithm with generalized averaged Gauss quadrature for complex full BSE. All of the four desired features listed in the beginning of this section are satisfied.

18

**Algorithm 5** The Lanczos algorithm for estimating the optical absorption spectrum for complex full BSE.

---

**Input:** A definite Bethe–Salpeter Hamiltonian matrix $H \in \mathbb{C}^{2n \times 2n}$, an optical transition vector $d \in \mathbb{R}^n$, a broadening factor $\sigma > 0$, the number of Lanczos steps $k$, and a set of frequencies $\{\omega_i\}_{i=1}^m$.

**Output:** The estimated absorption spectrum $\epsilon_\sigma(\omega)$ sampled at $\omega_i$ (for $1 \leq i \leq m$).

1: Perform $k$ Lanczos steps in $\Omega$-inner product with starting vector $[d^\mathsf{H}, \overline{d}^\mathsf{H}]^\mathsf{H}$ using Algorithm 4.

2: Formulate $\widehat{T}_k$ as defined in (13).

3: Compute the spectral decomposition $\widehat{T}_k = \widehat{S}_k \operatorname{diag}\left\{ \widehat{\theta}_1^2, \ldots, \widehat{\theta}_{2k-1}^2 \right\} \widehat{S}_k^\mathsf{H}$, where $\widehat{S}_k^\mathsf{H} \widehat{S}_k = I_{2k-1}$ and $\widehat{\theta}_{2k-1} \geq \cdots \geq \widehat{\theta}_2 > 0$.

4: Evaluate

$$
\epsilon_\sigma(\omega_i) = \operatorname{Re}\left( d^\mathsf{H} A d + d^\mathsf{H} B \overline{d} \right) \sum_{\substack{j=1 \\ \widehat{\theta}_j > 0}}^{2k-1} |\widehat{S}_k(1,j)|^2 \frac{g_\sigma(\omega_i - \widehat{\theta}_j) - g_\sigma(\omega_i + \widehat{\theta}_j)}{\widehat{\theta}_j}
$$

for $i = 1, \ldots, m$.

---

## 4.4 Connection with other Lanczos procedures

In this subsection, we establish the connection among several variants of the Lanczos procedures. The comparison includes the Lanczos procedures we have discussed in Sections 4.1 and 4.2, as well as that proposed in [32] and [11]. The connection with a variant of the symplectic Lanczos procedure from [34] is also discussed.

**Lanczos procedures for real BSE** In [32, Section 3], a Lanczos procedure that produces

$$
\hat{U}_k^\mathsf{T} \hat{V}_k = I_k, \qquad K \hat{U}_k = \hat{V} \hat{T}_k, \qquad M \hat{V}_k = \hat{U}_k \hat{D}_k \tag{31}
$$

is studied for real BSE, where $\hat{T}_k$ is symmetric tridiagonal, and $\hat{D}_k \succ 0$ is diagonal. By rescaling $\hat{U}_k$, $\hat{V}_k$ and $\hat{T}_k$ in (31) as

$$
U_k = \hat{U}_k \hat{D}_k^{1/2}, \qquad V_k = \hat{V}_k \hat{D}_k^{1/2}, \qquad T_k = \hat{D}_k^{1/2} \hat{T}_k \hat{D}_k^{1/2},
$$

we obtain

$$
U_k^\mathsf{T} V_k = I_k, \qquad K U_k = V T_k, \qquad M V_k = U_k,
$$

which is identical to the Lanczos procedure (17) in the $M$-inner product. As the rescaling is invertible, (31) and (17) are mathematically equivalent. Since there is no need to keep an additional diagonal matrix, (17) is slightly simpler compared to (31).

If both $H$ and the optical transition vector $d$ are real, the Lanczos procedure (23) simplifies to
$$KMU_k = U_k T_k + \beta_k u_{k+1} e_k^{\mathsf{H}}.$$

The orthogonality condition (24) becomes
$$V_k = MU_k, \qquad V_k^{\mathsf{H}} U_k = I_k, \tag{32}$$

or simply $U_k^{\mathsf{H}} M U_k = I_k$. Thus (23) and (17) are identical for real BSE. In the computation of the absorption spectrum for real BSE, (28) and (30) also reduce to (18) and (19), respectively. Therefore, Algorithm 5 can be regarded as a generalization of Algorithm 3 to complex BSE.

**Lanczos procedures for complex BSE**   In [10], a Lanczos procedure defined in terms of the $\Omega$-inner product, which produces
$$H\tilde{Q}_k = \tilde{Q}_k \tilde{T}_k + \tilde{\beta}_k \tilde{q}_{k+1} e_k^{\mathsf{H}}, \qquad \tilde{Q}_{k+1}^{\mathsf{H}} \Omega \tilde{Q}_{k+1} = I_{k+1}, \tag{33}$$

is presented. However, the projected symmetric tridiagonal matrix $\tilde{T}_k$ does not necessarily have a real spectrum that is symmetric with respect to the origin. Thus (33) is not structured preserving in general. In a subsequent paper [11], it was proposed that a structured starting vector $\tilde{q}_1 \in \mathcal{U}_0$ should be used in (33). With such a structured starting vector, it can be shown that $\tilde{T}_k$ is a real tridiagonal matrix whose diagonal entries are zeros. In addition the nonzero eigenvalues of $\tilde{T}_k$ appear in pairs $\pm\theta$. Hence (33) with $\tilde{q}_1 \in \mathcal{U}_0$ can be regarded as structure preserving. In the following we shall show that this Lanczos procedure is mathematically equivalent to (23).

We have shown that the real symmetric tridiagonal matrix $T_k$ in (20) and (23) is positive definite and componentwise nonnegative. Therefore it admits a Cholesky decomposition $T_k = L_k L_k^{\mathsf{H}}$ where
$$L_k = \text{tridiag} \left\{ \begin{matrix} & 0 & \cdots & & 0 & \\ \tilde{\beta}_1 & \cdots & & \cdots & & \tilde{\beta}_{2k-1} \\ & \tilde{\beta}_2 & \cdots & & \tilde{\beta}_{2k-2} & \end{matrix} \right\}$$

is a bidiagonal lower triangular matrix, which is also componentwise nonnegative. Multiply $\text{diag}\{I_k, L_k\}^{-\mathsf{H}}$ from the right to (23) yields
$$H \begin{bmatrix} U_k & V_k L_k^{-\mathsf{H}} \\ \overline{U}_k & -\overline{V_k L_k^{-\mathsf{H}}} \end{bmatrix} = \begin{bmatrix} U_k & V_k L_k^{-\mathsf{H}} \\ \overline{U}_k & -\overline{V_k L_k^{-\mathsf{H}}} \end{bmatrix} \begin{bmatrix} 0 & L_k \\ L_k^{\mathsf{H}} & 0 \end{bmatrix} + \beta_k \begin{bmatrix} u_{k+1} \\ \overline{u}_{k+1} \end{bmatrix} \begin{bmatrix} 0 \\ L_k^{-1} e_k \end{bmatrix}^{\mathsf{H}}.$$

Notice that $L_k^{-1} e_k$ is parallel to $e_k$. By setting
$$\tilde{U}_k = \frac{1}{\sqrt{2}} U_k, \qquad \tilde{V}_k = \frac{1}{\sqrt{2}} V_k L_k^{-\mathsf{H}}, \qquad \tilde{\beta}_{2k} = \sqrt{2}\, e_k L_k^{-1} e_k, \tag{34}$$

we arrive at a Lanczos procedure of the form

$$H \begin{bmatrix} \tilde{U}_k & \tilde{V}_k \\ \overline{\tilde{U}}_k & -\overline{\tilde{V}}_k \end{bmatrix} = \begin{bmatrix} \tilde{U}_k & \tilde{V}_k \\ \overline{\tilde{U}}_k & -\overline{\tilde{V}}_k \end{bmatrix} \begin{bmatrix} 0 & L_k \\ L_k^{\mathsf{H}} & 0 \end{bmatrix} + \tilde{\beta}_{2k} \begin{bmatrix} \tilde{u}_{k+1} \\ \overline{\tilde{u}}_{k+1} \end{bmatrix} e_{2k}^{\mathsf{H}}. \tag{35}$$

Let

$$\tilde{q}_{2j-1} = \begin{bmatrix} \tilde{u}_{2j-1} \\ \overline{\tilde{u}}_{2j-1} \end{bmatrix}, \qquad \tilde{q}_{2j} = \begin{bmatrix} \tilde{v}_{2j} \\ -\overline{\tilde{v}}_{2j} \end{bmatrix}.$$

Applying the permutation matrix $[e_1, e_{k+1}, e_2, e_{k+2}, \ldots, e_k, e_{2k}]$ from the right to (35) yields

$$H\tilde{Q}_{2k} = \tilde{Q}_{2k}\tilde{T}_{2k} + \tilde{\beta}_k \tilde{q}_{2k+1} e_{2k}^{\mathsf{H}},$$

where

$$\tilde{T}_{2k} = \text{tridiag} \left\{ \begin{matrix} & \tilde{\beta}_1 & \tilde{\beta}_2 & \cdots & & \tilde{\beta}_{2k-2} & \tilde{\beta}_{2k-1} & \\ 0 & & 0 & \cdots & \cdots & & 0 & & 0 \\ & \tilde{\beta}_1 & \tilde{\beta}_2 & \cdots & & \tilde{\beta}_{2k-2} & \tilde{\beta}_{2k-1} & \end{matrix} \right\}.$$

To obtain the orthogonality condition in terms of $\tilde{Q}_k$, we multiply (35) from left by

$$\begin{bmatrix} \tilde{U}_k & \tilde{V}_k \\ \overline{\tilde{U}}_k & -\overline{\tilde{V}}_k \end{bmatrix}^{\mathsf{H}} C_n.$$

Using (25) and simple algebraic manipulation, we obtain

$$\begin{bmatrix} \tilde{U}_k & \tilde{V}_k \\ \overline{\tilde{U}}_k & -\overline{\tilde{V}}_k \end{bmatrix}^{\mathsf{H}} \Omega \begin{bmatrix} \tilde{U}_k & \tilde{V}_k \\ \overline{\tilde{U}}_k & -\overline{\tilde{V}}_k \end{bmatrix} = I_{2k}.$$

Thus we have derived (33) from (23), assuming the number of Lanczos steps in (33) is even. As the transformation (34) is invertible, the two Lanczos procedures are mathematically equivalent.

The Lanczos procedure (33) can be used to approximate the absorption spectrum as follows:

$$\begin{aligned} \epsilon_\sigma(\omega) &= d_r^{\mathsf{H}} g_\sigma(\omega I_{2n} - H) d_l \\ &\approx d_r^{\mathsf{H}} \tilde{Q}_{2k} g_\sigma(\omega I_{2n} - \tilde{T}_{2k}) \tilde{Q}_{2k}^{\mathsf{H}} \Omega d_l \\ &= \frac{1}{2} \|d_l\|_\Omega^2 e_1^{\mathsf{H}} g_\sigma(\omega I_{2n} - \tilde{T}_{2k}) \tilde{T}_{2k}^{-1} e_1. \end{aligned} \tag{36}$$

The derivation of the last step requires similar effort compared to the proof of Theorem 3. The expression (36) is also mathematically equivalent to (28). The main difference between them is that the spectral decomposition of $\tilde{T}_{2k}$ instead of that of $T_k$ is needed. However, we remark that there exist subtle differences when the technique generalized averaged Gauss

21

quadrature is adopted. A direct application of generalized averaged Gauss quadrature replaces $\tilde{T}_{2k}$ by a $(4k-1) \times (4k-1)$ tridiagonal matrix

$$
\widehat{\tilde{T}}_{2k} = \text{tridiag} \left\{
\begin{array}{ccccccccc}
\tilde{\beta}_1 & \cdots & & \tilde{\beta}_{2k-1} & \tilde{\beta}_{2k} & \tilde{\beta}_{2k-2} & \cdots & & \tilde{\beta}_1 \\
0 & \cdots & & \cdots & 0 & 0 & \cdots & & \cdots & 0 \\
\tilde{\beta}_1 & \cdots & & \tilde{\beta}_{2k-1} & \tilde{\beta}_{2k} & \tilde{\beta}_{2k-2} & \cdots & & \tilde{\beta}_1
\end{array}
\right\}.
$$

The positive eigenvalues of $\widehat{\tilde{T}}_{2k}$ are not quite the same as the those of $\widehat{T}_k^{1/2}$, although the number of positive Gauss nodes in the generalized averaged Gauss quadrature is $2k-1$ for both case. We shall see from the numerical experiments the generalized averaged Gauss quadrature based on $\widehat{\tilde{T}}_{2k}$ is in general slightly worse than that based on $\widehat{T}_k$ in terms of accuracy.

We remark that in the discussion above we always assume that an even number of Lanczos steps is performed in (33). In fact, for an odd number of Lanczos steps, $\tilde{T}_{2k+1}$ always has a zero eigenvalue. In the view of Gauss quadrature for estimating the absorption spectrum, such a zero eigenvalue is not a very useful Gauss quadrature node because $\epsilon_\sigma(0) = 0$ is known trivially. Therefore, an even number of Lanczos steps should be performed when computing (33). Similarly, the zero eigenvalue of $\widehat{\tilde{T}}_{2k}$ is not very helpful. Thus we only consider the $2k-1$ positive eigenvalues of $\widehat{\tilde{T}}_{2k}$ to be useful in the generalized averaged Gauss quadrature.

**Connection with symplectic Lanczos procedure** We have shown that our new Lanczos procedure (23) is essentially equivalent to the one proposed in [11], and both are equivalent to (17) and the one in [32] when applying to real BSE. There exists other equivalent formulations. We present these formulations in this section, and exploit more properties of the Lanczos procedure.

Let

$$
\tilde{X}_k = \frac{U_k + V_k}{2}, \qquad \tilde{Y}_k = \frac{\overline{U}_k - \overline{V}_k}{2}, \qquad \tilde{A}_k = \frac{I_k + T_k}{2}, \qquad \tilde{B}_k = \frac{I_k - T_k}{2}.
$$

Then we reformulate (23) as

$$
H \begin{bmatrix} \tilde{X}_k & \overline{\tilde{Y}}_k \\ \tilde{Y}_k & \overline{\tilde{X}}_k \end{bmatrix} = \begin{bmatrix} \tilde{X}_k & \overline{\tilde{Y}}_k \\ \tilde{Y}_k & \overline{\tilde{X}}_k \end{bmatrix} \begin{bmatrix} \tilde{A}_k & \tilde{B}_k \\ -\tilde{B}_k & -\tilde{A}_k \end{bmatrix} + \frac{1}{2}\beta_k \begin{bmatrix} \tilde{x}_{k+1} & \overline{\tilde{y}}_{k+1} \\ \tilde{y}_{k+1} & \overline{\tilde{x}}_{k+1} \end{bmatrix} \begin{bmatrix} 0 & e_k^{\mathsf{H}} \\ 0 & e_k^{\mathsf{H}} \end{bmatrix}. \tag{37}
$$

The orthogonality condition (25) becomes

$$
\left( C_n \begin{bmatrix} \tilde{X}_k & \overline{\tilde{Y}}_k \\ \tilde{Y}_k & \overline{\tilde{X}}_k \end{bmatrix} C_k \right)^{\mathsf{H}} \begin{bmatrix} \tilde{X}_k & \overline{\tilde{Y}}_k \\ \tilde{Y}_k & \overline{\tilde{X}}_k \end{bmatrix} = \begin{bmatrix} \tilde{X}_k & -\overline{\tilde{Y}}_k \\ -\tilde{Y}_k & \overline{\tilde{X}}_k \end{bmatrix}^{\mathsf{H}} \begin{bmatrix} \tilde{X}_k & \overline{\tilde{Y}}_k \\ \tilde{Y}_k & \overline{\tilde{X}}_k \end{bmatrix} = I_{2k}. \tag{38}
$$

Although (23) is derived from (20), which uses the $\Omega$-inner product, the equivalent formulation (37) is a Lanczos procedure in the $C$-inner product. As a result, the projected matrix is a $2k \times 2k$ BSH matrix. As we have discussed in Section 2, the eigenvectors of $H$ are orthogonal in both the $\Omega$-inner product and the $C$-inner product. This suggests that our Lanczos procedure largely preserves properties of $H$. As a byproduct of this observation, we obtain the Cauchy interlacing property as stated in Theorem 4, which provides an estimate on the location of quadrature nodes in the Gauss quadrature. This can be viewed as a generalization of [32, Lemma 3.5]. A proof of Theorem 4 can be found in [28].

**Theorem 4.** *Let $T_k$ be defined as in (20) and suppose that the eigenvalues of $H$ and $T_k$ are $\pm\lambda_1$, $\pm\lambda_2$, ..., $\pm\lambda_n$, and $\theta_1^2$, $\theta_2^2$, ..., $\theta_k^2$, respectively, with $0 < \lambda_1 \le \lambda_2 \le \cdots \le \lambda_n$, $0 < \theta_1 \le \theta_2 \le \cdots \le \theta_k$. Then, under the assumption given in Theorem 2, we have*

$$\lambda_i \le \theta_i \le \lambda_{n-k+i}, \qquad (1 \le i \le k).$$

It has been shown in [27] that the matrix

$$i Q_n^{\mathsf{H}} H Q_n = J_n \begin{bmatrix} \mathrm{Re}(A+B) & \mathrm{Im}(A-B) \\ -\mathrm{Im}(A+B) & \mathrm{Re}(A-B) \end{bmatrix} =: J_n \tilde{M}$$

is a real Hamiltonian matrix with $\tilde{M} \succ 0$, where

$$J_n = \begin{bmatrix} 0 & I_n \\ -I_n & 0 \end{bmatrix}, \qquad Q_n = \frac{1}{\sqrt{2}} \begin{bmatrix} I_n & -iI_n \\ I_n & iI_n \end{bmatrix}.$$

Using the same unitary transformation, the Lanczos factorization (37) becomes

$$J_n \tilde{M} \begin{bmatrix} \mathrm{Re}(U_k) & \mathrm{Im}(V_k) \\ -\mathrm{Im}(U_k) & \mathrm{Re}(V_k) \end{bmatrix} = \begin{bmatrix} \mathrm{Re}(U_k) & \mathrm{Im}(V_k) \\ -\mathrm{Im}(U_k) & \mathrm{Re}(V_k) \end{bmatrix} \begin{bmatrix} 0 & T_k \\ -I_k & 0 \end{bmatrix} + [\text{rank } 1], \qquad (39)$$

and the orthogonality condition (38) becomes

$$\begin{bmatrix} \mathrm{Re}(U_k) & \mathrm{Im}(V_k) \\ -\mathrm{Im}(U_k) & \mathrm{Re}(V_k) \end{bmatrix}^{\mathsf{H}} J_n \begin{bmatrix} \mathrm{Re}(U_k) & \mathrm{Im}(V_k) \\ -\mathrm{Im}(U_k) & \mathrm{Re}(V_k) \end{bmatrix} = J_k. \qquad (40)$$

Since (40) indicates that the Lanczos vectors associated with $J_n \tilde{M}$ are symplectic, the Lanczos factorization (39) yields in fact a symplectic Lanczos procedure (see, e.g., [2, 4, 34]) for the real Hamiltonian matrix $J_n \tilde{M}$. Such a variant of symplectic Lanczos procedures has been discussed in [34]. Therefore, (23) can also be interpreted as a variant of symplectic Lanczos procedure. Finally we remark that for the purpose of computing the absorption spectrum, the starting vector in (39) should be chosen parallel to $\left[ \mathrm{Re}(d)^{\mathsf{H}}, -\mathrm{Im}(d)^{\mathsf{H}} \right]^{\mathsf{H}}$ if (39) is adopted.

23

## 4.5   Structure preserving Lanczos algorithm with paired starting vectors

Besides several equivalent structure preserving Lanczos procedures, there are also other structure preserving Lanczos procedures with paired starting vectors. Actually, when $\|u_1\|_2 \neq \|v_1\|_2$, Lanczos procedures of the form

$$H \begin{bmatrix} U_k & \overline{V}_k \\ V_k & \overline{U}_k \end{bmatrix} = \begin{bmatrix} U_k & \overline{V}_k \\ V_k & \overline{U}_k \end{bmatrix} \begin{bmatrix} A_k & B_k \\ -\overline{B}_k & -\overline{A}_k \end{bmatrix} + \begin{bmatrix} u_{k+1} & \overline{v}_{k+1} \\ v_{k+1} & \overline{u}_{k+1} \end{bmatrix} \begin{bmatrix} \beta_k e_k^{\mathsf{H}} & 0 \\ 0 & -\overline{\beta}_k e_k^{\mathsf{H}} \end{bmatrix} \tag{41}$$

can be constructed, where $A_k$ and $B_k$ are tridiagonal, and the orthogonality condition on the Lanczos vectors is either[2]

$$\begin{bmatrix} u_i & \overline{v}_i \\ v_i & \overline{u}_i \end{bmatrix}^{\mathsf{H}} \Omega \begin{bmatrix} u_j & \overline{v}_j \\ v_j & \overline{u}_j \end{bmatrix} = \delta_{ij} I_2 \tag{42}$$

or

$$C_2 \begin{bmatrix} u_i & \overline{v}_i \\ v_i & \overline{u}_i \end{bmatrix}^{\mathsf{H}} C_n \begin{bmatrix} u_j & \overline{v}_j \\ v_j & \overline{u}_j \end{bmatrix} = \begin{bmatrix} u_i & -\overline{v}_i \\ -v_i & \overline{u}_i \end{bmatrix}^{\mathsf{H}} \begin{bmatrix} u_j & \overline{v}_j \\ v_j & \overline{u}_j \end{bmatrix} = \delta_{ij} I_2. \tag{43}$$

As the eigenvalues of the projected matrix

$$H_k = \begin{bmatrix} A_k & B_k \\ -\overline{B}_k & -\overline{A}_k \end{bmatrix}$$

occur in pairs $\pm\theta$, (41) is regarded as structure preserving. In fact, from the discussion in the previous subsection, we also see that the condition (43) for BSH matrices is equivalent to the symplecticity condition for real Hamiltonian matrices.

When estimating the absorption spectrum using (41), we use $u_1 = d$, $v_1 = 0$ as the starting vectors because $d_l$ does not satisfy the condition $\|u_1\|_2 \neq \|v_1\|_2$. If (42) is used, the computed absorption spectrum is not guaranteed to be real. Hence, we do not consider Lanczos procedure (41) in the $\Omega$-inner product for computing the absorption spectrum. If the orthogonality condition (43) is adopted, the projected matrix $H_k$ is a definite BSH matrix. Let the spectral decomposition of $H_k$ be

$$H_k = \begin{bmatrix} S_1 & \overline{S}_2 \\ S_2 & \overline{S}_1 \end{bmatrix} \begin{bmatrix} \Theta & 0 \\ 0 & -\Theta \end{bmatrix} \begin{bmatrix} S_1 & -\overline{S}_2 \\ -S_2 & \overline{S}_1 \end{bmatrix}^{\mathsf{H}},$$

where $\Theta = \operatorname{diag}\{\theta_1, \ldots, \theta_k\} \succ 0$. It can be shown that

$$\epsilon_\sigma(\omega) \approx \|d\|_2^2 (e_1 - e_{k+1})^{\mathsf{H}} g_\sigma(\omega I - H_k)(e_1 + e_{k+1})$$

$$= \|d\|_2^2 \sum_{j=1}^{k} |S_1(1,j) - S_2(1,j)|^2 \big[ g_\sigma(\omega - \theta_j) - g_\sigma(\omega + \theta_j) \big]. \tag{44}$$

---

[2] If (42) is used, orthogonalization within each two dimensional subspace is required.

This formulation possesses the second and third features listed in the beginning of this section. However, theoretically Lanczos procedure in the $C$-inner product may sometimes break down due to $C$-neutral vectors.[3] Such a breakdown is not a lucky breakdown. It is also not very clear how to incorporate the technique of generalized averaged Gauss quadrature in (44).

We remark that in general (44) is not as good as (28) even if generalized averaged Gauss quadrature is not used. Our numerical experiments suggest that (44) typically requires about twice as many as Lanczos steps to achieve the same accuracy level compared to (28).[4] A brief explanation is that for the same number of Lanczos steps $k$, $H$ has been raised to the power $H^{2k}$ in (23), while $H$ has only been raised to the power $H^k$ in (41). A higher polynomial degree potentially provides better approximation quality.

# 5 Computational examples

In this section we present several examples to demonstrate the accuracy and efficiency of the Lanczos algorithm for computing the optical absorption spectrum. We implemented the Lanczos algorithms in the BerkeleyGW [7] software package. All tests were performed on the Linux cluster Edison at the National Energy Research Scientific Computing Center (NERSC).[5] Each computational node on Edison consists of 64 GB DDR3 1866 MHz memory and two sockets, with a 12-core Intel "Ivy Bridge" processor at 2.4 GHz on each socket. The computational nodes are connected by a Cray Aries network with Dragonfly topology, with 23.7 TB/s global bandwidth. Our tests make use of 10 computational nodes and 24 MPI processes per node. The Fortran 90 implementation of algorithms is compiled by the Intel Fortran compiler, and linked with the Cray LibSci and Cray MPI libraries. No multithreading feature is utilized.

For our calculations, we use a benchmark system consisting of a single-wall $(8, 0)$ carbon nanotube with 32 atoms, 128 electrons, and 64 Kohn–Sham spin-degenerate bands in the unit cell. As depicted in Figure 2, this system is periodic along the "c" axis, but confined along the other directions labeled by the axes "a" and "b", which makes this an interesting benchmark system. In particular, as we will discuss, the TDA may or may not be a good approximation depending on the direction of optical excitation in this particular system.

In general, crystal states can be written in a Bloch form as $\Psi_{nk}(r) = e^{ik \cdot r} u_{nk}(r)$, where $n$ is a band index, $k$ is a $k$-point, and $u_{nk}(r)$ is a cell-periodic complex-valued function. Because "c" is the only periodic direction, we only need to sample $k$-points along that axis. When solving the BSE, we include $n_v = 10$ valence states, $n_c = 12$ conduction states, and $n_k = 256$ $k$-points, so that $n = n_v n_c n_k = 30,720$, and we picked $g_\sigma$ as a Gaussian function with $\sigma = 100$ meV. The matrices $A$ and $B$ are both dense. We did not perform a

---

[3]A $C$-neutral vector is a vector $v \in \mathbb{C}^{2n}$ which satisfies $v^{\mathsf{H}} C_n v = 0$.

[4]A similar behavior has been observed in [33] when solving the linear response eigenvalue problem.

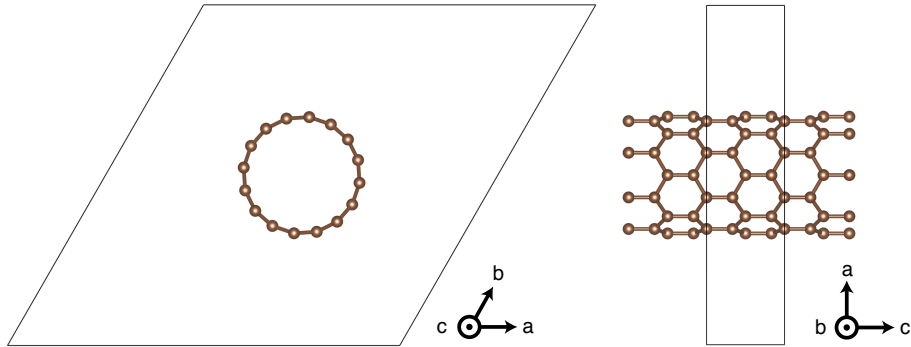[5]http://www.nersc.gov/users/computational-systems/edison/

Figure 2: Single-wall (8,0) carbon nanotube benchmark system. The black region represents the unit cell of the system, which is periodic along the "c" axis.

systematic convergence test with respect to the number of conducting bands. However, the use of $n_c = 12$ conducting bands already produces main physical features in the absorption spectrum also observed when a larger number of conduction bands are used.

**Full BSE vs. TDA.** In our first experiment, we calculate the absorption spectrum for two different directions for the light polarizations using both full BSE and TDA solvers. To exclude other sources of errors, we fully diagonalize the matrices $H$ and $A$, with dimensions 61,440 and 30,720, respectively. We can see from Figure 3 that even within the same system, the TDA can either be a valid approximation or give a qualitatively wrong absorption spectrum depending on the polarization direction of light. When the polarization of the optical excitation is along the "c" axis, which is a direction along which the system is periodic, the TDA is a good approximation for low-energy optical spectrum. However, if the light polarization is along any confined direction spanned by the "a" and "b" axes, a large difference between the two spectra can be observed. This can be understood from a large exciton–plasmon hybridization which couples to light polarized along the confined direction, and which can not be well-described within the TDA [10].

Thus this example confirms the necessity of developing full BSE solvers for absorption spectrum calculation. In the subsequent tests the light polarizations is chosen to be perpendicular to the tube so that using a full BSE solver is necessary.

**Effectiveness of the Lanczos algorithm.** In Figure 4 we plot the approximate absorption spectra obtained by running 32 steps of different variants of the Lanczos algorithm. We use the result obtained from full diagonalization as the "exact" solution to measure the accuracy of these Lanczos algorithms. The paired Lanczos algorithm described in Section 4.5 (Figure 4(a), abbreviated as PL) is clearly worse than the one with a single structured starting vector (Figure 4(b), abbreviated as SVL). The Lanczos algorithm pro-
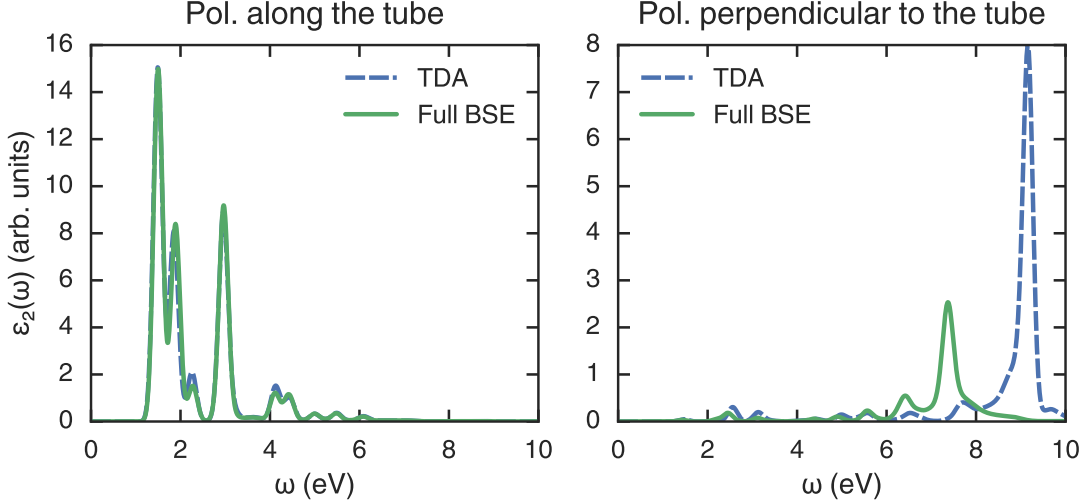
26

Figure 3: The absorption spectrum for a single-wall $(8,0)$ carbon nanotube for two different directions for the light polarizations. The label "Full BSE" refers to spectrum obtained from solving the full BSH in (7), and "TDA" refers to the spectrum obtained within the Tamm–Dancoff approximation.

posed in [11] (abbreviated as GMG) is equivalent to ours with a single structured starting vector and is hence omitted here. The technique of generalized averaged Gauss quadrature (abbreviated as GAGQ) clearly improves the accuracy of the Lanczos algorithm. With such a small number of Lanczos steps, our Lanczos algorithm with generalized averaged Gauss quadrature (i.e., Algorithm 5) already produces very satisfactory result. Though not very clear from this figure, our Algorithm 5 (Figure 4(d)) is slightly better than the one proposed in [11] with generalized averaged Gauss quadrature (Figure 4(c)) in this example.

To measure the accuracy of approximate absorption spectrum, we introduce the concept of *angle* between two functions as follows. Let $\xi(\omega)$ and $\zeta(\omega)$ be sufficiently smooth functions of $\omega$ over an interval $I$. Then the angle between $\xi(\omega)$ and $\zeta(\omega)$ is defined as

$$\angle(\xi, \eta) = \arccos \frac{\langle \xi, \zeta \rangle}{\sqrt{\langle \xi, \xi \rangle \langle \zeta, \zeta \rangle}}, \tag{45}$$

where

$$\langle \xi, \zeta \rangle = \int_I \xi(\omega) \overline{\zeta(\omega)} \, d\omega$$

is the usual $L^2$-inner product. The angle $\angle(\xi, \eta)$ is in fact the principal angle (also known as canonical angle) between two subspaces, $\mathrm{span}\,\{\xi(\omega)\}$ and $\mathrm{span}\,\{\eta(\omega)\}$, of $L^2(I)$. A small angle between two functions implies similar shapes of their curves. This allows us to measure the error of the approximate absorption spectrum compared to the "exact" one in
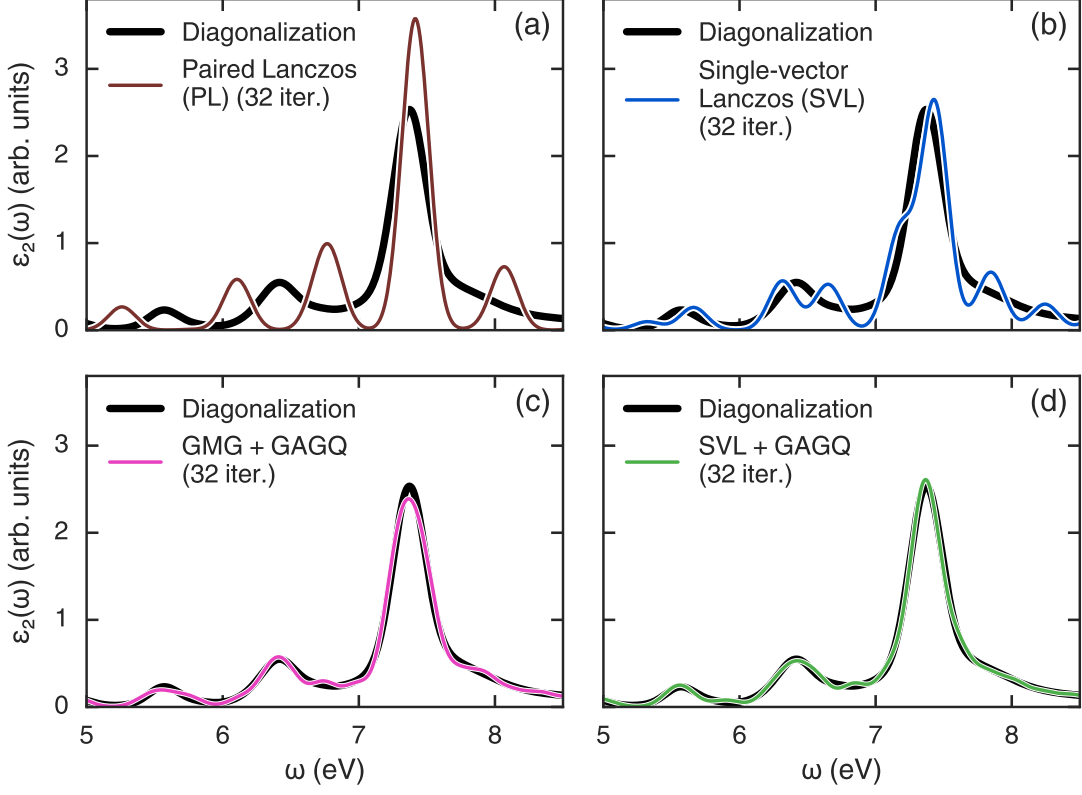
27

Figure 4: Comparison of the absorption spectra obtained from different variants of the Lanczos algorithm with the spectrum obtained from the full diagonalization of the BSH matrix.

terms of the angle between them. This measure is similar to the cross-correlation measure between two curves. In Figure 5 we plot the errors of different variants of the Lanczos algorithm, with the integrals in (45) approximated by rectangular rules using the sampling points of $\epsilon_2(\omega)$'s. It confirms our observation from Figure 4, not only for a single snapshot after 32 Lanczos steps, but also consistently throughout the whole iterative procedure. The difference between the two different variants with generalized averaged Gauss quadrature becomes more clear in Figure 5. Overall Algorithm 5 is better than the variant from [11] combined with generalized averaged Gauss quadrature. We remark that there are about 10% cases in this example involving nonpositive definite $\widehat{T}_k$ in Algorithm 5. Figure 5 shows that dropping the nonpositive eigenvalue of $\widehat{T}_k$ does not harm the accuracy.

It takes 62 iterations and 4.1 seconds for Algorithm 5 to achieve the accuracy level $10^{-3}$ (in terms of angles), which is more than sufficient for practical use. This is over 500 times faster compared to full diagonalization (2125.8 seconds). If the multiplications of $A$ and $B$
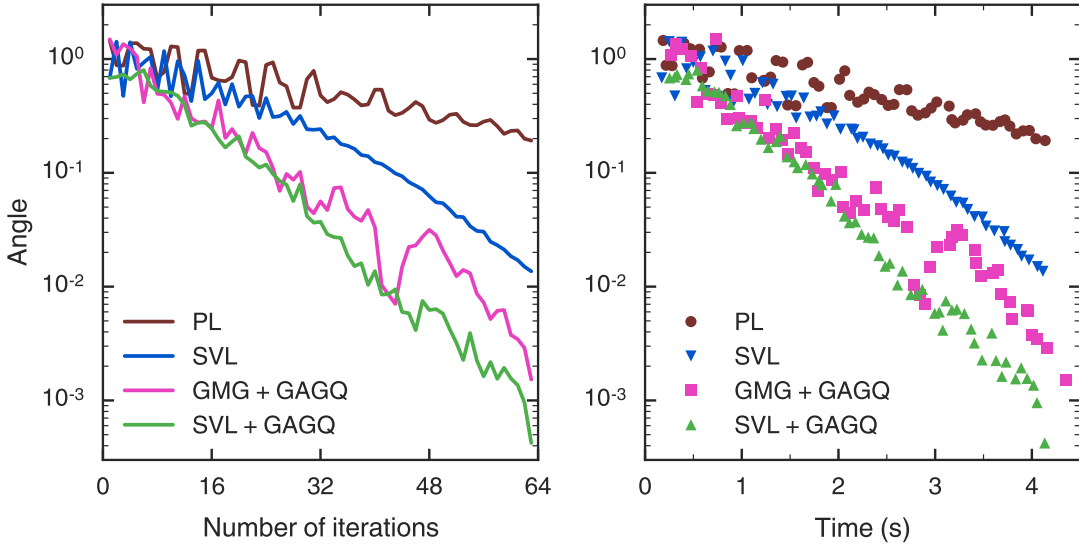
Figure 5: The convergence history of different variants of the Lanczos algorithm. The error is measured by the angle (45) between the approximate absorption spectrum and the one obtained from full diagonalization. The algorithms used here are the same as those in Figure 4.

with vectors can be implemented more efficient by further exploiting the structures of $A$ and $B$, the improvement is expected to be more significant.

In our test, the number of Lanczos steps is always prescribed by the user. We remark that it is possible to instead specify the desired accuracy in the input, and automatically determine the required number of Lanczos steps in the calculation. One strategy proposed in [19] is to estimate the error using the difference between the results obtained with and without generalized averaged Gauss quadrature. However, since this strategy relies on the result without generalized averaged Gauss quadrature, which is in a relatively low accuracy as we have shown in Figure 5, the estimate is in general too pessimistic. A better strategy is to use the difference between two consecutive iterations (i.e., $(k-1)$th and $k$th steps) instead in the stopping criterion.

**Systems with real-space inversion symmetry.** Our last example uses another system which has real-space inversion symmetry. When a system has real-space and time-reversal symmetry, the wave functions in reciprocal space, and thus the BSH matrix, can be written as real numbers [7]. We use bulk silicon for this benchmark, with $n_v = 4$, $n_c = 6$, and $n_k = 1,000$, so that the dimension of the $A$ and $B$ blocks of BSH is $n = 24,000$. We also use a Gaussian broadening in this system, but with $\sigma = 150$ meV. Since the BSH matrix is real, both Algorithms 3 and 5 are applicable, and are identical as discussed in Section 4.
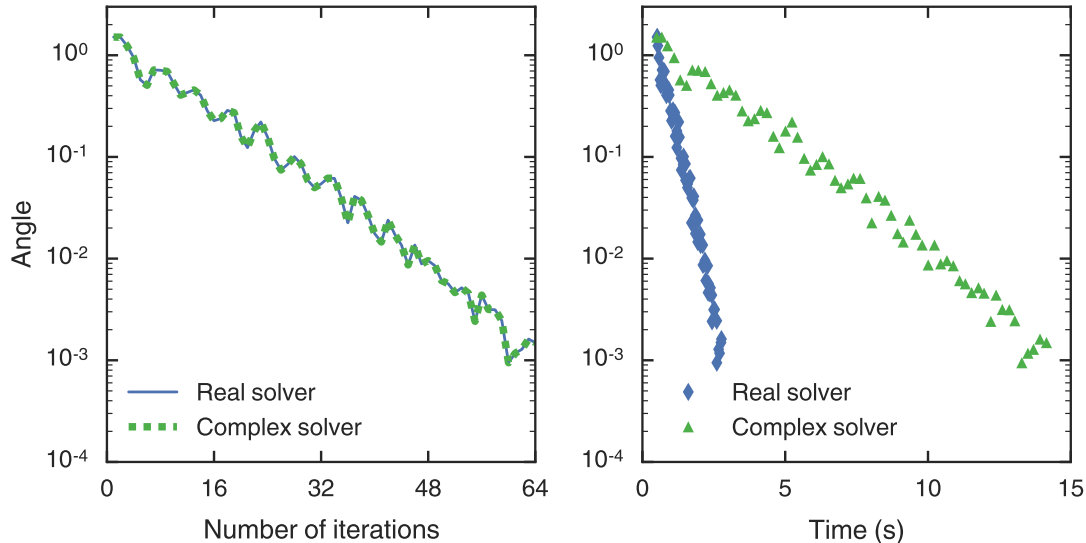
29

Figure 6: The convergence history of Algorithms 3 and 5 for bulk silicon. The error is measured by the angle (45) between the approximate absorption spectrum and the one obtained from full diagonalization.

Our experiment confirms the theoretical prediction. In Figure 6 we plot the convergence history (in terms of angles) of the two algorithms. The curves of convergence history are indeed on the top of each other in the left plot. As for the execution time, the real solver is faster than the complex one, due to some additional operations involving the imaginary parts in the complex solver when applying to real matrices.

# 6 Concluding remarks

In this paper we developed a simple structure preserving Lanczos procedure for definite Bethe–Salpeter Hamiltonian matrices and combined it with the recently developed technique of generalized averaged Gauss quadrature to estimate the optical absorption spectrum. Our Lanczos procedure possesses several attractive features, such as no serious breakdown, and preserving nonnegativity of the absorption spectrum. The use of alternative inner products based on the orthogonalities of the eigenvectors plays a key role in preserving the structure. By some theoretical analysis we established the equivalence between our Lanczos procedure with several existing Lanczos procedures in the literature, including the ones in [11, 32] for random phase approximation, and one variant of symplectic Lanczos procedure in [34]. Numerical experiments demonstrate that the Lanczos algorithm can provide accurate approximation of the absorption spectrum with a relatively small number of Lanczos steps. In addition, the technique of generalized averaged Gauss

30

quadrature largely improves the accuracy of the Lanczos algorithm. When this technique is applied, our Lanczos algorithm is more efficient and more accurate compared to other variants.

In this work the blocks $A$ and $B$ in the BSH matrix $H$ are formed as dense matrices. However, the Lanczos algorithm does not require these matrices to be explicitly formed. An implicit representation that allows to perform matrix–vector multiplication suffices. Efficient ways of constructing and applying the BSH matrix have been described in [16, 17, 22, 23, 24]. The development of other approximation and compression strategies is planned as future work.

## Acknowledgments

## References

[1] Z. Bai and R.-C. Li. Minimization principles for the linear response eigenvalue problem II: Computation. *SIAM J. Matrix Anal. Appl.*, 34(2):392–416, 2013.

[2] P. Benner and H. Fassbender. An implicitly restarted symplectic Lanczos method for the Hamiltonian eigenvalue problem. *Linear Algebra Appl.*, 263:75–111, 1997.

[3] P. Benner, H. Fassbender, and C. Yang. Some remarks on the complex $J$-symmetric eigenproblem. Preprint MPIMD/15-12, Max Planck Institute Magdeburg, 2015. Available from http://www.mpi-magdeburg.mpg.de/preprints/.

[4] P. Benner and A. Salam. The symplectic Lanczos process for Hamiltonian-positive matrices, 2011 (in preparation).

[5] J. Brabec, L. Lin, M. Shao, N. Govind, Y. Saad, C. Yang, and E. G. Ng. Efficient algorithms for estimating the absorption spectrum within linear response TDDFT. *J. Chem. Theory Comput.*, 11(11):5197–5208, 2015.

[6] S. M. Dancoff. Non-adiabatic meson theory of nuclear forces. *Phys. Rev.*, 78(4):382–385, 1950.

[7] J. Deslippe, G. Samsonidze, D. A. Strubbe, M. Jain, M. L. Cohen, and S. G. Louie. BerkeleyGW: A massively parallel computer package for the calculation of the quasiparticle and optical properties of materials and nanostructures. *Comput. Phys. Commun.*, 183(6):1269–1289, 2012.

[8] G. H. Golub and G. Meurant. Matrices, moments and quadrature. In D. F. Griffith and G. A. Watson, editors, *Numerical Analysis 1993*, volume 303 of *Pitman Research Notes in Mathematics*, pages 105–156, Essex, UK, 1994. Longman.

[9] G. H. Golub and G. Meurant. *Matrices, Moments and Quadrature with Applications*. Princeton University Press, Princeton, NJ, USA, 2010.

[10] M. Grüning, A. Marini, and X. Gonze. Exciton–plasmon states in nanoscale materials: Breakdown of the Tamm–Dancoff approximation. *Nano Lett.*, 9(8):2820–2824, 2009.

[11] M. Grüning, A. Marini, and X. Gonze. Implementation and testing of Lanczos-based algorithms for random-phase approximation eigenproblems. *Comput. Mater. Sci.*, 50:2148–2156, 2011.

[12] R. Haydock. The recursive solution of the Schrödinger equation. *Comput. Phys. Commun.*, 20(1):11–16, 1980.

[13] R. Haydock. The recursive solution of the Schrödinger equation. *Solid State Phys.*, 35:215–294, 1980.

[14] R. Haydock, V. Heine, and M. J. Kelly. Electronic structure based on the local atomic environment for tight-binding bands. *J. Phys. C: Solid State Phys.*, 5:2845–2858, 1972.

[15] R. Haydock, V. Heine, and M. J. Kelly. Electronic structure based on the local atomic environment for tight-binding bands: II. *J. Phys. C: Solid State Phys.*, 8:2591–2605, 1975.

[16] Y. Ping, D. Rocca, and G. Galli. Electronic excitations in light absorbers for photo-electrochemical energy conversion: first principles calculations based on many body perturbation theory. *Chem. Soc. Rev.*, 42:2437–2469, 2013.

[17] Y. Ping, D. Rocca, D. Lu, and G. Galli. *Ab initio* calculations of absorption spectra of semiconducting nanowires within many-body perturbation theory. *Phys. Rev. B*, 85:035316, 2012.

[18] D. Y. Qiu, F. H. da Jornada, and S. G. Louie. Optical spectrum of mos2: Many-body effects and diversity of exciton states. *Phys. Rev. Lett.*, 111(21):216805, 2013.

[19] L. Reichel, M. M. Spalević, and T. Tang. Generalized averaged Gauss quadrature rules for the approximation of matrix functionals. *BIT Numer. Math.*, 56(3):1045–1067, 2016.

[20] D. Rocca, Z. Bai, R.-C. Li, and G. Galli. A block variational procedure for the iterative diagonalization of non-Hermitian random-phase approximation matrices. *J. Chem. Phys*, 136(3):034111, 2012.

[21] D. Rocca, R. Gebauer, Y. Saad, and S. Baroni. Turbo charging time-dependent density-functional theory with Lanczos chains. *J. Chem. Phys*, 128(15):154105, 2008.

[22] D. Rocca, D. Lu, and G. Galli. Ab initio calculations of optical absorption spectra: Solution of the Bethe–Salpeter equation within density matrix perturbation theory. *J. Chem. Phys*, 133(16):164109, 2010.

[23] D. Rocca, Y. Ping, R. Gebauer, and G. Galli. Solution of the Bethe–Salpeter equation without empty electronic states: Application to the absorption spectra of bulk systems. *Phys. Rev. B*, 85:045116, 2012.

[24] D. Rocca, M. Vörös, A. Gali, and G. Galli. Ab initio optoelectronic properties of silicon nanoparticles: Excitation energies, sum rules, and Tamm–Dancoff approximation. *J. Chem. Theory Comput.*, 10(8):3290–3298, 2014.

[25] M. Rohlfing and S. G. Louie. Electron–hole excitations and optical spectra from first principles. *Phys. Rev. B*, 62(8):4927–4944, 2000.

[26] E. E. Salpeter and H. A. Bethe. A relativistic equation for bounded-state problems. *Phys. Rev.*, 84(6):1232–1242, 1951.

[27] M. Shao, F. H. da Jornada, C. Yang, J. Deslippe, and S. G. Louie. Structure preserving parallel algorithms for solving the Bethe–Salpeter eigenvalue problem. *Linear Algebra Appl.*, 488:148–167, 2016.

[28] M. Shao and C. Yang. Properties of definite Bethe–Salpeter eigenvalue problems. In *Proceedings of EPASA 2015 and EPASA 2014*, Lecture Notes in Computational Science and Engineering, Berlin, Germany, 2016 (to appear). Springer-Verlag.

[29] M. M. Spalević. On generalized averaged Gaussian formulas. *Math. Comp.*, 76(259):1483–1492, 2007.

[30] G. Strinati. Application of the Green's functions method to the study of the optical properties of semiconductors. *La Rivista del Nuovo Cimento*, 11(12):1–86, 1988.

[31] I. Y. Tamm. Relativistic interaction of elementary particles. *J. Phys. (USSR)*, 9:449–460, 1945.

[32] Z. Teng and R.-C. Li. Convergence analysis of Lanczos-type methods for the linear response eigenvalue problem. *J. Comput. Appl. Math.*, 247:17–33, 2013.

[33] E. Vecharynski, J. Brabec, M. Shao, N. Govind, and C. Yang. Efficient block preconditioned eigensolvers for linear response time-dependent density functional theory. 2015 (submitted).

[34] D. S. Watkins. On Hamiltonian and symplectic Lanczos processes. *Linear Algebra Appl.*, 385:23–45, 2004.

[35] R. Zimmermann. Influence of the non-Hermitian splitting terms on exciontic spectra. *Phys. Stat. Sol.*, 41:23–43, 1970.