

Error control in scientific modelling (MATH 500, Herbst)

Project 2: Band structures with guaranteed error bars

To be handed in via moodle by 12.01.2024

```
1 begin
2   using Brillouin
3   using DFTK
4   using LinearAlgebra
5   using Plots
6   using PlutoTeachingTools
7   using PlutoUI
8   using Printf
9   using LaTeXStrings
10  using Statistics
11  using Measurements
12  using DoubleFloats
13
14 end
```

Code paths for generic floating-point types activated in DFTK. Remember to add 'using GenericLinearAlgebra' to your user script. See https://docs.dftk.org/stable/examples/arbitrary_floattype/ for details.

Selection deleted

Table of Contents

Project 2: Band structures with guaranteed error bars

Introduction

Task 1: Cohen-Bergstresser Hamiltonians

Familiarisation and Bauer-Fike estimates

Task 2: Γ -point calculations

Band structure computations

Task 3: Bands with Bauer-Fike estimates

Kato-Temple estimates

Task 4: Simple and non-guaranteed Kato-Temple estimates

Task 5: A Schur-based estimate for the gap

Task 6: Band structure with guaranteed bounds

Adaptive diagonalisation techniques

Task 7: Developing a discretisation-adaptive LOBPCG

GTH pseudopotentials

Task 8: An error indicator for GTH pseudopotentials

Selection deleted

Introduction

In this project we will compute approximations of parts of the eigenspectra of a specific class of Hamiltonians

$$H = -\frac{1}{2}\Delta + V.$$

Precisely we will consider simple potentials V , which are periodic with respect to some lattice \mathbb{L} . For this setting we will further develop mathematically guaranteed estimates for the total numerical error of our eigenvalue computations, such that the exact answer is provably within our error bound. Following the lecture we will use Bloch-Floquet theory to reformulate this problem in terms of Bloch fibers

$$H_k = \frac{1}{2}(-i\nabla_x + k)^2$$

where the k -points are taken from the Brillouin zone Ω^* . Of main interest and the focus of our work is the so-called band structure, that is the dependency of the few lowest eigenvalues of the H_k on k . At the end of this project you will have coded up a routine to compute the band structure of simple periodic potentials V in a way that the error is fully tracked and mathematically guaranteed error bars can be annotated to the computed band structure.

For our numerical computation in this project we will mostly employ the density-functional toolkit (DFTK), which has already been installed alongside this notebook. You can find the source code on dftk.org and documentation as well as plenty of usage examples on docs.dftk.org. In particular an alternative introduction to periodic problems, which focuses more on the numerics than the mathematics is given on https://docs.dftk.org/stable/guide/periodic_problems/.

Selection deleted

As example systems in this work we will only consider face-centred cubic crystals in the diamond structure. As the name suggests this family of crystals includes diamond, but also materials such as bulk silicon, germanium or tin. For these the lattice

$$\mathbb{L} = \mathbb{Z}\mathbf{a}_1 + \mathbb{Z}\mathbf{a}_2 + \mathbb{Z}\mathbf{a}_3$$

can be constructed from the unit cell vectors

$$\mathbf{a}_1 = \frac{\mathbf{a}}{2} \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix} \quad \mathbf{a}_2 = \frac{\mathbf{a}}{2} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \quad \mathbf{a}_3 = \frac{\mathbf{a}}{2} \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}$$

where \mathbf{a} is the lattice constant, which differs between the materials. Keeping the notation from the lecture we will use Ω to refer to the unit cell

$$\Omega = \{ \mathbf{x} \in \mathbb{R}^3 \mid |\mathbf{x} - \mathbf{R}| > |\mathbf{x}| \quad \forall \mathbf{R} \in \mathbb{L} \setminus \{0\} \},$$

and employ \mathbb{L}^* to refer to the reciprocal lattice with vectors $\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3$ with $\mathbf{a}_i \cdot \mathbf{b}_j = 2\pi\delta_{ij}$. Ω^* denotes the first Brillouin zone (unit cell of \mathbb{L}^*).

We introduce the normalised plane waves

$$e_G(\mathbf{x}) = \frac{e^{iG \cdot \mathbf{x}}}{\sqrt{|\Omega|}} \quad \forall G \in \mathbb{L}^*,$$

which allow to express the Fourier series of the potential as

$$V(\mathbf{x}) = \frac{1}{\sqrt{|\Omega|}} \sum_{G \in \mathbb{L}^*} \hat{V}(G) e^{iG \cdot \mathbf{x}} = \sum_{G \in \mathbb{L}^*} \hat{V}(G) e_G(\mathbf{x}).$$

To simplify our notation we will assume that G -vectors always run over the countably infinite reciprocal lattice \mathbb{L}^* without making explicit reference to this set.

An early model to describe these materials theoretically is the Cohen-Bergstresser model [[CB1966](#)]. The special and simplifying property of this model is that only a relatively small number of Fourier coefficients is non-zero. In particular we can define a cutoff \mathcal{E}_V such that

Selection deleted

$$V(\mathbf{x}) = \sum_{|G| < \sqrt{2\mathcal{E}_V}} \hat{V}(G) e_G(\mathbf{x}),$$

where each $\hat{V}(G)$ is finite.

[[CB1966](#)]:

M. L. Cohen and T. K. Bergstresser Phys. Rev. **141**, 789 (1966) DOI [10.1103/PhysRev.141.789](https://doi.org/10.1103/PhysRev.141.789)

Task 1: Cohen-Bergstresser Hamiltonians

- (a) Using the results of the lecture show that H is self-adjoint on $L^2(\mathbb{R}^3)$ with the domain $D(H) = H^2(\mathbb{R}^3)$ if V is the Cohen-Bergstresser potential. *Hint:* Show and use that the Cohen-Bergstresser potential is bounded and \mathbb{L} -periodic.

Selection deleted

(a) Solution:

A function is called bounded if there exists a constant C such that:

$$|V(x)| \leq C \quad \forall x \in \mathbb{R}^3.$$

Applying triangular inequality and using the fact that $|e^{iG \cdot x}| = 1$ as a modulus of a complex number in a polar form we have:

$$\begin{aligned} |V(x)| &\leq \frac{1}{\sqrt{|\Omega|}} \sum_{|G|<\sqrt{2\mathcal{E}_V}} |\hat{V}(G)| |e^{iG \cdot x}| = \\ &= \sum_{|G|<\sqrt{2\mathcal{E}_V}} \frac{1}{\sqrt{|\Omega|}} |\hat{V}(G)| = C \quad \forall x \in \mathbb{R}^3, \end{aligned}$$

where $|\Omega|$ is a constant since it is the volume of the parallelepiped spanned by three vectors that form a unit cell. Only a small number of Fourier coefficients is non-zero because there is a finite number of vectors G which is from the reciprocal lattice \mathbb{L}^* and also inside the circle with the radius $\sqrt{2\mathcal{E}_V}$. Since $\hat{V}(G)$ is given to be finite we can conclude that C is indeed a constant.

Therefore, the Cohen-Bergstresser potential is bounded.

On the other hand, from the periodicity of the complex exponential $e^{iG \cdot x}$ on the lattice \mathbb{L} follows the \mathbb{L} -periodicity of $e_G(x)$. Since $V(x)$ is a finite sum of \mathbb{L} -periodic functions, it follows that $V(x)$ is itself \mathbb{L} -periodic.

Moreover, function $V(x) \in L^2(\Omega)$ which means that $\int_{\Omega} |V(x)|^2 dx$ is finite. Using Hölder's inequality, we have:

$$\begin{aligned} \int_{\Omega} |V(x)|^{3/2} dx &\leq \left(\int_{\Omega} (|V(x)|^{3/2})^{4/3} dx \right)^{3/4} \left(\int_{\Omega} 1^4 dx \right)^{1/4} = \\ &= \left(\int_{\Omega} |V(x)|^2 dx \right)^{3/4} \left(\int_{\Omega} 1^4 dx \right)^{1/4} < \infty, \end{aligned}$$

where $(\int_{\Omega} 1^4 dx)^{1/4}$ is finite as a measure of the unit cell Ω and $(\int (|V(x)|^2) dx)^{3/4}$ is also finite by the assumption that $V(x)$ is in L^2 .
Selection deleted

Therefore, $V(x) \in L_{per}^{3/2}(\Omega)$ and by applying theorem 10.1 from the lectures the operator H is self adjoint.

A consequence of (a) is that the Bloch-Floquet transformation yields Bloch fibers

$$H_k = T_k + V \quad \text{where} \quad T_k = \frac{1}{2}(-i\nabla_x + k)^2$$

which are self-adjoint on $L^2_{\text{per}}(\Omega)$ with $D(H_k) = H^2_{\text{per}}(\Omega)$. The natural inner product for our problem is thus the $\langle \cdot | \cdot \rangle_{L^2_{\text{per}}(\Omega)}$ inner product, which is the one we will employ unless a different one is specified.

For a particular \mathbf{k} -point we can thus discretise the problem using the plane-wave basis

$$\mathbb{B}_k^\mathcal{E} = \left\{ e_G \mid G \in \mathbb{L}^* \text{ and } \frac{1}{2}|G + k|^2 < \mathcal{E} \right\}$$

where \mathcal{E} is a chosen cutoff.

(b) Show that

$$\langle e_G | e_{G'} \rangle = \delta_{GG'}$$

and that

$$\langle e_G | T_k e_{G'} \rangle = \delta_{GG'} \frac{1}{2}|G + k|^2.$$

From this conclude that for any $e_G \in \mathbb{B}_k^\mathcal{E}$ we have

$$|\Delta G| > \sqrt{2\mathcal{E}_V} \implies \langle e_{G+\Delta G} | H e_G \rangle = 0, \quad (1)$$

i.e. that each plane wave only couples via the Hamiltonian with these plane waves, that differ in the wave vector by no more than $\sqrt{2\mathcal{E}_V}$.

Selection deleted

Solution (b):

We have

$$\langle e_G | e_{G'} \rangle = \int_{\Omega} e_G^*(x) e_{G'}(x) dx = \frac{1}{|\Omega|} \int_{\Omega} e^{i(G' - G) \cdot x} dx$$

In case when $G = G'$ we have:

$$\int_{\Omega} e^{i(G' - G) \cdot x} dx = \int_{\Omega} 1 dx = |\Omega|$$

In case of $G \neq G'$ denoting $(G - G') = G''$ which is also an element of the reciprocal lattice \mathbb{L}^* :

$$\int_{\Omega} e^{iG'' \cdot x} dx = \int_{\Omega} (\cos G'' \cdot x + i \sin G'' \cdot x) dx = \int_{\Omega} \cos G'' \cdot x dx.$$

Imaginary part is equal to zero since sine is an odd function. On the other hand, since we can represent the vectors x and G'' as follows:

$$x = \alpha_1 a_1 + \alpha_2 a_2 + \alpha_3 a_3$$

and

$$G'' = \beta_1 b_1 + \beta_2 b_2 + \beta_3 b_3,$$

we have

$$\cos G'' \cdot x = \cos(2\pi(\alpha_1 \beta_1 + \alpha_2 \beta_2 + \alpha_3 \beta_3)) = 0$$

since $a_i \cdot b_j = 2\pi \delta_{ij}$. Finally, we can see that

$$\langle e_G | e_{G'} \rangle = \frac{1}{|\Omega|} \delta_{G,G'} \cdot |\Omega|.$$

Considering

$$\begin{aligned} \text{Selection deleted } \langle e_G | T_k e_{G'} \rangle &= \int_{\Omega} e^{-iG \cdot x} \cdot \frac{1}{2} (-i\nabla_x + k)^2 e^{iG' \cdot x} dx \\ &= \frac{1}{2} \int_{\Omega} e^{-iG \cdot x} (-i\nabla_x + k) \left(-i\nabla_x e^{iG' \cdot x} + k e^{iG' \cdot x} \right) dx \\ &= \frac{1}{2} |G + k|^2 \langle e_G | e_{G'} \rangle = \frac{1}{2} \delta_{GG'} |G + k|^2 \end{aligned}$$

Finally, considering

$$H = -\frac{1}{2}\Delta + V.$$

We have

$$\begin{aligned} \langle e_{G+\Delta G} | H e_G \rangle &= -\frac{1}{2}|G|^2 \langle e_{G+\Delta G} | e_G \rangle + \langle e_{G+\Delta G} | V e_G \rangle = \\ &= 0 + \int_{\Omega} \sum_{|G'| < \sqrt{2\mathcal{E}_V}} \hat{V}(G') e^{i \cdot G' \cdot x} e^{-i \Delta G \cdot x} dx = 0 \end{aligned}$$

using the fact that $|\Delta G| > \sqrt{2\mathcal{E}_V}$ and that $e_G \in \mathbb{B}_k^{\mathcal{E}}$.

Familiarisation and Bauer-Fike estimates

Work through https://docs.dftk.org/stable/guide/periodic_problems/ as well as <https://docs.dftk.org/stable/guide/tutorial/> to get some basic familiarity with DFTK.

Following loosely [another documented example](#), this code sets up a calculation of silicon using the Cohen-Bergstresser potential and $\mathcal{E} = 10$.

```

1 begin
2   Si = ElementCohenBergstresser(:Si)
3   atoms = [Si, Si]
4   positions = [ones(3)/8, -ones(3)/8]
5   lattice = Si.lattice_constant / 2 .* [[0 1 1.]; [1 0 1.]; [1 1 0.]]
6 end;

```

Selection deleted

```

PlaneWaveBasis discretization:
  architecture          : DFTK.CPU()
  num. mpi processes   : 1
  num. julia threads   : 4
  num. blas threads    : 1
  num. fft threads     : 1

  Ecut                  : 10.0 Ha
  fft_size              : (24, 24, 24), 13824 total points
  kgrid type            : Monkhorst-Pack
  kgrid                 : [1, 1, 1]
  num. irredu. kpoints  : 1

Discretized Model(custom, 3D):
  lattice (in Bohr)    : [0           , 5.13061    , 5.13061
                         [5.13061   , 0           , 5.13061
                         [5.13061   , 5.13061    , 0
  unit cell volume      : 270.11 Bohr³

  atoms                 : Si₂
  atom potentials       : DFTK.ElementCohenBergstresser(Si)
                           DFTK.ElementCohenBergstresser(Si)

  num. electrons         : 8
  spin polarization      : none
  temperature            : 0 Ha

  terms                 : Kinetic()
                           DFTK.AtomicLocal()

```

```
1 begin
2     model = Model(lattice, atoms, positions; terms=[Kinetic(), AtomicLocal()])
3     basis_one = PlaneWaveBasis(model; Ecut=10.0, kgrid=(1, 1, 1));
4 end
```

The `kgrid` parameter selects the \mathbf{k} -point mesh $\mathbb{K} \subset \overline{\Omega^*}$, which is employed by determining the number of \mathbf{k} -points in each space dimension. In this case only a single \mathbf{k} -point is used, namely the origin of $\overline{\Omega^*}$:

```
[KPoint([      0,      0,      0], spin = 1, num. G vectors = 411)]  
1 basis_one.kpoints
```

The origin of the Brillouin zone has a special name and is usually called Γ -point. Before treating the case of band structure computations (where \mathbf{k} is varied), we stick to using only a single \mathbf{k} -point.

Each KPoint in DFTK automatically stores a representation of the basis $\mathbb{B}_k^{\mathcal{E}}$. If you are curious, the **G-Selection** of the respective plane waves e_G can be looked at using the functions

```
[StaticArraysCore.SVector{3, Float64}: [0.0, 0.0, 0.0], StaticArraysCore.SVector{3, Float64}: [0.0, 0.0, 0.0], StaticArraysCore.SVector{3, Float64}: [0.0, 0.0, 0.0]]
```

Similarly a discretised representation of \mathbf{H} restricted to the fibres matching the chosen \mathbf{k} -point mesh can be obtained using

```
ham_one =
Hamiltonian(PlaneWaveBasis discretization:
    architecture           : DFTK.CPU())
, [DftHami
◀ ━━━━━━ ━━━━━━ ━━━━▶
1 ham_one = Hamiltonian(basis_one)
```

Let us denote by $P_k^{\mathcal{E}}$ the projection into the plane-wave basis $\mathbb{B}_k^{\mathcal{E}}$ and further by $H_k^{\mathcal{E}\mathcal{E}} \equiv P_k^{\mathcal{E}} H_k P_k^{\mathcal{E}}$ the discretised fibers with elements

$$(H_k^{\mathcal{E}\mathcal{E}})_{GG'} = (P_k^{\mathcal{E}} H_k P_k^{\mathcal{E}})_{GG'} = \langle e_G, H_k e_{G'} \rangle \quad \text{for } G, G' \in \mathbb{B}_k^{\mathcal{E}}.$$

Representations of $H_k^{\mathcal{E}\mathcal{E}}$ in DFTK are available via indexing of Hamiltonian objects. E.g. `ham[1]` corresponds to the fiber of `basis.kpoints[1]`, `ham[2]` to `basis.kpoints[2]` and so on. This can also be seen from comparing their sizes with the size of $\mathbb{B}_k^{\mathcal{E}}$ (available via the `G_vectors_cart`):

`DFTK.DftHamiltonianBlock`

```
1 begin
2     @show length(G_vectors_cart(basis_one, basis_one.kpoints[1]))
3     @show size(ham_one[1])
4     @show typeof(ham_one[1])
5 end
```

```
length(G_vectors_cart(basis_one, basis_one.kpoints[1])) = 411
size(ham_one[1]) = (411, 411)
typeof(ham_one[1]) = DFTK.DftHamiltonianBlock
```

?

These `DftHamiltonianBlock` objects very much behave like matrices, e.g. they can be multiplied with vectors, but also used in iterative solvers.

The `diagonalize_all_kblocks` function does some optimisations and book-keeping to solve for the `n_bands` lowest eigenpairs at each H_k . Here we solve for 6 eigenpairs:

```
(λ = [[-0.079228, 0.3843, 0.3843, 0.3843, 0.510143, 0.510143]], X = [411×6 Matrix{ComplexF
-0 503153+0 806491
◀ ━━━━━━ ━━━━▶
1 begin
2     n_bands = 6
3     eigres_one = diagonalize_all_kblocks(DFTK.lobpcg_hyper, ham_one, n_bands)
4 end
Selection deleted
```

The returned object `eigres_one` now contains the eigenvalues in `band_data.λ` and the Bloch waves in `band_data.X`. Again both are a vector along the k -points, further `band_data.X[ik]` is of size `length(G_vectors_cart(basis, kpoint)) × n_bands`:

(411, 6)

```

1 begin
2   @show size(eigres_one.X[1])
3   @show (length(G_vectors_cart(basis_one, basis_one.kpoints[1])), n_bands)
4 end

```

size(eigres_one.X[1]) = (411, 6)
 $(\text{length}(\text{G_vectors_cart}(\text{basis_one}, \text{basis_one}.kpoints[1])), \text{n_bands}) = (411, 6)$?

As a result if we wanted to compute at each \mathbf{k} -point the residual

$$\mathbf{H}_k^{\mathcal{E}\mathcal{E}} \tilde{\mathbf{X}}_{kn} - \tilde{\lambda}_{kn} \tilde{\mathbf{X}}_{kn} \quad (2)$$

after we have been able to diagonalise $\mathbf{H}_k^{\mathcal{E}\mathcal{E}}$, we can simply iterate as such:

```

1 for (ik, kpt) in enumerate(basis_one.kpoints)
2   hamk = ham_one[ik]
3   λk   = eigres_one.λ[ik]
4   Xk   = eigres_one.X[ik]
5
6   residual_k = hamk * Xk - Xk * Diagonal(λk)
7   println(ik, " ", norm(residual_k))
8 end

```

1 1.906458902168061e-6 ?

Task 2: Γ -point calculations

(a) For the case of employing only a single \mathbf{k} -point (`kgrid = (1, 1, 1)`) vary \mathcal{E} (i.e. `Ecut`) between 5 and 30. Taking $\mathcal{E} = 80$ as a reference, plot the convergence of the first two eigenpairs of \mathbf{H}_k at the Γ point.

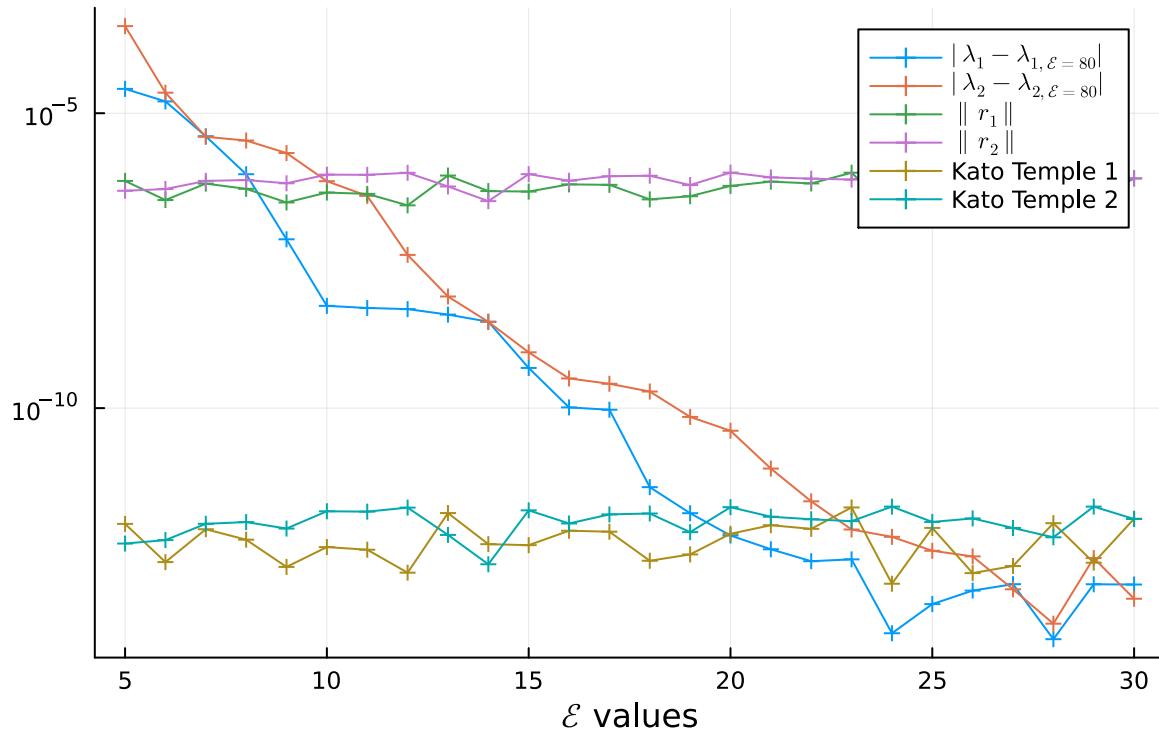
Selection deleted

```
9257×2 Matrix{ComplexF64}:
 -0.229462-0.922464im 7.86019e-8+2.26517e-8im
 -0.025416-0.102175im -0.0820593+0.238735im
 -0.00235226-0.00945631im -0.00077991+0.00226897im
 4.47873e-5+0.000180067im -0.000226507+0.000658962im
 3.26618e-7+1.31129e-6im 1.46728e-5-4.26949e-5im
 2.83319e-7+1.13659e-6im -2.04392e-7+5.95927e-7im
 -6.02277e-10-2.78884e-9im 2.36527e-8-6.86782e-8im
 :
 -3.61888e-8-1.45086e-7im 6.02275e-9-3.6302e-7im
 -7.97043e-8-3.22212e-7im -4.34045e-6+1.35597e-5im
 -3.88321e-5-0.000156105im 8.52411e-5-2.79917e-5im
 -2.73819e-5-0.000110068im 0.00180911+0.00496245im
 -0.00247903-0.009966im 0.0435335+0.00199638im
 0.025416+0.102175im 0.262493+0.254536im
```

```
1 begin
2
3   Ecuts = 5:30
4   Ecut_ref = 80
5   n_bands_2a = 2
6   ik = 1
7
8   eigenvalues_2a = zeros(length(Ecuts), n_bands_2a)
9   eigenvec1 = []
10  eigenvec2 = []
11
12  for (i, Ecut) in enumerate(Ecuts)
13    basis_2a = PlaneWaveBasis(model; Ecut=Ecut, kgrid=(1, 1, 1))
14    ham_2a = Hamiltonian(basis_2a)
15    eigres_2a = diagonalize_all_kblocks(DFTK.lobpcg_hyper, ham_2a, n_bands_2a)
16
17    eigenvalues_2a[i, :] = eigres_2a.λ[ik][1:n_bands_2a]
18    Xk = eigres_2a.X[ik]
19    push!(eigenvec1, Xk[1])
20    push!(eigenvec2, Xk[2])
21
22  end
23
24  basis_r = PlaneWaveBasis(model; Ecut=Ecut_ref, kgrid=(1, 1, 1))
25  ham_r = Hamiltonian(basis_r)
26  eigres_r = diagonalize_all_kblocks(DFTK.lobpcg_hyper, ham_r, n_bands_2a)
27  λ_ref = eigres_r.λ[ik]
28  X_ref = eigres_r.X[ik]
29
30 end
31
```

Selection deleted

Convergence of the first two eigenvalues



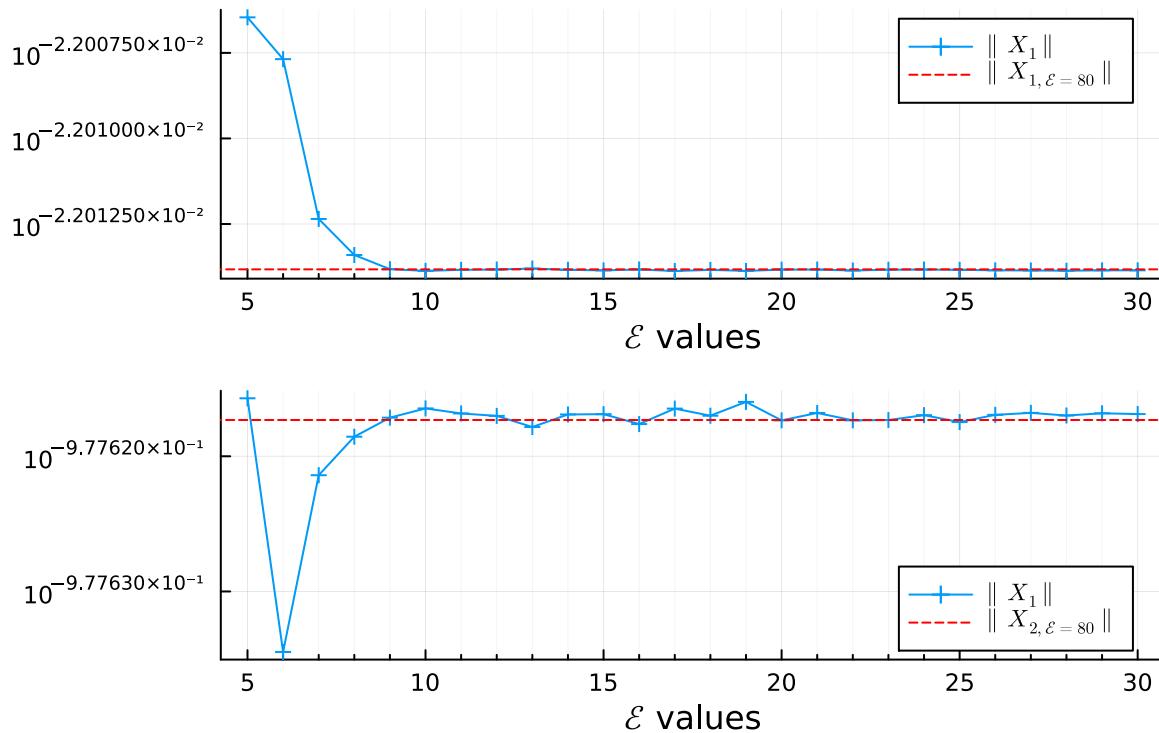
```

1 begin
2     p = plot(yaxis=:log, xlabel=L"\mathcal{E} values")
3     λ1 = abs.(eigenvalues_2a[:, 1] .- λ_ref[1])
4     λ2 = abs.(eigenvalues_2a[:, 2] .- λ_ref[2])
5
6
7     plot!(p, Ecuts, λ1, shape=:cross, label=L"|\lambda_1-\lambda_1,\mathcal{E}=80|")
8     plot!(p, Ecuts, λ2, shape=:cross, label=L"|\lambda_2-\lambda_2,\mathcal{E}=80|")
9
10    plot!(p, Ecuts, ρ_bauer_fike[:, 1], shape=:cross, label=L"\parallel r_1 \parallel")
11    plot!(p, Ecuts, ρ_bauer_fike[:, 2], shape=:cross, label=L"\parallel r_2 \parallel")
12
13    plot!(p, Ecuts, ρ_kato_temple[:, 1], shape=:cross, label="Kato Temple 1")
14    plot!(p, Ecuts, ρ_kato_temple[:, 2], shape=:cross, label="Kato Temple 2")
15
16    title!(p, "Convergence of the first two eigenvalues")
17 end

```

Selection deleted

Convergence of the first two eigenvectors



```

1 begin
2     p_x= plot(layout=(n_bands_2a, 1), xlabel=L"\mathcal{E} values", yscale=:log10,
3             minorgrid=true)
4     X1=norm.(eigenvec1)
5     X2=norm.(eigenvec2)
6
7     plot!(p_x, Ecuts, X1, shape=:cross, subplot=1, label=L"\parallel X_1 \parallel")
8     hline!([norm.(X_ref[1])], line=:dash, color=:red, subplot=1, label=L"\parallel X_{1,\mathcal{E}=80}\parallel")
9
10    plot!(p_x, Ecuts, X2, shape=:cross, subplot=2, label=L"\parallel X_1 \parallel")
11    hline!([norm.(X_ref[2])], line=:dash, color=:red, subplot=2, label=L"\parallel X_{2,\mathcal{E}=80}\parallel")
12
13    title!(p_x, subplot=1, "Convergence of the first two eigenvectors")
14
15 end

```

Selection deleted

We want to compare this convergence behaviour with a first estimate of the eigenvalue discretisation error based on the Bauer-Fike bound. Given an approximate eigenpair $(\tilde{\lambda}_{kn}, \tilde{X}_{kn})$ of the fiber H_k we thus need access to the residual

$$r_{kn} = H_k \tilde{X}_{kn} - \tilde{\lambda}_{kn} \tilde{X}_{kn}.$$

If we performed our computation using the plane-wave basis $\mathbb{B}_k^{\mathcal{E}}$, then by construction

$$P_k^{\mathcal{E}} \tilde{X}_{kn} = \tilde{X}_{kn},$$

such that

$$\begin{aligned} P_k^{\mathcal{E}} r_{kn} &= P_k^{\mathcal{E}} H_k P_k^{\mathcal{E}} P_k^{\mathcal{E}} \tilde{X}_{kn} - \tilde{\lambda}_{kn} P_k^{\mathcal{E}} \tilde{X}_{kn} \\ &= H_k^{\mathcal{E}\mathcal{E}} \tilde{X}_{kn} - \tilde{\lambda}_{kn} \tilde{X}_{kn}, \end{aligned}$$

which is exactly the quantity we computed in (2) above. In other words $P_k^{\mathcal{E}} r_{kn}$ is the residual corresponding to the iterative diagonalisation, thus leading to the **algorithm error**.

However, we are also interested in the **discretisation error**, which we obtain from the missing residual term in (2), namely

$$Q_k^{\mathcal{E}} r_{kn} = (1 - P_k^{\mathcal{E}}) r_{kn} = Q_k^{\mathcal{E}} H_k P_k^{\mathcal{E}} \tilde{X}_{kn} = H_k^{\mathcal{E}^\perp \mathcal{E}} \tilde{X}_{kn}$$

where $Q_k^{\mathcal{E}} = (1 - P_k^{\mathcal{E}})$ is the orthogonal projector to $P_k^{\mathcal{E}}$ and we introduced the notation

$$H_k^{\mathcal{E}^\perp \mathcal{E}} = Q_k^{\mathcal{E}} H_k P_k^{\mathcal{E}}.$$

While computing $P_k^{\mathcal{E}} r_{kn}$ is generally easy (we just did it above in 5 lines), obtaining $Q_k^{\mathcal{E}} r_{kn}$ is impossible for general potentials V : since $H_k P_k^{\mathcal{E}} \tilde{X}_{kn}$ can have support anywhere on $L^2_{\text{per}}(\Omega)$, estimating $Q_k^{\mathcal{E}} r_{kn}$ usually requires making a mathematical statement about the behaviour of H_k on all of $L^2_{\text{per}}(\Omega)$. Estimating discretisation errors is thus substantially more challenging than estimating algorithm errors.

(b) In our case of Cohen-Bergstesser Hamiltonians we are more lucky. Using (1) show that there exists a cutoff $\mathcal{F} > \mathcal{E}$, such that

Selection deleted

$$r_{kn} = P_k^{\mathcal{F}} r_{kn} = H_k^{\mathcal{F}\mathcal{F}} \tilde{X}_{kn} - \tilde{\lambda}_{kn} \tilde{X}_{kn}.$$

Having diagonalised the Hamiltonian using cutoff \mathcal{E} we can thus obtain the *full residual* by a computation of the Hamiltonian-vector product using an elevated cutoff \mathcal{F} .

Solution (b)

We need to find such cutoff $\mathcal{F} > \mathcal{E}$ that

$$r_{kn} = P_k^{\mathcal{F}} r_{kn} \iff Q_k^{\mathcal{F}} r_{kn} = 0 \iff H_k^{\mathcal{F}^\perp \mathcal{E}} \tilde{X}_{kn} = 0,$$

where

$$\left(H_k^{\mathcal{F}^\perp \mathcal{E}} \right)_{GG'} = \langle e_G, H_k e_{G'} \rangle$$

If $e_G \in (\mathbb{B}_k^{\mathcal{F}})^\perp$ and $e_{G'} \in \mathbb{B}_k^{\mathcal{E}}$ which means that

$$\frac{1}{2}|G + k|^2 < \mathcal{F} \quad \text{and} \quad \frac{1}{2}|G' + k|^2 < \mathcal{E},$$

from which it follows that

$$e_G \in (\mathbb{B}_k^{\mathcal{F}})^\perp \iff \frac{1}{2}|G + k|^2 \geq \mathcal{F} \quad (*).$$

Considering $e_G \in \mathbb{B}_k^{\mathcal{F}}$ and $e_{G'} \in \mathbb{B}_k^{\mathcal{E}}$ and applying the reverse triangular inequality we have

$$\begin{aligned} |G - G'| &= |G + k - G' - k| \geq ||G + k| - |G' + k|| \geq \\ &\geq \sqrt{2\mathcal{F}} - \sqrt{2\mathcal{E}} \end{aligned}$$

From (1) we can see that for any $e_G \in \mathbb{B}_k^{\mathcal{F}}$

$$|\Delta G| > \sqrt{2\mathcal{E}_V} \implies \langle e_{G+\Delta G} | H_k e_G \rangle = 0.$$

To have $\langle e_{G+\Delta G} | H_k e_G \rangle \neq 0$ and since $\mathbb{B}_k^{\mathcal{E}} \subset \mathbb{B}_k^{\mathcal{F}}$ for $e_G \in \mathbb{B}_k^{\mathcal{F}}$ and $e_{G'} \in \mathbb{B}_k^{\mathcal{E}}$ we need

$$|\Delta G| = |G - G'| \leq \sqrt{2\mathcal{E}_V}$$

From both inequalities we have

$$\sqrt{\mathcal{F}} - \sqrt{\mathcal{E}} = \sqrt{\mathcal{E}_V}.$$

Therefore, if we take

Selection deleted

$$\mathcal{F} = \left(\sqrt{\mathcal{E}_V} + \sqrt{\mathcal{E}} \right)^2.$$

it will follow that

$$\begin{aligned} \left(H_k^{\mathcal{F}^\perp \mathcal{E}} \right)_{GG'} &= \langle e_G, H_k e_{G'} \rangle = \langle e_G, V e_{G'} \rangle = \\ &= \delta_{GG'} \frac{1}{2} |G + k|^2 + \int_{\Omega} \sum_{|G''| < \sqrt{2\mathcal{E}_V}} \hat{V}(G'') e^{i \cdot G'' \cdot x} e^{i(G-G') \cdot x} dx = 0 \end{aligned}$$

since $G' \neq G$ from (*).

(c) In DFTK nobody stops you from using multiple bases of different size. Moreover, having obtained a set of bloch waves `X_small` using `basis_small` you can obtain its representation on the bigger `basis_large` using the function

```
x_large = transfer_blochwave(X_small, basis_small, basis_large)
```

Using this technique you can compute the application of the Hamiltonian using a bigger basis (just as `Hamltionian(basis_large) * X_large`). Use this setup to vary \mathcal{F} and use this to estimate \mathcal{E}_V numerically. Check your estimate for various values of \mathcal{E} to ensure it is consistent. Rationalise your results by taking a look at the [Cohen Bergstresser implementation in DFTK](#).

Solution (c):

Selection deleted

```

1 begin
2
3     E_list = 5:10
4     F_found_list = []
5     for E in E_list
6
7         local basis_small = PlaneWaveBasis(model; Ecut=E, kgrid=(1, 1, 1))
8         local ham = Hamiltonian(basis_small)
9         local eigres = diagonalize_all_kblocks(DFTK.lobpcg_hyper, ham, 2)
10
11    F_list = E:5:E
12    local ik = 1
13
14    res_norms = []
15
16    for (i, F) in enumerate(F_list)
17        basis_large = PlaneWaveBasis(model; Ecut=F, kgrid=(1, 1, 1))
18        ham_large = Hamiltonian(basis_large)
19        eigres_large = diagonalize_all_kblocks(DFTK.lobpcg_hyper, ham_large, 2)
20
21        Xk_large = transfer_blochwave(eigres.X, basis_small, basis_large)
22
23        residual = ham_large[ik] * Xk_large[ik] - Xk_large[ik] *
24            Diagonal(eigres.λ[ik])
25        push!(res_norms, norm(residual))
26
27        if i > 1 && abs(res_norms[i] - res_norms[i-1]) < 1e-8
28            push!(F_found_list, F)
29            break
30        end
31    end
32
33 end
34 end
35

```

Ev_estimated = 2.1720066733367176

```

1 Ev_estimated = mean((sqrt.(F_found_list) - sqrt.(E_list)).^2)
```

Since the reciprocal lattice constant $b = \frac{2\pi}{a}$ and from the Cohen Bergstresser implementation in DFTK the potential is calculated for the sum over G such that

Selection deleted

$$|G|^2 < 11 \left(\frac{2\pi}{a} \right)^2 \approx 4.124$$

by definition given above

$$|G|^2 < 2\mathcal{E}_V \implies \mathcal{E}_V > 2.062$$

2.0621726989864846

```
1 0.5 * 11 * (2 * π / Si.lattice_constant) ^ 2
```

On the other hand, G' such that

$$|G'|^2 > 11 \left(\frac{2\pi}{a} \right)^2$$

are not included while computing V meaning that

$$|G'|^2 > 2\mathcal{E}_V$$

giving an upper bound for \mathcal{E}_V .

2.24965

```
1 begin
2     l_bound = 11 * (2π / Si.lattice_constant) ^ 2
3     G_norms = (norm.(G_vectors_cart(basis_one, basis_one.kpoints[1]))) .^ 2
4     0.5 * sort([g for g in round.(G_norms, sigdigits=5) if g > l_bound])[1]
5 end
```

Therefore,

$$2.06 < \mathcal{E}_V < 2.25$$

which is in line with the value estimated above.

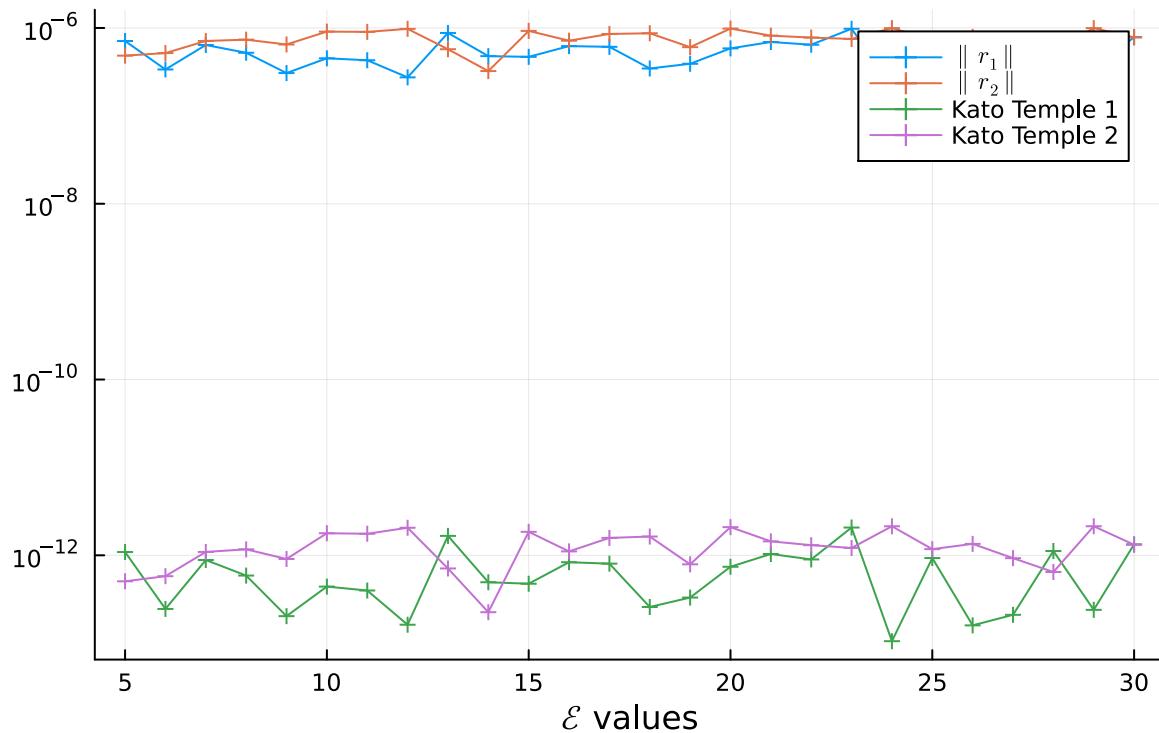
(d) Based on the Bauer-Fike bound estimate the algorithm and arithmetic error for the first two eigenpairs at the Γ point and for using cutoffs between $\mathcal{E} = 5$ and $\mathcal{E} = 30$. Add these estimates to your plot in **(a)**. What do you observe regarding the tightness of the bound ?

Selection deleted

```
1 begin
2     local n_bands = 2
3     ρ_bauer_fike = zeros(length(Ecuts), n_bands)
4
5     for (i, Ecut) in enumerate(Ecuts)
6         basis_2a = PlaneWaveBasis(model; Ecut=Ecut, kgrid=(1, 1, 1))
7         ham_2a = Hamiltonian(basis_2a)
8         eigres_2a = diagonalize_all_kblocks(DFTK.lobpcg_hyper, ham_2a, n_bands_2a)
9
10    eigenvalues_2a[i, :] = eigres_2a.λ[ik][1:n_bands_2a]
11    Xk = eigres_2a.X[ik]
12
13
14    residual_k = ham_2a[ik] * Xk - Xk * Diagonal(eigres_2a.λ[ik])
15    ρ_bauer_fike[i, :] = norm.(eachcol(residual_k))
16
17 end
18
19
20 end
```

Selection deleted

Bauer Fike and Kato Temple bounds



```

1 begin
2
3     local p= plot(yaxis=:log, xlabel=L"\mathcal{E} values")
4
5     plot!(p, Ecuts, ρ_bauer_fike[:, 1], shape=:cross, label=L"\parallel r_1 \parallel")
6     plot!(p, Ecuts, ρ_bauer_fike[:, 2], shape=:cross, label=L"\parallel r_2 \parallel")
7
8     plot!(p, Ecuts, ρ_kato_temple[:, 1], shape=:cross, label="Kato Temple 1")
9     plot!(p, Ecuts, ρ_kato_temple[:, 2], shape=:cross, label="Kato Temple 2")
10
11    title!(p, "Bauer Fike and Kato Temple bounds")
12 end

```

We can notice that Bauer Fike bounds for both eigenvalues are not very tight but for the first \mathcal{E} we can see that discretization error dominates knowing that Bauer Fike dives a bound for arithmetic and algorithmic errors.

Selection deleted

Band structure computations

Next we turn our attention to band structure computations. Recall that the n -th band is the mapping from \mathbf{k} to the n -th eigenvalue $\lambda_{\mathbf{k}n}$ of the Bloch fiber $\mathbf{H}_\mathbf{k}$. For standard lattices — like the diamond-FCC lattices we consider — there are tabulated standard 1D paths through the Brillouin zone Ω^* , that should be used to maximise the physical insight one may gain from the computation.

In DFTK this path can be automatically determined. The following function provides a corresponding list of \mathbf{k} -points to consider (the warnings can be ignored in this function call):

kpath =

```
Brillouin.KPaths.KPathInterpolant{3}: [StaticArraysCore.SVector{3, Float64}: [0.0, 0.0, 0  
1 kpath = interpolate(irrfbz_path(model); density=15)
```

The density parameter in the above call determines how many \mathbf{k} -points are considered along the path. A larger density leads to more points.

Based on this data we can use the `compute_bands` function to diagonalise all \mathbf{H}_k along the path using a defined density and cutoff `Ecut`. Internally this function uses `diagonalize_all_kblocks`, so the returned data structure is alike.

`compute_bands_auto` (generic function with 1 method)

```
1 function compute_bands_auto(model; Ecut, kwargs...)
2     basis = PlaneWaveBasis(model; Ecut, kgrid=(1, 1, 1))
3     compute_bands(basis, kpath; show_progress=false, kwargs...)
4 end
```

band_data =

```
(basis = PlaneWaveBasis discretization:  
    architecture      : CPU()  
  
1 band_data = compute_bands_auto(model; n_bands=6, tol=1e-3, Ecut=10)
```

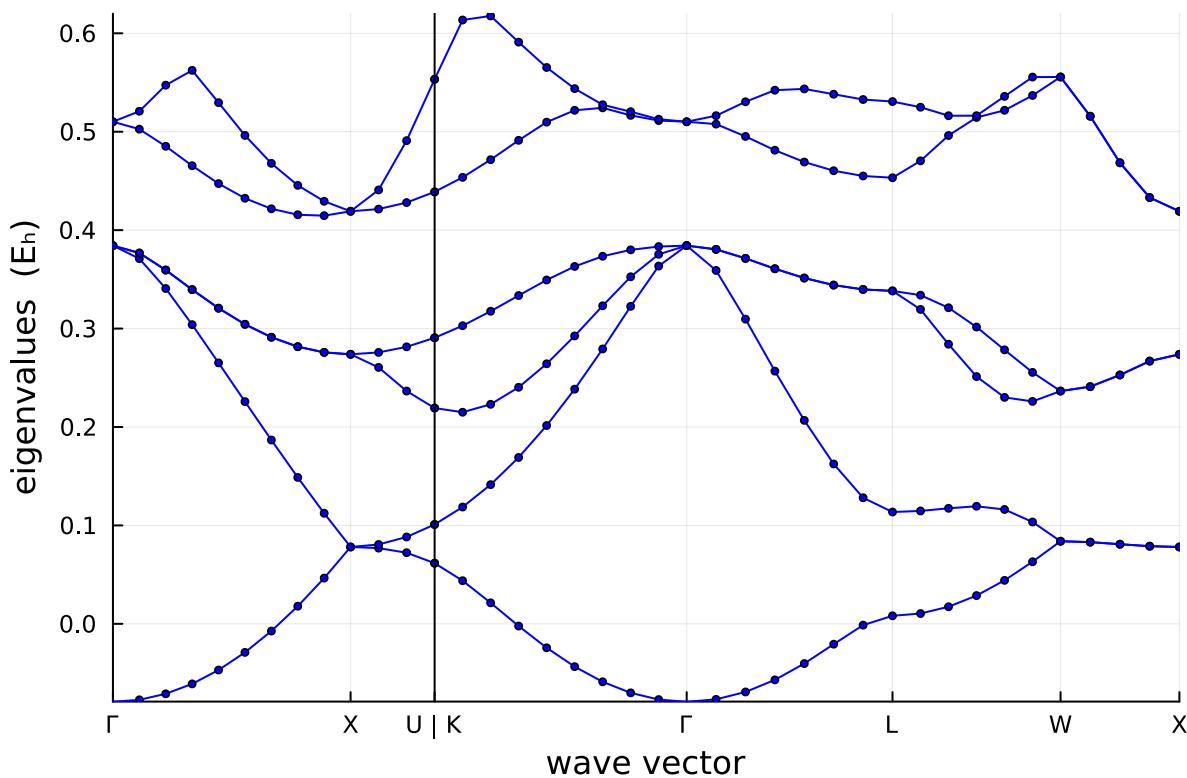
For example we obtain the eigenvalues of the 3rd k -point

```
Selection deleted  
KPoint([ 0.111,      0,  0.111], spin = 1, num. G vectors = 399)  
1 band_data.basis.kpoints[3]
```

using the expression

```
[ -0.0711114, 0.34071, 0.359448, 0.359448, 0.485213, 0.547298]
```

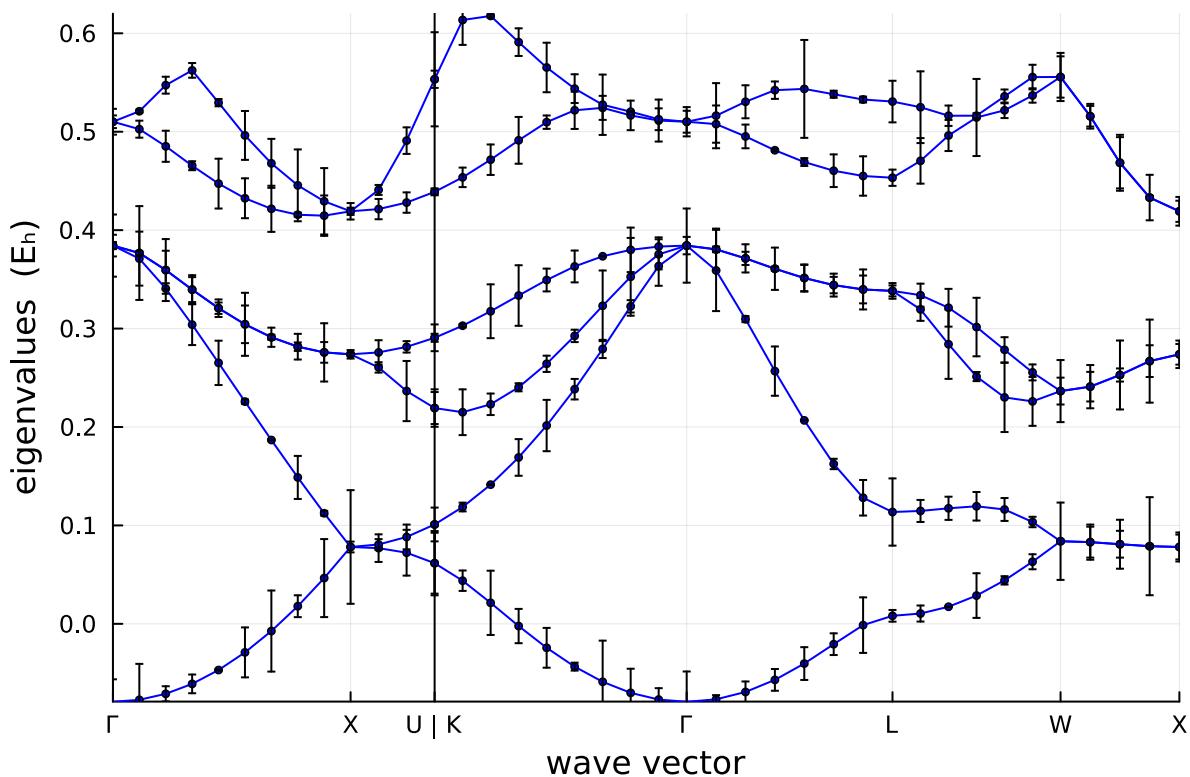
Once the data has been computed, a nice bandstructure plot can be produced using



```
1 DFTK.plot_band_data(kpath, band_data)
```

Error bars for indicating eigenvalue errors can also be easily added:

Selection deleted



```

1 begin
2     λerror = [0.02 * abs.(randn(size(λk))) for λk in band_data.λ] # dummy data
3     data_with_errors = merge(band_data, (; λerror))
4     DFTK.plot_band_data(kpath, data_with_errors)
5 end

```

Task 3: Bands with Bauer-Fike estimates

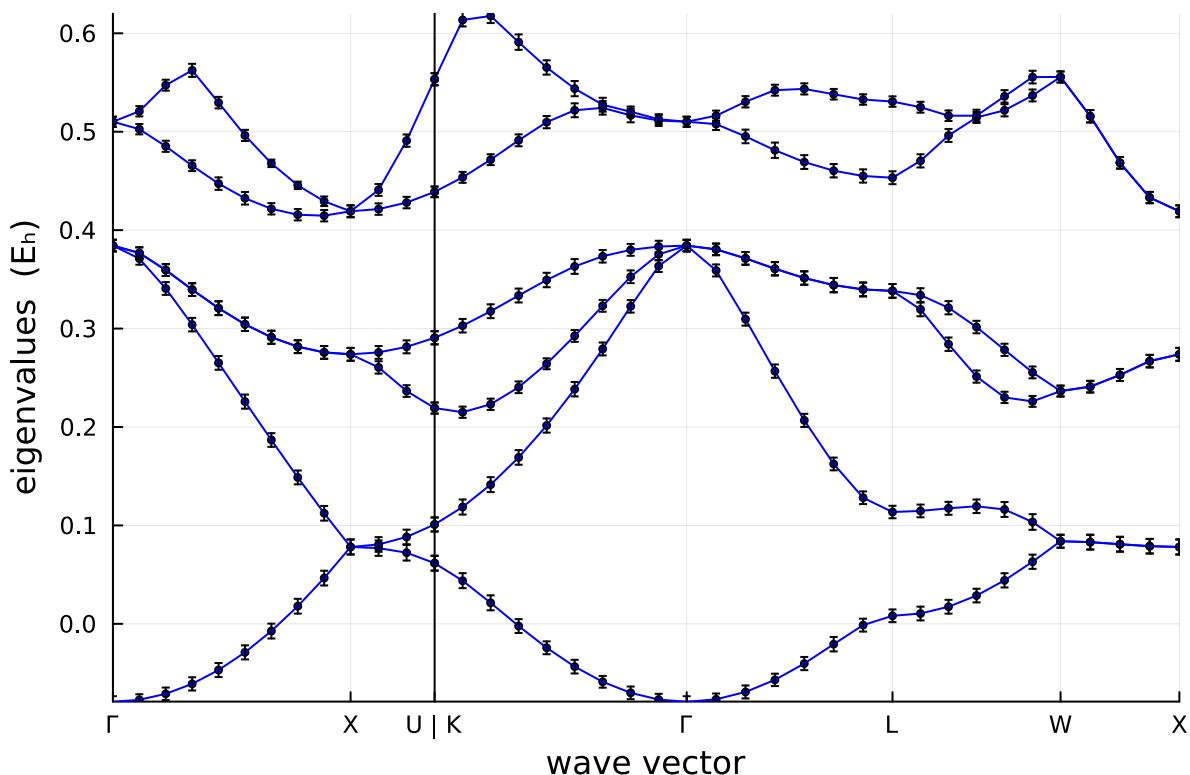
For $\mathcal{E} = 7$ a k -point density of 15 and the 6 lowest eigenpairs plot a band structure in which you annotate the bands with error bars estimated from Bauer-Fike. To compute the Hamiltonian-vector product employing a larger basis set with cutoff `Ecut_large`, employ the following piece of code, which forwards the k -Point coordinates from one basis to another:

```

basis_change_Ecut (generic function with 1 method)
1 function basis_change_Ecut(basis_small, Ecut_large)
2     PlaneWaveBasis(basis_small.model, Ecut_large,
3                      basis_small.kcoords_global,
4                      basis_small.kweights_global)
5 end

```

Selection deleted



```

1 begin
2     local Ecut = 7
3     Fcut = (sqrt(Ev_estimated) + sqrt(Ecut)) ^ 2
4     local band_data = compute_bands_auto(model; n_bands=6, tol=1e-3, Ecut=Ecut)
5     local basis_large = basis_change_Ecut(band_data.basis, Fcut)
6
7     local X_large = transfer_blochwave(band_data.X , band_data.basis, basis_large)
8     local ham = Hamiltonian(basis_large)
9
10    ρ_bauer_fike3 = []
11
12    for (ik, kpt) in enumerate(band_data.basis.kpoints)
13        residual_k = ham[ik] * X_large[ik] - X_large[ik] * Diagonal(band_data.λ[ik])
14        push!(ρ_bauer_fike3, norm.(eachcol(residual_k)))
15    end
16
17    band_data3_bf=merge(band_data,(;λerror=ρ_bauer_fike3))
18    DFTK.plot_band_data(kpath, band_data3_bf)
19
20 end

```

Selection deleted

Kato-Temple estimates

Following up on Task 2 (d) we want to improve the bounds using Kato-Temple estimates. Similar to our previous developments in this regard the challenging ingredient is estimating a lower bound to the gap. If λ_{kn} denotes the eigenvalue closest to the approximate eigenvalue $\tilde{\lambda}_{kn}$, we thus need a lower bound to

$$\delta_{kn} = \min_{s \in \sigma(H_k) \setminus \{\lambda_{kn}\}} |s - \tilde{\lambda}_{kn}|.$$

The usual approach we employed so far (see e.g. Sheet 6) was to assume that our obtained approximations $\tilde{\lambda}_{kn}$ for $n = 1, \dots, N$ did not miss any of the exact eigenvalues. With this assumption we would employ a combination of Kato-Temple and Bauer-Fike bound to obtain a lower bound to δ_{kn} .

Already for estimating the eigenvalue error of iterative eigensolvers, this can be a far-fetched assumptions as we saw in previous exercises. For estimating the discretisation error where the discretisation by nature of projecting into a finite-dimensional space inevitably changes the properties of the spectrum, this is even more far fetched. In fact in practice it is not rare that an unsuited basis changes the order of eigenpairs, completely neglects certain eigenpairs etc. Since the goal of *a posteriori* estimation is exactly to detect such cases, we will in Task 5 also develop techniques based on weaker assumptions.

Before doing so, we first extend the usual technique to infinite dimensions.

Selection deleted

Task 4: Simple and non-guaranteed Kato-Temple estimates

- (a)** Notice that the Kato-Temple theorem requires the targeted eigenvalue λ_{kn} to be isolated, see the conditions of Theorem 9.9. and 9.8. in the notes. Considering our setting of periodic Schrödinger operators in L^2 . Looking at the spectral properties of H and H_k argue why there is even hope that Kato-Temple can be employed in our setting?
- (b)** Assume that our numerical procedure (discretisation + diagonalisation) for all considered cutoffs is sufficiently good, such that no eigenpairs are missed. Employ the previously discussed (e.g. Sheet 6) technique of combining Kato-Temple and Bauer-Fike estimates to obtain an improved estimate for the error in the eigenvalue. Add your estimate for the first eigenvalue at the Γ -point to your plot in Task 2 (d). What do you observe ?
- (c)** Perform a band structure computation at cutoff $\mathcal{E} = 7$ annotated with the error estimated using the estimate developed in (b). For approximate eigenvalues where you cannot obtain a Kato-Temple estimate (e.g. degeneracies), fall back to a Bauer-Fike estimate. Play with the cutoff. What do you observe ? Do the error bars correspond to the expectations of a variational convergence ? Apart from the probably unjustified assumption in (b), what is the biggest drawback of the Kato-Temple estimate ?
-

Solution (a):

From the notes, we know that eigenvalues of H_k are well separated. And that the spectrum of H is a union of spectrums of H_k . As discussed in the lecture an energy band gives us a well-defined minimum and maximum which can be employed for using the Kato-Temple theorem.

Solution (b):

Selection deleted

get_KT_bound (generic function with 1 method)

```

1 function get_KT_bound(computed_values, res_norms)
2     n = size(res_norms)[1]
3
4     δ_list = zeros(n)
5     δ_list[1] = abs(computed_values[1] - computed_values[2]) -
6         res_norms[2]
7     δ_list[end] = abs(computed_values[end] - computed_values[end-1]) -
8         res_norms[end-1]
9     for i in 2:Int64(n - 1)
10        δ_left = abs(computed_values[i] - computed_values[i-1]) -
11            res_norms[i-1]
12        δ_right = abs(computed_values[i] - computed_values[i+1]) -
13            res_norms[i+1]
14        δ_list[i] = min(δ_left, δ_right)
15    end
16
17    δ = [max(0, i) for i in δ_list]
18    error_Kato_Temple = res_norms .^2 ./ δ
19    error_Kato_Temple
20
21 end

```

```

1 begin
2     local ik = 1
3     ρ_kato_temple = zeros(length(Ecuts), n_bands_2a)
4
5     for (i, Ecut) in enumerate(5:30)
6
7         # Kato-Temple
8         ρ_kato_temple[i, :] = get_KT_bound(eigenvalues_2a[i, :], ρ_bauer_fike[i, :])
9
10    end
11 end
12
13

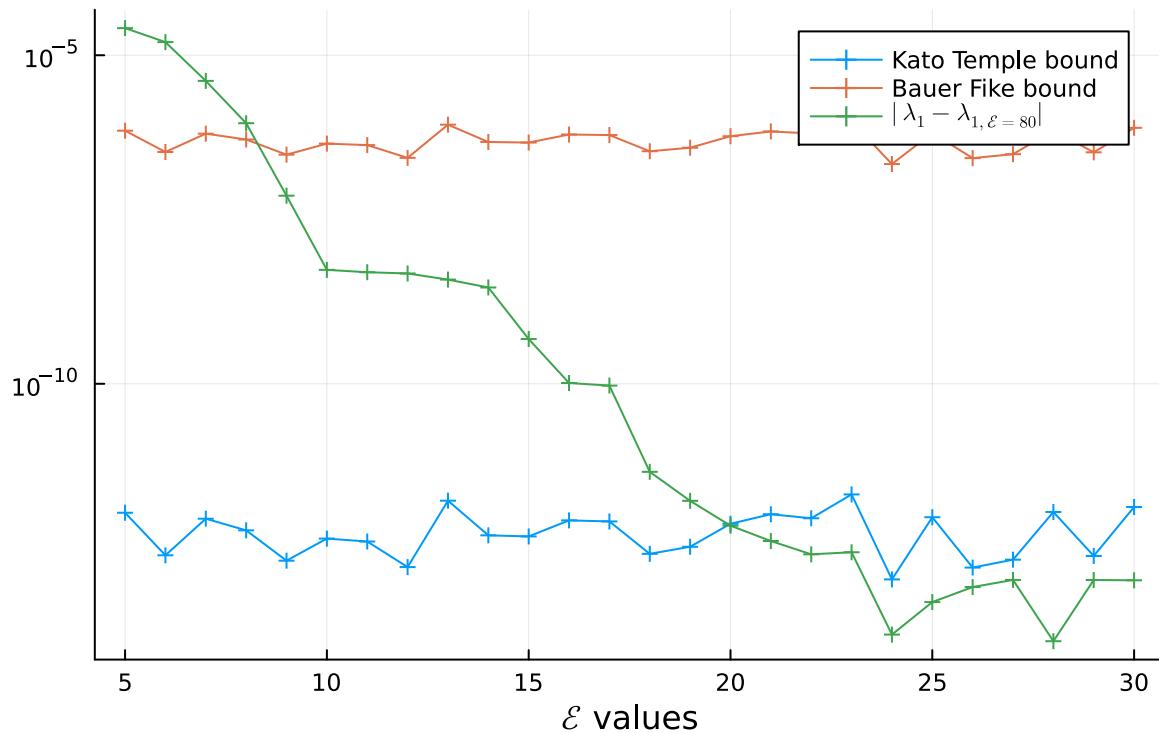
```

[1.08969e-12, 2.45311e-13, 8.84173e-13, 5.88207e-13, 2.03148e-13, 4.40665e-13, 3.9757e-13,

1 ρ_kato_temple[:, 1]

Selection deleted

Kato-Temple and Bauer Fike for the first eigenvalue



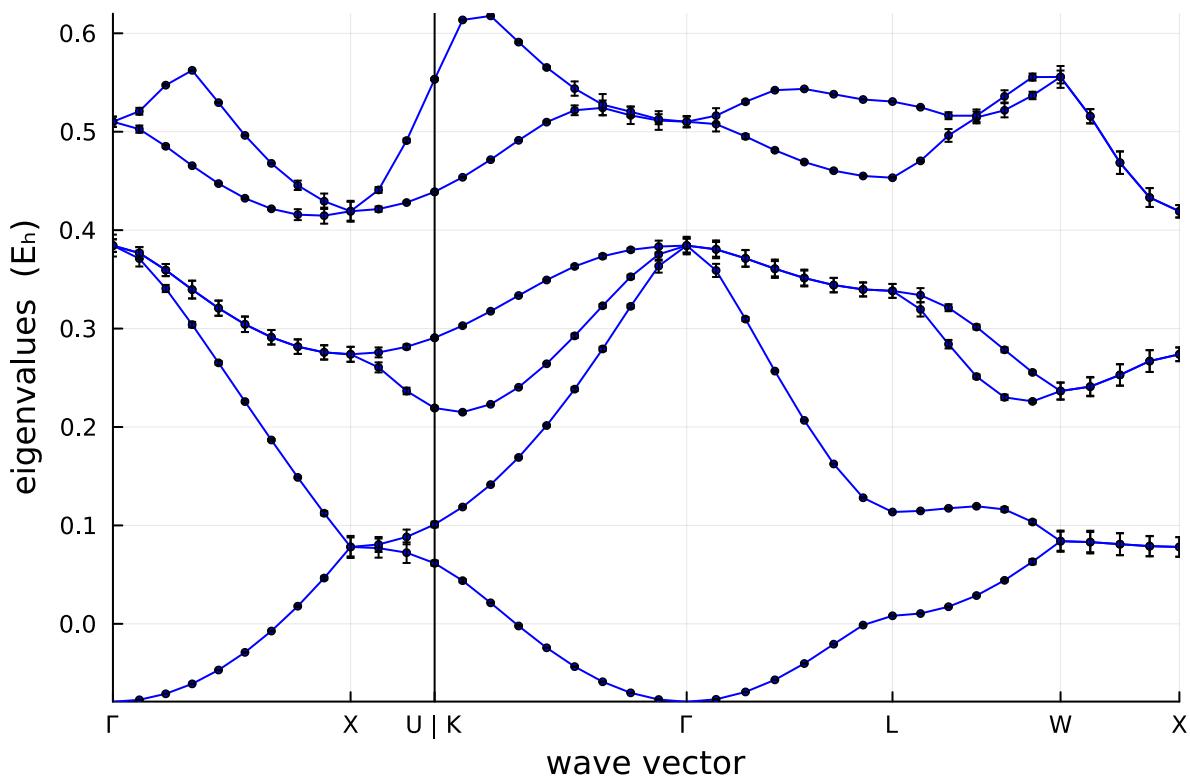
```

1 begin
2     local p = plot(yaxis=:log, xlabel=L"\mathcal{E} values")
3
4
5     plot!(p, Ecuts, rho_kato_temple[:, 1], shape=:cross, label="Kato Temple bound")
6     # plot!(p, Ecuts, rho_kato_temple[2], shape=:cross, label="Kato Temple eval2")
7
8     plot!(p, Ecuts, rho_bauer_fike[:, 1], shape=:cross, label="Bauer Fike bound")
9     # # plot!(p, Ecuts, rho_bauer_fike[:, 2], shape=:cross, label="Bauer Fike eval2")
10
11    plot!(p, Ecuts, lambda_1, shape=:cross, label=L"|\lambda_1 - \lambda_{1,\mathcal{E}=80}|")
12    # # plot!(p, Ecuts, lambda_2, shape=:cross, label=L"|\lambda_2 - \lambda_{2,\mathcal{E}=80}|")
13
14    title!(p, "Kato-Temple and Bauer Fike for the first eigenvalue")
15
16 end

```

In this case the Kato-Temple bound is tighter than its Bauer-Fike bound which makes Kato-Temple Selection deleted our task. Moreover, having a bound for algorithmic and arithmetic error we can conclude depending on the value \mathcal{E} where discretization error dominates in our approximate solution.

Solution (c):



```

1 begin
2     local Ecut = 7
3     local Fcut = (sqrt(Ev_estimated) + sqrt(Ecut)) ^ 2
4     local band_data = compute_bands_auto(model; n_bands=6, tol=1e-2, Ecut=Ecut)
5     local basis_large = basis_change_Ecut(band_data.basis, Fcut)
6
7     local X_large = transfer_blochwave(band_data.X , band_data.basis, basis_large)
8     local ham = Hamiltonian(basis_large)
9
10    ρ_bauer_fike4 = []
11    ρ_kato_temple4 = []
12    ρ_best = []
13
14    for (ik, kpt) in enumerate(band_data.basis.kpoints)
15        residual_k = ham[ik] * X_large[ik] - X_large[ik] * Diagonal(band_data.λ[ik])
16        r_norms = norm.(eachcol(residual_k))
17        kt_bound = get_KT_bound(band_data.λ[ik], r_norms)
18
19        push!(ρ_bauer_fike4, r_norms)
20        push!(ρ_kato_temple4, kt_bound)
21
22        push!(ρ_best, min.(r_norms, kt_bound))
Selection deleted
23
24
25    band_data4_bf=merge(band_data,(;λerror=ρ_best))
26    DFTK.plot_band_data(kpath, band_data4_bf)
27
28 end
29

```

Analyzing the plots from (c) and task 3 we observe that the error bars derived from the Kato-Temple estimate are tighter compared to those obtained only from the Bauer-Fike estimate. At the same time, the error bars tend to decrease with increasing cutoff, which would align with the expectations of smaller discretization errors leading to more accurate results. The main limitation is that we can apply the Kato-Temple theorem only when the targeted eigenvalue is isolated which is not always the case.

Task 5: A Schur-based estimate for the gap

Based on the definition of the basis cutoffs \mathcal{E} and \mathcal{F} and the respective projectors $P_k^{\mathcal{E}}, P_k^{\mathcal{F}}$ into these bases as well as $Q_k^{\mathcal{E}}$ and $Q_k^{\mathcal{F}}$ as the projectors to the complements of the bases, we can identify the following orthogonal decompositions of the Hilbert space $L^2_{\text{per}}(\Omega)$:

$$\begin{aligned} id &= (P_k^{\mathcal{E}} + Q_k^{\mathcal{E}}) \\ &= (P_k^{\mathcal{E}} + P_k^{\mathcal{F}} Q_k^{\mathcal{E}} + Q_k^{\mathcal{F}} Q_k^{\mathcal{E}}) \\ &= (P_k^{\mathcal{E}} + P_k^{\mathcal{R}} + Q_k^{\mathcal{F}}) \end{aligned}$$

where we defined $P_k^{\mathcal{R}} \equiv P_k^{\mathcal{F}} Q_k^{\mathcal{E}}$. Using the notations

$$\begin{aligned} H_k^{\mathcal{EF}} &= P_k^{\mathcal{E}} H_k P_k^{\mathcal{F}} & H_k^{\mathcal{RE}} &= P_k^{\mathcal{R}} H_k P_k^{\mathcal{E}} \\ H_k^{\mathcal{EE}^\perp} &= P_k^{\mathcal{E}} H_k Q_k^{\mathcal{E}} & H_k^{\mathcal{EF}^\perp} &= P_k^{\mathcal{E}} H_k Q_k^{\mathcal{F}} \end{aligned}$$

and so on, we can decompose the Hamiltonian fiber into blocks:

$$H_k = \begin{pmatrix} H_k^{\mathcal{EE}} & H_k^{\mathcal{EE}^\perp} \\ H_k^{\mathcal{E}^\perp \mathcal{E}} & \textcolor{red}{H_k^{\mathcal{E}^\perp \mathcal{E}^\perp}} \end{pmatrix} = \begin{pmatrix} H_k^{\mathcal{EE}} & H_k^{\mathcal{ER}} & H_k^{\mathcal{EF}^\perp} \\ H_k^{\mathcal{RE}} & \textcolor{red}{H_k^{\mathcal{RR}}} & H_k^{\mathcal{RF}^\perp} \\ H_k^{\mathcal{F}^\perp \mathcal{E}} & \textcolor{red}{H_k^{\mathcal{F}^\perp \mathcal{R}}} & H_k^{\mathcal{F}^\perp \mathcal{F}^\perp} \end{pmatrix}.$$

We note that our iterative diagonalisation as part of the band structure calculation indeed gives us access to a few eigenpairs $(\tilde{\lambda}_{kn}, \tilde{X}_{kn})$ of $H_k^{\mathcal{EE}}$. Moreover due to Courant-Fisher $\lambda_{kn} \leq \tilde{\lambda}_{kn}$, i.e. an upper bound to λ_{kn} is readily at hand. Suppose now we additionally had a lower bound μ_{kn} , i.e.

$$\mu_{kn} \leq \lambda_{kn} \leq \tilde{\lambda}_{kn}$$

then

$$\text{Selection deleted} \quad \delta_{kn} \geq \min \left(\mu_{kn} - \tilde{\lambda}_{k,n-1}, \mu_{k,n+1} - \tilde{\lambda}_{k,n} \right) \quad (3).$$

Finding such a suitable μ_{kn} provably is the goal of this task.

(a) Show that for Cohen-Bergstresser Hamiltonians and an appropriate choice of \mathcal{F} depending on \mathcal{E} (Task 2 (c)), the following structure of the Hamiltonian is obtained:

$$H_k = \begin{pmatrix} H_k^{\mathcal{EE}} & V_k^{\mathcal{ER}} & 0 \\ V_k^{\mathcal{RE}} & H_k^{\mathcal{RR}} & V_k^{\mathcal{RF}^\perp} \\ 0 & V_k^{\mathcal{F}^\perp\mathcal{R}} & H_k^{\mathcal{F}^\perp\mathcal{F}^\perp} \end{pmatrix},$$

where the use of V instead of H indicates that the kinetic term does not contribute to the respective block and $\mathbf{0}$ indicates an all-zero block.

Solution (a):

Using the results from task 2 (b), namely how \mathcal{F} is defined, we can conclude that indeed

$$H_k^{\mathcal{EF}^\perp} = H_k^{\mathcal{EF}^\perp} = 0$$

Furthermore,

$$(T_k^{\mathcal{RE}})_{GG'} = \langle e_G | T_k e_{G'} \rangle = \delta_{GG'} \frac{1}{2} |G + k|^2 = 0$$

since

$$\mathcal{E} < \frac{1}{2} |G + k|^2 < \mathcal{F} \quad \text{and} \quad \frac{1}{2} |G' + k|^2 < \mathcal{E}.$$

This means that the kinetic terms

$$T_k^{\mathcal{RE}} = T_k^{\mathcal{ER}} = 0$$

Analogically,

$$(T_k^{\mathcal{F}^\perp\mathcal{R}})_{GG'} = \langle e_G | T_k e_{G'} \rangle = \delta_{GG'} \frac{1}{2} |G + k|^2 = 0$$

where

$$\mathcal{F} < \frac{1}{2} |G + k|^2 \quad \text{and} \quad \mathcal{E} < \frac{1}{2} |G' + k|^2 < \mathcal{F}.$$

Selection deleted

Therefore,

$$T_k^{\mathcal{F}^\perp\mathcal{R}} = T_k^{\mathcal{RF}^\perp} = 0$$

We focus on the splitting of H_k into four blocks. An important realisation for our purpose is now that if $H_k - \mu_{kn}$ has exactly $n - 1$ negative eigenvalues, then we necessarily have $\mu_{kn} < \lambda_{kn}$. In the following we will thus develop conditions on a parameter μ , such that $H_k - \mu$ has a provable number of negative eigenvalues.

A first step is the Haynsworth inertia additivity formula [[Hay1968](#)]. For a matrix

$$H_k - \mu = \begin{pmatrix} H_k^{\mathcal{E}\mathcal{E}} - \mu & V_k^{\mathcal{E}\mathcal{E}^\perp} \\ V_k^{\mathcal{E}^\perp\mathcal{E}} & H_k^{\mathcal{E}^\perp\mathcal{E}^\perp} - \mu \end{pmatrix}.$$

this theorem states that

$$N(H_k - \mu) = N(H_k^{\mathcal{E}\mathcal{E}} - \mu) + N(S_\mu)$$

where $N(\mathcal{A})$ denotes the number of negative eigenvalues of the operator \mathcal{A} and S_μ is the Schur complement

$$S_\mu = (H_k^{\mathcal{E}^\perp\mathcal{E}^\perp} - \mu) - V_k^{\mathcal{E}^\perp\mathcal{E}} (H_k^{\mathcal{E}\mathcal{E}} - \mu)^{-1} V_k^{\mathcal{E}\mathcal{E}^\perp}.$$

[[Hay1968](#)]:

E. V. Haynsworth. Linear Algebra Appl. 1, 73 (1968)

Selection deleted

(b) Using this statement explain why a $\mu \in (\tilde{\lambda}_{k,n-1}, \tilde{\lambda}_{kn})$ such that $S_\mu \geq 0$ is a guaranteed lower bound to λ_{kn} .

(c) We proceed to obtain conditions that ensure $S_\mu \geq 0$. Let $X = L^2_{\text{per}}(\Omega) \setminus \text{span}(\mathbb{B}_k^\mathcal{E})$. Recall that $S_\mu \geq 0$ exactly if

$$\langle x, S_\mu x \rangle \geq 0 \quad \forall x \in X$$

Prove the following statements:

- First $\forall x \in X$

$$\langle x | (H_k^{\mathcal{EE}} - \mu)x \rangle \geq \mathcal{E} - \|V_k^{\mathcal{E}^\perp \mathcal{E}^\perp}\|_{\text{op}} - \mu,$$

where $\|\cdot\|_{\text{op}}$ is the standard Hilbert operator norm

$$\|\mathcal{A}\|_{\text{op}} = \sup_{0 \neq \varphi \in L^2_{\text{per}}(\Omega)} \frac{\langle \varphi, \mathcal{A}\varphi \rangle}{\langle \varphi, \varphi \rangle}$$

- Second, given a full eigendecomposition $H_k^{\mathcal{EE}} = \tilde{X}_k \tilde{\Lambda}_k \tilde{X}_k^H$ show that for $\mu \in I_n$ with the open interval

$$I_n = \left(\frac{1}{2} (\tilde{\lambda}_{k,n-1} + \tilde{\lambda}_{kn}), \tilde{\lambda}_{kn} \right)$$

we have $\forall x \in X$:

$$\begin{aligned} \langle x | V_k^{\mathcal{E}^\perp \mathcal{E}} (H_k^{\mathcal{EE}} - \mu)^{-1} V_k^{\mathcal{E} \mathcal{E}^\perp} x \rangle &\leq \left\| (V_k^{\mathcal{RE}} \tilde{X}_k) (\tilde{\Lambda}_k - \mu)^{-1} (V_k^{\mathcal{RE}} \tilde{X}_k)^H \right\|_{\text{op}} \\ &\leq \frac{\|V_k\|_{\text{op}}^2}{\tilde{\lambda}_{kn} - \mu} \end{aligned}$$

- Combining both show that

$$S_\mu \geq \mathcal{E} - \mu - l_1^V - \frac{(l_1^V)^2}{\tilde{\lambda}_{kn} - \mu} \quad \text{where} \quad l_1^V = \sum_{|G| < \sqrt{2\mathcal{E}_V}} |\hat{V}(G)|. \quad (4)$$

Selection deleted

You may want to use Young's inequality, that is

$$\|V_k\|_{\text{op}} \leq \sum_{G \in \mathbb{L}^*} |\hat{V}(G)|.$$

Solution (c: 1):

First of all

$$H_k^{\mathcal{E}^\perp \mathcal{E}^\perp} = Q_k^\mathcal{E} H_k Q_k^\mathcal{E} = Q_k^\mathcal{E} T_k Q_k^\mathcal{E} + Q_k^\mathcal{E} V Q_k^\mathcal{E}$$

where the elements of the kinetic term are defined as the following

$$(T_k^{\mathcal{E}^\perp \mathcal{E}^\perp})_{GG'} = \langle e_G | T_k e_{G'} \rangle = \frac{1}{2} \delta_{GG'} |G + k|^2$$

It means that $T_k^{\mathcal{E}^\perp \mathcal{E}^\perp}$ is diagonal and we have:

$$\langle x | T_k^{\mathcal{E}^\perp \mathcal{E}^\perp} x \rangle = \frac{1}{2} \delta_{GG'} |G + k|^2 \langle x | x \rangle \geq \mathcal{E} \langle x | x \rangle.$$

On the other hand,

$$|\langle x | V_k^{\mathcal{E}^\perp \mathcal{E}^\perp} x \rangle| \leq \|x\| \|V_k^{\mathcal{E}^\perp \mathcal{E}^\perp} x\|$$

so that

$$-\|x\| \|V_k^{\mathcal{E}^\perp \mathcal{E}^\perp} x\| \leq \langle x | V_k^{\mathcal{E}^\perp \mathcal{E}^\perp} x \rangle \leq \|x\| \|V_k^{\mathcal{E}^\perp \mathcal{E}^\perp} x\|$$

Therefore,

$$\langle x | V_k^{\mathcal{E}^\perp \mathcal{E}^\perp} x \rangle \geq -\|x\| \|V_k^{\mathcal{E}^\perp \mathcal{E}^\perp} x\| \geq -\|x\| \|V_k^{\mathcal{E}^\perp \mathcal{E}^\perp}\|_{\text{op}} \|x\| = -\|V_k^{\mathcal{E}^\perp \mathcal{E}^\perp}\|_{\text{op}} \langle x | x \rangle$$

In the end, putting it all together and using Cauchy–Schwarz inequality, we have

$$\langle x | (H_k^{\mathcal{E}^\perp \mathcal{E}^\perp} - \mu) x \rangle \geq (\mathcal{E} - \|V_k^{\mathcal{E}^\perp \mathcal{E}^\perp}\|_{\text{op}} - \mu) \langle x | x \rangle \geq (\mathcal{E} - \|V_k^{\mathcal{E}^\perp \mathcal{E}^\perp}\|_{\text{op}} - \mu)$$

since vector x are normalized.

Selection deleted

Solution (c: 2):

Using a full eigendecomposition $H_k^{\mathcal{EE}} = \tilde{X}_k \tilde{\Lambda}_k \tilde{X}_k^H$ we can express $(H_k^{\mathcal{EE}} - \mu)^{-1}$ as following

$$(H_k^{\mathcal{EE}} - \mu)^{-1} = \tilde{X}_k (\tilde{\Lambda}_k - \mu)^{-1} \tilde{X}_k^H,$$

and applying the same logic as above the inner product can be bounded above by the operator norm:

$$\langle x | V_k^{\mathcal{E}^\perp \mathcal{E}} (H_k^{\mathcal{EE}} - \mu)^{-1} V_k^{\mathcal{E} \mathcal{E}^\perp} x \rangle \leq \|V_k^{\mathcal{E}^\perp \mathcal{E}} \tilde{X}_k (\tilde{\Lambda}_k - \mu)^{-1} \tilde{X}_k^H V_k^{\mathcal{E} \mathcal{E}^\perp}\|_{\text{op}}$$

Since we work with Cohen-Bergstresser Hamiltonians and \tilde{X}_k are unitary, applying the sub-multiplicative property of the operator norm we get:

$$\begin{aligned} \|V_k^{\mathcal{E}^\perp \mathcal{E}} \tilde{X}_k (\tilde{\Lambda}_k - \mu)^{-1} \tilde{X}_k^H V_k^{\mathcal{E} \mathcal{E}^\perp}\|_{\text{op}} &= \| (V_k^{\mathcal{RE}} \tilde{X}_k) (\tilde{\Lambda}_k - \mu)^{-1} (V_k^{\mathcal{RE}} \tilde{X}_k)^H \|_{\text{op}} \leq \\ &\leq \|V_k^{\mathcal{RE}}\|_{\text{op}} \|(\tilde{\Lambda}_k - \mu)^{-1}\|_{\text{op}} \|V_k^{\mathcal{RE}}\|_{\text{op}} \end{aligned}$$

Since $V_k^{\mathcal{RE}}$ is part of the Hamiltonian and the operator norm of the $H_k^{\mathcal{EE}}$ is bounded by its spectral radius:

$$\|V_k^{\mathcal{RE}}\|_{\text{op}} \leq \|H_k^{\mathcal{EE}}\|_{\text{op}}$$

Finally, considering the definition of the interval I_n , we can conclude that $\mu < \tilde{\lambda}_{kn}$. Therefore,

$$\| (V_k^{\mathcal{RE}} \tilde{X}_k) (\tilde{\Lambda}_k - \mu)^{-1} (V_k^{\mathcal{RE}} \tilde{X}_k)^H \|_{\text{op}} \leq \frac{\|V_k\|_{\text{op}}^2}{\tilde{\lambda}_{kn} - \mu}$$

Selection deleted

Solution (c: 3):

We will start with the Young's inequality for the operator norm:

$$\|V_k\|_{\text{op}} \leq \sum_{|G| < \sqrt{2\mathcal{E}_V}} |\hat{V}(G)| = l_1^V.$$

Then for the first inequality:

$$\left\langle x \mid (H_k^{\mathcal{E}^\perp \mathcal{E}^\perp} - \mu)x \right\rangle \geq \mathcal{E} - \|V_k^{\mathcal{E}^\perp \mathcal{E}^\perp}\|_{\text{op}} - \mu \geq \mathcal{E} - l_1^V - \mu$$

For the second inequality:

$$\left\langle x \mid V_k^{\mathcal{E}^\perp \mathcal{E}} (H_k^{\mathcal{E} \mathcal{E}} - \mu)^{-1} V_k^{\mathcal{E} \mathcal{E}^\perp} x \right\rangle \leq \frac{\|V_k\|_{\text{op}}^2}{\tilde{\lambda}_{kn} - \mu} \leq \frac{(l_1^V)^2}{\tilde{\lambda}_{kn} - \mu}$$

Finally,

$$\begin{aligned} \langle x, S_\mu x \rangle &= \left\langle x, \left((H_k^{\mathcal{E}^\perp \mathcal{E}^\perp} - \mu) - V_k^{\mathcal{E}^\perp \mathcal{E}} (H_k^{\mathcal{E} \mathcal{E}} - \mu)^{-1} V_k^{\mathcal{E} \mathcal{E}^\perp} \right) x \right\rangle = \\ &= \left\langle x, (H_k^{\mathcal{E}^\perp \mathcal{E}^\perp} - \mu) x \right\rangle - \left\langle x, V_k^{\mathcal{E}^\perp \mathcal{E}} (H_k^{\mathcal{E} \mathcal{E}} - \mu)^{-1} V_k^{\mathcal{E} \mathcal{E}^\perp} x \right\rangle. \\ &\geq \mathcal{E} - l_1^V - \mu - \frac{(l_1^V)^2}{\tilde{\lambda}_{kn} - \mu} \end{aligned}$$

Selection deleted

The strategy to determine μ is thus as follows:

1. Compute I_1^V . For this we need a Fourier representation of V . A real space representation is available as `Vreal = DFTK.total_local_potential(ham)` where `ham` is a `DFTK.Hamiltonian`, which can be transformed into Fourier space using a fast-fourier transform `fft(basis, kpoint, Vreal)` where `basis` is the `ham.basis` and `kpoint` is the `Kpoint` object matching the currently processed \mathbf{k} . Note that this assumes that $\mathcal{E}_V < \mathcal{E}$.
2. Using the lower bound (4) find the largest $\mu \in I_n$ for which it can be ensured that $S_\mu \geq 0$. Since we are not using the exact S_μ here, but only a lower bound, this step can fail (i.e. both possible solutions μ are outside of I_n). In this case we are unable to obtain a lower bound for $\tilde{\lambda}_{kn}$ and thus unable to obtain a Kato-Temple estimate.
3. If 2. goes through we set $\mu_{kn} = \mu$ and thus have guaranteed bounds $\mu_{kn} \leq \lambda_{kn} \leq \tilde{\lambda}_{kn}$.

(d) Focusing on the first two eigenpairs of the Γ -point of the Cohen-Bergstresser model, plot both the guaranteed lower bound μ_{kn} from above procedure as well as the non-guaranteed bound of Task 4 as you increase \mathcal{E} . Take the computation of $\tilde{\lambda}_{nk}$ for $\mathcal{E} = 80$ as the reference. Which bound is sharper? While it should be emphasised that better guaranteed bounds than our development are possible (see for example [HLC2020]), can you comment on the advantages and disadvantages of guaranteed error bounds?

[HLC2020]:

M. F. Herbst, A. Levitt and E. Cancès. Faraday Discuss., **224**, 227 (2020). DOI [10.1039/DoFD00048E](https://doi.org/10.1039/DoFD00048E)

Selection deleted

```
MethodError: no method matching fft(::DFTK.Planewavebasis{Float64, Float64, DFTK.CPU, Array{StaticArraysCore.SVector{3, Int64}, 3}, Array{StaticArraysCore.SVector{3, Float64}, 3}, Vector{StaticArraysCore.SVector{3, Int64}}}, ::DFTK.Kpoint{Float64, Vector{StaticArraysCore.SVector{3, Int64}}}, ::Array{ComplexF64, 4})
Closest candidates are:
fft(::DFTK.Planewavebasis, ::DFTK.Kpoint, !Matched::AbstractArray{T, 3} where T; kwargs...)
@ DFTK C:\Users\User\.julia\packages\DFTK\HG5Ae\src\fft.jl:114
fft(!Matched::AbstractArray{<:Real}, ::Any)
@ AbstractFFTs C:\Users\User\.julia\packages\AbstractFFTs\4iQz5\src\definitions.jl:214
fft(!Matched::AbstractArray{<:Complex{<:Union{Integer, Rational}}}, ::Any)
@ AbstractFFTs C:\Users\User\.julia\packages\AbstractFFTs\4iQz5\src\definitions.jl:216
...
1. top-level scope @ [ Local: 8 [inlined]
```

```
1 begin
2     local Ecut = 5
3     local basis = Planewavebasis(model; Ecut=Ecut, kgrid=(1, 1, 1))
4     local ham = Hamiltonian(basis)
5     local Vreal = ComplexF64.(DFTK.total_local_potential(ham))
6
7     kpoint = basis.kpoints[1]
8     Vfourier = fft(ham.basis, kpoint, Vreal); #kpoint,
9
10    l1_V = sum(abs.(Vfourier))
11
12    # eigres = diagonalize_all_kblocks(DFTK.lobpcg_hyper, ham, n_bands=6)
13    # λkn = eigres.λ[1][1:2]
14
15 end
```

Task 6: Band structure with guaranteed bounds

Based on the results so far compute a band structure with guaranteed error bars for $\mathcal{E} = 7$ a k -point density of 15 and the 6 first bands (i.e. the 6 lowest eigenpairs). For each eigenvalue:

- Estimate the **algorithm** and **arithmetic error** by re-computing the in-basis residual $P_k^{\mathcal{E}} r_{kn}$ using Double64. For this you need to build a basis and a Hamiltonian that employes Double64 as the working precision. To change DFTK's internal working precision follow the [Arbitrary floating-point types](#) documentation page. From a basis_double64 using Double64 precision a ~~Selection deleted~~ corresponding Hamiltonian is obtained by Hamiltonian(basis_double64) as shown before.
Hint: Before computing products Hamiltonian(basis_double64) * x, ensure that x has been converted to Double64 as well.
- Estimate the **discretisation error** using both the Bauer-Fike estimates of Tasks 2 & 3 as well as the Kato-Temple estimates of Task 5. In the band structure annotate the tightest bound that is available to you.

1 Enter cell code...

Selection deleted

UndefVarError: `basisD64` not defined

1. top-level scope @ Local: 11

```

1 begin
2
3     λ_conv_double64=[]
4     X_conv_double64=[]
5     res_norm_double64=[]
6     ham_list_double64=[]
7     ρ_bauer_fike_double64 = Vector{Vector{Float64}}(){}
8     println(ρ_bauer_fike_double64)
9
10
11    ham_double64 = Hamiltonian(basisD64)
12    eigres_double64 = diagonalize_all_kblocks(DFTK.lobpcg_hyper, ham_double64, 6)
13
14    #large_base3=basis_change_Ecut(band_data_double64.basis,20)
15    #ham_large3=Hamiltonian(large_base3)
16    #eigres_large_3 = diagonalize_all_kblocks(DFTK.lobpcg_hyper, ham_large3, 6)
17
18
19    for (ik, kpt) in enumerate(bands_double64.basis.kpoints)
20
21
22
23
24        hamk = ham_double64[ik]
25        #hamk_large= ham_large3[ik]
26
27        λk = eigres_double64.λ[ik]
28        Xk = eigres_double64.X[ik]
29
30        #λk_large = eigres_large_3.λ[ik]
31        #Xk_large = eigres_large_3.X[ik]
32
33        residual_k = hamk * Xk - Xk * Diagonal(λk)
34        column_norms = [norm(residual_k[:,i]) for i in 1:6]
35
36        print(size(residual_k))
37        println(norm.(residual_k))
38        println(typeof(norm.(residual_k)))
39        println(size(norm.(residual_k)))
40
41        push!(ham_list_double64,hamk)
Selection deleted
42        push!(ρ_bauer_fike_double64,column_norms)
43        push!(λ_conv_double64,λk)
44        push!(X_conv_double64,Xk)
45        push!(res_norm_double64,norm(residual_k))
46    end
47    println(ρ_bauer_fike_double64)
48
49 end

```



Adaptive diagonalisation techniques

The main purpose of *a posteriori* error bounds as we have developed in the previous tasks is to understand the error in the current computational setup. A good error bound is sharp, cheap to compute and ideally guaranteed, i.e. the actual error is always smaller. One can thus think of the error bars as safety checks. If one is not yet satisfied with the current quality of results, one can always choose a more accurate numerical setup and re-compute. For example if the estimated discretisation error is too large, one just increases \mathcal{E} until one is satisfied.

As simple as this procedure is, it has one notable drawback: It leads to a nested set of loops, where for each attempted \mathcal{E} (outer loop), we iteratively diagonalise the H_k using LOBPCG (inner loop). A natural question is thus, whether one can fuse these two loops, i.e. adapt \mathcal{E} while the LOBPCG is converging, such that at each step the algorithm error and discretization error are balanced.

To make this idea more clear, let us recall a basic LOBPCG implementation, which in this case only performs a diagonalisation of the Hamiltonian fiber at the gamma point of a passed model:

Selection deleted

nonadaptive_lobpcg (generic function with 1 method)

```

1  function nonadaptive_lobpcg(model::Model{T}, Ecut, n_bands;
2                                maxiter=100, tol=1e-6, verbose=false) where {T}
3      kgrid = (1, 1, 1) #  $\Gamma$  point only
4      basis = PlaneWaveBasis(model; Ecut, kgrid)
5      ham = Hamiltonian(basis)
6      hamk = ham[1] # Select  $\Gamma$  point
7      prec = PreconditionerTPA(hamk) # Initialise preconditioner
8      X = DFTK.random_orbitals(hamk.basis, hamk.kpoint, n_bands)
9
10     converged = false
11     λ = NaN
12     residual_norms = NaN
13     residual_history = []
14
15     P = zero(X)
16     R = zero(X)
17     for i in 1:maxiter
18         if i > 1
19             Z = hcat(X, P, R)
20         else
21             Z = X
22         end
23         Z = Matrix(qr(Z).Q) # QR-based orthogonalisation
24
25         # Rayleigh-Ritz
26         HZ = hamk * Z
27         λ, Y = eigen(Hermitian(Z' * HZ))
28         λ = λ[1:n_bands]
29         Y = Y[:, 1:n_bands]
30         new_X = Z * Y
31
32         # Compute residuals and convergence check
33         R = HZ * Y - new_X * Diagonal(λ)
34         residual_norms = norm.(eachcol(R))
35         push!(residual_history, residual_norms)
36         verbose && @printf "%3i %8.4g %8.4g\n" i λ[end] residual_norms[end]
37         if maximum(residual_norms) < tol
38             converged = true
39             X .= new_X
40             break
41         end
42
43         # Precondition and update
Selection deleted
44         DFTK.precondprep!(prec, X)
45         ldiv!(prec, R)
46         P .= X - new_X
47         X .= new_X
48
49         # Additional step:
50         # Move to larger basis ?
51     end
52
53     (; λ, X, basis, ham, converged, residual_norms, residual_history)

```

This implementation can be easily compared to the `DFTK.lobpcg_hyper` function of DFTK. For example:

```
[-6.4615e-14, -2.65676e-13, -8.01026e-14, -6.45428e-13, -1.68643e-13, -4.84579e-12]
```

```

1 let
2   Ecut      = 10
3   n_bands   = 6
4   basis     = PlaneWaveBasis(model; Ecut=10, kgrid=(1, 1, 1))
5   ham       = Hamiltonian(basis)
6   res_dftk = diagonalize_all_kblocks(DFTK.lobpcg_hyper, ham, n_bands)
7   res_dftk_λ_Γ = res_dftk.λ[1] # Extract eigenvalues at Γ point
8
9   # Notice: This function always only computes one k-Point
10  res_nonadaptive = nonadaptive_lobpcg(model, Ecut, n_bands)
11
12  res_nonadaptive.λ - res_dftk_λ_Γ
13 end

```

In order to make this routine discretisation-adaptive we need to add an additional step as indicated in the above source code: After we have checked for convergence and as we are updating `X`, `P` and `R` to their new values, we might additionally want to increase the discretisation basis to employ a new, increased cutoff value $\mathcal{E} + \Delta$ to discretise the Hamiltonian. Using

```
transfer_blochwave(X, basis_small, basis_small.kpoints[1],
                     basis_large, basis_large.kpoints[1])
```

we can then transfer the state vectors `X`, `P` and `R` to this basis.

This means while the iterative eigensolver is converging we also improve the discretisation basis by successively increasing \mathcal{E} . Each LOBPCG iteration may thus run using a different discretisation basis.

Selection deleted

Task 7: Developing a discretisation-adaptive LOBPCG

Assume the LOBPCG currently runs at cutoff \mathcal{F} . Due to the properties of the Cohen-Bergstresser Hamiltonian (Task 2 (b)) there is a smaller cutoff $\mathcal{E} < \mathcal{F}$ such that (Task 5 (a))

$$Q_k^{\mathcal{F}} H P_k^{\mathcal{E}} = 0.$$

(a) Prove that if $\|Q_k^{\mathcal{E}} \tilde{X}_{kn}\| = 0$, then

$$|\lambda_{kn} - \tilde{\lambda}_{kn}| \leq \|P_k^{\mathcal{F}} r_{kn}\|,$$

i.e. that the total error is driven to zero as the iterative diagonalisation employing a cutoff \mathcal{F} is converging. Based on this argue why a large value for $\|Q_k^{\mathcal{E}} \tilde{X}_{kn}\|$ relative to $\|P_k^{\mathcal{F}} r_{kn}\|$ provides a reasonable indicator to decide when to refine the discretisation.

(b) Following the idea of (a) code up an adaptive LOBPCG by monitoring the ratio

$$\frac{\|Q_k^{\mathcal{E}} \tilde{X}_{kn}\|}{\|P_k^{\mathcal{F}} r_{kn}\|} \quad (5).$$

Note, that one way to compute $Q_k^{\mathcal{E}} \tilde{X}_{kn}$ is as $\tilde{X}_{kn} - P_k^{\mathcal{E}} \tilde{X}_{kn}$ where $P_k^{\mathcal{E}}$ can be obtained by transferring from the large basis to the small basis and back to the large basis, i.e.

```
X_small = transfer_blochwave(X_large,
                               basis_large, basis_large.kpoints[1],
                               basis_small, basis_small.kpoints[1])
P_X      = transfer_blochwave(X_small,
                               basis_small, basis_small.kpoints[1],
                               basis_large, basis_large.kpoints[1])
```

Whenever (5) becomes too large, you should switch to a finer discretisation basis, e.g. by switching from a cutoff \mathcal{E} to $\mathcal{E} + \Delta$. Experiment with this setup to find what are good ratios (5) to indicate switching and good values for Δ . Some strategy for exploration:

- Assume we eventually want to target a cutoff $\mathcal{F}^{\text{final}} = 80$. Save the convergence history for running the full LOBPCG at exactly this cutoff. This is your reference.
- With your adaptive methods you should try to loose as little as possible in the rate of convergence compared to the reference convergence profile, but try to use as small values of \mathcal{F} for as many iterations as you can.
- To develop your approach, first only compute a single eigenvalue (`n_bands = 1`) and then start considering more than one. Note, that you will need to adapt (5) and take appropriate maxima / minima over all computed eigenpairs to decide when to switch to the next bigger basis.
- Be creative and experiment. In this Task there is no "best" solution.

Solution (a):

Using the definition of \mathcal{F} from task 2 (b), triangular inequality, and assuming $\tilde{\mathbf{X}}_{kn}$ to be normalized we have

$$\begin{aligned} |\lambda_{kn} - \tilde{\lambda}_{kn}| &= |\lambda_{kn} - \tilde{\lambda}_{kn}| \|\tilde{\mathbf{X}}_{kn}\| \leq \\ &\leq \|r_{kn}\| + \|H_k \tilde{\mathbf{X}}_{kn} - \lambda_{kn} \tilde{\mathbf{X}}_{kn}\| \approx \\ &\approx \|r_{kn}\| = \|P_k^{\mathcal{F}} r_{kn}\|, \end{aligned}$$

where

$$r_{kn} = H_k \tilde{\mathbf{X}}_{kn} - \tilde{\lambda}_{kn} \tilde{\mathbf{X}}_{kn}.$$

If the ratio (5) is large it indicates that the approximate eigenvector $\tilde{\mathbf{X}}_{kn}$ has significant components outside the subspace spanned by $P_k^{\mathcal{E}}$ meaning that the chosen cutoff \mathcal{E} may be insufficient to capture the eigenvector.

Solution (b):

Selection deleted

adaptive_lobpcg (generic function with 1 method)

```

1  function adaptive_lobpcg(model::Model{T}, Ecut, n_bands;
2                                maxiter=100, tol=1e-6, verbose=false, threshold=10)
3      where {T}
4          kgrid = (1, 1, 1) #  $\Gamma$  point only
5          basis = PlaneWaveBasis(model; Ecut, kgrid)
6          ham = Hamiltonian(basis)
7          hamk = ham[1] # Select  $\Gamma$  point
8          prec = PreconditionerTPA(hamk) # Initialise preconditioner
9          X = DFTK.random_orbitals(hamk.basis, hamk.kpoint, n_bands)
10
11     converged = false
12     λ = NaN
13     residual_norms = NaN
14     residual_history = []
15     Ecut_new = Ecut
16
17     basis_prev = basis
18
19     P = zero(X)
20     R = zero(X)
21     for i in 1:maxiter
22         if i > 1
23             Z = hcat(X, P, R)
24         else
25             Z = X
26         end
27         Z = Matrix(qr(Z).Q) # QR-based orthogonalisation
28
29         # Rayleigh-Ritz
30         HZ = hamk * Z
31         λ, Y = eigen(Hermitian(Z' * HZ))
32         λ = λ[1:n_bands]
33         Y = Y[:, 1:n_bands]
34         new_X = Z * Y
35
36         # Compute residuals and convergence check
37         R = HZ * Y - new_X * Diagonal(λ)
38         residual_norms = norm.(eachcol(R))
39         push!(residual_history, residual_norms)
40         verbose && @printf "%3i %8.4g %8.4g\n" i λ[end] residual_norms[end]
41         if maximum(residual_norms) < tol
42             converged = true
43             X .= new_X
44             break
45         end
46
47         # Precondition and update
48         DFTK.precondprep!(prec, X)
49         ldiv!(prec, R)
50         P .= X - new_X
51         X .= new_X
52

```

```

53      # Additional step:
54      basis_small = basis_change_Ecut(basis_prev, Ecut_new - 3)
55      X_E = transfer_blochwave_kpt(X,
56                                      basis_prev, basis_prev.kpoints[1],
57                                      basis_small, basis_small.kpoints[1])
58      P_X = transfer_blochwave_kpt(X_E,
59                                      basis_small, basis_small.kpoints[1],
60                                      basis_prev, basis_prev.kpoints[1])
61
62
63      # Computing the ratio
64      ratio = norm(X - P_X) / norm(residual_norms)
65
66      # Moving to a larger basis
67      if ratio > threshold
68
69          Ecut_new = Ecut_new + 2
70          println("i = ", i, ", Ecut = ", Ecut_new)
71          basis_new = PlaneWaveBasis(model; Ecut=Ecut_new, kgrid)
72
73          X = transfer_blochwave_kpt(X, basis_prev, basis_prev.kpoints[1],
74                                      basis_new, basis_new.kpoints[1])
75          P = transfer_blochwave_kpt(P, basis_prev, basis_prev.kpoints[1],
76                                      basis_new, basis_new.kpoints[1])
77          R = transfer_blochwave_kpt(R, basis_prev, basis_prev.kpoints[1],
78                                      basis_new, basis_new.kpoints[1])
79
80          hamk = Hamiltonian(basis_new)[1]
81          prec = PreconditionerTPA(hamk)
82          basis_prev = basis_new
83
84      end
85
86  end
87
88  (; λ, X, basis, ham, converged, residual_norms, residual_history)
89 end

```

```
[-0.079228, 0.384299, 0.384299, 0.384299, 0.510143, 0.510143]
```

```
1 begin
2     local Ecut      = 5
3     local n_bands   = 6
4
5     res_adaptive = adaptive_lobpcg(model, Ecut, n_bands, threshold=0.2)
6     res_adaptive.λ
7 end
```

```
i = 1, Ecut = 7
i = 6, Ecut = 9
i = 9, Ecut = 11
i = 11, Ecut = 13
i = 13, Ecut = 15
i = 20, Ecut = 17
i = 26, Ecut = 19
i = 29, Ecut = 21
i = 31, Ecut = 23
i = 36, Ecut = 25
i = 39, Ecut = 27
i = 41, Ecut = 29
```

(?)

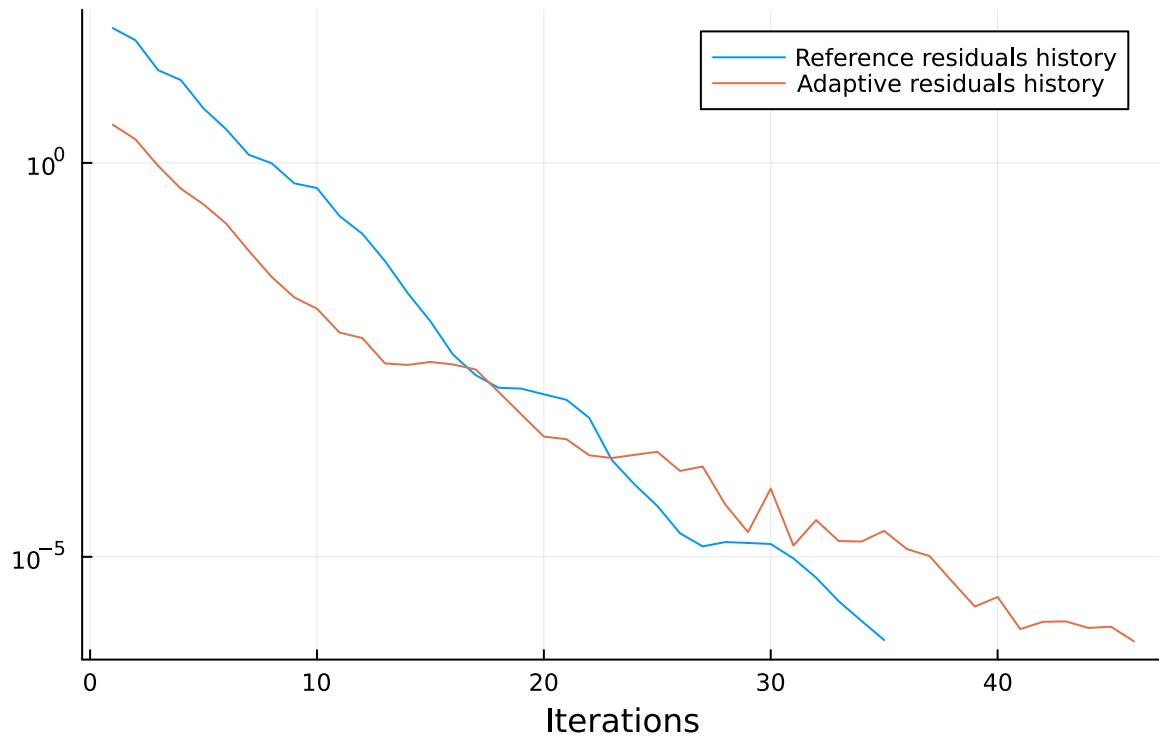
```
[-0.079228, 0.384299, 0.384299, 0.384299, 0.510143, 0.510143]
```

```
1 begin
2     local Ecut      = 80
3     local n_bands   = 6
4
5     res_nonadaptive = nonadaptive_lobpcg(model, Ecut, n_bands)
6     res_nonadaptive.λ
7 end
```

```
[-0.0791827, 0.389126, 0.389126, 0.389126, 0.512338, 0.512338]
```

```
1 begin
2     local Ecut      = 3
3     local n_bands   = 6
4
5     res_nonadaptive5 = nonadaptive_lobpcg(model, Ecut, n_bands)
6     res_nonadaptive5.λ
7 end
```

Rate of convergence



```

1 begin
2   local p = plot(yaxis=:log, xlabel="Iterations")
3   plot!(1:size(res_nonadaptive.residual_history)[1], norm.
4     (res_nonadaptive.residual_history), label="Reference residuals history")
5   plot!(1:size(res_adaptive.residual_history)[1], norm.
6     (res_adaptive.residual_history), label="Adaptive residuals history")
7   title!(p, "Rate of convergence")
8 end

```

Experimenting with various thresholds for the ratio given in equation (5) and values of Δ we can see that the rate of convergence using the adaptive method, with an appropriately set threshold and Δ , can effectively match the reference convergence profile. Moreover, the developed approach has better computational efficiency by using smaller cutoffs during the early steps in the algorithm.

GTH pseudopotentials

To close off this project we will now consider slightly more realistic potentials than the Cohen-Bergstesser model. We will look at the local part of the Goedecker-Teter-Hutter pseudopotentials [GTH96]. Neglecting some minor details involving the structure of the system, these potentials can be understood in Fourier space as

$$\hat{V}(p) = P(p) \exp(-p^2/\sigma^2),$$

i.e. a strongly decaying Gaussian times a polynomial prefactor. Instead of just a few non-zero Fourier coefficients like in the Cohen-Bergtresser case, we thus have the weaker setting, that the Fourier coefficients only strongly decay as $|G| \rightarrow \infty$. Unlike our development in Task 2 (b), we thus cannot find an elevated cutoff $\mathcal{F} > \mathcal{E}$, such that the residual can be computed exactly. However, due to the strong decay the truncated residual $P_k^{\mathcal{F}} r_{kn}$ of Tasks 2 becomes a better and better approximation to r_{kn} the larger \mathcal{E} and \mathcal{F} are taken. Thus the ideas of Tasks 2 can still be employed to compute *approximate* errors for this GTH potential — though without any guarantees that the actual error is smaller than these estimates.

[GTH96]:

S. Goedecker, M. Teter, and J. Hutter Phys. Rev. B **54**, 1703 (1996). DOI [10.1103/PhysRevB.54.1703](https://doi.org/10.1103/PhysRevB.54.1703)

Task 8: An error indicator for GTH pseudopotentials

A setup for the GTH pseudopotential model, discretised for a provided cutoff \mathcal{E} and using only the Γ -point is given by the function

`make_gth_basis` (generic function with 1 method)

```

1 function make_gth_basis(Ecut)
2     Si = ElementPsp(:Si, psp=load_psp("hgh/lda/si-q4"))
3     atoms = [Si, Si]
4     positions = [ones(3)/8, -ones(3)/8]
5     lattice = 5 .* [[0 1 1.]; [1 0 1.]; [1 1 0.]]
6     model = Model(lattice, atoms, positions; terms=[Kinetic(), AtomicLocal()])
7     PlaneWaveBasis(model; Ecut, kgrid=(1, 1, 1))
8 end

```

Assume a secondary cutoff $\mathcal{F} > \mathcal{E}$, which is large enough that the error bounds developed in Task 2 are sufficiently good to be useful. In other words we will assume that the residual computed in $\mathbb{B}_k^{\mathcal{F}}$, i.e.

$$P_k^{\mathcal{F}} r_{kn} = H_k^{\mathcal{FF}} \tilde{X}_{kn} - \tilde{\lambda}_{kn} \tilde{X}_{kn},$$

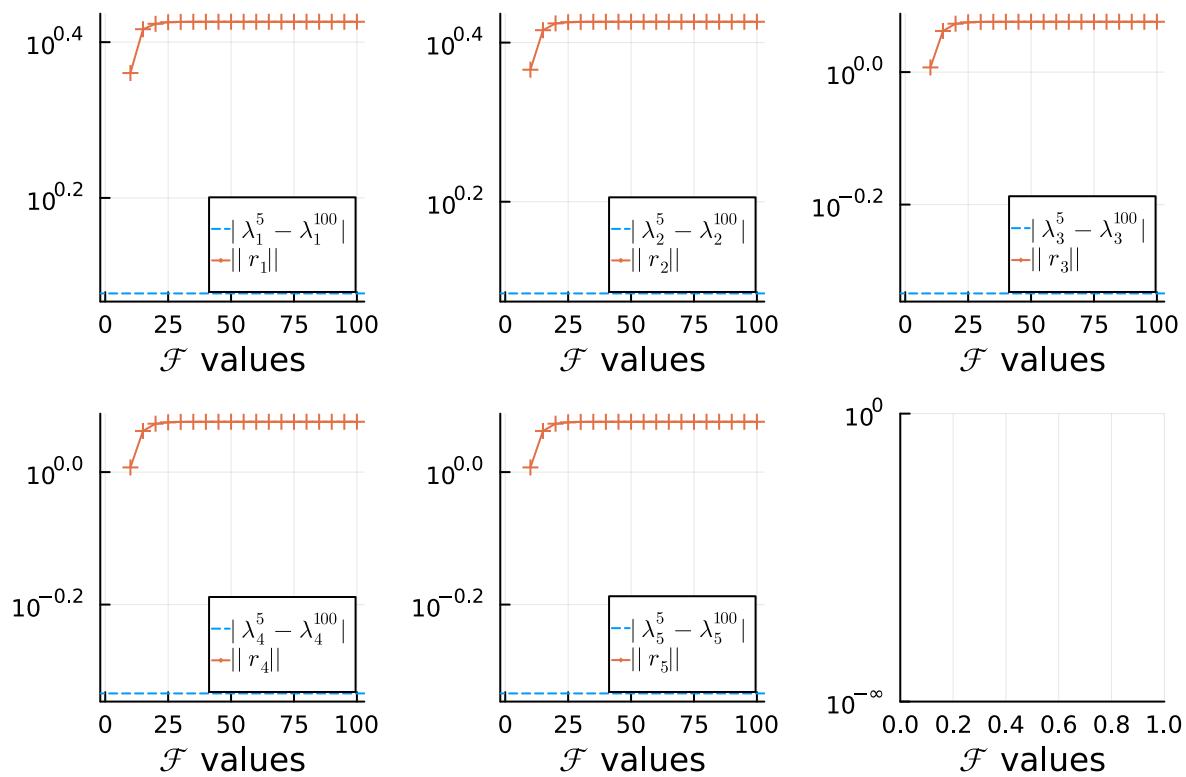
is a good estimate of the true residual.

Focusing on the first **5** eigenvalues of the Γ point of the GTH Hamiltonian, investigate the Bauer-Fike error bounds based on this approximate residual. Investigate in particular:

- (a)** Take $\mathcal{E} = 100$ as a reference and plot the error bound for $\mathcal{E} = 5$ and for varying \mathcal{F} between **10** and **100**. Does the bound always seem to hold? Does it converge with \mathcal{F} ?
- (b)** For a few offsets Δ between **5** and **30** set $\mathcal{F} = \mathcal{E} + \Delta$ and consider the **2nd** and **3rd** eigenvalue. Plot both the obtained error estimate as well as the error of the eigenvalue as you vary \mathcal{E} between **5** and **50**. Again take $\mathcal{E} = 100$ as a reference. What offset Δ would you recommend based on this investigation?

Solution (a):

```
1 begin
2     n_bands_8a = 5
3     local ik = 1
4
5     Fcuts = 10:5:100
6     ρ_BF = zeros(length(Fcuts), n_bands_8a)
7
8     basis_100 = make_gth_basis(100)
9     ham_100 = Hamiltonian(basis_100)
10    eigres_100 = diagonalize_all_kblocks(DFTK.lobpcg_hyper, ham_100, n_bands_8a)
11    λ_100 = eigres_100.λ[ik]
12
13    basis_5 = make_gth_basis(5)
14    ham_5 = Hamiltonian(basis_5)
15    eigres_5 = diagonalize_all_kblocks(DFTK.lobpcg_hyper, ham_5, n_bands_8a)
16    λ_5 = eigres_5.λ[ik]
17
18    r_ref = abs.(λ_5 - λ_100)
19
20    for (i, Fcut) in enumerate(Fcuts)
21        basis = make_gth_basis(Fcut)
22        ham = Hamiltonian(basis)
23        X_F = transfer_blochwave(eigres_5.X, basis_5, basis)
24
25        residual_k = ham[ik] * X_F[ik] - X_F[ik] * Diagonal(eigres_5.λ[ik])
26        ρ_BF[i, :] = norm.(eachcol(residual_k))
27
28    end
29
30 end
31
```



```

1 begin
2
3     local p = plot(yaxis=:log, layout=(2, 3),
4                         empty_at = 6, xlabel=L"\mathcal{F} values")
5
6     for i in 1:n_bands_8a
7
8         hline!([r_ref[i]], subplot=i, line=:dash,
9             label=string(latexstring("|\lambda^5_{\$}(i) - \lambda^{100}_{\$}(i)|")))
10        plot!(p, Fcuts, ρ_BF[:, i], subplot=i, shape=:cross,
11              label=string(latexstring("||r_{\$}(i)||")))
12
13    end
14    plot!(legend=:bottomright)
15 end

```

No strict ticks found

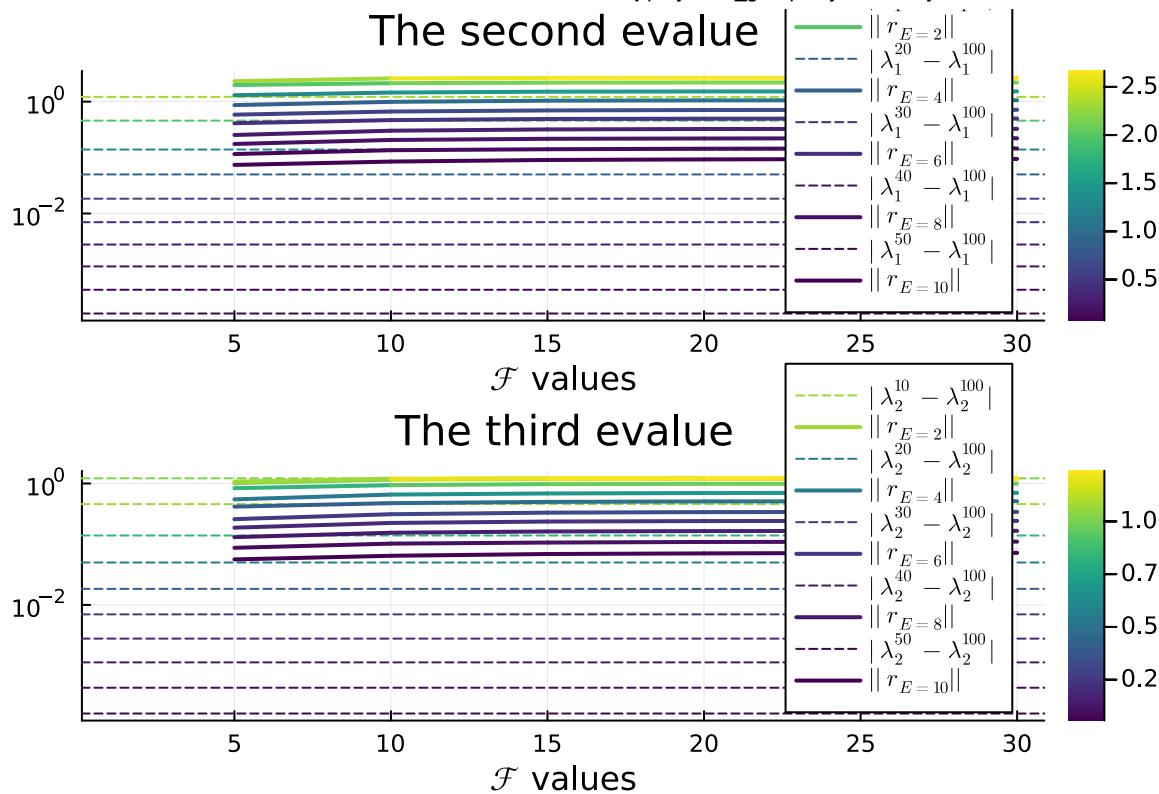
No strict ticks found

Increasing \mathcal{F} we can see that all bounds are converging and always are bigger than the reference value.

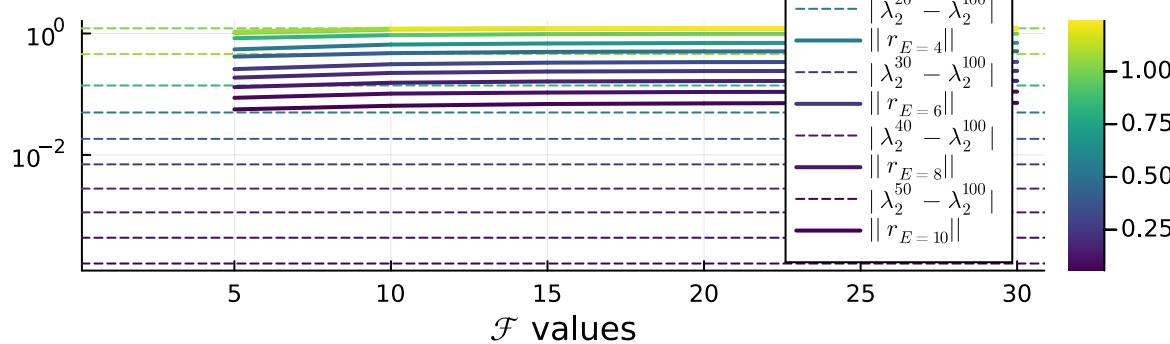
Solution (b):

```
1 begin
2     Delta_range = 5:5:30
3     E_range= 5:5:50
4     local n_bands = 3
5     e_BF      = zeros(length(E_range), length(Delta_range), 2)
6     e_ref     = zeros(length(E_range), 2)
7
8     for (i, E) in enumerate(E_range)
9         basis_E      = make_gth_basis(E)
10        ham_E       = Hamiltonian(basis_E)
11        eigres_E    = diagonalize_all_kblocks(DFTK.lobpcg_hyper, ham_E, n_bands)
12
13        e_ref[i, :] = abs.(eigres_E.λ[ik][2:3] - λ_100[2:3])
14
15        for (j, Δ) in enumerate(Delta_range)
16            F = E + Δ
17            basis      = make_gth_basis(F)
18            ham       = Hamiltonian(basis)
19            X_F       = transfer_blochwave(eigres_E.X, basis_E, basis)
20
21            res_k = ham[ik] * X_F[ik] - X_F[ik] * Diagonal(eigres_E.λ[ik])
22            e_BF[i, j, :] = norm.(eachcol(res_k))[2:3]
23        end
24    end
25 end
```

The second eval



The third eval



```

1 begin
2     gradient = [:blue, :red]
3     local p = plot(yaxis=:log, layout=(2, 1), xlabel=L"\$\\mathcal{F} \$ values")
4
5     for i in 1:2
6
7         for (j, E) in enumerate(E_range)
8
9             if E % 10 == 0
10                label_h = string(latexstring("|\lambda^{$(E)}_{-$(i)} - \lambda^{100}_{-$(i)}|"))
11                label_r = string(latexstring("||r_{-E = $(j)}||"))
12            else
13                label_h = ""
14                label_r = ""
15            end
16
17            hline!([e_ref[j]], subplot=i, line=:dash,
18                  line_z=e_BF[j, :, i],
19                  color=cgrad(:viridis),
20                  label=label_h
21        )
22
23            plot!(p, Delta_range, e_BF[j, :, i],
24                  subplot=i, linewidth = 2,
25                  line_z=e_BF[j, :, i],
26                  color=cgrad(:viridis),
27                  label=label_r
28        )
29    end
30 end
31
32 plot!(legend=:bottomright)

```

```
33     plot!(title="The second evalve", subplot=1)
34     plot!(title="The third evalve", subplot=2)
35
36 end
```