

# Projectability disentanglement for accurate and automated electronic-structure Hamiltonians

Junfeng Qiao\*

*Theory and Simulations of Materials (THEOS), and National Centre for Computational Design and Discovery of Novel Materials (MARVEL), École Polytechnique Fédérale de Lausanne, 1015 Lausanne, Switzerland*

Giovanni Pizzi and Nicola Marzari

*Theory and Simulations of Materials (THEOS), and National Centre for Computational Design and Discovery of Novel Materials (MARVEL), École Polytechnique Fédérale de Lausanne, 1015 Lausanne, Switzerland and Laboratory for Materials Simulations (LMS), Paul Scherrer Institut (PSI), CH-5232 Villigen PSI, Switzerland*

(Dated: July 6, 2023)

Maximally-localized Wannier functions (MLWFs) are a powerful and broadly used tool to characterize the electronic structure of materials, from chemical bonding to dielectric response to topological properties. Most generally, one can construct MLWFs that describe isolated band manifolds, e.g. for the valence bands of insulators, or entangled band manifolds, e.g. in metals or describing both the valence and the conduction manifolds in insulators. Obtaining MLWFs that describe a target manifold accurately and with the most compact representation often requires chemical intuition and trial and error, a challenging step even for experienced researchers and a roadblock for automated high-throughput calculations. Here, we present a powerful approach that automatically provides MLWFs spanning the occupied bands and their natural complement for the empty states, resulting in Wannier Hamiltonian models that provide a tight-binding picture of optimized atomic orbitals in crystals. Key to the success of the algorithm is the introduction of a projectability measure for each Bloch state onto atomic orbitals (here, chosen from the pseudopotential projectors) that determines if that state should be kept identically, discarded, or mixed into a disentangling algorithm. We showcase the accuracy of our method by comparing a reference test set of 200 materials against the selected-columns-of-the-density-matrix algorithm, and its reliability by constructing Wannier Hamiltonians for 21,737 materials from the Materials Cloud.

## I. INTRODUCTION

In periodic crystals, the electronic structure is usually described using one-particle Bloch wavefunctions. While choosing a basis set that is also periodic to describe these wavefunctions can often be beneficial, an alternative approach is to adopt localized orbitals in real space. One such choice of orbitals are Wannier functions (WFs), that can be obtained by Fourier transforming the periodic wavefunctions from reciprocal to real space. WFs are not unique, as they depend on the choice of the gauge (i.e., the choice of the phases of the wavefunctions) at each point in the Brillouin zone (BZ). Maximally-localized Wannier functions (MLWFs) [1–4] are obtained by a gauge choice that is optimized to provide the most localized set of WFs, i.e., those that minimize the sum of their quadratic spread in real space [1]. Having a very localized representation of the electronic structure not only provides an insightful analysis of chemical bonding in solids, but also brings a formal connection between the MLWF centers and the modern theory of electric polarization [5]. Moreover, the real-space locality of MLWF allows for accurate and fast interpolation of physical operators [6], enabling calculations of material properties that require dense samplings of the BZ, such as Fermi

surface, orbital magnetization [7], anomalous Hall conductivity [8, 9], and spin Hall conductivity [10], to name a few. Practically, one obtains MLWFs starting from a set of Bloch wavefunctions, calculated e.g., from density-functional theory (DFT). Often, these Bloch states are projected onto some localized orbitals (usually chosen by the user) to generate initial guesses for MLWFs. In an insulator, by minimizing the spread functional [1] which measures localization, one can obtain a set of MLWFs, i.e., “Wannierize” a material. The Wannierization contains an additional disentanglement step [2] if the target Bloch states are not isolated from other band manifolds. For such entangled bands—metals or the conduction bands of insulators—one needs to first identify the relevant Bloch states that will be used to construct MLWFs, and then mix or “disentangle” these from all the Bloch states [2]. Practically, the choices for the initial projections and states to be disentangled substantially influence the shape and the quality of the final MLWFs.

In recent years, a lot of effort has been devoted to obtaining high-quality MLWFs and automate the Wannierization procedure. Focus of the research can be categorized into the following classes: (a) Novel minimization algorithms, such as: the symmetry-adapted WF method that adds constraints to impose the symmetries of the resulting WFs [11]; the simultaneous diagonalization algorithm that directly minimizes the spread functional for an isolated (or “ $\Gamma$ -only”) system [12]; the partly-occupied WF method, where the total spread is directly

---

\* junfeng.qiao@epfl.ch

minimized in one step [13, 14], rather than performing a two-step minimization for its gauge-invariant and gauge-dependent parts as in the standard procedure [2]; or the variational formulation, that combines single-step optimization with manifold optimization to make the minimization algorithm more robust [15]; (b) new forms for the spread functional, such as the selectively localized WFs (SLWFs) for which only a subset of WFs of interest are localized and a penalty term is added to constrain the position of the WF centers [16], or the spread-balanced WF method, that adds a penalty term to distribute the spread as uniformly as possible among all WFs [17]; (c) targeting a subset of orbitals, e.g. SLWF for a subset of MLWFs [16] or the optimized projection functions method where starting projections for the Wannierization are generated from a larger group of initial ones [18]; (d) matrix manifold algorithms instead of projection methods to construct a smooth gauge in a non-iterative way [19, 20]; (e) basis-vector decomposition of the density matrix, e.g. the selected columns of the density matrix (SCDM) algorithm [21, 22], that starts from the density matrix of the system and uses QR decomposition with column pivoting (QRCP) to automatically generate an optimal set of basis vectors from the columns of the density matrix.

At the same time, high-throughput (HT) calculations have become increasingly popular for materials discovery and design. Calculations and results managed by workflow engines are collected into databases of original calculations, such as the Materials Project [23], AFLOW [24], OQMD [25], CMR [26], and the Materials Cloud [27], or aggregated, as in NOMAD [28]. Thanks to recent research advances on Wannierization algorithms, it starts now to be possible to run HT Wannierizations for many materials and generate tight-binding (TB) models that reliably describe their physics. So far, several attempts have been made in this direction. Gresch *et al.* [29] gathered 195 Wannier TB Hamiltonians and applied post-processing symmetrization to study strained III-V semiconductor materials. Vitale *et al.* [30] implemented the SCDM algorithm and designed a protocol to determine automatically the remaining free parameters of the algorithm; this protocol, implemented into automated workflows, was verified to work well for band interpolations on a set of 200 structures (metals, or valence and conduction bands of insulators) and 81 insulators (valence bands only). Garrity and Choudhary [31] accumulated a Wannier TB Hamiltonian database of 1771 materials using the standard hydrogenic orbital projections. However, there are still several challenges for an accurate and automated HT Wannierization, some of which might be more relevant depending on the research goal and the specific property to compute: MLWFs should be able to faithfully represent the original band structure, often (e.g., for transport properties) at least for those bands close to the Fermi energy; MLWFs should resemble the physically intuitive atomic orbitals for solids that would enter into Bloch sums; the algorithm should be fully and

reliably automated and the implementation should be efficient for HT calculations.

To overcome the challenges mentioned above, in this paper we present a new methodology for automated Wannierization. First, we choose physically-inspired orbitals as initial projectors for MLWFs, that is, the pseudo-atomic orbitals (PAOs) from pseudopotentials [32]. Then, for each state  $|n\mathbf{k}\rangle$  ( $n$  is the band index,  $\mathbf{k}$  is the Bloch quasi-momentum) we decide if it should be dropped, kept identically, or thrown into the disentanglement algorithm depending on the value of its projectability onto the chosen set of PAOs, replacing the standard disentanglement and frozen manifolds based only on energy windows. This approach naturally and powerfully targets the TB picture of atomic orbitals in crystals, as it will also become apparent from our results. Moreover, we fully automate this approach and implement it in the form of open-source `AiiDA` [33–35] workflows. To assess its effectiveness and precision, we compare the quality of the band interpolation and the locality of the Wannier Hamiltonians generated with the present approach, which we name as projectability-disentangled Wannier functions (PDWFs), with the results from the SCDM algorithm [30]. Statistics from 200 materials demonstrate that PDWFs are more localized and more atomic-like, and the band interpolation is accurate at the meV scale. Furthermore, to demonstrate the reliability and automation of our method and workflows, we carry out a large-scale high-throughput Wannierization of 21,737 materials from the Materials Cloud [27, 36].

To set the context for the following paragraphs, here we briefly summarize the notations for WFs; a detailed description can be found in Refs. [1–3]. WFs  $|w_{n\mathbf{R}}\rangle$  are unitary transformations of Bloch wavefunctions  $|\psi_{m\mathbf{k}}\rangle$ , given by

$$|w_{n\mathbf{R}}\rangle = \frac{V}{(2\pi)^3} \int_{\text{BZ}} d\mathbf{k} e^{-i\mathbf{k}\cdot\mathbf{R}} \sum_{m=1}^{J \text{ or } J_{\mathbf{k}}} |\psi_{m\mathbf{k}}\rangle U_{m\mathbf{n}\mathbf{k}}, \quad (1)$$

where  $\mathbf{k}$  and  $\mathbf{R}$  are the Bloch quasi-momentum in the BZ and a real-space lattice vector, respectively;  $m$  is the band index, and  $n$  is the Wannier-function index (running from 1 to the number of WFs  $J$ ). For an isolated group of bands,  $J$  is equal to the number of bands, and the  $U_{m\mathbf{n}\mathbf{k}}$  are unitary matrices; for entangled bands, the number of bands considered at each  $k$ -point is  $J_{\mathbf{k}} \geq J$ , and the  $U_{m\mathbf{n}\mathbf{k}}$  are semi-unitary rectangular matrices. MLWFs are the minimizers of the quadratic spread functional [1]

$$\Omega = \sum_{n=1}^J \left[ \langle w_{n\mathbf{0}} | \mathbf{r}^2 | w_{n\mathbf{0}} \rangle - |\langle w_{n\mathbf{0}} | \mathbf{r} | w_{n\mathbf{0}} \rangle|^2 \right]. \quad (2)$$

Since Eq. (2) is a minimization problem with multiple local minima, initial guesses for  $U_{m\mathbf{n}\mathbf{k}}$  substantially influence the optimization path and the final minimum obtained. In order to target the most localized and chemically appealing solution, Marzari and Vanderbilt [1] used

hydrogenic wavefunctions  $|g_n\rangle$  (i.e., analytic solutions of the isolated hydrogenic Schrödinger equation) to provide a set of sensible initial guesses  $|\phi_{n\mathbf{k}}\rangle$ , after projection on the space defined by the relevant Bloch states:

$$|\phi_{n\mathbf{k}}\rangle = \sum_{m=1}^{J \text{ or } J_{\mathbf{k}}} |\psi_{m\mathbf{k}}\rangle \langle \psi_{m\mathbf{k}} | g_n \rangle. \quad (3)$$

The projection matrices  $A_{m\mathbf{n}\mathbf{k}} = \langle \psi_{m\mathbf{k}} | g_n \rangle$ , after Löwdin orthonormalization [37], form the initial guesses for  $U_{m\mathbf{n}\mathbf{k}}$ . We underline that while the gauge of Bloch wavefunctions  $|\psi_{m\mathbf{k}}\rangle$  is arbitrary, Eq. (3) is invariant to such gauge freedom: suppose  $|\psi'_{i\mathbf{k}}\rangle$  are also solutions of the electronic structure problem, then  $|\psi'_{i\mathbf{k}}\rangle$  are related to  $|\psi_{m\mathbf{k}}\rangle$  by some unitary matrices  $|\psi'_{i\mathbf{k}}\rangle = \sum_m |\psi_{m\mathbf{k}}\rangle U_{mi\mathbf{k}}$ ; thus  $|\phi_{n\mathbf{k}}\rangle = \sum_m |\psi_{m\mathbf{k}}\rangle \langle \psi_{m\mathbf{k}} | g_n \rangle = \sum_m \sum_i |\psi'_{i\mathbf{k}}\rangle U_{im\mathbf{k}}^* U_{mi\mathbf{k}} \langle \psi'_{i\mathbf{k}} | g_n \rangle = \sum_i |\psi'_{i\mathbf{k}}\rangle \langle \psi'_{i\mathbf{k}} | g_n \rangle$  does not depend on the gauge of Bloch wavefunctions, where superscript  $*$  denotes conjugate transpose. For entangled bands, the “standard” disentanglement approach [2] uses energy windows to choose the disentanglement and frozen manifolds: (a) an (outer) disentanglement window that includes a large set of Bloch states, which can be mixed together to obtain a smaller disentangled manifold; (b) an (inner) frozen window that specifies a smaller set of Bloch states (often states around Fermi energy) which are kept unchanged in the final disentangled manifold.

Since in the following sections the present results are compared with SCDM, we also summarize the SCDM procedure here. The SCDM method [21] starts from the real-space density matrix  $\langle \mathbf{r} | P_{\mathbf{k}} | \mathbf{r}' \rangle$  where  $P_{\mathbf{k}} = \sum_{m=1}^{J_{\mathbf{k}}} |\psi_{m\mathbf{k}}\rangle \langle \psi_{m\mathbf{k}}|$ , and uses QR factorization with column pivoting (QRCP) to decompose  $\langle \mathbf{r} | P_{\mathbf{k}} | \mathbf{r}' \rangle$  into a set of localized real-space orbitals, thanks to the near-sightedness principle [38, 39] stating that the matrix elements  $\langle \mathbf{r} | P_{\mathbf{k}} | \mathbf{r}' \rangle$  decay exponentially with the distance between two points  $\mathbf{r}$  and  $\mathbf{r}'$  in insulating systems. While storing the full  $\langle \mathbf{r} | P_{\mathbf{k}} | \mathbf{r}' \rangle$  is memory intensive (it has size  $N_{\mathbf{r}} \times N_{\mathbf{r}}$ , where  $N_{\mathbf{r}}$  is the number of real-space grid points), one can equivalently decompose the matrix formed by the real-space Bloch wavefunctions  $\Psi_{\mathbf{k}}^* = [\psi_{1\mathbf{k}}, \dots, \psi_{J_{\mathbf{k}}\mathbf{k}}]^*$ , which has a smaller size  $J_{\mathbf{k}} \times N_{\mathbf{r}}$ . For periodic systems, often the choice of columns in the QRCP algorithm can be performed using the wavefunctions at the  $\Gamma$  point only ( $\Psi_{\Gamma}$ ) [40], and the same column selection is then used for all other  $k$ -points. For entangled bands, since the density matrix is not continuous across the  $k$ -points, one can construct a quasi-density matrix (or equivalently a matrix of wavefunctions)  $\sum_{m=1}^{J_{\mathbf{k}}} |\psi_{m\mathbf{k}}\rangle f(\varepsilon_{m\mathbf{k}}) \langle \psi_{m\mathbf{k}}|$ , where  $f(\varepsilon_{m\mathbf{k}})$  is a smooth function of the energy eigenvalues  $\varepsilon_{m\mathbf{k}}$ , specifying the target energy window for the constructed MLWFs. Often the complementary error function  $\frac{1}{2} \operatorname{erfc}(\frac{\varepsilon - \mu}{\sigma})$  is chosen as  $f(\varepsilon)$ , and the choice of  $\mu$  and  $\sigma$  determines the shape of MLWFs, as well as band-interpolation quality. Using projectability, defined later in Eq. (5),  $\mu$  and  $\sigma$  can be automatically chosen, thus automating the Wannier-

ization process [30].

## II. RESULTS AND DISCUSSIONS

### A. Pseudo-atomic-orbital projections

In addition to the hydrogenic orbitals discussed above, alternative starting guesses for the Wannierization can be used. For instance, in pseudopotential plane-wave methods, PAOs are localized orbitals originating from the pseudopotential generation procedure [32]. In this procedure, for each element, atomic wavefunctions of an isolated atom are pseudized to remove the radial nodes and are localized functions around the atom; spherical harmonics with well-defined angular-momentum character ( $s$ ,  $p$ ,  $d$ , or  $f$ ) are chosen for their angular dependency. Then, the PAOs are summed over lattice points with appropriate phases to obtain Bloch sums, Fourier transformed to a plane-wave basis, Löwdin-orthonormalized, and finally taken as the projectors for initial projections. PAOs are commonly used for analyzing the orbital contributions to band structures, as the basis set for non-iterative construction of TB Hamiltonians [32], or as projectors in DFT+Hubbard calculations [41].

In order to understand the contribution of each orbital  $|g_n\rangle$  to a Bloch state  $|\psi_{m\mathbf{k}}\rangle$ , we define a measure of projectability as the square of the inner product between  $|\psi_{m\mathbf{k}}\rangle$  and  $|g_n\rangle$ :

$$p_{nm\mathbf{k}} = |\langle g_n | \psi_{m\mathbf{k}} \rangle|^2; \quad (4)$$

the projectability of  $|\psi_{m\mathbf{k}}\rangle$  onto all PAOs is then defined as

$$p_{m\mathbf{k}} = \sum_n p_{nm\mathbf{k}}. \quad (5)$$

If the projectors  $|g_n\rangle$  are complete for  $|\psi_{m\mathbf{k}}\rangle$ , then  $p_{m\mathbf{k}} = 1$ . The band projectability is a very useful criterion to identify the orbital character of the bands; this is exemplified in Fig. 1a, where we show the projectability of the bands of graphene onto  $2s$  and  $2p$  PAOs for carbon. It is immediately apparent how one can easily identify states in the conduction manifold that have a strong  $2p$  and  $2s$  component.

Compared with the hydrogenic projections, which is the method used by default in `Wannier90` [4] and its interface code to `Quantum ESPRESSO` [42] (called `pw2wannier90.x`), PAOs are better adapted to each element since they come exactly from the pseudopotential used in the actual solid-state calculation. Moreover, in pseudopotentials with semicore states, the PAOs for semicores are nodeless and those for valence wavefunctions have at least one radial node (so as to be orthogonal to the semicore states with same angular momentum); thus band projectability can clearly differentiate semicore from valence, making PAOs more convenient than the hydrogenic orbitals, for which the user would need to

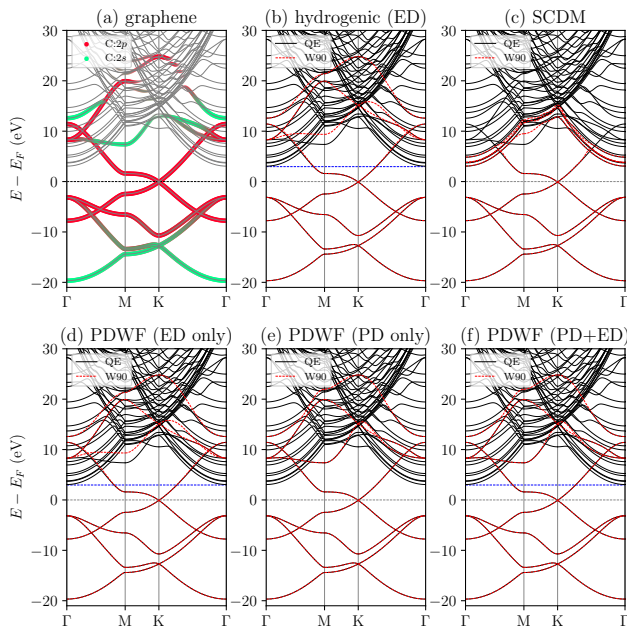


FIG. 1. **Comparisons of graphene band structures interpolated using different methods.** (a) DFT band structure, shown as grey lines. The colored dots represent the projectabilities onto carbon  $2s$  (green) and  $2p$  (red) orbitals. The size of each dot is proportional to the total projectability  $p_{m\mathbf{k}}$  of the band  $m$  at  $k$ -point  $\mathbf{k}$ ; see Eq. (5). For a detailed plot of total projectability, see Fig. S1. Comparisons of the original and the Wannier-interpolated bands for (b) hydrogenic projections with energy disentanglement (ED), (c) SCDM, (d) PAO with ED, (e) PAO with projectability disentanglement (PD), and (f) PAO with PD+ED. The Fermi energy  $E_F$  (horizontal black dashed line) is at zero; the horizontal blue dashed line denotes the top of the inner energy window, where applicable.

manually set the correct radial functions for both semi-core and valence projectors. For these reasons, we use in this work the PAOs as initial and more accurate projections. If needed, higher energy orbitals not included in the pseudopotential file can be constructed, for example, using solutions of Schrödinger equation under confinement potential [43, 44] (see also discussion in Section II F).

## B. Projectability disentanglement

As mentioned, the standard disentanglement approach selects the disentanglement and frozen manifolds via two energy windows [2]. We refer to this as energy disentanglement (ED). However, since bands have dispersions across the BZ, a fixed window for all  $k$ -points might not be an optimal choice. Taking the graphene band structure (Fig. 1a) as an example, the bands with large projectability are mixed with many free-electron bands

with zero projectability (grey bands in the conduction region). In this case, one is faced with several options for the outer and inner energy windows, each with different shortcomings: (a) If the inner window includes free-electron bands, the final MLWFs are mixtures of  $2s$ ,  $2p$  atomic orbitals and free-electron bands, delocalizing the resulting MLWFs; (b) if the outer window excludes both the free-electron bands and the atomic-orbital states inside free-electron bands, the WFs lack the anti-bonding part of the bonding/anti-bonding closure [13], again degrading the localization of WF; (c) if the upper bound of the inner window is set to its maximal allowed value, i.e. the blue dashed line positioned at the minimum of free-electron bands in Fig. 1b, and all the DFT eigenstates are included in the outer window, the disentanglement algorithm [2] will extract an optimally smooth manifold, at the expense of decreasing the chemical representability of the atomic-orbital bands in the free-electron region; in other words, the MLWFs obtained lose the information of the TB atomic orbitals in this chemical environment (see Fig. 1b).

The graphene case highlights the limitations of the standard ED. Instead, we propose here to select the disentanglement and frozen manifolds based on the projectability  $p_{m\mathbf{k}}$  of each state on the chosen PAOs (i.e., states are selected irrespective of their energy, but rather based on their chemical representativeness). Specifically, we select states based on two thresholds  $p_{\min}$  and  $p_{\max}$ : (a) If  $p_{m\mathbf{k}} < p_{\min}$ , the state  $\psi_{m\mathbf{k}}$  is discarded. (b) If  $p_{m\mathbf{k}} \geq p_{\max}$ , the state  $\psi_{m\mathbf{k}}$  is kept identically. Crucially, all states for which  $p_{\min} \leq p_{m\mathbf{k}} < p_{\max}$  are thrown in the disentanglement algorithm. Optimal numerical values for  $p_{\min}$  and  $p_{\max}$  are discussed later. In the case of graphene,  $p_{\max}$  identifies the fully atomic-orbital states inside the free-electron bands, while  $p_{\min}$  removes the fully free-electron bands from the disentanglement process, preventing the mixing of atomic and free-electron states. The two thresholds  $p_{\min}$  and  $p_{\max}$  constitute the parameters of the disentanglement process, replacing the four defining energy windows (the lower and upper bounds of the outer and inner energy windows). We note that projectability disentanglement is different from partly-occupied WF [13, 14] in that the latter uses an energy window to select frozen states and minimizes the total spread functional directly, while projectability disentanglement selects the localized states using projectability instead of a constant energy window across  $k$ -points. In fact, one can combine projectability disentanglement with a variational formulation [15] to construct MLWFs by minimizing directly the total spread functional.

Ideally, if PAOs were always a complete set to describe valence and near-Fermi-energy conduction bands, the PD would select the most relevant Bloch states and accurately interpolate these DFT bands. However, since the PAOs are fixed orbitals from isolated single-atom calculations for each element, if the chemical environment in the crystal structure is significantly different from that of pseudopotential generation, then the total projectability

$p_{m\mathbf{k}}$  might be smaller than 1 for bands around the conduction band minimum (CBM) or even for valence bands. In such cases, one solution is to increase the number of PAOs, i.e., adding more projectors with higher angular momentum, as we will discuss in Section II F. However, since one almost always wants to correctly reproduce valence bands (plus possibly the bottom of the conduction) but at the same time keep the Wannier Hamiltonian small for computational reasons, we suggest to additionally freeze all the states that sit below the Fermi energy in metals (or below the CBM for insulators) and also those a few eV above (typically, 2 eV or so). Such a combination of PD+ED gives accurate interpolation of bands below and around the Fermi energy (or band edges for insulators), as well as maximally restoring the atomic-orbital picture.

We stress here that, even if we call the resulting Wannier functions PDWFs for clarity, our optimal suggestion is to always also freeze the states in the energy window mentioned above, as we discuss in the next sections.

### C. Comparison

We choose four prototypical materials to discuss the present method: graphene, silicon, copper, and strontium vanadate ( $\text{SrVO}_3$ ). Graphene is a difficult case where atomic-orbital states highly mix with free-electron bands; silicon tests the Wannierization of both valence and conduction bands of an insulator; copper is a test on a metal; and  $\text{SrVO}_3$  represents the class of (metallic) perovskites. We compare the shapes, centers, and spreads of the resulting MLWFs using the five methods mentioned earlier: hydrogenic projection with ED (i.e., the standard approach), SCDM, PAO projection with ED, PAO projection with PD, and PAO projection with PD+ED.

#### 1. Graphene

The original and interpolated band structures for the five methods discussed are shown in Figs. 1b to 1f. The blue dashed lines in Figs. 1b, 1d and 1f indicate the top of the inner energy window, which is set optimally (and manually) to just below the free-electron bands, to freeze as much as possible the atomic-orbital states but exclude any free-electron state. For PD and PD+ED, we choose  $p_{\max} = 0.85$  and  $p_{\min} = 0.02$  (we will discuss later on the choice of these thresholds). Comparing Fig. 1d and Fig. 1b, one sees that ED produces similar bands irrespective of using hydrogenic or PAO projection. However, as shown in Fig. 2 (first and third row), the MLWFs for the two cases fall into slightly different minima: MLWFs from hydrogenic projection with ED are  $p_z$  and hybridized  $s \pm p$  orbitals pointing towards the center of the hexagon, while MLWFs from PAO with ED are  $p_z$ ,  $p_x$ , and  $s \pm p_y$ . This is due to the fact that the PAO projections guide the minimization towards spherical har-

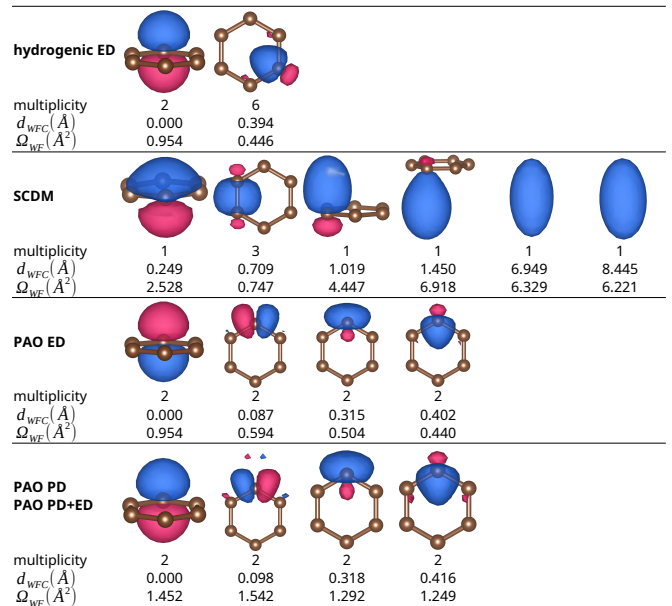


FIG. 2. **Graphene MLWFs: shapes, centers, and spreads obtained using different methods.**  $d_{\text{WFC}}$  is the distance of the WF center from the nearest-neighbor atom, and  $\Omega_{\text{WF}}$  is the MLWF spread. The multiplicity is the number of equivalent MLWFs, i.e. having the same  $d_{\text{WFC}}$ ,  $\Omega_{\text{WF}}$ , and shape, but different orientations.

monics, while the hydrogenic projections are farther away from such local minimum and the optimization algorithm happens to escape and converge to a better minimum. A possible future work is to introduce more advanced optimization algorithms to improve the convergence of maximal localization. Both the PAO with PD and PAO with PD+ED cases reach the same set of MLWFs,  $p_z$ ,  $p_x$ , and  $s \pm p_y$ , but with larger spreads than the PAO with ED, since the PD and PD+ED freeze more states, giving thus less freedom for maximal localization. Nevertheless, the interpolated bands of the PAO with PD and PAO with PD+ED cases can much better reproduce the atomic-orbital states inside the free-electron bands. Finally, compared to other cases, SCDM includes some free-electron bands, some of which can be even reproduced by the Wannier interpolation. However, in order to follow those free-electron bands, abrupt changes of character and band derivative are needed in the conduction band. As required by Nyquist–Shannon sampling theorem[45], this results in a denser  $\mathbf{k}$ -space sampling needed to obtain a good interpolation quality. Moreover, the MLWFs are much more delocalized and do not resemble atomic orbitals: as shown in Fig. 2, the last two MLWFs for SCDM are floating away from the graphene 2D lattice, blurring the TB picture of atomic orbitals in solids.

## 2. Silicon

The SCDM method obtains four front-bonding and four back-bonding MLWFs, while all other cases lead to atom-centered  $s$  and  $p$  MLWFs, as shown in Fig. S3. While overall the SCDM bands (Fig. 3c) seem to reproduce relatively better the higher conduction bands, they fail to correctly reproduce the bottom of the conduction band near the X point, induce more wiggles around X and W, and have much larger spreads. Due to the low projectability of Bloch states around X ( $p_{m\mathbf{k}}$  around 0.83), the CBM is not correctly reproduced in the PAO with PD, as these are not frozen in PD with the current choice of  $p_{\max} = 0.95$  and  $p_{\min} = 0.01$ . To explicitly freeze the CBM,  $p_{\max}$  would need to be lowered below 0.83. However, such kind of decrease will also result in freezing some high-energy conduction bands, degrading the localization. PD+ED overcomes this by explicitly freezing the near-Fermi-energy and low-projectability states at the CBM, but still only freezing those atomic-orbital states in the high-energy conduction bands that possess high projectability (see Fig. 3f), thus improving band interpolation. We note that the lower projectability of silicon CBM is intrinsic to the material—its CBM also includes  $3d$  character. Therefore, by adding  $d$  PAOs, the CBM projectability increases (from 0.83 to 0.99) and one can restore a high-quality band-structure interpolation within the PD method: as shown in Fig. 3e, the low-energy conduction bands are correctly reproduced once we regenerate a silicon pseudopotential including  $3d$  PAOs. Therefore, PD is sufficient to obtain an accurate band interpolation if enough PAOs are included (we will also discuss this later in Section IIF). For completeness, we show the SCDM interpolation using the regenerated pseudopotential in Fig. 3c: the added  $d$  PAOs help select a larger manifold thanks to the increased projectability, enabling SCDM to reproduce higher conduction bands, as well as fixing the wrong interpolation at the W point. Moreover, additional PAOs can also benefit ED, since the frozen window can be enlarged to reproduce more states. In general, adding more PAOs improves interpolation quality in cases where the target bands have low projectability, at the price of increased computational cost. PD+ED is a better option for reaching a good interpolation accuracy while keeping the size of the corresponding TB model small.

## 3. Copper and $\text{SrVO}_3$

Results for copper and  $\text{SrVO}_3$  are only shown in the SI (Figs. S4, S6, S7 and S9), since the conclusions are the same: PD+ED consistently provides the best interpolation quality among all methods we consider, while not requiring to increase the size of the Hamiltonian model, and results in WFs that resemble atomic orbitals or their hybridization.

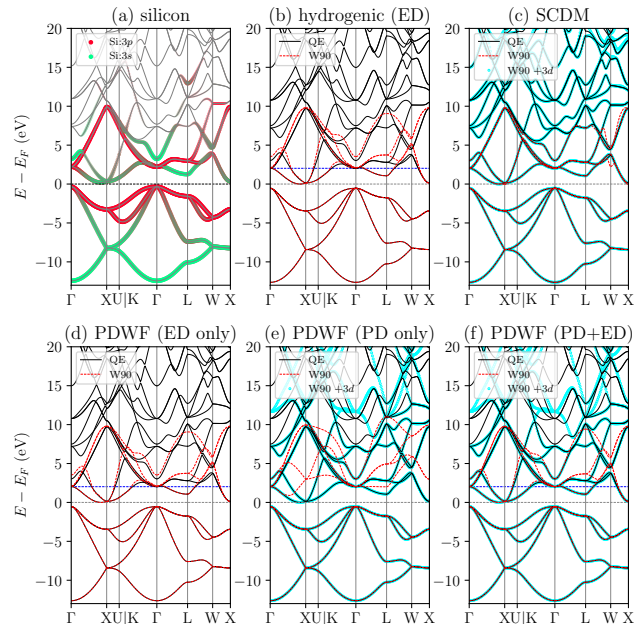


FIG. 3. **Comparisons of silicon band structures interpolated using different methods.** (a) DFT band structure, shown as grey lines. The colored dots represent the projectabilities of silicon  $3s$  (green) and  $3p$  (red) orbitals. The size of the dot is proportional to the total projectability  $p_{m\mathbf{k}}$  of the band  $m$  at  $k$ -point  $\mathbf{k}$ . For a detailed plot of total projectability, see Fig. S2. Comparisons of the original and the Wannier-interpolated bands for (b) hydrogenic projections with ED, (c) SCDM, (d) PAO with ED, (e) PAO with PD, and (f) PAO with PD+ED. The CBM (horizontal black dashed line) is at zero; the horizontal blue dashed line denotes the top of the inner energy window, i.e., CBM + 2eV, where applicable. Note in (c), (e), and (f), the cyan lines with circle markers show the interpolated bands obtained including also  $3d$  orbitals, and consequently increasing the dimensionality of the disentangled manifold. These additional states are beneficial because of the presence of an intrinsic  $d$  component at the bottom of the conduction manifold, and lead to more accurate band interpolations.

## D. High-throughput verification

In this section we discuss the applicability of the present PDWF method to obtain, in a fully automated way and without user input, WFs for any material. In order to assess quantitatively its performance, we compare it to SCDM, that can also be fully automated (see Ref. [30]).

In all results that follow, we exclude semicore orbitals in both methods, since these low-energy states correspond to almost flat bands and do not play any role in the chemistry of the materials. We compare quantitatively the band interpolation quality between the two methods and the corresponding WF centers and spreads on the 200-structure set used in Ref. [30] for both oc-

cupied and unoccupied bands, totalling 6818 MLWFs for each method. In accordance with Refs. [30, 46], the band interpolation quality is measured by the average band distance,

$$\eta_\nu = \sqrt{\frac{\sum_{n\mathbf{k}} \tilde{f}_{n\mathbf{k}} (\epsilon_{n\mathbf{k}}^{\text{DFT}} - \epsilon_{n\mathbf{k}}^{\text{Wan}})^2}{\sum_{n\mathbf{k}} \tilde{f}_{n\mathbf{k}}}}, \quad (6)$$

and the max band distance,

$$\eta_\nu^{\text{max}} = \max_{n\mathbf{k}} \left( \tilde{f}_{n\mathbf{k}} \left| \epsilon_{n\mathbf{k}}^{\text{DFT}} - \epsilon_{n\mathbf{k}}^{\text{Wan}} \right| \right), \quad (7)$$

where  $\tilde{f}_{n\mathbf{k}} = \sqrt{f_{n\mathbf{k}}^{\text{DFT}}(E_F + \nu, \sigma) f_{n\mathbf{k}}^{\text{Wan}}(E_F + \nu, \sigma)}$  and  $f(E_F + \nu, \sigma)$  is the Fermi-Dirac distribution. Here  $E_F + \nu$  and  $\sigma$  are fictitious Fermi levels and smearing widths which we choose for comparing a specific range of bands. Since the Wannier TB model describes the low-energy valence electrons, it is expected that the band interpolation deviates from the original in the higher conduction band region. Therefore, the higher  $\nu$  is, the larger  $\eta_\nu$  is expected to be. In the following paragraphs, we will use  $\eta_0$  and  $\eta_2$  to compare bands below  $E_F$  and  $E_F + 2\text{eV}$ , respectively;  $\sigma$  is always fixed at  $0.1\text{eV}$ .

In the supplementary information Section S8, we provide comparisons between the Wannier-interpolated bands and the DFT bands for both PDWF and SCDM, their respective band distances, and the Hamiltonian decay plots for each of the 200 materials. We discuss these properties in the following.

### 1. Projectability thresholds and automation

For PDWF, we set the maximum of the inner window to the Fermi energy + 2eV for metals, or to the CBM + 2eV for insulators, to fully reproduce states around Fermi energy or the band edges. We also specify the two additional parameters  $p_{\text{min}}$  and  $p_{\text{max}}$ . From our tests, in most cases  $p_{\text{max}} = 0.95$  and  $p_{\text{min}} = 0.01$  already produce very good results. However, since chemical environments vary across different crystal structures, the two parameters are not universal and influence the quality of band interpolation. Figure 4 shows the variation of band distances w.r.t.  $p_{\text{min}}$  and  $p_{\text{max}}$  for several materials. For  $\text{Al}_3\text{V}$  (Figs. 4a and 4b),  $\eta_0$  and  $\eta_2$  reach a minimum at two different sets of parameters, i.e.,  $p_{\text{max}} = 0.99$ ,  $p_{\text{min}} = 0.01$  and  $p_{\text{max}} = 0.97$ ,  $p_{\text{min}} = 0.01$ , respectively. In some cases, the variation of  $\eta$  w.r.t.  $p_{\text{max}}$  and  $p_{\text{min}}$  can be non-monotonic and display multiple local minima: For instance, in  $\text{Au}_2\text{Ti}$  (Fig. 4c) at  $p_{\text{min}} = 0.01$ ,  $\eta_2$  decreases from  $p_{\text{max}} = 0.90$  to 0.95 but increases from  $p_{\text{max}} = 0.95$  to 0.98 and finally reaches a local minimum at  $p_{\text{max}} = 0.99$ . In other cases,  $\eta$  can be quite stable and largely independent of the parameters: e.g., for  $\text{Ba}_6\text{Ge}_{10}$  (Fig. 4d),  $\eta_2$  reaches the same minimum for  $p_{\text{max}} = 0.99$  to 0.88.

Therefore, we implement an iterative optimization workflow to automatically find the optimal values for

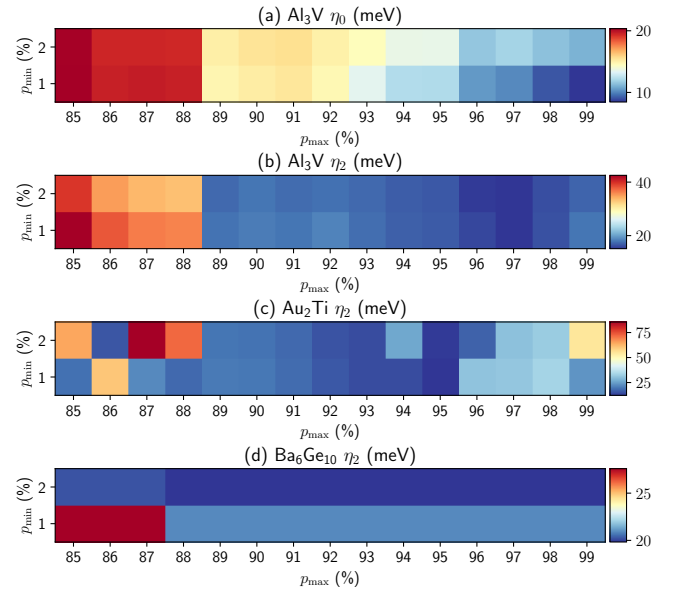


FIG. 4. **Quality of band interpolations: band distances for different choices of  $p_{\text{min}}$  and  $p_{\text{max}}$ .** (a)  $\eta_0$  of  $\text{Al}_3\text{V}$ , (b)  $\eta_2$  of  $\text{Al}_3\text{V}$ , (c)  $\eta_2$  of  $\text{Au}_2\text{Ti}$ , and (d)  $\eta_2$  of  $\text{Ba}_6\text{Ge}_{10}$ . Note the color scale is different for each plot.

$p_{\text{max}}$  and  $p_{\text{min}}$ , in order to fully automate the Wannierization procedure. The workflow is released as part of the `aiana-wannier90-workflows` package [47]. First, we run a QE band structure workflow to get the reference DFT bands for calculating  $\eta_2$ ; in addition, the DFT bands are also used to calculate the band gap of the material. Second, we run an optimization workflow with the following settings: The maximum of the inner window is set to Fermi energy + 2eV for metals and CBM + 2eV for insulators, respectively;  $p_{\text{max}}$  and  $p_{\text{min}}$  are set to the defaults of 0.95 and 0.01, respectively. Third, if the average band distance  $\eta_2$  is less than a threshold (set to 10 meV here), the workflow stops; otherwise, the workflow iterates on a mesh of  $p_{\text{max}}$  and  $p_{\text{min}}$ , i.e.  $p_{\text{max}}$  decreasing from 0.99 to 0.80 with step size -0.01, and  $p_{\text{min}} = 0.01$  or 0.02, until  $\eta_2 \leq \text{threshold}$ . If  $\eta_2$  is still larger than the threshold after exhausting all the parameter combinations, the workflow will output the minimum- $\eta_2$  calculation.

### 2. Band distance

To compare quantitatively the band interpolation quality of SCDM and PDWF, we Wannierize the 200 structures mentioned earlier and calculate their band distances with respect to the corresponding DFT bands. We choose  $\eta_2$  and  $\eta_2^{\text{max}}$  to compare near-Fermi-energy bands. The histograms of the band distances for the 200 structures are shown in Fig. 5. To directly compare SCDM and PDWF, the mean and median value

of  $\eta$  of the 200 calculations are shown as vertical lines in each panel. For PDWF, the mean  $\eta_2$  is 4.231 meV, to be compared with 11.201 meV for SCDM. For  $\eta_2^{\max}$  (that is a more stringent test of the quality of interpolation) the PDWF method also performs better, with a  $\eta_2^{\max} = 36.743$  meV vs. 84.011 meV for SCDM. We can also observe this trend in Fig. 5: For  $\eta_2$  and  $\eta_2^{\max}$ , the PDWF histogram bins are much more clustered towards  $\eta = 0$ . Note that in the cumulative histograms of  $\eta_2$ , at  $\eta = 20$  meV, the PDWF cumulative count is closer to the total number of calculations (200). This indicates that the PDWF has a higher success rate in reducing the interpolation error below 20 meV. Similarly, for  $\eta_2^{\max}$ , PDWF has a higher success rate in reducing the interpolation error under 100 meV (to get a better overview of  $\eta$  and  $\eta^{\max}$ , we further show the same histograms of  $\eta$  in a wider range 0 meV to 100 meV, and  $\eta^{\max}$  in range 0 meV to 500 meV, in Figs. S11 and S12). To reduce the effect of major outliers, we can also compare the interpolation accuracy of successful calculations, i.e., excluding the outlier calculations which have significantly large band distances. As shown in Table S1, the  $\eta_2^{\leq 20}$ , i.e., the average of all the calculations for which  $\eta_2 \leq 20$  meV, indicates that PDWF (2.922 meV) is twice as good as SCDM (5.280 meV), and also has a higher success rate: for  $\eta_2^{\leq 20}$ ,  $193/200 = 96.5\%$  of the structures have  $\eta_2 \leq 20$  meV, while for SCDM it is  $183/200 = 91.5\%$ . More details are listed in Table S1.

In summary, PDWF provides more accurate and robust interpolations, especially for bands around the Fermi energy or the band gap edges, which are the most relevant bands for many applications. Last but not least, a higher energy range can be accurately interpolated by increasing the number of PAOs (see Section IIF).

### 3. MLWF centers

Since we are aiming at restoring a tight-binding atomic-orbital picture with PDWF, we compare the distance of the WF centers from the nearest-neighboring (NN) and next-nearest-neighboring (NNN) atoms, again both for SCDM and PDWF. For each method, we compute  $d_{\text{NN}}$  and  $d_{\text{NNN}}$ , i.e., the average distance of all the 6818 MLWFs from the respective NN and NNN atoms. If  $d_{\text{NN}}$  is 0, then the atomic-orbital picture is strictly preserved. However, this is unlikely to happen since there is no constraint on the WF centers during both the disentanglement and the localization, and the final PDWFs, resembling atomic orbitals, are optimized according to the chemical environment. Still, if a WF center is much closer to the NN atom than to the NNN atom, then one can still assign it to the NN atom, preserving the atomic-orbital picture. Figure 6 shows the histograms for  $d_{\text{NN}}$  and  $d_{\text{NNN}}$  for the two methods. The PDWF average  $d_{\text{NN}} = 0.43$  Å is smaller than the SCDM  $d_{\text{NN}} = 0.53$  Å, and correspondingly the PDWF  $d_{\text{NNN}} = 2.19$  Å is instead larger than the SCDM  $d_{\text{NNN}} = 2.11$  Å.

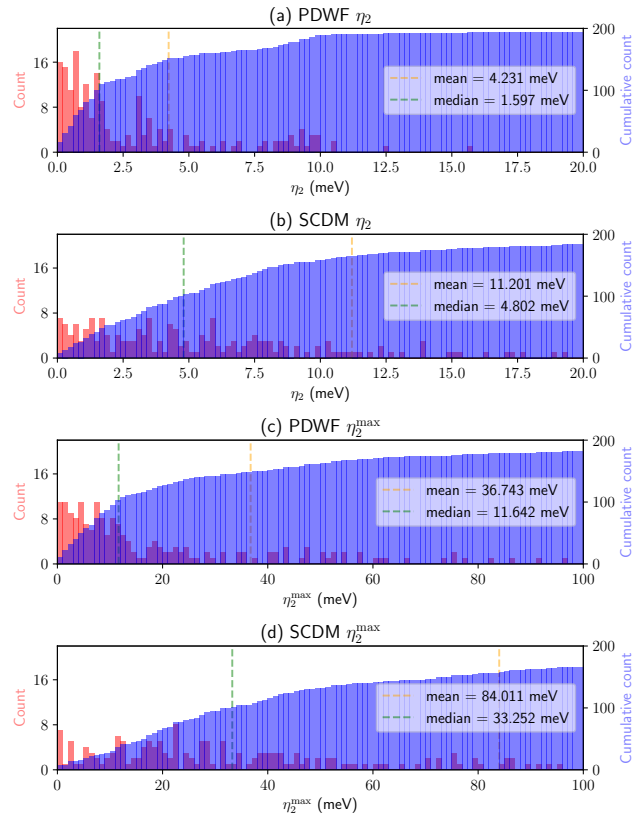


FIG. 5. **Histogram (red) and cumulative histogram (blue) of the band distances  $\eta_2$  and  $\eta_2^{\max}$  for 200 reference structures.** (a)  $\eta_2$  of PDWF, (b)  $\eta_2$  of SCDM, (c)  $\eta_2^{\max}$  of PDWF, and (d)  $\eta_2^{\max}$  of SCDM. The orange (green) vertical line is the mean (median) of the band distance for the 200 structures; their values are shown in the right of each panel; PDWF provides approximately an improvement by a factor of 3.

This can also be observed in Fig. 6: The overlap of the  $d_{\text{NN}}$  and  $d_{\text{NNN}}$  histograms is smaller for PDWF than for SCDM. To further understand the overlaps, we plot the histogram of the ratio  $d_{\text{NN}}/d_{\text{NNN}}$  of each MLWF in the insets of Fig. 6. For a MLWF, if  $d_{\text{NN}}/d_{\text{NNN}} = 1$ , then the MLWF is a bonding orbital centered between two atoms; while if  $d_{\text{NN}}/d_{\text{NNN}} \ll 1$ , then it can be regarded as an (almost) atomic orbital. The histogram of the ratio of SCDM has a long tail extending towards 1.0, i.e., there are a large number of SCDM MLWFs sitting close to bond centers; on the contrary, the vast majority of the PDWF MLWFs are closer to the NN atom.

We can further compare the effect of maximal localization on the WF centers. The WFs from the projection matrices  $A_{m\mathbf{nk}}$  are strictly atom-centered, i.e.  $d_{\text{NN}} = 0$ . The inset of Fig. S13a shows the histogram of the initial WFs, i.e., after disentanglement and before maximal localization, and the final MLWFs, i.e., after maximal localization, for PDWF. If one chooses  $d_{\text{NN}} \leq 0.1$  Å as the criterion for atom-centered MLWFs, then  $5594/6818 =$



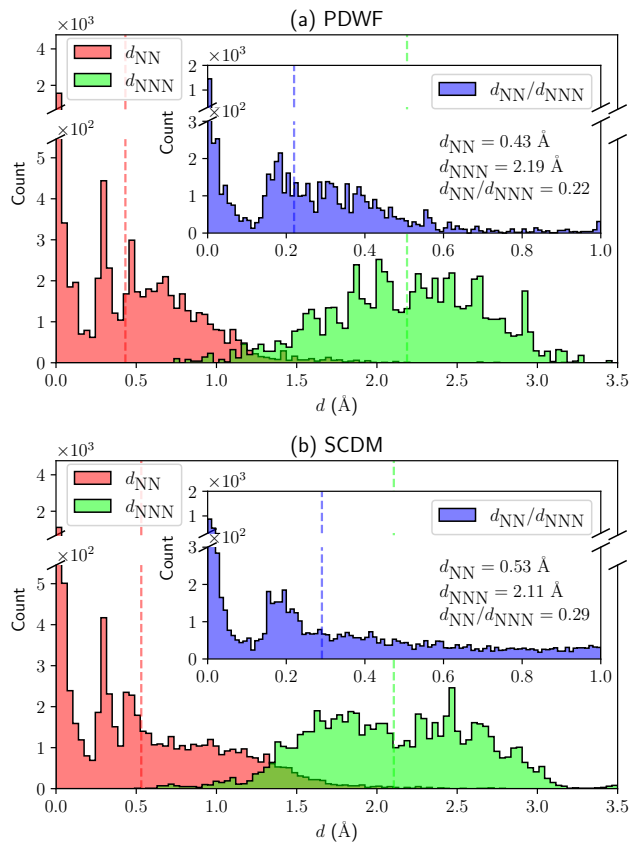


FIG. 6. Histogram of the distances of the WF centers from the NN atom (red,  $d_{NN}$ ) and NNN atom (green,  $d_{NNN}$ ), for 200 reference structures. (a) PDWF and (b) SCDM. The inset of each panel shows the histogram of the ratio of  $d_{NN}/d_{NNN}$ . The numbers in the lower right of each inset are the averages over all the 6818 MLWFs; PDWF provides MLWFs that are both closer to the NN atom and further away from the NNN atom.

82.0% of the initial WFs and  $2045/6818 = 30.0\%$  of the final MLWFs are atom-centered. The disentanglement and maximal localization improve the band interpolation, but since there is no constraint on the WF center in the spread functional Eq. (2), many of the final MLWF centers are not atom-centered. As a comparison, for SCDM,  $955/6818 = 14.0\%$  of the initial WFs and  $1823/6818 = 26.7\%$  of the final MLWFs are atom-centered. For completeness, the statistics and histograms of initial and final  $d_{NN}$ ,  $d_{NNN}$ , and  $d_{NN}/d_{NNN}$  are shown in Table S2 and Fig. S13.

In summary, for PDWF, most of the initial WFs (after disentanglement and before maximal localization) are atom-centered; many drift a bit away from atom centers during the localization, but the MLWFs are still much closer to the NN than to NNN atoms. For SCDM, most of the initial WFs are away from atom centers, and maximal localization pushes some of the WFs back to atoms, but

there is still a large number of MLWFs for which an atom representing the WF center cannot be clearly identified. To exactly fix the MLWFs to atomic positions, one needs to add constraints to the spread functional [16], at the cost of potentially having worse interpolators. However, this is beyond the scope of the current work, and here we rely on the atom-centered PAO projectors to guide the MLWFs towards the atomic positions, so that the final MLWFs are optimally localized and atom-centered.

#### 4. MLWF spreads

Next, we investigate the spread distributions of SCDM and PDWF. Usually, we want localized MLWFs to restore the TB atomic orbitals. Figure 7 shows the histograms of the spread distributions for the two methods. The SCDM spreads have a long tail extending over  $10 \text{ \AA}^2$  in Fig. 7b, due to its inclusion of free-electron states in the density matrix, thus resulting in more delocalized MLWFs as discussed earlier (see e.g. Fig. 2). On the contrary, the PDWF selects and freezes atomic-orbital states from the remaining bands, leading to much more localized MLWFs, thus much more clustered in a narrow range of  $0 \text{ \AA}^2$  to  $4 \text{ \AA}^2$ , and already at  $5 \text{ \AA}^2$  the cumulative histogram almost reaches the total number of MLWFs (see Fig. 7a). This can be interpreted as follows: The PAO initial projections guide the spread minimization toward the (local) minimum resembling spherical harmonics, whereas the SCDM-decomposed basis vectors are designed to be mathematical objects spanning as much as possible the density matrix, but result in WFs for which it is harder to assign definite orbital characters.

We can further compare the average initial (after disentanglement but before maximal localization) and final (after disentanglement and maximal localization) spreads between the two methods, as shown in Table S3 and corresponding histograms in Fig. S14. Maximal localization is needed to bring SCDM spreads, from the initial  $\Omega^i = 30.82 \text{ \AA}^2$  to the final  $\Omega^f = 3.54 \text{ \AA}^2$ ; For PDWF, the initial  $\Omega^i = 2.72 \text{ \AA}^2$  is already excellent, and much better than the final  $\Omega^f$  for SCDM; localization then brings it to an optimal  $\Omega^f = 1.41 \text{ \AA}^2$ .

#### 5. Hamiltonian decay

Finally, we compare the decay length of the Wannier gauge Hamiltonian between the two methods in Fig. 8. Thanks to the localization of MLWFs, the expectation values of quantum mechanical operators in the MLWF basis, such as the Hamiltonian  $H(\mathbf{R})$ , decay rapidly with respect to the lattice vector  $\mathbf{R}$  (exponentially in insulators[48, 49] and properly disentangled metals). To compare this decay for the Hamiltonian matrix elements, we approximate the Frobenius norm of the Hamiltonian

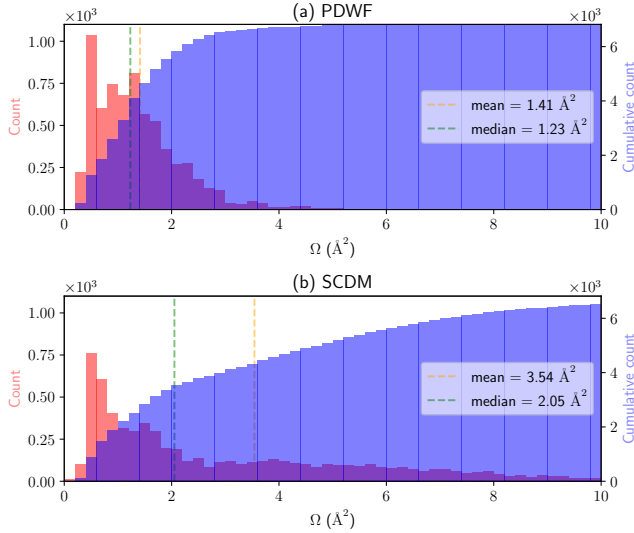


FIG. 7. **Histogram (red) and cumulative histogram (blue) of WF spreads for 200 reference structures.** (a) PDWF and (b) SCDM. The orange (green) vertical line is the mean (median) spread of the 6818 MLWFs, their values are shown in the right of each panel. The long tail of MLWF spreads obtained with SCDM is absent in PDWF.

as

$$\|H(\mathbf{R})\| = \|H(\mathbf{0})\| \exp\left(-\frac{\|\mathbf{R}\|}{\tau}\right), \quad (8)$$

where  $\tau$  measures the decay length. Then  $\tau$  is fitted by least squares to the calculated  $\|H(\mathbf{R})\|$ ; as shown in Fig. 8a, the Hamiltonian of PDWF decays faster than SCDM for  $\text{Br}_2\text{Ti}$ , which is selected here to represent the general trend between PDWF and SCDM Hamiltonians. Fig. 8b shows the histogram of  $\tau$  for the 200 materials; the mean  $\tau$  are 2.266 Å for PDWF and 2.659 Å for SCDM, respectively, indicating that the PDWF Hamiltonian decays faster than SCDM, consistent with the better band interpolation of PDWF discussed in Fig. 5.

### E. High-throughput Wannierization

Based on the above verification, we run a HT Wannierization using PDWF for 21,737 materials, selected from the non-magnetic materials of the MC3D database [36]. Figure 9 shows the band distance histograms for  $\eta_2$  and  $\eta_2^{\max}$ . Overall, the statistics follow the same trend as the 200 materials set in Fig. 5: the average  $\eta_2$  and average  $\eta_2^{\max}$  are 3.685 meV and 42.768 meV, respectively. Note in Fig. 9a the  $\eta_2$  is not truncated at 10 meV, but rather due to the automated optimization workflow: results that have  $\eta_2$  larger than a threshold (10 meV) are further optimized with respect to  $p_{\min}$  and  $p_{\max}$ , thus improving the average band distance  $\eta_2$ . In Table S4

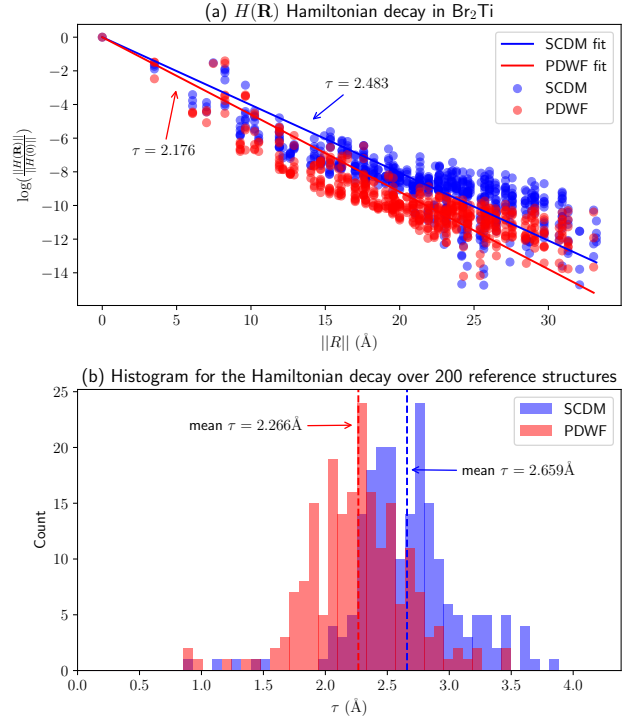


FIG. 8. **Exponential decay of the Hamiltonian  $H(\mathbf{R})$  in the basis of MLWFs.** (a) Exponential-form fitting of Frobenius norm of the Hamiltonian  $\|H(\mathbf{R})\|$  w.r.t. to the 2-norm of lattice vector  $\|\mathbf{R}\|$  for the case of  $\text{Br}_2\text{Ti}$ , for PDWF (red) and SCDM (blue). The  $\tau$  reported are the fitted decay lengths of the PDWF and SCDM Hamiltonians, respectively. (b) Histogram of decay lengths  $\tau$  for the 200 reference materials, obtained using PDWF (red) and SCDM (blue). The vertical lines indicate the mean  $\tau$  of PDWF and SCDM, respectively.

we show several other statistics for the band distances. The excellent interpolation quality of PDWF can be assessed, for instance, from the number of systems with  $\eta_2 \leq 20$  meV, that are  $\approx 97.8\%$  of all the calculations (21259/21737); the corresponding bands distance calculated on these 21259 calculations is  $\eta_2^{\leq 20} = 2.118$  meV. This remarkable result show how automated and reliable Wannierizations can now be deployed automatically both for individual calculation and for HT application.

### F. Additional PAOs for high-energy high-accuracy interpolation

Based on the HT Wannierization results, one can identify cases where the interpolation quality can be further improved by increasing the number of PAOs. Typically, the number of PAOs is determined during pseudopotential generation, and they are usually the orbitals describing low-energy valence electrons. In some cases, the bonding/anti-bonding combinations of these PAOs

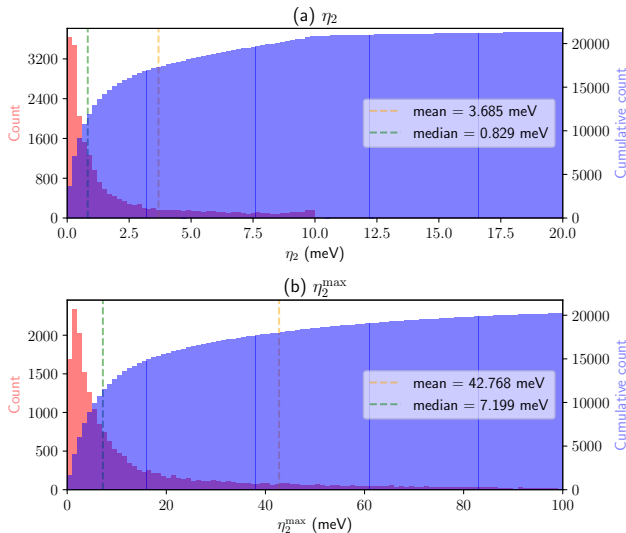


FIG. 9. **Histogram (red) and cumulative histogram (blue) of the PDWF band distances for 21,737 non-magnetic structures obtained from the materials cloud MC3D database [36].** (a) Average band distance  $\eta_2$  and (b) max band distance  $\eta_2^{\max}$ . The orange (green) vertical line is the mean (median) of the band distance for the 21,737 structures; their values are shown in the right of each panel.

are not sufficient to span the space of target conduction bands, leading to a loss of interpolation quality. We use silicon as an example to illustrate the difficulties of accurately describing its CBM [50], which is not located at any high-symmetry  $k$ -point, but along the  $\Gamma$ -X line. The common choice of one  $s$  and three  $p$  hydrogenic or PAOs projectors per atom results in oscillations in the Wannier-interpolated bands at the meV level. To remedy this, one can use a larger set of PAOs, e.g., by regenerating a silicon pseudopotential including  $d$  PAOs as discussed in Section IIC2. However, generating a new pseudopotential requires extensive testing and validation, therefore another solution could be using a set of PAOs different from the pseudopotential ones. To compare this second approach, we test here also PAOs obtained from the `OpenMX` code [44], and Wannierize silicon using one  $s$ , three  $p$ , and five  $d$  PAOs per atom using ED. This provides a much better description of the CBM, as shown in Fig. S17. Moreover, the additional  $d$  orbitals allow to raise the inner energy window and better reproduce a larger number of conduction bands, as shown in Fig. S18, which might be beneficial for some applications. For completeness, we also show the WF spreads and shapes of  $d$  orbitals in Fig. S19. However, there are some caveats to this approach. When using external PAOs, ideally one should generate them using the same pseudization scheme as the pseudopotentials used in the DFT calculations. The PAOs from `OpenMX` are instead generated using a different scheme, resulting in

lower projectabilities (smaller than one even for the valence bands, as shown in Fig. S21). In such case, PD cannot reproduce the original bands (see Fig. S20b), thus ED (with a higher inner energy window) is needed to obtain accurate interpolation (see Fig. S18d). In comparison, the pseudopotential PAOs which we regenerated with  $3d$  orbitals (as discussed in Section IIC2) are better projectors for the wavefunctions. Indeed, the first 12 bands have projectabilities almost equal to 1, and as a consequence PD itself already provides accurate band interpolation (all the low-energy conduction states are frozen since their projectabilities are high, see Fig. S20a). Moreover, we mention that when adding additional projectors one needs to make sure that they have the correct number of radial nodes: e.g., the gold pseudopotential from SSSP [46] contains  $5s + 5p$  semicore states, and  $6s + 5d$  orbitals for valence electrons. If one wants to add an additional  $6p$  orbital, it is important to ensure that the  $6p$  orbital has one radial node, such that it is orthogonal to the nodeless  $5p$  semicore state; Otherwise, the Bloch wavefunctions would project onto the  $5p$  semicore state, and PD would only disentangle the  $5p$  semicore states instead of the  $6p$  orbitals contributing to bands above the Fermi energy. In summary, including more projectors can further improve the interpolation quality, but at the expense of increasing the number of orbitals in the model. The combination of PD and ED enables to improve the interpolation quality of low-projectability states while keeping the TB model size small. Automatic checks could be implemented in the future in the `AiiDA` workflows to detect whether the projectability drops below a certain threshold, and in that case either raise a warning or automatically add more projectors.

### III. CONCLUSIONS

We present an automated method for the automated, robust, and reliable construction of tight-binding models based on MLWFs. The approach applies equally well to metals, insulators and semiconductors, providing in all cases atomic-like orbitals that span both the occupied states, and the empty ones whose character remains orbital-like and not free-electron-like. The method is based on the band projectability onto pseudo-atomic orbitals to select which states are kept identically, dropped, or passed on to the established disentanglement procedure. We augment such projectability-based selection with an additional energy window to guarantee that all states around the Fermi level or the conduction band edge are well reproduced, showing that such a combination enables accurate interpolation even when minimal sets of initial atomic orbitals are chosen. This results in compact Wannier tight-binding models that provide accurate band interpolations while preserving the picture of atomic orbitals in crystals. We refer to the method collectively as projectability-disentangled Wannier functions (PDWF).

The Wannierization process is implemented as fully automated `AiiDA` workflows. We compare PDWFs with the other method that is also fully automated, namely SCDM. We show with a detailed study of 200 structures that PDWFs lead to more accurate band interpolations (with errors with respect to the original bands at the meV scale), and are more atom-centered and more localized than those originating from SCDM. The high accuracy in band interpolations, the target atomic orbitals obtained, and the low computational cost make PDWFs an ideal choice for automated or high-throughput Wannierization, which we demonstrate by performing the Wannierization of 21,737 non-magnetic structures from the Materials Cloud MC3D database.

#### IV. METHODS

We implement the PAO projection in the `pw2wannier90.x` executable inside `Quantum ESPRESSO` (QE) [42, 51]; the PD and PD+ED methods are implemented on top of the `Wannier90` code [4]. In terms of the practical implementation, computing PAO projections is more efficient in both computational time and memory than the SCDM QR decomposition with column pivoting (QRCP) algorithm, since the  $A_{m\mathbf{k}}$  matrices (i.e., the inner products of Bloch wavefunctions with PAOs) can be evaluated in the plane-wave  $G$  vector space, rather than requiring a Fourier transform and decomposition of very large real-space wavefunction matrices. Furthermore, since the HT Wannierization can be computationally intensive, we implement a “ $k$ -pool parallelization strategy” inside `pw2wannier90.x`, similarly to the main `pw.x` code of QE, to efficiently utilize many-core architectures by parallelizing over “pools” of processors for the almost trivially-parallel computations at each  $k$ -point. Test results show that  $k$ -pool parallelization significantly improves the efficiency of `pw2wannier90.x` (benchmarks are shown in Fig. S10).

The DFT calculations are carried out using QE, with the SSSP efficiency (version 1.1, PBE functional) library [46] for pseudopotentials and its recommended energy cutoffs. The HT calculations are managed with the `AiiDA` infrastructure [33–35] which submits QE and `Wannier90` calculations to remote clusters, parses, and stores the results into a database, while also orchestrating all sequences of simulations and workflows. The automated `AiiDA` workflows are open-source and hosted on `GitHub` [47]. The workflows accept a crystal structure as input and provide the Wannier-interpolated band structure, the real-space MLWFs, and a number of additional quantities as output. Semicore states from pseudopotentials are automatically detected and excluded from the Wannierizations, except for a few cases where some semicore states overlap with valence states; in such cases, all the semicore states are Wannierized, otherwise the band interpolation quality would be degraded, especially for

SCDM. A regular  $k$ -point mesh is used for the Wannier calculations, with a  $k$ -point spacing of  $0.2 \text{ \AA}^{-1}$ , as selected by the protocol in Vitale *et al.* [30]. MLWFs are rendered with VESTA [52]. Figures are generated by `matplotlib` [53].

#### V. DATA AVAILABILITY

All data generated for this work can be obtained from the Materials Cloud Archive (<https://doi.org/10.24435/materialscloud:v4-e9>).

#### VI. CODE AVAILABILITY

All codes used for this work are open-source; the latest stable versions can be downloaded at <http://www.wannier.org/>, <https://www.quantum-espresso.org/>, <https://www.aiida.net/>, and <https://github.com/aiidateam/aiida-wannier90-workflows>.

The modifications to the codes mentioned above implemented for this work will become available in the next releases of `Quantum ESPRESSO` (`pw2wannier90.x`) and `Wannier90`.

#### VII. ACKNOWLEDGEMENTS

We acknowledge financial support from the NCCR MARVEL (a National Centre of Competence in Research, funded by the Swiss National Science Foundation, grant No. 205602), the Swiss National Science Foundation (SNSF) Project Funding (grant 200021E.206190 “FISH4DIET”). The work is also supported by a pilot access grant from the Swiss National Supercomputing Centre (CSCS) on the Swiss share of the LUMI system under project ID “PILOT MC EPFL-NM 01”, a CHRONOS grant from the CSCS on the Swiss share of the LUMI system under project ID “REGULAR MC EPFL-NM 02”, and a grant from the CSCS under project ID s0178.

#### VIII. AUTHOR CONTRIBUTIONS

J. Q. implemented and tested the PDWF method. N. M. suggested to use projectability thresholds. G. P. and N. M. supervised the project. All authors analyzed the results and contributed to writing the manuscript.

#### IX. COMPETING INTERESTS

The authors declare that there are no competing interests.

- [1] N. Marzari and D. Vanderbilt, Maximally localized generalized Wannier functions for composite energy bands, *Phys. Rev. B* **56**, 12847 (1997).
- [2] I. Souza, N. Marzari, and D. Vanderbilt, Maximally localized Wannier functions for entangled energy bands, *Phys. Rev. B* **65**, 035109 (2001).
- [3] N. Marzari, A. A. Mostofi, J. R. Yates, I. Souza, and D. Vanderbilt, Maximally localized Wannier functions: Theory and applications, *Rev. Mod. Phys.* **84**, 1419 (2012).
- [4] G. Pizzi, V. Vitale, R. Arita, S. Blügel, F. Freimuth, G. Géranton, M. Gibertini, D. Gresch, C. Johnson, T. Koretsune, J. Ibañez-Azpiroz, H. Lee, J.-M. Lihm, D. Marchand, A. Marrazzo, Y. Mokrousov, J. I. Mustafa, Y. Nohara, Y. Nomura, L. Paulatto, S. Poncé, T. Ponweiser, J. Qiao, F. Thöle, S. S. Tsirkin, M. Wierzbowska, N. Marzari, D. Vanderbilt, I. Souza, A. A. Mostofi, and J. R. Yates, Wannier90 as a community code: new features and applications, *J. Phys.: Condens. Matter* **32**, 165902 (2020).
- [5] R. Resta and D. Vanderbilt, *Theory of Polarization: A Modern Approach* (Springer, 2007) pp. 31–68.
- [6] Y.-S. Lee, M. B. Nardelli, and N. Marzari, Band Structure and Quantum Conductance of Nanostructures from Maximally Localized Wannier Functions: The Case of Functionalized Carbon Nanotubes, *Phys. Rev. Lett.* **95**, 076804 (2005).
- [7] M. G. Lopez, D. Vanderbilt, T. Thonhauser, and I. Souza, Wannier-based calculation of the orbital magnetization in crystals, *Phys. Rev. B* **85**, 014435 (2012).
- [8] X. Wang, J. R. Yates, I. Souza, and D. Vanderbilt, Ab initio calculation of the anomalous Hall conductivity by Wannier interpolation, *Phys. Rev. B* **74**, 195118 (2006).
- [9] J. R. Yates, X. Wang, D. Vanderbilt, and I. Souza, Spectral and Fermi surface properties from Wannier interpolation, *Phys. Rev. B* **75**, 195121 (2007).
- [10] J. Qiao, J. Zhou, Z. Yuan, and W. Zhao, Calculation of intrinsic spin Hall conductivity by Wannier interpolation, *Phys. Rev. B* **98**, 214402 (2018).
- [11] R. Sakuma, Symmetry-adapted Wannier functions in the maximal localization procedure, *Phys. Rev. B* **87**, 235109 (2013).
- [12] F. Gygi, J.-L. Fattebert, and E. Schwegler, Computation of Maximally Localized Wannier Functions using a simultaneous diagonalization algorithm, *Comput. Phys. Commun.* **155**, 1 (2003).
- [13] K. S. Thygesen, L. B. Hansen, and K. W. Jacobsen, Partly Occupied Wannier Functions, *Phys. Rev. Lett.* **94**, 026405 (2005).
- [14] K. S. Thygesen, L. B. Hansen, and K. W. Jacobsen, Partly occupied Wannier functions: Construction and applications, *Phys. Rev. B* **72**, 125119 (2005).
- [15] A. Damle, A. Levitt, and L. Lin, Variational Formulation for Wannier Functions with Entangled Band Structure, *Multiscale Model. Simul.* **17**, 167 (2019).
- [16] R. Wang, E. A. Lazar, H. Park, A. J. Millis, and C. A. Marianetti, Selectively localized Wannier functions, *Phys. Rev. B* **90**, 165125 (2014).
- [17] P. F. Fontana, A. H. Larsen, T. Olsen, and K. S. Thygesen, Spread-balanced Wannier functions: Robust and automatable orbital localization, *Phys. Rev. B* **104**, 125140 (2021).
- [18] J. I. Mustafa, S. Coh, M. L. Cohen, and S. G. Louie, Automated construction of maximally localized Wannier functions: Optimized projection functions method, *Phys. Rev. B* **92**, 165134 (2015).
- [19] É. Cancès, A. Levitt, G. Panati, and G. Stoltz, Robust determination of maximally localized Wannier functions, *Phys. Rev. B* **95**, 075114 (2017).
- [20] D. Gontier, A. Levitt, and S. Siraj-dine, Numerical construction of Wannier functions through homotopy, *J. Math. Phys.* **60**, 031901 (2019).
- [21] A. Damle, L. Lin, and L. Ying, Compressed Representation of Kohn–Sham Orbitals via Selected Columns of the Density Matrix, *J. Chem. Theory Comput.* **11**, 1463 (2015).
- [22] A. Damle and L. Lin, Disentanglement via Entanglement: A Unified Method for Wannier Localization, *Multiscale Model. Simul.* **16**, 1392 (2018).
- [23] A. Jain, S. P. Ong, G. Hautier, W. Chen, W. D. Richards, S. Dacek, S. Cholia, D. Gunter, D. Skinner, G. Ceder, and K. A. Persson, Commentary: The Materials Project: A materials genome approach to accelerating materials innovation, *APL Mater.* **1**, 011002 (2013).
- [24] S. Curtarolo, W. Setyawan, G. L. W. Hart, M. Jahnatek, R. V. Chepulskii, R. H. Taylor, S. Wang, J. Xue, K. Yang, O. Levy, M. J. Mehl, H. T. Stokes, D. O. Demchenko, and D. Morgan, AFLOW: An automatic framework for high-throughput materials discovery, *Comput. Mater. Sci.* **58**, 218 (2012).
- [25] J. E. Saal, S. Kirklin, M. Aykol, B. Meredig, and C. Wolverton, Materials Design and Discovery with High-Throughput Density Functional Theory: The Open Quantum Materials Database (OQMD), *JOM* **65**, 1501 (2013).
- [26] Computational Materials Repository, <https://cmr.fysik.dtu.dk/>, [Online; accessed 2023-02-28].
- [27] L. Talirz, S. Kumbhar, E. Passaro, A. V. Yakutovich, V. Granata, F. Gargiulo, M. Borelli, M. Uhrin, S. P. Huber, S. Zoupanos, C. S. Adorf, C. W. Andersen, O. Schütt, C. A. Pignedoli, D. Passerone, J. VandeVondele, T. C. Schulthess, B. Smit, G. Pizzi, and N. Marzari, Materials Cloud, a platform for open computational science, *Sci. Data* **7**, 299 (2020).
- [28] C. Draxl and M. Scheffler, NOMAD: The FAIR concept for big data-driven materials science, *MRS Bulletin* **43**, 676 (2018).
- [29] D. Gresch, Q. Wu, G. W. Winkler, R. Häuselmann, M. Troyer, and A. A. Soluyanov, Automated construction of symmetrized Wannier-like tight-binding models from ab initio calculations, *Phys. Rev. Mater.* **2**, 103805 (2018).
- [30] V. Vitale, G. Pizzi, A. Marrazzo, J. R. Yates, N. Marzari, and A. A. Mostofi, Automated high-throughput Wannierisation, *npj Comput. Mater.* **6**, 66 (2020).
- [31] K. F. Garrity and K. Choudhary, Database of Wannier tight-binding Hamiltonians using high-throughput density functional theory, *Sci. Data* **8**, 106 (2021).
- [32] L. A. Agapito, S. Ismail-Beigi, S. Curtarolo, M. Fornari, and M. B. Nardelli, Accurate tight-binding Hamiltonian matrices from ab initio calculations: Minimal basis sets, *Phys. Rev. B* **93**, 035104 (2016).

- [33] G. Pizzi, A. Cepellotti, R. Sabatini, N. Marzari, and B. Kozinsky, AiiDA: automated interactive infrastructure and database for computational science, *Comput. Mater. Sci.* **111**, 218 (2016).
- [34] S. P. Huber, S. Zoupanos, M. Uhrin, L. Talirz, L. Kahle, R. Häuselmann, D. Gresch, T. Müller, A. V. Yakutovich, C. W. Andersen, F. F. Ramirez, C. S. Adorf, F. Gargiulo, S. Kumbhar, E. Passaro, C. Johnston, A. Merkys, A. Cepellotti, N. Mounet, N. Marzari, B. Kozinsky, and G. Pizzi, AiiDA 1.0, a scalable computational infrastructure for automated reproducible workflows and data provenance, *Sci. Data* **7**, 300 (2020).
- [35] M. Uhrin, S. P. Huber, J. Yu, N. Marzari, and G. Pizzi, Workflows in AiiDA: Engineering a high-throughput, event-based engine for robust and modular computational workflows, *Comput. Mater. Sci.* **187**, 110086 (2021).
- [36] Materials Cloud three-dimensional crystals database (MC3D), <https://www.materialscloud.org/discover/mc3d/dashboard/ptable>, [Online; accessed 2023-02-28].
- [37] P.-O. Löwdin, On the Non-Orthogonality Problem Connected with the Use of Atomic Wave Functions in the Theory of Molecules and Crystals, *J. Chem. Phys.* **18**, 365 (1950).
- [38] E. Prodan and W. Kohn, Nearsightedness of electronic matter, *Proceedings of the National Academy of Sciences* **102**, 11635 (2005).
- [39] M. Benzi, P. Boito, and N. Razouk, Decay Properties of Spectral Projectors with Applications to Electronic Structure, *SIAM Rev.* **55**, 3 (2013).
- [40] A. Damle, L. Lin, and L. Ying, SCDM-k: Localized orbitals for solids via selected columns of the density matrix, *J. Comput. Phys.* **334**, 1 (2017).
- [41] R. Mahajan, I. Timrov, N. Marzari, and A. Kashyap, Importance of intersite Hubbard interactions in  $\beta$ -MnO<sub>2</sub>: A first-principles DFT+U+V study, *Phys. Rev. Materials* **5**, 104402 (2021).
- [42] P. Giannozzi, O. Baseggio, P. Bonfà, D. Brunato, R. Car, I. Carnimeo, C. Cavazzoni, S. de Gironcoli, P. Delugas, F. F. Ruffino, A. Ferretti, N. Marzari, I. Timrov, A. Urru, and S. Baroni, Quantum ESPRESSO toward the exascale, *J. Chem. Phys.* **152**, 154105 (2020).
- [43] T. Ozaki, Variationally optimized atomic orbitals for large-scale electronic structures, *Phys. Rev. B* **67**, 155108 (2003).
- [44] T. Ozaki and H. Kino, Numerical atomic basis orbitals from H to Kr, *Phys. Rev. B* **69**, 195113 (2004).
- [45] S. N. Alan Oppenheim, Alan Willsky, *Signals and Systems* (Prentice Hall, 1997) p. 957.
- [46] G. Prandini, A. Marrazzo, I. E. Castelli, N. Mounet, and N. Marzari, Precision and efficiency in solid-state pseudopotential calculations, *npj Comput. Mater.* **4**, 72 (2018).
- [47] aiiida-wannier90-workflows: A collection of advanced automated workflows to compute Wannier functions using AiiDA and the Wannier90 code, <https://github.com/aiidateam/aiida-wannier90-workflows>, [Online; accessed 2023-02-28].
- [48] C. Brouder, G. Panati, M. Calandra, C. Mourougane, and N. Marzari, Exponential Localization of Wannier Functions in Insulators, *Phys. Rev. Lett.* **98**, 046402 (2007).
- [49] G. Panati and A. Pisante, Bloch Bundles, Marzari-Vanderbilt Functional and Maximally Localized Wannier Functions, *Commun. Math. Phys.* **322**, 835 (2013).
- [50] S. Poncé, F. Macheda, E. R. Margine, N. Marzari, N. Bonini, and F. Giustino, First-principles predictions of Hall and drift mobilities in semiconductors, *Phys. Rev. Res.* **3**, 043022 (2021).
- [51] P. Giannozzi, S. Baroni, N. Bonini, M. Calandra, R. Car, C. Cavazzoni, D. Ceresoli, G. L. Chiarotti, M. Cococcioni, I. Dabo, A. D. Corso, S. de Gironcoli, S. Fabris, G. Fratesi, R. Gebauer, U. Gerstmann, C. Gougoussis, A. Kokalj, M. Lazzeri, L. Martin-Samos, N. Marzari, F. Mauri, R. Mazzarello, S. Paolini, A. Pasquarello, L. Paulatto, C. Sbraccia, S. Scandolo, G. Sclauzero, A. P. Seitsonen, A. Smogunov, P. Umari, and R. M. Wentzcovitch, QUANTUM ESPRESSO: a modular and open-source software project for quantum simulations of materials, *J. Phys.: Condens. Matter* **21**, 395502 (2009).
- [52] K. Momma and F. Izumi, VESTA 3 for three-dimensional visualization of crystal, volumetric and morphology data, *J. Appl. Crystallogr.* **44**, 1272 (2011).
- [53] J. D. Hunter, Matplotlib: A 2D Graphics Environment, *Comput. Sci. Eng.* **9**, 90 (2007).