

StochasticGW-GPU: Rapid Quasi-Particle Energies for Molecules beyond 10,000 Atoms

Phillip S. Thomas,¹ Minh Nguyen, Dimitri Bazile, Tucker Allen, Barry Y. Li, Wenfei Li, Mauro Del Ben, Jack Deslippe, and Daniel Neuhauser*



Cite This: *J. Chem. Theory Comput.* 2026, 22, 3960–3970



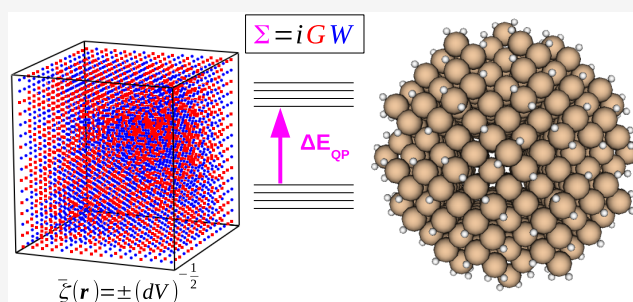
Read Online

ACCESS |

Metrics & More

Article Recommendations

ABSTRACT: *StochasticGW* is a code for computing accurate quasi-particle (QP) energies of molecules and material systems in the *GW* approximation. *StochasticGW* utilizes the stochastic Resolution of the Identity (sROI) technique to enable a massively parallel implementation with computational costs that scale semilinearly with system size, allowing the method to access systems with tens of thousands of electrons. We introduce a new implementation, *StochasticGW-GPU*, for which the main bottleneck steps have been ported to GPUs and give substantial performance improvements over previous versions of the code. We showcase the new code by computing band gaps of hydrogenated silicon clusters (Si_xH_y) containing up to 10,001 atoms and 35,144 electrons, and we obtain individual QP energies with a statistical precision of better than ± 0.03 eV with times-to-solution of less than 1 h.



INTRODUCTION

In recent years, predicting electronic properties of materials from first principles has become a key step in the materials design process, greatly reducing laboratory time and costs by directing synthetic efforts toward the most promising material candidates for a given application. Properties of interest, including band gaps, ionization potentials, and optical spectra, can be computed via electronic structure methods implemented in commercially available and open source software. For excited states, post-Hartree–Fock methods, including multireference configuration interaction^{1,2} and equation-of-motion coupled-cluster methods,^{3,4} while being gold standards for accuracy, are only applicable to small molecules since the computational cost of these methods grows steeply with the number of electrons. Due to their more favorable scaling, density functional theory (DFT)-based methods⁵ have become the industry standards for predicting ground state energies of large molecules and materials;⁶ however, their accuracy is poor when used to predict quasi-particle (QP) energies corresponding to excited states.^{7–9} Excited-state methods that can be applied on top of a DFT starting point, such as time-dependent (TD)-DFT,¹⁰ the *GW* approximation,^{9,11,12} including its extensions using perturbation theory,¹³ and the Bethe–Salpeter Equation (BSE) approach,¹⁴ provide superior accuracy compared to DFT, but they are expensive to apply, limiting excited state calculations to systems containing $\sim 10,000$ electrons.^{15–20}

The *GW* method has emerged as a robust and routinely used tool for computing QP energies of material systems,^{11,21,22} and *GW* implementations are now found in many quantum chemistry/materials software packages.^{23–36} Here, one approximates the self-energy operator, Σ , which embodies the many-body electron exchange and correlation effects, as the product of the single-particle Green’s function, G , and the screened Coulomb interaction, W , i.e., $\Sigma = iGW$. In common practice, one initiates a *GW* calculation by first solving the Kohn–Sham equation using a DFT method of choice to generate the starting orbitals and energies,

$$\left[-\frac{1}{2}\nabla^2 + V_{\text{ion}} + V_H + V_{\text{XC}}^{\text{KS}} \right] \phi_k^{\text{KS}} = \epsilon_k^{\text{KS}} \phi_k^{\text{KS}} \quad (1)$$

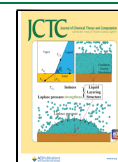
where V_{ion} , V_H , and V_{XC} are the ionic, Hartree, and exchange–correlation potential terms, respectively, and ϕ_k^{KS} and ϵ_k^{KS} are the k -th orbital and energy eigenpair. To obtain QP wave functions and energies, one starts by setting $\phi_k^{\text{QP}} = \phi_k^{\text{KS}}$ and $\epsilon_k^{\text{QP}} = \epsilon_k^{\text{KS}}$ and then solves the analogous Dyson equation,

Received: January 21, 2026

Revised: April 2, 2026

Accepted: April 2, 2026

Published: April 7, 2026



$$\left[-\frac{1}{2}\nabla^2 + V_{\text{ion}} + V_H + \Sigma(\varepsilon_k^{\text{QP}}) \right] \phi_k^{\text{QP}} = \varepsilon_k^{\text{QP}} \phi_k^{\text{QP}} \quad (2)$$

for ϕ_k^{QP} and $\varepsilon_k^{\text{QP}}$ to self-consistency.^{37–39} For many practical applications, it is sufficient to solve the equation in a single pass, possibly from a preoptimized starting point.⁴⁰ This is referred to as the G_0W_0 approximation, and this is what we use throughout the paper with the zero subscript omitted for clarity.

Evaluating the self-energy operator is costly and can be tackled by one of two strategies broadly defined as “deterministic” and “stochastic”. In deterministic *GW*, the overall cost is dominated by computing the inverse dielectric ε^{-1} and Σ operator matrix elements, requiring one to evaluate many integrals and summations over valence-conduction orbital pairs; this formally scales as $O(N_e^4)$ for an N_e -electron molecule or periodic system. Considerable efforts have been directed toward improving this scaling: one can achieve $O(N_e^3 \log N_e)$ complexity by employing, for example, the space-time formulation and using the fast Fourier transform (FFT) to transform to and from the real space.^{41–44} Interpolative density fitting⁴⁵ also achieves cubic complexity, potentially with smaller prefactors than the real space-time methods. The stochastic pseudobands approach⁴⁶ can be used to reduce the size of the valence space needed to converge matrix elements of Σ even further, decreasing the overall scaling to $O(N_e^{2.4})$.

The developments described above have spurred increasing interest in performing large-scale *GW* calculations,^{16,19,20} and several massively parallel deterministic *GW* implementations have been benchmarked. Zhang et al. have recently demonstrated a portable GPU implementation in the Berkeley*GW* code capable of scaling efficiently to entire exascale architectures, achieving excellent time to solution (on the order of minutes) for the computation of quasi-particle (QP) energies in semiconductor systems containing up to 17,574 atoms in the simulation cell.²⁰ Yu and Govoni computed states of an interface model of Si and Si₃N₄ with up to 2376 atoms and 10,368 electrons on 10,368 V100 GPUs in ~578 min (summed total of *wstat* and *wfreq* steps) using the GPU-enabled WEST code.¹⁸ Wu et al. reported calculations on 13,824-atom, 13,824-electron LiH supercells on 4,49,280 SW26010Pro cores in 285 s using a massively parallel version of PWDFIT.¹⁹ Very recently, Vetsch et al. have performed nonequilibrium Green’s function calculations on hydrogen-passivated silicon nanoribbons with up to 42,240 atoms on 37,600 MI250X GPUs in 42 s per iteration using a novel self-consistent *GW* algorithm with domain decomposition, implemented in their QuaTrEx code.⁴⁷

For systems containing thousands of atoms or more, one can evaluate the self-energy operator using a stochastic *GW* formulation at greatly reduced cost, as detailed by our previous works.^{40,48–53} Here, we briefly summarize the main features of the method. First, we evaluate the self-energy operator in the time domain to exploit the direct product computation of $\Sigma(t)$ from Green’s function G and screened Coulomb potential W ; we Fourier transform the resulting $\Sigma(t)$ to $\Sigma(\omega)$ only in the final stage of the calculation. Second, we invoke the stochastic Resolution of Identity (sRoI)^{48,50,54} and define random orbital functions to use as bases for evaluating the Green’s function G and effective polarization W . We then computed the expectation values of these operators using real-time propagation and accumulated statistical averages over products

of random samples. This is the main ingredient of stochastic *GW*, and it has the advantage of allowing one to decouple the spatial and time dependence in the six-dimensional integrals needed to evaluate $\Sigma(t)$.⁴⁸ As a result, instead of requiring the full space of occupied and unoccupied orbitals and energies $\{\phi_\nu, \varepsilon_k\}$ (which typically number in the tens of thousands for a thousand-atom molecule), we, in effect, evaluate G and W using compact sets of stochastic linear combinations of the occupied or unoccupied orbital space. An additional benefit of sRoI is that computations in the stochastic bases can be done independently, enabling the critical path of the calculation to be made embarrassingly parallelizable. Third, we incorporate sparse stochastic compression^{50,55–60} in our stochastic time-dependent Hartree propagation algorithm^{48,61} for evaluating W . This enables computing components of W over a collection of randomly chosen short segments without needing a full spatial grid, reducing storage costs. Finally, instead of projecting each stochastic sample onto the full set of occupied orbitals $\{\phi_k^{\text{KS}}\}_{\text{occ}}$ from the preliminary DFT calculation (resulting in substantial I/O and computational costs), we filter these samples to generate occupied stochastic orbitals. We construct a filter from a Chebyshev polynomial expansion of the Kohn–Sham Hamiltonian.⁶² While the filtering approach has the disadvantage of requiring many terms to produce a sharp cutoff at the Fermi energy, we recently found that the expansion length of the filter can be greatly decreased⁵³ by relaxing the expansion to have zero weights inside the band gap (where no states are present); this, in turn, reduces the number of matrix-vector products needed to prepare the occupied stochastic orbitals.

The above framework enables a near-linear $O(N_e \log N_e)$ scaling stochastic *GW* algorithm, with costs dominated by performing FFTs on the spatial grid. While development of stochastic algorithms for computing electronic properties has lagged behind that of deterministic ones,²⁰ for QP energies stochastic *GW* is well suited for handling large systems at a much reduced computational cost compared to deterministic *GW*. Some large-scale stochastic calculations have been performed: Vlcek et al.⁵⁰ computed HOMO–LUMO gaps of Γ -point Diamond and Silicon supercells containing up to 2744 atoms and 10976 electrons in under 2000 core hours on an HPC cluster containing 144 nodes with 1728 Intel Xeon E5–2680v3@2.5 GHz processors. More recently, Brooks et al.¹⁷ used *StochasticGW* to compute twist-induced localized Moire states of bilayer phosphorene sheets containing up to 2708 atoms and 13,540 electrons. Both of these calculations were performed with a CPU-only version of the code without gapped filtering. In this paper, we assimilate the ideas described above in a new GPU-accelerated version of *StochasticGW*, which we showcase by computing QP energies of clusters containing upward of 10,001 atoms and 35,144 electrons on ~1000 GPUs with times-to-solution of less than 1 h.

IMPLEMENTATION DETAILS

Algorithmic Overview

Our stochastic *GW* implementation is similar to that described previously⁵⁰ and with the inclusion of the gapped filtering⁵³ technique. Here, we only summarize the key components of the algorithm; see the earlier works for a more detailed explanation of the methodology. A block diagram, shown in Figure 1, depicts the major portions of the code.

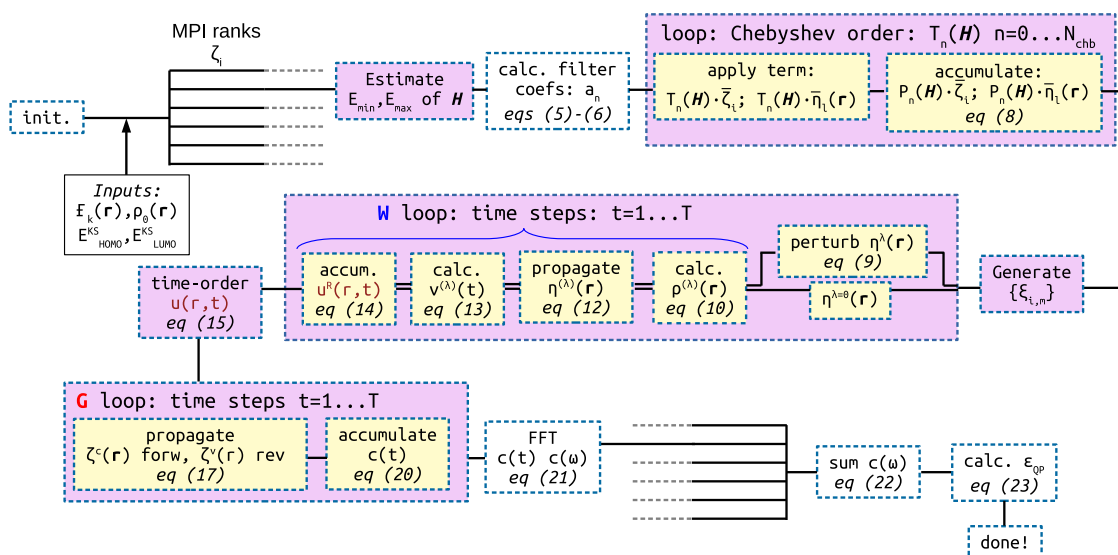


Figure 1. Block diagram of the main steps of the StochasticGW algorithm. Each MPI rank performs the same operations on a different set of data (see the text for details). Steps enclosed in shaded boxes were ported to GPUs.

The *StochasticGW* code requires as inputs: (1) coordinates of the atoms, (2) a pseudopotential for each atomic type, (3) the ground state density $\rho_0(\mathbf{r})$, (4) estimates of the energies of the highest occupied and lowest unoccupied molecular orbitals ($E_{\text{HOMO}}^{\text{KS}}$ and $E_{\text{LUMO}}^{\text{KS}}$, respectively), and (5) a spatial orbital ϕ_k for which the quasi-particle energy ϵ_k^{QP} is desired. Items (1)–(3) are needed to construct the Kohn–Sham Hamiltonian internally in *StochasticGW*; item (4) defines the cutoff region of the gapped filter. We obtain items (3)–(5) from a preliminary DFT calculation.

We begin by constructing a spectral filter to apply the Heaviside operator. The Heaviside operator (Θ) can be applied as an N -degree polynomial P_N in the Hamiltonian \hat{H} , where the individual terms are Chebyshev polynomials of the first kind $T_n(\hat{H})$:

$$\begin{aligned} \Theta(\mu - \hat{H}) &\simeq P_N(\hat{H}) \\ &= \sum_{n=0}^{N_{\text{chb}}} a_n T_n(\hat{H}), \end{aligned} \quad (3)$$

where μ is the chemical potential, a_n are expansion coefficients, and N_{chb} is the maximum degree of Chebyshev polynomial needed to build the filter. One requires the upper and lower spectral bounds of the Kohn–Sham Hamiltonian \hat{H}^{KS} to shift-and-scale the eigenvalue spectrum into the interval $[-1, 1]$; there are various schemes to obtain these bounds, but we find that one of the simplest, a shifted power iteration, works well for this purpose. Estimating the spectral bounds is needed if the upper bound is not provided by the preliminary DFT calculation or if one wants to perform stochastic GW calculations on a grid interpolated from the grid used in the initial DFT step. In the gapped-filtering method, rather than applying the cutoff of filter Θ at a specific value of μ , we instead apply it over the band gap containing μ , i.e., over $E_{\text{HOMO}}^{\text{KS}} \leq \mu \leq E_{\text{LUMO}}^{\text{KS}}$. Next, we need filter expansion coefficients a_n . Instead of using the canonical Chebyshev weight function $w(x) = (1 - x^2)^{-1/2}$, we define a modified function that has zero weighting inside the band gap,

$$\tilde{w}(x) = w(x)(\Theta(x_{\text{H}} - x) + \Theta(x - x_{\text{L}})) \quad (4)$$

where x_{H} and x_{L} are $E_{\text{HOMO}}^{\text{KS}}$ and $E_{\text{LUMO}}^{\text{KS}}$ mapped into $[-1, 1]$, respectively. We compute the coefficients by minimizing the functional

$$J = \int_{-1}^1 \left| \Theta(x - x_0) - \sum_{n=0}^{N_{\text{chb}}} a_n T_n(x) \right| \tilde{w}(x) dx \quad (5)$$

As described in more detail in ref 53, the N_{chb} filter coefficients a_n can be computed by solving the linear system $\mathbf{M}\mathbf{a} = \mathbf{b}$ with elements

$$M_{i,n} = \int_{-1}^1 \tilde{w}(x) T_i(x) T_n(x) dx \quad (6)$$

$$b_i = \int_{-1}^{x_{\text{H}}} \tilde{w}(x) T_i(x) dx \quad (7)$$

After obtaining the filter, we generate N_{ζ} random “white-noise” start orbitals, $|\zeta_i(\mathbf{r}, t=0)\rangle$, for the stochastic realization of G ; we dub these the Monte Carlo (MC) samples. For each, we also generate N_{η} additional white-noise orbitals, $|\bar{\eta}_{i,l}(\mathbf{r}, t=0)\rangle$, needed to calculate the action of the time-dependent effective interaction operator $W(t)$ on a vector related to each $|\zeta_i\rangle$. We apply the Heaviside filter to both sets of $\{|\zeta_i\rangle\}$ and $\{|\bar{\eta}_{i,l}\rangle\}$ orbitals in order to project them onto random linear combinations of the occupied orbital subspace $\{\phi_k^{\text{KS}}\}_{\text{occ}}$:

$$\begin{aligned} \zeta_i &= \sum_{n=0}^{N_{\text{chb}}} a_n T_n(\hat{H}) \bar{\zeta}_i \\ \eta_{i,l} &= \sum_{n=0}^{N_{\text{chb}}} a_n T_n(\hat{H}) \bar{\eta}_{i,l} \end{aligned} \quad (8)$$

Subsequently, we evaluated the diagonal time-dependent self-energy matrix element, $\langle \phi_k | \Sigma(t) | \phi_k \rangle$, in two phases. In the first, we use linear-response time-dependent Hartree⁴⁸ to compute the action of the retarded polarization interaction W^{R} on the occupied states. Algorithmically, we create a perturbed copy of each $\eta_{i,l}$,

$$\eta_{i,l}^{\lambda}(\mathbf{r}, t=0) = e^{-i\lambda v_{\text{pert}}(\mathbf{r})} \eta_{i,l}(\mathbf{r}) \quad (9)$$

where $v_{\text{pert}}(\mathbf{r}) = \int v(\mathbf{r}, \mathbf{r}') \bar{\zeta}(\mathbf{r}') \phi(\mathbf{r}') d\mathbf{r}'$ and λ has a small value of $10^{-4} E_h^{-1}$.

In the W^R propagation phase, at each time step, we first evaluate the time-dependent density,

$$\rho_i^{(\lambda)}(\mathbf{r}, t) = C_{\text{norm}} \frac{2}{N_\eta} \sum_{l \leq N_\eta} |\eta_{i,l}^{(\lambda)}(\mathbf{r}, t)|^2 \quad (10)$$

where the superscripted (λ) indicates performing the action for both perturbed and unperturbed copies of the $\{\eta_{i,l}\}$, and C_{norm} is a constant which normalizes the density to the number of electrons. We then compute the difference between time-dependent Hartree potentials at times t and $t = 0$ from their respective densities,

$$v_{H_i}^{(\lambda)}(\mathbf{r}, t) \equiv v_{H_i}^{(\lambda)}[\rho_i^{(\lambda)}(t)](\mathbf{r}) - v_{H_i}^{(\lambda)}[\rho_i^{(\lambda)}(t=0)](\mathbf{r}) \quad (11)$$

and we propagate both perturbed and unperturbed $\{\eta_{i,l}\}$ in time under the action of the time-dependent Hamiltonian,

$$|\eta_{i,l}^{(\lambda)}(t + dt)\rangle = e^{-iH^{(\lambda)}(t)dt} |\eta_{i,l}^{(\lambda)}(t)\rangle \quad (12)$$

To compute the causal response function $u_i^R(\mathbf{r}, t)$, we must first project the time-dependent Hartree potential into the sparse stochastic basis $\{\xi_{i,m}\}$,

$$v_{i,m}^{(\lambda)}(t) = \langle \xi_{i,m}(\mathbf{r}) | v_{H_i}^{(\lambda)}(\mathbf{r}, t) \rangle \quad (13)$$

from which we then accumulate the causal response function, $u_i^R(\mathbf{r}, t)$, as the difference of the perturbed and unperturbed time-dependent potentials, summed over the set of $\{\xi_{i,m}\}$,

$$u_i^R(\mathbf{r}, t) = \frac{1}{N_\xi} \sum_{m \leq N_\xi} \xi_{i,m}(\mathbf{r}) \frac{v_{i,m}^{(\lambda)}(t) - v_{i,m}^{\lambda=0}(t)}{\lambda} \quad (14)$$

Once the propagation is complete, we time-order⁶³ the accumulated $u_i^R(\mathbf{r}, t)$ to produce the effective polarization potential $u_i(\mathbf{r}, t)$:

$$\begin{aligned} u_i^R(\mathbf{r}, t) &\rightarrow u_i^R(\mathbf{r}, \omega) = \int_0^\infty e^{-1/2\gamma^2 t^2} e^{i\omega t} u_i^R(\mathbf{r}, t) dt \\ \rightarrow u_i(\mathbf{r}, \omega) &= \begin{cases} (u_i^R(\mathbf{r}, \omega))^* & \omega < 0 \\ u_i^R(\mathbf{r}, \omega) & \omega \geq 0 \end{cases} \\ \rightarrow u_i(\mathbf{r}, t) &= \frac{1}{2\pi} \int_{-\infty}^\infty e^{-i\omega t} u_i(\mathbf{r}, \omega) d\omega, \end{aligned} \quad (15)$$

where γ is a damping factor.

In the second phase, we evaluate the action of the Green's function iG on each $\bar{\zeta}_i$. We first define the occupied and unoccupied components of each stochastic sample, ζ_i^v and ζ_i^c , at time $t = 0$,

$$\begin{aligned} \zeta_i^v(\mathbf{r}) &= P_N(\hat{H}) \bar{\zeta}_i(\mathbf{r}) \\ \zeta_i^c(\mathbf{r}) &= (I - P_N(\hat{H})) \bar{\zeta}_i(\mathbf{r}) = \bar{\zeta}_i(\mathbf{r}) - \zeta_i^v(\mathbf{r}) \end{aligned} \quad (16)$$

where we obtained ζ_i^v by applying the Chebyshev filter in eq 8. We then propagate the occupied component backward in time while simultaneously propagating the unoccupied component forward in time:

$$\begin{aligned} \zeta_i^v(\mathbf{r}, t - dt) &= -e^{-iH_0 dt} \zeta_i^v(\mathbf{r}, t) \\ \zeta_i^c(\mathbf{r}, t + dt) &= e^{-iH_0 dt} \zeta_i^c(\mathbf{r}, t) \end{aligned} \quad (17)$$

The stochastic realization of Green's function is

$$iG(\mathbf{r}, \mathbf{r}', t) = \frac{1}{N_\zeta} \sum_{\zeta} \zeta(\mathbf{r}, t) \bar{\zeta}(\mathbf{r}') \quad (18)$$

Having obtained $u(\mathbf{r}, t)$ and $\zeta(\mathbf{r}, t)$ from the first and second phases, respectively, the diagonal self-energy matrix element for orbital ϕ_k becomes

$$\begin{aligned} \langle \phi_k | \Sigma(t) | \phi_k \rangle &= \int \int \phi_k(\mathbf{r}) iG(\mathbf{r}, \mathbf{r}', t) W(\mathbf{r}, \mathbf{r}', t) \phi_k(\mathbf{r}) d\mathbf{r} d\mathbf{r}' \\ &= \frac{1}{N_\zeta} \sum_{\zeta} \int \int \phi_k(\mathbf{r}) \zeta(\mathbf{r}, t) W(\mathbf{r}, \mathbf{r}', t) \bar{\zeta}(\mathbf{r}') d\mathbf{r} d\mathbf{r}' \\ &= \frac{1}{N_\zeta} \sum_{\zeta} \int \phi_k(\mathbf{r}) \zeta(\mathbf{r}, t) u(\mathbf{r}, t) d\mathbf{r} \end{aligned} \quad (19)$$

Numerically, instead of evaluating eq 19 directly, we compute a time correlation function associated with each individual sample $\zeta_i(\mathbf{r}, t)$,

$$c_i(t) \equiv \int \phi_k(\mathbf{r}) \zeta_i(\mathbf{r}, t) u_i(\mathbf{r}, t) d\mathbf{r} \quad (20)$$

and we apply a Fourier transform to obtain the frequency-dependent form,

$$c_i(\omega) = \int_{-\infty}^\infty e^{-1/2\gamma^2 t^2} e^{i\omega t} c_i(t) dt \quad (21)$$

Finally, we accumulate the frequency-resolved self-energy matrix element by summing over the N_ζ samples

$$\langle \phi_k | \Sigma(\omega) | \phi_k \rangle = \frac{1}{N_\zeta} \sum_{i=1}^{N_\zeta} c_i(\omega) \quad (22)$$

and we obtain quasi-particle energy, ϵ_k^{QP} , by solving

$$\epsilon_k^{QP} = \epsilon_k^{KS} + \langle \phi_k | X + \Sigma(\omega = \epsilon_k^{QP}) - V_{XC} | \phi_k \rangle \quad (23)$$

where X is the sROI^{48,50} realization of the Fock exchange operator in the basis of $\{\eta_{i,l}\}$ and all other quantities have been previously defined.

Scaling of the Method

A key aim of our stochastic GW formulation is to achieve computational scaling that grows slowly, ideally linearly, with respect to the system size. The most numerically intensive portions of the algorithm apply matrix-vector products repeatedly to the set of $\{\eta_{i,l}\}$ during the filtering and propagation phases. Here, one applies either the Kohn–Sham Hamiltonian, \hat{H}^{KS} , or the evolution operator, $e^{-i\hat{H}(t)\Delta t}$, to a set of vectors with each having length $N_g = N_x N_y N_z$. We apply matrix-vector products in a Fourier grid representation, whereby FFT pairs switch between position and momentum representations. Applying the individual kinetic and potential energy operators scales as $O(N_g)$, but the cost of each Hamiltonian/evolution operation is dominated by the $O(N_g \log_2 N_g)$ FFT cost. Thus, to filter the $N_\zeta N_\eta$ starting orbitals, one applies a length N_{chb} filter at a cost of $O(N_\zeta N_\eta N_{chb} N_g \log_2 N_g)$ operations. Likewise, propagating the full set of $\{\eta_{i,l}\}$ for N_τ time steps has a cost scaling as $O(N_\zeta N_\eta N_\tau N_g \log_2 N_g)$. Accumulating $u^R(\mathbf{r}, t)$ costs

$(N_{\xi}N_{\zeta}N_{\tau}N_{\eta}f_g)$, where f_g is the fractional length of each of the fragmented stochastic functions $\{\xi_{i,m}\}$ relative to the full spatial grid length, N_g .

For tackling quasi-particle energies of large molecules, it is important to consider the dependence of each parameter on system size. The number of MC samples, N_{ζ} , and the number of occupied stochastic orbitals, N_{η} , determine the statistical accuracy of the QP energies and do not increase with system size (N_{η} actually decreases with increasing system size due to self-averaging). The number of time steps, N_{τ} , determines the energy resolution of $\Sigma(\omega)$ and is also independent of the system size. The number of grid points, N_g , while cubic in dimension ($N_g = N_x N_y N_z$), grows linearly overall with system size due to spatial packing of atoms in 3-dimensional space. For accumulating $u^R(\mathbf{r}, t)$, the statistical error does increase with the ratio $\frac{N_g}{N_{\xi}}$. This means that one must simultaneously

increase the number of stochastic segments N_{ξ} as the grid size N_g increases to prevent growth of errors. However, in the sparse stochastic basis, the cost increase from requiring larger N_{ξ} can be offset by decreasing the fractional length f_g of each segment $\{\xi_{i,m}\}$ (i.e., by using more ξ vectors but making them “sparser”).⁵⁰ Finally, the number of Chebyshev coefficients, N_{chb} , needed to fit the gapped filter depends on the width of the band gap relative to the spectral width of the Kohn–Sham Hamiltonian. The value of N_{chb} needed to accurately fit the filter does increase with system size due to the larger spectral width of \hat{H}^{KS} , but this can be mitigated by applying a kinetic energy cutoff. In summary, as long as care is taken to manage the growth of the N_{ξ} and N_{chb} parameters accordingly, one can achieve near-linear scaling with system size in stochastic GW calculations.

GPU Implementation

The original `StochasticGW` code (through v.2.0) is written in Fortran 90 and parallelized using the Message Passing Interface (MPI). A key feature of `StochasticGW` is that the N_{ζ} Monte Carlo samples can be processed independently of one another, resulting in an embarrassing parallelism over large portions of the algorithm. Additionally, the code contains an option to extend the MPI-level parallelism over the N_{η} occupied stochastic functions at the cost of an additional call to `mpi_allreduce()` at each time step (needed to compute the time-dependent density $\rho(\mathbf{r}, t)$). In the original implementation, operations over grid points are performed in serial. For systems containing thousands of atoms or more, the grids are large enough that these operations become significant serial bottlenecks, which motivated us to develop a GPU port to handle them in parallel.

In the GPU implementation, we retain the idea of processing each of the N_{ζ} MC samples with a separate MPI rank, but the N_{η} occupied stochastic functions per sample reside on the same MPI rank so that MPI calls are not needed at each time step to evaluate the time-dependent density $\rho(\mathbf{r}, t)$. The GPU code utilizes kernels written using OpenACC directives and calls to specialized libraries (cuRAND and cuFFT) when needed. To maximize efficiency, attention must be given to minimize the amount of data transferred between the host CPU and each GPU and to organize the computational workload to expose as much parallelism to the GPU as possible. To this end, we performed several optimizations.

First, we structured the arrays containing the stochastic orbitals so that each kernel can process the orbitals in a single-

instruction multiple-data (SIMD) fashion. Many of the operations in the filtering and propagation cycles, such as applying the kinetic and potential energy operators, are simple element-wise array multiplications, which are highly vectorizable on GPU hardware. In each case that follows, we offload the arrays once onto a single GPU, retrieving the result only after the full set of filtering or propagation iterations. For the filtering cycle, this means that on each MPI rank, we pack the $\bar{\zeta}_i$ associated with rank i along with its set of $\{\bar{\eta}_{i,l}\}; l = 1 \dots N_{\eta}$ orbitals into an array of size $(N_g \times N_{\eta} + 1)$. For the propagation cycle involving the set of $\{\eta_{i,l}\}$, the perturbed and unperturbed copies can be processed in parallel, so we pack both copies into an array of size $(N_g \times N_{\eta} \times 2)$. The MC sample ζ_i cannot be propagated in parallel with the $\{\eta_{i,l}\}$ here since we accumulate $c_i(t)$ (which requires effective polarization potential $u(\mathbf{r}, t)$ in eq 20) in tandem with the ζ_i propagation. However, since reverse-time propagation of ζ_i and forward-time propagation of its orthogonal complement are operationally identical (other than a difference in sign), we can pack these functions into an array of size $(N_g \times 2)$ and achieve a parallel performance boost for propagation of ζ_i as well.

Not all operations in the filtering and propagation steps are trivial to vectorize. Normalizations appear periodically in each of the filtering, propagation, and spectral estimation stages; each requires the summation over values defined over N_g grid points. For instance, for normalizations performed in the $\{\eta_{i,l}\}$ propagation cycle, at most only $2N_{\eta}$ operations can be performed in parallel instead of $2N_{\eta}N_g$. For the largest systems in this work, $N_{\eta} \sim 8$, while $N_g \sim 1.6 \times 10^8$, meaning that the benefits of having many parallel threads are largely lost in each normalization kernel call. To solve this, we divided the N_g grid points into short segments of length L and perform a reduction ($O(\log_2 L)$ operations) on each segment. This optimization partitions the array into blocks that fit into the L1 cache and allows us to parallelize sums over grid points over $\frac{N_g}{L}$ threads at the cost of having to perform an atomic add by each thread after the sum over each segment has been accumulated. The optimal value of L is hardware-dependent; on NVIDIA A100 GPUs, we achieved the best performance with $L \sim 256$. In this manner, we achieve an overall parallelism of up to $2N_{\eta} \frac{N_g}{L}$ threads in normalization calls.

Second, the main computations needed to accumulate $u(\mathbf{r}, t)$ have also been ported to the GPU. We generate the $\{\xi_{i,m}\}$ basis via calls to the cuRAND library, and we compute the overlaps $\langle \xi | u^R(t) \rangle$ on-the-fly during the $\{\eta_{i,l}\}$ propagation phase in a segmented fashion similar to the one described above for normalization. Here, we multiply two arrays of sizes $(N_{\xi} \times N_g f_g)$ and $(N_g f_g \times 2)$ parallelized over $2N_{\xi} \frac{N_g f_g}{L}$ threads, where each ξ function is processed in segments of length $L = 32$, and the factor of 2 again arises from performing the unperturbed- η and perturbed- η propagations in parallel. Finally, we perform the time-ordering operation to convert the resulting $u^R(\mathbf{r}, t)$ to $u(\mathbf{r}, t)$ by calling the cuFFT library before and after an OpenACC kernel was used for performing the complex conjugation step.

Utilities

The newest (3.0) version of `StochasticGW` is freely available⁶⁴ on GitHub and includes several utilities to aid researchers in preparing inputs for the code:

The `dft2sgw` utility reads and preprocesses results from a preliminary DFT calculation. This utility requires a DFT output file and a set of `.cube` files as input; `dft2sgw` prepares an input file, `sgwinp.txt`, containing atomic coordinates, HOMO and LUMO energies (for gapped filtering), the spatial charge density, and a requested set of orbitals for the system of interest. `dft2sgw` also has a functionality, enabled via the `FFTW`⁶⁵ library, to up- or down-sample the orbital/density spatial grid from the preliminary DFT calculation in case a different grid for the stochastic *GW* step is desired. The utility currently supports Quantum ESPRESSO,^{25,66} the Real-Space Multigrid (RMG),^{67,68} DFT code, and CP2K³⁵ (but note that `StochasticGW` does not yet include pseudopotential support for CP2K).

`StochasticGW` also includes two utilities, `plotfilter.py` and `plotorbital.py`, which use the `Matplotlib`⁶⁹ python package to generate plots related to stochastic *GW* calculations:

The `plotfilter.py` utility plots the filter and depicts the log of the magnitudes of the filter coefficients and is useful for checking the quality of the Chebyshev expansion of the filter.

The `plotorbital.py` utility visualizes the atomic coordinates, spatial orbitals, and charge density contained in `sgwinp.txt`; this feature allows one to quickly select orbitals of interest for a subsequent `StochasticGW` calculation.

NUMERICAL EXPERIMENTS

We now test our implementation of `StochasticGW` by computing QP energies of a series of nonperiodic nanoclusters, $\text{Si}_{293}\text{H}_{172}$, $\text{Si}_{703}\text{H}_{300}$, $\text{Si}_{5031}\text{H}_{1172}$, $\text{Si}_{7745}\text{H}_{1572}$, $\text{Si}_{8381}\text{H}_{1620}$. We constructed each cluster from a uniformly expanded silicon superlattice of size $15 \times 15 \times 15$ using the experimental unit cell parameter for silicon ($a = b = c = 10.26$ Bohr) corresponding to the diamond cubic structure with an eight-atom unit cell.⁷⁰ We then shifted the coordinate origin to the geometric center of the superlattice and applied a spherical truncation, retaining only Si atoms within 20–70 Bohr of the origin. The truncated cluster is passivated with hydrogen atoms to saturate dangling bonds, and the resulting structure is relaxed to its equilibrium geometry using the MMFF94 force field⁷¹ as implemented in Open Babel software.⁷² Figure 2 depicts the largest cluster, $\text{Si}_{8381}\text{H}_{1620}$.

We performed the initial periodic DFT calculations to generate the orbitals and charge densities using the RMG^{67,68} DFT code. The DFT Hamiltonian uses the GGA PBE exchange-correlation functional with Troullier-Martins⁷³ norm-conserving pseudopotentials. For each system, we performed the calculation on the Γ k-point in a primitive cubic cell with isotropic sampling. Each cluster is separated from its periodic image by a vacuum layer of 11–17 bohr. The initial DFT step provides the energy estimates $E_{\text{HOMO}}^{\text{KS}}$ and $E_{\text{LUMO}}^{\text{KS}}$ used to define the gapped filter for the *GW* step. Cell and grid parameters, along with HOMO and LUMO energies, are given in Table 1.

We then used `StochasticGW` to compute ϵ_{QP} for the HOMO and LUMO orbitals of each system. The Kohn–Sham Hamiltonian in `StochasticGW` uses the same pseudopotentials and grids as the previous DFT step; here, we employ

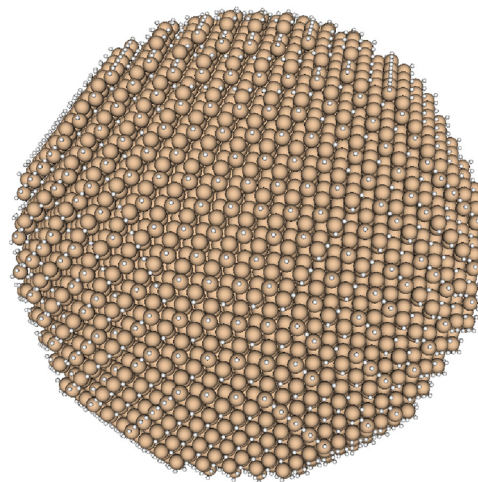


Figure 2. $\text{Si}_{8381}\text{H}_{1620}$ cluster. Silicon and hydrogen atoms are shown as brown and white spheres, respectively.

the PBE functional⁷⁴ as implemented in the LibXC⁷⁵ library and apply an energy cutoff of 28 hartree to the kinetic energy operator. This cutoff value is more conservative than the default kinetic energy cutoffs selected by RMG-DFT for each grid 1. In the filtering step, for the largest system we studied, $\text{Si}_{8381}\text{H}_{1620}$, the energy difference $E_{\text{LUMO}}^{\text{KS}} - E_{\text{HOMO}}^{\text{KS}}$ is $\sim 0.11\%$ of the full spectral range of \hat{H}^{KS} . Even though the cutoff is spread over the full band gap, it is still sharp enough to require 8192 Chebyshev terms to reduce the Gibbs oscillations to negligible levels outside of the band gap (Figure 3). Although this filter length is an order of magnitude larger than that used in our recent calculation on the naphthalene molecule ($N_{\text{chb}} = 450$),⁵³ it is still less than the lengths required in earlier calculations performed on much smaller systems without gapped filtering ($N_{\text{chb}} \sim 18,000$).⁴⁹ The oscillations present inside of the band gap (inset in Figure 3) arise due to the near singularity of matrix \mathbf{M} in eq 7, but these features are confined to the band gap where there are no eigenstates, and therefore they do not affect the occupied space projection.

For each calculation, we averaged 1024 Monte Carlo samples, which is sufficient to achieve a statistical error of better than 0.03 eV for all QP energies. The number of time steps, N_{τ} , is controlled internally by the energy-broadening parameter, γ , which we apply when Fourier transforming the self-energy element from the time domain to the frequency domain,

$$\langle \phi_k | \Sigma(\omega) | \phi_k \rangle = \int \langle \phi_k | \Sigma(t) | \phi_k \rangle e^{-\gamma^2 t^2 / 2} e^{i\omega t} dt \quad (24)$$

We use a time step size of $\Delta t = 0.05 E_{\text{F}}^{-1} \hbar$ in the split-operator propagation of the orbitals. The number of time steps to obtain a desired energy resolution is $N_{\tau} \approx \frac{3}{\gamma \Delta t}$; for all calculations in this work, we set $\gamma = 0.06 E_{\text{F}} \hbar^{-1}$, which yields a propagation length of $N_{\tau} = 1000$ time steps over 50 atomic time units. Numbers of occupied stochastic orbitals (N_{η}) and segmented stochastic functions (N_{ϵ}), along with the fractional grid lengths (f_g) for the latter, are chosen at values similar to those in previous works.^{17,49,50} Input parameters are summarized in Table 2. All calculations were run on 256 GPU nodes of NERSC-Perlmutter; each node contained one AMD EPYC 7763 processor running at 2.5 GHz and 4 NVIDIA A100 GPUs.

Table 1. Details of Preliminary DFT Calculations on Each Cluster, Including Numbers of Electrons (N_e), Points in the Spatial Grid (N_g), Grid Spacing (Δ_g , Bohr), Kinetic Energy Cutoffs (E_{cut}^k , Hartrees), along with Resulting HOMO and LUMO Energies and Band Gaps (eV)

| system | N_e | N_g | Δ_g | E_{cut}^k | $E_{\text{HOMO}}^{\text{KS}}$ | $E_{\text{LUMO}}^{\text{KS}}$ | KS band gap |
|--------------------------------------|--------|------------------|------------|--------------------|-------------------------------|-------------------------------|-------------|
| Si ₂₉₃ H ₁₇₂ | 1344 | 128 ³ | 0.4429 | 25.2 | -3.259 | -1.526 | 1.733 |
| Si ₇₀₅ H ₃₀₀ | 3120 | 128 ³ | 0.5167 | 18.5 | -1.928 | -0.457 | 1.471 |
| Si ₅₀₃₁ H ₁₁₇₂ | 21,296 | 256 ³ | 0.5000 | 19.7 | -1.705 | -0.776 | 0.928 |
| Si ₇₇₄₅ H ₁₅₇₂ | 32,552 | 256 ³ | 0.5400 | 16.9 | -0.986 | -0.121 | 0.865 |
| Si ₈₃₈₁ H ₁₆₂₀ | 35,144 | 256 ³ | 0.5600 | 15.7 | -1.095 | -0.245 | 0.851 |

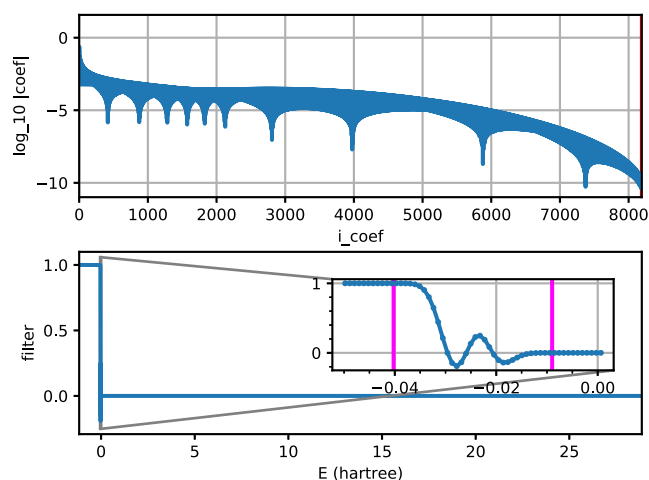


Figure 3. (Top) Plot of the log of the absolute magnitudes of Chebyshev coefficients used to construct the gapped filter for the Si₈₃₈₁H₁₆₂₀ cluster. (Bottom) Reconstructed filter, where the inset shows an expansion of the region of the band gap. The purple vertical lines in the inset indicate the positions of $E_{\text{HOMO}}^{\text{KS}}$ and $E_{\text{LUMO}}^{\text{KS}}$.

Table 2. Parameters Used in Stochastic GW Calculations

| description | parameter | value |
|---|--------------------|--------|
| number of Monte Carlo samples | N_c | 1024 |
| number of occupied stochastic orbitals | N_η | 8 |
| number of segmented stochastic functions | N_g | 10,000 |
| grid fraction of each segmented function | f_g | 0.003 |
| number of Chebyshev polynomials in filter | N_{chb} | 8192 |
| damping parameter (Hartrees) | γ | 0.06 |
| kinetic energy cutoff (Hartrees) | E_{cut}^k | 28.0 |

Figure 4 plots the QP energies of the HOMO and LUMO and their difference for each system; these values are also listed

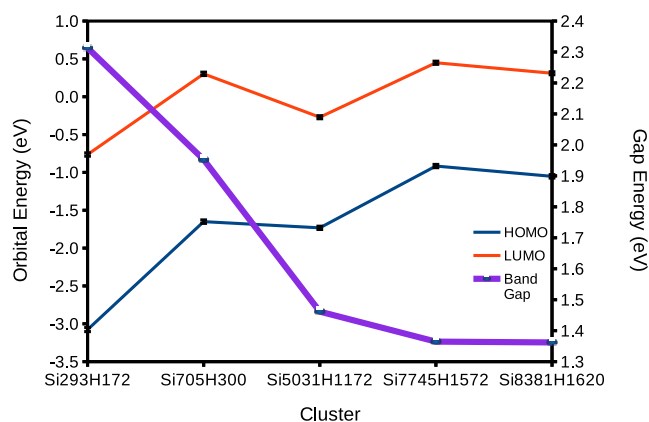


Figure 4. QP orbital energies and band gaps are given for each cluster.

in Table 3. The statistical errors in the MC energies are shown as error bars in the HOMO and LUMO traces and are small compared with the magnitudes of the energies. Moreover, comparing the band gaps across the five clusters, the band gaps show convergent behavior toward ~ 1.36 eV, suggesting that the largest clusters are approaching the bulk limit for our choice of density functional and pseudopotential.

Table 3 also lists the wall times for all calculations. The two smaller clusters have comparable times of 800 ± 40 s, and the larger three clusters have timings of 2700 ± 120 s, where the main factor behind the difference is the size of the spatial grid (128^3 vs 256^3 for the smaller and larger clusters, respectively). For a given spatial grid, one expects an increase in runtime for the larger systems for two reasons: first, the potential energy terms containing the pseudopotential contribution must be applied via atomic operations on grid points where pseudopotentials for neighboring atoms overlap. Second, the higher density of states in larger systems causes the spectral range estimation to converge more slowly. However, these differences are smaller than the variation in performance in MPI and I/O operations that occur over the large scale of these runs.

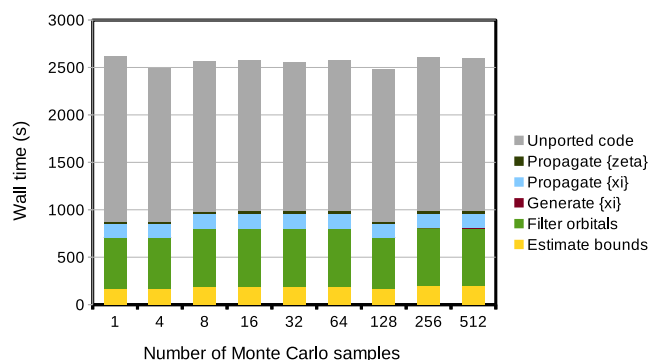
We also performed a set of tests to measure the efficiency of parallelizing over Monte Carlo samples for the HOMO calculation of the largest Si₈₃₈₁H₁₆₂₀ system (Figure 5). Here, the number of samples, N_c , is set equal to both the number of MPI ranks and the number of GPUs; therefore, this test is a measure of the weak scaling of the code. Note that the total runtimes are all ~ 2500 s, similar to the full 1024-sample run, demonstrating nearly ideal scaling with number of samples.

From Figure 5, one can see that the portions of the code that have been ported to GPUs collectively account for $\sim 38\%$ of the total runtime. The remaining, unported portion of the code includes I/O operations, preparation of the grid representation of the Hamiltonian, constructing the $\{\zeta\}$ and $\{\xi\}$ orbitals prior to filtering, solving the linear system to generate the filter coefficients,⁵³ and collecting and postprocessing the samples to produce the final QP energies.

We also measured the speedup factors for the individual GPU-ported sections of the code relative to their CPU counterparts. A full-scale run of the original CPU code is not possible for the larger clusters due to the wall time limit on NERSC-Perlmutter, so we instead performed several comparative tests focusing on individual portions of the code. These tests are summarized in Table 4, where we compared the runtimes of each section on a single NVIDIA A100 GPU with those on a single AMD EPYC 7763 core. While these tests compare time-to-solution for a single MC sample, the timings would be similar for a full $N_c = 1024$ -sample run with each sample running on a separate MPI rank. We also measured the GPU memory cost per MPI rank in single MC sample runs with all GPU-ported sections of the code

Table 3. Results of Stochastic GW Calculations on Each Cluster, Including Quasi-Particle Energies and Band Gaps (eV) and Calculation Wall Times (s)

| system | $E_{\text{HOMO}}^{\text{QP}}$ | $E_{\text{LUMO}}^{\text{QP}}$ | QP band gap | $t_{\text{HOMO}}^{\text{wall}}$ | $t_{\text{LUMO}}^{\text{wall}}$ |
|--------------------------------------|-------------------------------|-------------------------------|-------------|---------------------------------|---------------------------------|
| Si ₂₉₃ H ₁₇₂ | -3.077 ± 0.027 | -0.764 ± 0.022 | 2.313 | 836 | 770 |
| Si ₇₀₅ H ₃₀₀ | -1.650 ± 0.023 | 0.302 ± 0.022 | 1.953 | 835 | 788 |
| Si ₅₀₃₁ H ₁₁₇₂ | -1.732 ± 0.021 | -0.271 ± 0.020 | 1.462 | 2609 | 2617 |
| Si ₇₇₄₅ H ₁₅₇₂ | -0.916 ± 0.022 | 0.449 ± 0.021 | 1.365 | 2702 | 2688 |
| Si ₈₃₈₁ H ₁₆₂₀ | -1.052 ± 0.023 | 0.310 ± 0.029 | 1.362 | 2669 | 2812 |

**Figure 5.** Wall times spent in each portion of the code for calculations on the HOMO state of Si₈₃₈₁H₁₆₂₀, with different numbers of Monte Carlo samples.**Table 4. Timings and Speedups of GPU Portions of StochasticGW Relative to the CPU Portions, for Calculation on the HOMO State of Si₈₃₈₁H₁₆₂₀**

| portion | t_{CPU} (s) | t_{GPU} (s) | $t_{\text{CPU}}/t_{\text{GPU}}$ |
|---|----------------------|----------------------|---------------------------------|
| propagate $\{\zeta\}$ | 4947 | 31 | 160 |
| propagate $\{\eta\}$ | 37,711 | 153 | 246 |
| generate $\{\xi\}$ | 662 | 0.08 | 8764 |
| filter $\{\zeta\}, \{\eta\}$ | 9796 | 199 | 49 |
| estimate $[E_{\text{min}}, E_{\text{max}}]$ | 26,364 | 191 | 138 |
| unported code | 1237 | 1 | 1 |
| total (incl. unported) | 81,292 | 1811 | 45 |

activated. For these tests, performed on the Si₂₉₃H₁₇₂ and Si₈₃₈₁H₁₆₂₀ systems, we measured values of 2.3 and 18 GB, respectively.

We ran each test for the HOMO state of Si₈₃₈₁H₁₆₂₀ on a single MC sample with other parameters the same as in Table 2 except that here we decreased the filter length by a factor of 16 (to $N_{\text{chb}} = 512$) to reduce the CPU time needed for this step. As Table 4 shows, the GPU filtering step achieves a $\sim 50\times$ speedup over its CPU counterpart. The propagation and spectral estimation steps achieve even higher speedups of 150–250 \times . This is due not only to porting the routines to GPUs but also to optimizations that were not present in the CPU code, such as premultiplying potential energy factors before offloading them to the GPU. The step to generate the $\{\xi\}$ segments, while requiring much less time than the propagation and filtering steps, exhibited the largest performance improvement resulting from replacing the serial calls to the KISS random number generator with calls to the cuRAND library. Finally, the last two rows of Table 4 list the timings of the unported code and the sum of all timings, including unported portions of the code, showing that the overall speedup of the GPU implementation of StochasticGW is $\sim 45\times$ that of the CPU code.

CONCLUSION

In this work, we describe a new implementation of the StochasticGW code. Our code utilizes the stochastic Resolution of the Identity (sROI)^{48,50} technique, which allows one to decouple the main steps of the GW method into independent, statistical operations that can be performed massively in parallel. In deterministic GW methods, the cost is dominated by computing matrix elements over indices representing the occupied and unoccupied orbitals. In contrast, in the stochastic method, the cost depends on operations over the full spatial grid and accumulating a sufficient number of Monte Carlo samples to achieve a desired statistical accuracy. Therefore, compared with deterministic GW, the cost of stochastic GW grows much more slowly with respect to the size of the molecule or material system of interest.

Motivated by the large-scale parallelism available in modern GPU hardware, we ported the major computational motifs of the algorithm to GPUs. These include estimating the spectral width of the Kohn–Sham Hamiltonian, filtering the initial orbitals by projecting onto an occupied subspace, and propagating the orbitals under the influence of a time-dependent Hamiltonian. Each of these steps is applied via a sequence of vectorized OpenACC kernels and calls to GPU-optimized FFT libraries.

We showcased the GPU implementation by computing the quasi-particle energies of the HOMO and LUMO orbitals of five hydrogen-passivated silicon clusters. The band gaps show convergent behavior toward a bulk-like limit at ca. 1.36 eV. QP calculations on the largest system, Si₈₃₈₁H₁₆₂₀, with 10,001 atoms and 35,144 electrons, can be completed in only ~ 45 min with the workload partitioned with one MC sample per GPU. For this system, the GPU version of StochasticGW achieves roughly 45 \times speedup in time-to-solution relative to the CPU version over the entire execution of the code. This work opens the way for computing the QP energies of even larger systems.

AUTHOR INFORMATION

Corresponding Author

Daniel Neuhauser – Department of Chemistry and Biochemistry, University of California, Los Angeles, California 90095, United States; Email: dxn@ucla.edu

Authors

Phillip S. Thomas – National Energy Research Scientific Computing Center (NERSC), Lawrence Berkeley National Laboratory, Berkeley, California 94720, United States;

orcid.org/0009-0007-3794-0150

Minh Nguyen – Center for Computing Research, Sandia National Laboratories, Albuquerque, New Mexico 87123, United States

Dimitri Bazile – Department of Chemistry and Biochemistry, University of California, Los Angeles, California 90095, United States

Tucker Allen – Department of Chemistry and Biochemistry, University of California, Los Angeles, California 90095, United States; orcid.org/0000-0003-4764-5802

Barry Y. Li – Department of Chemistry and Biochemistry, University of California, Los Angeles, California 90095, United States; orcid.org/0000-0001-8469-6890

Wenfei Li – Department of Chemistry and Biochemistry, University of California, Los Angeles, California 90095, United States

Mauro Del Ben – Applied Mathematics and Computational Research Division, Lawrence Berkeley National Laboratory, Berkeley, California 94720, United States

Jack Deslippe – National Energy Research Scientific Computing Center (NERSC), Lawrence Berkeley National Laboratory, Berkeley, California 94720, United States

Complete contact information is available at:

<https://pubs.acs.org/10.1021/acs.jctc.6c00116>

Author Contributions

[†]P.T. is the first author.

Notes

The authors declare no competing financial interest.

ACKNOWLEDGMENTS

This work was supported through the Center for Computational Study of Excited State Phenomena in Energy Materials (C2SEPEM), which is funded by the U.S. Department of Energy, Office of Science, Basic Energy Sciences, Materials Sciences and Engineering Division, via contract no. DE-AC02-05CH11231, as part of the Computational Materials Sciences Program. Work in Daniel Neuhauser's group was supported by the National Science Foundation Grant No. CHE2245253. Computational resources were provided by the National Energy Research Scientific Computing Center (NERSC), a U.S. Department of Energy Office of Science User Facility operating under contract no. DE-AC02-05CH11231 using NERSC awards BES-ERCAP0032245 (project m2651) and DDR-NESAP-ERCAP0033923 (project m4941).

REFERENCES

- (1) Werner, H.-J.; Knowles, P. J. An efficient internally contracted multiconfiguration reference configuration interaction method. *J. Chem. Phys.* **1988**, *89*, 5803–5814.
- (2) Buenker, R. J.; Peyerimhoff, S. D.; Butscher, W. Applicability of the multi-reference double-excitation CI (MRD-CI) method to the calculation of electronic wavefunctions and comparison with related techniques. *Mol. Phys.* **1978**, *35*, 771–791.
- (3) Krylov, A. I. Equation-of-Motion Coupled-Cluster Methods for Open-Shell and Electronically Excited Species: The Hitchhikers Guide to Fock Space. *Annu. Rev. Phys. Chem.* **2008**, *59*, 433–462.
- (4) Stanton, J. F.; Bartlett, R. J. The equation of motion coupled-cluster method. A systematic biorthogonal approach to molecular excitation energies, transition probabilities, and excited state properties. *J. Chem. Phys.* **1993**, *98*, 7029–7039.
- (5) Hohenberg, P.; Kohn, W. Inhomogeneous Electron Gas. *Phys. Rev.* **1964**, *136*, No. B864.
- (6) Das, S.; Kanungo, B.; Subramanian, V.; Panigrahi, G.; Motamarri, P.; Rogers, D.; Zimmerman, P.; Gavini, V. et al. *Large-Scale Materials Modeling at Quantum Accuracy: Ab Initio Simulations of Quasicrystals and Interacting Extended Defects in Metallic Alloys*,

Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis; 2023; pp 1–12.

(7) Martin, R. M. *Electronic Structure: Basic Theory and Practical Methods*; Cambridge University Press, 2004; p 624.

(8) Dreizler, R. M.; Gross, E. K. U. *Density Functional Theory*; Springer Berlin Heidelberg, 1990.

(9) Aryasetiawan, F.; Gunnarsson, O. The GW method. *Rep. Prog. Phys.* **1998**, *61*, 237–312.

(10) Runge, E.; Gross, E. K. U. Density-Functional Theory for Time-Dependent Systems. *Phys. Rev. Lett.* **1984**, *52*, 997–1000.

(11) Hedin, L. New Method for Calculating the One-Particle Green's Function with Application to the Electron-Gas Problem. *Phys. Rev.* **1965**, *139*, A796–A823.

(12) Martin, R. M. *Interacting Electrons*; Reining, L.; Ceperley, D. M., Eds.; Cambridge University Press: Cambridge, 2016.

(13) Li, Z.; Antonius, G.; Wu, M.; da Jornada, F. H.; Louie, S. G. Electron-Phonon Coupling from Ab Initio Linear-Response Theory within the GW Method: Correlation-Enhanced Interactions and Superconductivity in $Ba_{1-x}K_xBiO_3$. *Phys. Rev. Lett.* **2019**, *122*, No. 186402.

(14) Onida, G.; Reining, L.; Rubio, A. Electronic excitations: density-functional versus many-body Green's-function approaches. *Rev. Mod. Phys.* **2002**, *74*, 601–659.

(15) Del Ben, M.; da Jornada, F. H.; Canning, A.; Wichmann, N.; Raman, K.; Sasanka, R.; Yang, C.; Louie, S. G.; Deslippe, J. Large-scale GW calculations on pre-exascale HPC systems. *Comput. Phys. Commun.* **2019**, *235*, 187–195.

(16) Del Ben, M.; Yang, C.; Li, Z.; Jornada, F. H. d.; Louie, S. G.; Deslippe, J. *Accelerating Large-Scale Excited-State GW Calculations on Leadership HPC Systems*, SC20: International Conference for High Performance Computing, Networking, Storage and Analysis; 2020; pp 1–11.

(17) Brooks, J.; Weng, G.; Taylor, S.; Vlcek, V. Stochastic many-body perturbation theory for Moire states in twisted bilayer phosphorene. *J. Phys.: Condens. Matter* **2020**, *32*, No. 234001.

(18) Yu, V. W.-z.; Govoni, M. GPU Acceleration of Large-Scale Full-Frequency GW Calculations. *J. Chem. Theory Comput.* **2022**, *18*, 4690–4707.

(19) Wu, W.; Zhou, Z.; Jiang, Q. et al. *Enabling 13K-Atom Excited-State GW Calculations via Low-Rank Approximations and HPC on the New Sunway Supercomputer*, SC24: International Conference for High Performance Computing, Networking, Storage and Analysis; 2024; pp 1–14.

(20) Zhang, B.; Weinberg, D.; Hsu, C.-E.; Altman, A. R.; Shi, Y.; White, J. B.; Vigil-Fowler, D.; Louie, S. G.; Deslippe, J. R.; da Jornada, F. H.; Li, Z.; Del Ben, M. *Advancing Quantum Many-Body GW Calculations on Exascale Supercomputing Platforms*, Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis; 2025; pp 48–59.

(21) Hybertsen, M. S.; Louie, S. G. First-Principles Theory of Quasiparticles: Calculation of Band Gaps in Semiconductors and Insulators. *Phys. Rev. Lett.* **1985**, *55*, 1418–1421.

(22) Hybertsen, M. S.; Louie, S. G. Electron correlation in semiconductors and insulators: Band gaps and quasiparticle energies. *Phys. Rev. B* **1986**, *34*, 5390–5413.

(23) Kresse, G.; Hafner, J. Ab initio molecular dynamics for liquid metals. *Phys. Rev. B* **1993**, *47*, 558–561.

(24) Kresse, G.; Hafner, J. Ab initio molecular-dynamics simulation of the liquid-metal-amorphous-semiconductor transition in germanium. *Phys. Rev. B* **1994**, *49*, 14251–14269.

(25) Giannozzi, P.; Baroni, S.; Bonini, N.; et al. QUANTUM ESPRESSO: a modular and open-source software project for quantum simulations of materials. *J. Phys.: Condens. Matter* **2009**, *21*, No. 395502.

(26) Gonze, X.; Amadon, B.; Anglade, P. M.; et al. ABINIT: First-principles approach to material and nanosystem properties. *Comput. Phys. Commun.* **2009**, *180*, 2582–2615.

- (27) Marini, A.; Hogan, C.; Gruning, M.; Varsano, D. Yambo: An ab initio tool for excited state calculations. *Comput. Phys. Commun.* **2009**, *180*, 1392–1403.
- (28) Blum, V.; Gehrke, R.; Hanke, F.; Havu, P.; Havu, V.; Ren, X.; Reuter, K.; Scheffler, M. Ab initio molecular simulations with numeric atom-centered orbitals. *Comput. Phys. Commun.* **2009**, *180*, 2175–2196.
- (29) Deslippe, J.; Samsonidze, G.; Strubbe, D. A.; Jain, M.; Cohen, M. L.; Louie, S. G. BerkeleyGW: A massively parallel computer package for the calculation of the quasiparticle and optical properties of materials and nanostructures. *Comput. Phys. Commun.* **2012**, *183*, 1269–1289.
- (30) The Elk code. <http://elk.sourceforge.net/>.
- (31) Gulans, A.; Kontur, S.; Meisenbichler, C.; Nabok, D.; Pavone, P.; Rigamonti, S.; Sagmeister, S.; Werner, U.; Draxl, C. exciting: a full-potential all-electron package implementing density-functional theory and many-body perturbation theory. *J. Phys.:Condens. Matter* **2014**, *26*, No. 363202.
- (32) Govoni, M.; Galli, G. Large Scale GW Calculations. *J. Chem. Theory Comput.* **2015**, *11*, 2680–2696.
- (33) Jacquemin, D.; Duchemin, I.; Blase, X. Benchmarking the Bethe-Salpeter Formalism on a Standard Organic Molecular Set. *J. Chem. Theory Comput.* **2015**, *11*, 3290–3304.
- (34) Bruneval, F.; Rangel, T.; Hamed, S. M.; Shao, M.; Yang, C.; Neaton, J. B. MOLGW 1: Many-body perturbation theory software for atoms, molecules, and clusters. *Comput. Phys. Commun.* **2016**, *208*, 149–161.
- (35) Kühne, T. D.; Iannuzzi, M.; Del Ben, M.; et al. CP2K: An electronic structure and molecular dynamics software package - Quickstep: Efficient and accurate electronic structure calculations. *J. Chem. Phys.* **2020**, *152*, No. 194103.
- (36) Schlipf, M.; Lambert, H.; Zibouche, N.; Giustino, F. SternheimerGW: A program for calculating GW quasiparticle band structures and spectral functions without unoccupied states. *Comput. Phys. Commun.* **2020**, *247*, No. 106856.
- (37) Bruneval, F.; Vast, N.; Reining, L. Effect of self-consistency on quasiparticles in solids. *Phys. Rev. B* **2006**, *74*, No. 045102.
- (38) van Schilfgaarde, M.; Kotani, T.; Faleev, S. Quasiparticle Self-Consistent GW Theory. *Phys. Rev. Lett.* **2006**, *96*, No. 226402.
- (39) Vlček, V.; Baer, R.; Rabani, E.; Neuhauser, D. Simple eigenvalue-self-consistent Δ GW. *J. Chem. Phys.* **2018**, *149*, No. 174107.
- (40) Allen, T.; Nguyen, M.; Neuhauser, D. GW with hybrid functionals for large molecular systems. *J. Chem. Phys.* **2024**, *161*, No. 164116.
- (41) Rieger, M. M.; Steinbeck, L.; White, I.; Rojas, H.; Godby, R. The GW space-time method for the self-energy of large systems. *Comput. Phys. Commun.* **1999**, *117*, 211–228.
- (42) Liu, P.; Kaltak, M.; Klimes, J.; Kresse, G. Cubic-scaling GW: Towards fast quasiparticle calculations. *Phys. Rev. B* **2016**, *94*, No. 165109.
- (43) Wilhelm, J.; Golze, D.; Talirz, L.; Hutter, J.; Pignedoli, C. A. Toward GW Calculations on Thousands of Atoms. *J. Phys. Chem. Lett.* **2018**, *9*, 306–312.
- (44) Kim, M.; Martyna, G. J.; Ismail-Beigi, S. Complex-time shredded propagator method for large-scale GW calculations. *Phys. Rev. B* **2020**, *101*, No. 035139.
- (45) Yeh, C.-N.; Morales, M. A. Low-Scaling Algorithms for GW and Constrained Random Phase Approximation Using Symmetry-Adapted Interpolative Separable Density Fitting. *J. Chem. Theory Comput.* **2024**, *20*, 3184–3198.
- (46) Altman, A. R.; Kundu, S.; da Jornada, F. H. Mixed Stochastic-Deterministic Approach for Many-Body Perturbation Theory Calculations. *Phys. Rev. Lett.* **2024**, *132*, No. 086401.
- (47) Vetsch, N.; Maeder, A.; Maillou, V.; Winka, A.; Cao, J.; Kwasniewski, G.; Deuschle, L.; Hoefler, T.; Ziogas, A. N.; Luisier, M. Ab-initio Quantum Transport with the GW Approximation, 42,240 Atoms, and Sustained Exascale Performance; 2025; pp 1–13.
- (48) Neuhauser, D.; Gao, Y.; Arntsen, C.; Karshenas, C.; Rabani, E.; Baer, R. Breaking the Theoretical Scaling Limit for Predicting Quasiparticle Energies: The Stochastic GW Approach. *Phys. Rev. Lett.* **2014**, *113*, No. 076402.
- (49) Vlček, V.; Rabani, E.; Neuhauser, D.; Baer, R. Stochastic GW Calculations for Molecules. *J. Chem. Theory Comput.* **2017**, *13*, 4997–5003.
- (50) Vlček, V.; Li, W.; Baer, R.; Rabani, E.; Neuhauser, D. Swift GW beyond 10,000 electrons using sparse stochastic compression. *Phys. Rev. B* **2018**, *98*, No. 075107.
- (51) Bradbury, N. C.; Nguyen, M.; Caram, J. R.; Neuhauser, D. Bethe–Salpeter equation spectra for very large systems. *J. Chem. Phys.* **2022**, *157*, No. 031104.
- (52) Bradbury, N. C.; Allen, T.; Nguyen, M.; Ibrahim, K. Z.; Neuhauser, D. Optimized attenuated interaction: Enabling stochastic Bethe-Salpeter spectra for large systems. *J. Chem. Phys.* **2023**, *158*, No. 154104.
- (53) Nguyen, M.; Neuhauser, D. Gapped-filtering for efficient Chebyshev expansion of the density projection operator. *Chem. Phys. Lett.* **2022**, *806*, No. 140036.
- (54) Hutchinson, M. A stochastic estimator of the trace of the influence matrix for laplacian smoothing splines. *Commun. Stat. - Simul. Comput.* **1990**, *19*, 433–450.
- (55) Bradbury, N. C.; Allen, T.; Nguyen, M.; Neuhauser, D. Deterministic/Fragmented-Stochastic Exchange for Large-Scale Hybrid DFT Calculations. *J. Chem. Theory Comput.* **2023**, *19*, 9239–9247.
- (56) Sereda, M.; Allen, T.; Bradbury, N. C.; Ibrahim, K. Z.; Neuhauser, D. Sparse-Stochastic Fragmented Exchange for Large-Scale Hybrid Time-Dependent Density Functional Theory Calculations. *J. Chem. Theory Comput.* **2024**, *20*, 4196–4204.
- (57) Bradbury, N. C.; Li, B. Y.; Allen, T.; Caram, J. R.; Neuhauser, D. No more gap-shifting: Stochastic many-body-theory based TDHF for accurate theory of polymethine cyanine dyes. *J. Chem. Phys.* **2024**, *161*, No. 141101.
- (58) Li, B. Y.; Duong, T.; Allen, T.; Bradbury, N. C.; Caram, J. R.; Neuhauser, D. Parameterized attenuated exchange for generalized TDHF@ v_W applications. *J. Chem. Phys.* **2025**, *163*, No. 034102.
- (59) Chen, K.; Li, B. Y.; Allen, T.; Neuhauser, D. Mixed Plane-wave and Localized Orbital Basis for Sparse-Stochastic Hybrid Time-Dependent Density Functional Theory. *J. Chem. Theory Comput.* **2025**, *21*, 8509–8517.
- (60) Allen, T.; Li, B. Y.; Duong, T.; Williams, K.; Neuhauser, D. Efficient plane-wave approach to generalized Kohn-Sham density functional theory of solids with mixed deterministic and stochastic exchange. *Phys. Rev. B* **2025**, *112*, No. 155104.
- (61) Baer, R.; Neuhauser, D. Real-time linear response for time-dependent density-functional theory. *J. Chem. Phys.* **2004**, *121*, 9803–9807.
- (62) Baer, R.; Head-Gordon, M. Chebyshev expansion methods for electronic structure calculations on large molecular systems. *J. Chem. Phys.* **1997**, *107*, 10003–10013.
- (63) Fetter, A. L.; Walecka, J. D. *Quantum Theory of Many Particle Systems*; McGraw-Hill: New York, 1971; p 299.
- (64) The StochasticGW code. <https://github.com/stochasticGW/stochasticGW>.
- (65) Frigo, M.; Johnson, S. The Design and Implementation of FFTW3. *Proc. IEEE* **2005**, *93*, 216–231.
- (66) Giannozzi, P.; Barone, O.; Bonfanti, P.; Brunato, D.; Car, R.; Carnimeo, I.; Cavazzoni, C.; de Gironcoli, S.; Delugas, P.; Ferrari Ruffino, F.; Ferretti, A.; Marzari, N.; Timrov, I.; Urru, A.; Baroni, S. Quantum ESPRESSO toward the exascale. *J. Chem. Phys.* **2020**, *152*, No. 154105.
- (67) Briggs, E. L.; Sullivan, D. J.; Bernholc, J. Real-space multigrid-based approach to large-scale electronic structure calculations. *Phys. Rev. B* **1996**, *54*, 14362–14375.
- (68) Hodak, M.; Wang, S.; Lu, W.; Bernholc, J. Implementation of ultrasoft pseudopotentials in large-scale grid-based electronic structure calculations. *Phys. Rev. B* **2007**, *76*, No. 085108.

(69) Hunter, J. D. Matplotlib: A 2D Graphics Environment. *Comput. Sci. Eng.* **2007**, *9*, 90–95.

(70) Hom, T.; Kiszzenik, W.; Post, B. Accurate lattice constants from multiple reflection measurements. II. Lattice constants of germanium silicon, and diamond. *J. Appl. Crystallogr.* **1975**, *8*, 457–458.

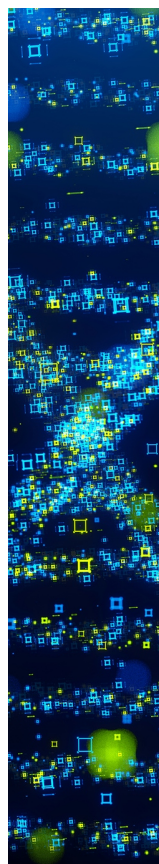
(71) Halgren, T. A. Merck molecular force field. I. Basis, form, scope, parameterization, and performance of MMFF94. *J. Comput. Chem.* **1996**, *17*, 490–519.

(72) O'Boyle, N. M.; Banck, M.; James, C. A.; Morley, C.; Vandermeersch, T.; Hutchison, G. R. Open Babel: An open chemical toolbox. *J. Cheminf.* **2011**, *3*, No. 33, DOI: [10.1186/1758-2946-3-33](https://doi.org/10.1186/1758-2946-3-33).

(73) Troullier, N.; Martins, J. L. Efficient pseudopotentials for plane-wave calculations. *Phys. Rev. B* **1991**, *43*, 1993–2006.

(74) Perdew, J. P.; Burke, K.; Ernzerhof, M. Generalized Gradient Approximation Made Simple. *Phys. Rev. Lett.* **1996**, *77*, 3865–3868.

(75) Lehtola, S.; Steigemann, C.; Oliveira, M. J.; Marques, M. A. Recent developments in libxc - A comprehensive library of functionals for density functional theory. *SoftwareX* **2018**, *7*, 1–5.



CAS BIOFINDER DISCOVERY PLATFORM™

STOP DIGGING THROUGH DATA —START MAKING DISCOVERIES

CAS BioFinder helps you find the
right biological insights in seconds

Start your search

