

**Gradient-descent optimization of fermion nodes in the diffusion Monte Carlo technique**John McFarland<sup>1,\*</sup> and Efstratios Manousakis<sup>1,2,†</sup><sup>1</sup>*Department of Physics and National High Magnetic Field Laboratory, Florida State University, Tallahassee, Florida 32306, USA*<sup>2</sup>*Department of Physics, National and Kapodistrian University of Athens, Panepistimioupolis, Zografos, 157 84 Athens, Greece*

(Received 4 November 2021; accepted 2 March 2022; published 17 March 2022)

We present a method for optimizing the location of the fermion ground-state nodes using a combination of diffusion Monte Carlo (DMC) and projected gradient descent (PGD). A PGD iteration shifts the parameters of an arbitrary node-fixing trial function in the opposite direction of the DMC energy gradient, while maintaining the cusp condition for atomic electrons. The energy gradient is calculated from DMC walker distributions by one of three methods we derive from an exact analytical expression. We combine our energy gradient calculation methods with different gradient-descent algorithms and a projection operator that maintains the cusp condition. We apply this stochastic PGD method to trial functions of Be, Li<sub>2</sub>, and Ne, all consisting of a single Slater determinant with randomized parameters, and find that the nodes dramatically improve to the same DMC energy as nodes optimized by variational Monte Carlo. Our method, therefore, departs from the standard procedure of optimizing the nodes with a non-DMC scheme such as variational Monte Carlo, density-functional theory, or configuration-interaction-based calculation, which do not directly minimize the DMC energy.

DOI: [10.1103/PhysRevA.105.032815](https://doi.org/10.1103/PhysRevA.105.032815)**I. INTRODUCTION**

Diffusion Monte Carlo (DMC) [1–3], also known as Green’s-function Monte Carlo or projector Monte Carlo, is a technique for projecting the many-body wave function to the ground state. It has been used to accurately tackle bosonic or nonfrustrated quantum-spin systems (see Ref. [4] for a list of references). In the case of fermionic systems, in general one has to restrict oneself to the fixed-node approximation, where it has been accurately applied to nuclear physics [5–7] when the required initial trial function was accurate enough. Node relaxation has led to accurate results for the electron gas in the continuum [2] and for electrons on a lattice [8–10]. There are variants of the method, such as the constraint path [11,12], which have been successfully applied to condensed-matter physics [12,13] problems.

The DMC method samples the ground-state wave function  $\Psi$  of  $N$  particles with walkers. These are points in a  $dN$ -dimensional space ( $d$  being the dimensionality of space) with a probability density equal to  $\Psi\Psi_g$ , where  $\Psi_g$  is a guiding function used for importance sampling. DMC projects  $\Psi$  to the ground state by propagating it along the imaginary-time axis, which is done by changing the position and weight of each walker in a way that samples the Green’s function (i.e., the matrix elements of the evolution operator). To prevent large variations in weights, walkers are regularly deleted, duplicated, or combined in a way that the initial and final  $\Psi\Psi_g$  are proportional.

For a fermionic  $\Psi$ , it is generally necessary to impose zero boundary conditions at an *a priori* chosen nodal surface. That

is, the nodal surface of  $\Psi$  is taken to be the nodal surface of a trial function  $\Psi_t$ , which is an approximate ground-state wave function produced by a non-DMC method. Typically  $\Psi_t$  is a product of a symmetric Jastrow function [14] factor, which describes correlations, and an antisymmetric factor, which is a single Slater determinant or a combination of Slater determinants [15] that describes the nodes. It is standard to set  $\Psi_t = \Psi_g$ , in which case nodal boundary conditions are naturally imposed by imaginary-time propagation alone. In the present paper we distinguish  $\Psi_t$  from  $\Psi_g$  because we often use a nodeless  $\Psi_g$ , in which case nodal boundary conditions are imposed by deleting those walkers that attempt to cross the nodes.

DMC error is improved with a  $\Psi_t$  that better approximates the exact ground-state wave function. The only source of error that cannot be eliminated with the DMC propagation is the fixed-node error, which is contained in just the antisymmetric part of  $\Psi_t$ . The other sources of error, namely, the finite time-step error, statistical error, and population control error, can be respectively controlled with a smaller time step, a greater sample number, and a larger walker population.

This makes the choice of the nodes of  $\Psi_t$  the fundamental approximation of DMC. The resulting fixed-node error typically ranges from about 82 to 435 meV per atom [16], and is generally controlled by increasing the complexity of  $\Psi_t$  and better optimizing its parameters. Because DMC energy is an upper bound to the true ground-state energy [17], a parameter optimization method that minimizes DMC energy is the most accurate. However, optimization methods in use will minimize some other quantity, such as the variational Monte Carlo (VMC) energy, the Kohn-Sham energy, or the local VMC energy variance [18,19].

The first attempt to optimize the parameters of  $\Psi_t$  using DMC walker distributions alone was by Reboredo *et al.* [20].

\*swqecs@gmail.com

†manousakis@gmail.com

They iteratively generated a new  $\Psi_t$  by projecting the coefficients of its determinants (or pfaffians) from the walker distribution of the previous  $\Psi_t$ . Mindful that it becomes expensive to evaluate a  $\Psi_t$  with an increasing number of determinants, they also proposed a method to reduce the complexity of  $\Psi_t$  by minimizing a cost function between the initial and reduced  $\Psi_t$  that heavily penalized changes of the nodes.

We propose to directly use the DMC energy as a ‘‘cost function’’ when optimizing the parameters of  $\Psi_t$ . This differs from using the VMC energy as a cost function [21]. Our method iteratively performs projected gradient descent (PGD) on  $\Psi_t$  by shifting its parameters in a direction roughly opposite to the energy gradient, while maintaining any constraints (the electron-nuclear cusp condition in our examples) by projecting the parameters back to the surface satisfying the constraint with each iteration. Gradient descent has an advantage for large parameter number  $N_p$  in that it does not require an estimate of  $N_p \times N_p$  matrices, which is required by many others, e.g., the linear method.

We calculate the gradient of each parameter iteration from walker samples using one of three methods that we label A, B, and C. The walker distributions are produced by DMC imaginary-time propagation with nodes fixed by the parameters. We experiment with several gradient-descent algorithms commonly used for machine learning. In order to keep our method self-reliant, during PGD we do not rely on a preoptimized Jastrow function, which does not affect the fixed-node error, although it does affect the other errors.

Our paper is organized as follows: In Sec. II we present the three methods for evaluating the derivative of the energy from the DMC walker distribution. In Sec. III we describe our implementation of gradient descent, including the DMC scheme used, practical issues, and the form of  $\Psi_t$  and  $\Psi_g$ . In Sec. IV we present tests on the accuracy and speed of the three methods used and results of PGD on trial functions of Be, Li<sub>2</sub>, Ne, and F<sub>2</sub>. In Sec. V, we discuss the advantages of our method and ways to improve it. We discuss parameter fluctuations in Appendix A, and we describe the gradient-descent algorithms used in Appendix B and we compare them in Appendix C.

## II. METHOD

A PGD iteration shifts the nodes in a direction expected to lower the DMC energy  $E$ . The nodes are determined by the parameters  $\theta_i$  of the antisymmetric part of  $\Psi_t$ , and the shift of the nodes is determined by the gradient  $\frac{\partial E}{\partial \theta_i}$ . The gradient of a nodal surface is calculated with walker samples produced by DMC imaginary-time propagation of  $\Psi\Psi_g$ , with  $\Psi$  the fixed-node ground state and  $\Psi_g$  a guiding function. Thus PGD iterations change the nodal surface, while DMC iterations generate data for the gradient of a given nodal surface.

Provided the  $\theta_i$  are not already at a local minimum,  $E$  will be lowered by the following change of parameters:

$$\theta_i \rightarrow \theta_i - a_i \frac{\partial E}{\partial \theta_i}, \quad (1)$$

provided that  $a_i$  is positive and sufficiently small. However, the presence of stochastic error in our gradient means that only the expectation of  $E$  will be lowered with sufficiently small  $a_i$ . If the  $\theta_i$  are constrained (e.g., from the cusp condition or

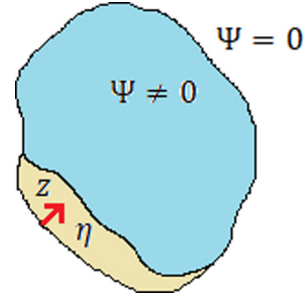


FIG. 1. Nodal pocket.

symmetry) we modify Eq. (1) to

$$\theta_i \rightarrow P(\theta_i - a_i \frac{\partial E}{\partial \theta_i}), \quad (2)$$

where  $P$  is the projection operator that moves the parameters to the nearest point on the manifold that satisfies the constraint. Some gradient-descent algorithms, e.g., Adam, simulate momentum with friction by replacing  $\frac{\partial E}{\partial \theta_i}$  with an exponential trailing average of current and past iterations.

### A. Derivative of the energy

Central to our method is calculating the DMC energy gradient in parameter space. Working in atomic units, we derive an expression for  $\frac{\partial E}{\partial \theta_i}$ , also derived by Berman [22]. To simplify the derivation, let us again set the ground state  $\Psi$  to zero except for one nodal pocket, as justified in Sec. II A. Now let us examine how the energy, given as

$$E = \frac{\langle \Psi | H | \Psi \rangle}{\langle \Psi | \Psi \rangle}, \quad (3)$$

changes due to an infinitesimal change  $\delta\theta_i$  in one of the parameters, which shifts the node and, thus, changes  $\Psi$  to  $\Psi + \delta\Psi$ . By doing that the change in energy  $\delta E$  is given by

$$\delta E = \frac{\langle \delta\Psi | (H - E) | \Psi \rangle}{\langle \Psi | \Psi \rangle} + \frac{\langle \Psi | (H - E) | \delta\Psi \rangle}{\langle \Psi | \Psi \rangle} + \frac{\langle \delta\Psi | (H - E) | \delta\Psi \rangle}{\langle \Psi | \Psi \rangle}. \quad (4)$$

We convert the above expression to an integral over the 3N dimensional position vector  $\mathbf{R}$  and integrate by parts to obtain

$$\delta E = \frac{\int d\mathbf{R} \ 2\delta\Psi(H - E)\Psi + \delta\Psi(H - E)\delta\Psi}{\langle \Psi | \Psi \rangle}. \quad (5)$$

The  $\delta\Psi(H - E)\Psi$  term of Eq. (5) is zero except at the original nodes, where the action of the Laplacian on a discontinuity of  $\nabla\Psi$  produces a delta function. The term  $\delta\Psi(H - E)\delta\Psi$  of the above equation is second order in  $\delta\theta_i$  except at the original and shifted nodes, where the action of the Laplacian on the discontinuities of  $\nabla\delta\Psi$  produces zeroth-order delta functions. Thus, we only need to consider the operation of the Laplacian on these discontinuities at the original and shifted nodes to evaluate Eq. (5) up to first order in  $\delta\theta_i$ .

Let us now choose a coordinate system with the shape of the nodes. Let  $z$  be the coordinate perpendicular to the node, zero at the initial node, and pointing towards the direction where  $\Psi \neq 0$  (see Fig. 1 for an illustration). Let  $\mathbf{A}$  be the

remaining  $3N - 1$  dimensional coordinates that parametrize the nodal surface at  $z = 0$ . Then  $\Psi$  to first order is

$$\Psi(\mathbf{A}, z) = g(\mathbf{A})\mathbf{R}(z), \quad (6)$$

where  $g(\mathbf{A}) \equiv |\nabla\Psi(\mathbf{A}, z = 0^+)|$ , and  $\mathbf{R}$  is the ramp function:

$$\mathbf{R}(z) = \begin{cases} z & \text{if } 0 \leq z \\ 0 & \text{if } z < 0 \end{cases}.$$

To express  $\delta\Psi$  near the node, let us define  $\eta(\mathbf{A})$  as the displacement of the node in the  $z$  direction that results from  $\delta\theta_i$ . Then, to first order in  $\Psi$ ,  $\delta\Psi$  is the difference between the displaced and original wave function, that is,

$$\delta\Psi(\mathbf{A}, z) = g(\mathbf{A})\mathbf{R}[z - \eta(\mathbf{A})] - g(\mathbf{A})\mathbf{R}(z). \quad (7)$$

The delta functions resulting from the Laplacian are then given by

$$\nabla^2\Psi(\mathbf{A}, z) = g(\mathbf{A})\delta(z), \quad (8)$$

$$\nabla^2\delta\Psi(\mathbf{A}, z) = g(\mathbf{A})\delta(z - \eta) - g(\mathbf{A})\delta(z). \quad (9)$$

Since only the delta functions contribute to first order, we leave only their contribution to Eq. (5), yielding

$$\delta E = -\frac{\int d\mathbf{R} \delta\Psi(\mathbf{A}, z)g(\mathbf{A})[\delta(z) + \delta(z - \eta)]}{2\langle\Psi|\Psi\rangle}. \quad (10)$$

We reduce this to an integral over the nodal surface by integrating over  $z$ . Absorbing the Jacobian determinant into  $d\mathbf{A}$  we obtain

$$\delta E = -\frac{\int_{\text{node}} d\mathbf{A} g(\mathbf{A})[\delta\Psi(\mathbf{A}, 0) + \delta\Psi(\mathbf{A}, \eta)]}{2\langle\Psi|\Psi\rangle}. \quad (11)$$

From Eq. (7) we find that  $\delta\Psi(\mathbf{A}, 0) + \delta\Psi(\mathbf{A}, \eta) = g(\mathbf{A})\eta(\mathbf{A})$  both when  $\eta(\mathbf{A}) > 0$  and when  $\eta(\mathbf{A}) < 0$ . This yields

$$\delta E = \frac{\int_{\text{node}} d\mathbf{A} \eta(\mathbf{A})g^2(\mathbf{A})}{2\langle\Psi|\Psi\rangle}. \quad (12)$$

By taking a derivative with respect to  $\theta_i$  we derive our analytical expression:

$$\frac{\partial E}{\partial\theta_i} = \frac{\int_{\text{node}} d\mathbf{A} |\nabla\Psi(\mathbf{A})|^2\partial_{\theta_i}\eta(\mathbf{A})}{2\langle\Psi|\Psi\rangle}. \quad (13)$$

### B. Practical methods to estimate the energy gradient

Starting from Eq. (13) we have derived three different equations for calculating the energy gradient from walker distributions, which we refer to as method A, B, and C. Since DMC usually samples terms linear in  $\Psi$ , but Eq. (13) contains two terms bilinear in  $\Psi$ , an approximation used by methods A and B replaces the bilinear terms with a mixed estimate of the true ground-state wave function and the trial function, i.e.,

$$\langle\Psi|\Psi\rangle \rightarrow \langle\Psi_t|\Psi_t\rangle, \quad |\nabla\Psi|^2 \rightarrow \nabla\Psi_t \cdot \nabla\Psi. \quad (14)$$

Then, using

$$\frac{\partial\eta}{\partial\theta_i} = -\frac{1}{|\nabla\Psi_t|} \frac{\partial\Psi_t}{\partial\theta_i}, \quad (15)$$

we approximate Eq. (13) with

$$\frac{\partial E}{\partial\theta_i} \approx -\frac{\int_{\text{node}} d\mathbf{A} |\nabla\Psi(\mathbf{A})|\partial_{\theta_i}\Psi_t(\mathbf{A})}{2\langle\Psi_t|\Psi_t\rangle}. \quad (16)$$

Unfortunately, this relies on the quality of  $\Psi_t$ , although this quality is improved with PGD. Method C does not make this approximation, and thus in principle requires no Jastrow optimization.

### C. Method A: The energy gradient from a standard walker distribution

Method A calculates Eq. (16) with a standard walker distribution equal to  $\Psi_t\Psi$ , which we think will make it easiest to implement into existing DMC code. It requires the evaluation of  $\Psi_t$ , its parameter derivative, and the parameter derivative of its  $3N$  dimensional Laplacian. To get Eq. (16) in a form where this is possible, we change the nodal integral to a volume integral. Using the  $z$  coordinate as defined in Sec. II and Gauss's theorem, we transform Eq. (16) to

$$\frac{\partial E}{\partial\theta_i} \approx -\frac{1}{2\langle\Psi_t|\Psi_t\rangle} \int_{\text{node}} d\mathbf{A} \lim_{\epsilon \rightarrow 0} \int_{-\epsilon}^{\epsilon} dz \frac{\partial\Psi_t}{\partial\theta_i} \nabla^2\Psi, \quad (17)$$

where we used the fact that  $|\nabla\Psi| = \frac{\partial\Psi}{\partial z}$  for  $z > 0$  and  $\Psi = 0$  for  $z < 0$ .

We replace  $\nabla^2$  with  $2(E - H)$  in Eq. (17) in order to increase the bounds of the  $z$  integral to all space, which is possible since  $(E - H)\Psi = 0$  everywhere except the nodes. We then integrate by parts to obtain

$$\frac{\partial E}{\partial\theta_i} \approx \frac{\int_{\text{vol}} d\mathbf{R} \Psi(H - E)\partial_{\theta_i}\Psi_t}{\int_{\text{vol}} d\mathbf{R} \Psi_t\Psi}. \quad (18)$$

Interestingly, one can also arrive at Eq. (18) by taking the derivative of the mixed estimator  $\frac{\langle\Psi|H|\Psi_t\rangle}{\langle\Psi|\Psi_t\rangle}$  with  $\Psi$  fixed.

Method A calculates the integrals of Eq. (18) from positions  $R_k$  and weights  $w_k$  of a walker distribution equal to  $\Psi_t\Psi$  with the equation

$$\frac{\partial E}{\partial\theta_i} \approx \frac{\sum_k w_k \Psi_t^{-1}(R_k)(H - E)\partial_{\theta_i}\Psi_t(R_k)}{\sum_k w_k}. \quad (19)$$

### D. Method B: The energy gradient from a nodal walker distribution

Method B requires the evaluation of  $\Psi_t$  and its parameter derivative. It samples  $|\nabla\Psi|$  of Eq. (16) directly from a walker distribution on the nodal surface, which we generate by recording walkers that cross the nodal surface and are thus deleted. We note that shifting crossed walker positions closer to the nodes using Newton's method noticeably improves performance. Unlike those of method A, the precrossed walkers cannot be guided by  $\Psi_t$ . This is partly because no walkers will cross the nodes of  $\Psi_t$  in the zero time-step limit. Instead we guide them with a nodeless function  $\Psi_g$ .

We wish to relate  $|\nabla\Psi|$  at the nodal surface to the number of walkers that cross the nodal surface per area  $dF/dA$  during time  $\Delta\tau$ . Let us assume a constant fixed-node walker distribution equal to  $\Psi_g\Psi$ . If we suddenly remove the nodal surface, there will be an initial increase of walker population resulting from walkers that would have crossed the nodal surface where it is present. Therefore the resulting population increase per nodal area to first order in time is  $dF/dA$ . Let us make  $\Delta\tau$  small enough that  $\Psi$  does not change at a distance  $\epsilon$  from the node; then, the population increase per nodal area from the

released nodes is

$$\frac{dF(A)}{dA} = \Psi_g(\mathbf{R}_A) \int_{-\epsilon}^{\epsilon} \int_0^{\Delta\tau} dz d\tau \frac{d\Psi(\mathbf{A}, z)}{d\tau}, \quad (20)$$

where  $\mathbf{R}_A = (\mathbf{A}, 0)$ . We can rewrite the above as

$$\frac{dF(A)}{dA} = \Psi_g(\mathbf{R}_A) \int_{-\epsilon}^{\epsilon} \int_0^{\Delta\tau} dz d\tau (-H)\Psi(\mathbf{A}, z). \quad (21)$$

We take the  $\epsilon \rightarrow 0$  limit of Eq. (21), which leaves only a kinetic contribution from the discontinuity of  $\nabla\Psi$ :

$$\frac{dF(A)}{dA} = \Psi_g(\mathbf{R}_A) \int_{-\epsilon}^{\epsilon} \int_0^{\Delta\tau} dz d\tau \frac{1}{2} \nabla^2 \Psi(\mathbf{A}, z). \quad (22)$$

Then, using

$$\begin{aligned} \nabla^2 \Psi(\mathbf{A}, \epsilon) &= \partial_z^2 \Psi(\mathbf{A}, \epsilon), \\ |\nabla \Psi(\mathbf{A}, \epsilon)| &= \partial_z \Psi(\mathbf{A}, \epsilon), \\ |\nabla \Psi(\mathbf{A}, -\epsilon)| &= 0, \end{aligned} \quad (23)$$

we write

$$\frac{dF(A)}{dA} = \frac{1}{2} \Delta\tau \Psi_g(\mathbf{R}_A) |\nabla \Psi(\mathbf{R}_A)|. \quad (24)$$

We can now use Eq. (24) to sample the  $|\nabla \Psi(\mathbf{R}_A)|$  term of Eq. (16), which yields the equation

$$\frac{\partial E}{\partial \theta_i} \approx - \frac{\sum_j w_j \Psi_g^{-1}(R_j) \partial_{\theta_i} \Psi_t(R_j)}{\Delta\tau \langle \Psi_t | \Psi \rangle}, \quad (25)$$

where  $R_j$  and  $w_j$  are, respectively, the positions and weights of walkers that crossed the nodes during the period  $\Delta\tau$ . To sample the denominator, we simply use

$$\Delta\tau \langle \Psi_t | \Psi \rangle = \int_0^{\Delta\tau} d\tau \sum_k \frac{\Psi_t(R_k, \tau)}{\Psi_g(R_k, \tau)} w_k(\tau), \quad (26)$$

with  $R_k$  and  $w_k$  the positions and weights of all walkers as a function of time  $\tau$ .

Method B therefore calculates the energy gradient with

$$\frac{\partial E}{\partial \theta_i} \approx - \frac{\sum_j w_j \Psi_g^{-1}(R_j) \partial_{\theta_i} \Psi_t(R_j)}{\int_0^{\Delta\tau} d\tau \sum_k w_k \Psi_g^{-1}(R_k) \Psi_t(R_k)}. \quad (27)$$

### E. Method C: The exact energy gradient

Method C does not use the approximation of Eq. (14). It samples one factor of  $|\nabla \Psi|$  in Eq. (13) with the same nodal walker distribution used by method B, and samples the other factor of  $|\nabla \Psi|$  using forward walking. Just as with method B, method C requires the evaluation of  $\Psi_t$  and its parameter derivative.

We start by writing the  $\Psi$  projector as a time evolution operator

$$\frac{|\Psi\rangle \langle \Psi|}{\langle \Psi | \Psi \rangle} = \lim_{\tau \rightarrow \infty} \exp[\tau(E-H)] \quad (28)$$

and contract the left side of Eq. (28) with  $\langle \Psi_t |$  and the right with  $|J\rangle_{\mathbf{R}}$  where

$$|J\rangle_{\mathbf{R}} \equiv \frac{1}{\Psi_t(\mathbf{R})} |\mathbf{R}\rangle, \quad (29)$$

resulting in

$$\frac{\Psi(\mathbf{R}) \langle \Psi_t | \Psi \rangle}{\Psi_t(\mathbf{R}) \langle \Psi | \Psi \rangle} = \lim_{\tau \rightarrow \infty} \langle \Psi_t | \exp[\tau(E-H)] | J \rangle_{\mathbf{R}}. \quad (30)$$

We recognize that

$$\langle \Psi_t | \exp[\tau(E-H)] | J \rangle_{\mathbf{R}} = \int d\mathbf{R}' \tilde{G}(\mathbf{R}', \mathbf{R}, \tau), \quad (31)$$

where  $\tilde{G}$ , a Green's function, is the evolution operator of  $\Psi_t \Psi$ . Because  $\Psi(\mathbf{R})/\Psi_t(\mathbf{R})$  approaches  $|\nabla \Psi(\mathbf{R}_A)|/|\nabla \Psi_t(\mathbf{R}_A)|$  as  $\mathbf{R}$  approaches the nodal surface (since  $\nabla = \partial_z$  at  $z \rightarrow 0^+$ ), at the nodal surface we have

$$|\nabla \Psi(\mathbf{R}_A)| = |\nabla \Psi_t(\mathbf{R}_A)| \frac{\langle \Psi | \Psi \rangle}{\langle \Psi_t | \Psi \rangle} \lim_{\tau \rightarrow \infty} \int d\mathbf{R}' \tilde{G}(\mathbf{R}', \mathbf{R}_A, \tau). \quad (32)$$

We now substitute one  $|\nabla \Psi|$  term of Eq. (13) with Eq. (32) and again use  $\frac{\partial \eta}{\partial \theta_i} = -\frac{1}{|\nabla \Psi_t|} \frac{\partial \Psi_t}{\partial \theta_i}$  to get

$$\begin{aligned} \frac{\partial E}{\partial \theta_i} &= -\frac{1}{2 \langle \Psi_t | \Psi \rangle} \int_{\text{node}} d\mathbf{A} \frac{\partial \Psi_t}{\partial \theta_i} |\nabla \Psi(\mathbf{R}_A)| \\ &\times \lim_{\tau \rightarrow \infty} \int d\mathbf{R}' \tilde{G}(\mathbf{R}', \mathbf{R}_A, \tau). \end{aligned} \quad (33)$$

Equation (33) is identical to Eq. (16) except the added factor  $\lim_{\tau \rightarrow \infty} \int d\mathbf{R}' \tilde{G}(\mathbf{R}', \mathbf{R}_A, \tau)$ . We sample this factor with the value  $\xi$  we define as the final weight of a walker that had an initial weight of 1 and an initial position of  $\mathbf{R}_A$ , and was then propagated for time  $\tau$  with  $\Psi_t$  as the guiding function, that is,

$$\langle \xi(\mathbf{R}_A, \tau) \rangle = \int d\mathbf{R}' \tilde{G}(\mathbf{R}', \mathbf{R}_A, \tau). \quad (34)$$

Instead of taking the  $\tau \rightarrow \infty$  limit, we cut the propagation time short at  $\tau_c$ . The value we choose should be long enough for  $\langle \xi(\mathbf{R}_A, \tau) \rangle$  to become roughly constant.

Placing a walker guided by  $\Psi_t$  at the nodes can be problematic since the local energy and raw drift velocity  $\mathbf{V} = \nabla \ln \Psi_t$  are divergent at the nodes when  $\Psi_t$  is used as the guiding function. When calculating  $\xi$ , we suppress the effects of both divergences in the standard way, where we modify the drift velocity and weight increase by multiplying them with  $\frac{-1 + \sqrt{1 + 2V^2 \delta\tau}}{V^2 \delta\tau}$ , where  $\delta\tau$  is the time step. To ensure walkers move towards the positive side of the node, we also multiply the drift velocity by  $\text{sign}(\Psi_t)$ .

Method C calculates Eq. (33) with node crossing walkers following the same logic we used for method B and Eq. (16), but with the extra factor  $\xi$ . We multiply Eq. (27) with  $\xi(R_j, \tau_c)$  to obtain our equation for method C:

$$\frac{\partial E}{\partial \theta_i} \approx - \frac{\sum_j w_j \xi(R_j, \tau_c) \Psi_g^{-1}(R_j) \partial_{\theta_i} \Psi_t(R_j)}{\int_0^{\Delta\tau} d\tau \sum_k w_k \Psi_g^{-1}(R_k) \Psi_t(R_k)}. \quad (35)$$

## III. IMPLEMENTATION

To test our method, we developed a DMC code with walker propagation that for the most part follows the prescription of Umrigar *et al.* [23]. In addition, we introduced our PGD iteration method in the code. The main steps are outlined next.

We will apply our technique on three systems: atomic Be, Li<sub>2</sub>, and atomic Ne.

### A. Trial and guiding function

The  $\Psi_t$  used for all of our calculations were of the well-known Slater-Jastrow form

$$\Psi_t = J\Psi_S, \quad \Psi_S = D^\uparrow D^\downarrow, \quad (36)$$

where  $D^\uparrow$  and  $D^\downarrow$  are spin-up and spin-down Slater determinants of single-particle orbitals  $\phi_a(r)$  described in the next paragraph.  $J$  is a Jastrow function:

$$J = \prod_{i<j} f_{ij}, \quad f_{ij} = e^{-\frac{a_{ij}r_{ij}}{1+br_{ij}}} \quad (37)$$

where  $r_{ij}$  is the electron-electron distance and  $a_{ij}$  is equal to 1/4 (or 1/2) when both  $i$  and  $j$  correspond to electrons of parallel (or antiparallel) spin projections [23].

The single-particle orbitals  $\phi$  of the Slater determinants are made of a basis of Slater functions, taking the form

$$\phi_a(r) = \sum_{b=1}^{N_{\text{basis}}} C_{ab} r_b^{n_b-1} e^{-\xi_b r_b} Y_{l_b m_b}(\hat{r}_b), \quad (38)$$

where  $r_b = |\mathbf{r} - \mathbf{R}_b|$ , and  $n_b$ ,  $l_b$ , and  $m_b$  are quantum numbers characterizing the  $b$  basis state which is centered at the atomic nuclear position  $\mathbf{R}_b$ . The parameters we optimize are  $C_{ab}$  and  $\xi_b$ , with  $\xi_b$  being shared by all orbitals.

Since  $J$  is nodeless, and, thus, does not affect the energy, our method cannot improve it. We therefore do not include it during PGD since we have no *a priori* knowledge of its parameters. Although there is a reason to expect that a quality  $J$  will improve the methods A and B, since they rely on a  $\Psi \rightarrow \Psi_t$  approximation, we do not notice such an improvement. Method C on the other hand makes no such approximation, so in principle, it does not require an optimized  $J$ . We only include  $J$  when evaluating the energy after optimization.

Our walker distributions for methods B and C use an electron-nuclear Jastrow function as a guiding function  $\Psi_g$ . It has the form

$$\Psi_g = \prod_{i,k} \exp\left(-\frac{Z_k |\mathbf{R}_k - \mathbf{r}_i|}{1 + |\mathbf{R}_k - \mathbf{r}_i|}\right), \quad (39)$$

with  $\mathbf{R}_k$  the nuclear coordinate,  $Z_k$  the nuclear charge, and  $\mathbf{r}_i$  the electron coordinate. We choose this  $\Psi_g$  instead of  $\Psi_g = 1$  to reduce the statistical, time-step, and population growth errors. This  $\Psi_g$  also results in many more walkers crossing the node per time, and the electrons are far less likely to ionize. Ionization can still be a problem with poor nodes; when this is the case, we surround the system with a potential barrier.

### B. Projection on the cusp-condition satisfying parameter space

The electron-nuclear cusp condition [24,25] is given by

$$S_k \equiv -\frac{1}{2} \frac{d \ln \overline{\Psi_t^2}}{dr} \Big|_{r=0}, \quad S_k = Z_k, \quad (40)$$

where  $r_k$  is the radial coordinate from nucleus  $k$ ,  $Z_k$  is the nuclear charge, and  $\overline{\Psi_t^2}$  is defined as the angular average of  $\Psi_t^2$  about  $r_k = 0$ . It is a necessary requirement to avoid a large

divergence of local energy near the nucleus, which gives rise to a large increase of all errors besides the fixed-node error. Since the cusps  $S_k$  of our  $\Psi_t$  depend on  $\Psi_S$  and its varied parameters  $\theta_i$ , we project  $\theta_i$  back to the manifold of the cusp condition after each iteration of gradient descent.

The cusp condition of  $\Psi_S$  is satisfied when the cusps  $S_{ka}$  of all single-particle orbitals  $\phi_a$  satisfy  $S_{ka} = Z_k$ . We assume the shift in parameters is small enough (due to small  $a_i$ ) that a linear approximation of the cusp  $\tilde{S}_{ka}$  can be made at the preprojected parameters  $\theta_{i0}$ :

$$\tilde{S}_{ka}(\theta_i) \equiv S_{ka}(\theta_{i0}) + \sum_i (\theta_i - \theta_{i0}) \frac{\partial S_{ka}(\theta_{i0})}{\partial \theta_i}. \quad (41)$$

$\theta_i$  is then shifted to the closest point satisfying

$$\tilde{S}_{ka}(\theta_i) = \tilde{S}_{ka}(\theta_{i0}) + c[Z_k - \tilde{S}_{ka}(\theta_{i0})], \quad (42)$$

where  $c$  is added for stability and is between zero and one (we used 0.5). We repeat this until  $S_{ka} - Z_k$  is within a threshold.

### C. Details of the PGD iterations

After each PGD iteration, we update the energy and effective time step [17] with data of that iteration. We use the mixed estimate

$$E = \frac{\langle \Psi_t | H | \Psi \rangle}{\langle \Psi_t | \Psi \rangle},$$

for method A, and for methods B and C we use the growth estimate

$$E = \Delta\tau^{-1} \ln \frac{\Sigma(0)}{\Sigma(\Delta\tau)},$$

where  $\Sigma(\tau)$  is the sum of the weights of all the walkers at imaginary time  $\tau$ , ignoring weight normalization.

To calculate the gradient, we propagate the walker distribution  $f$  until the number of samples (walker positions and weights) reaches a threshold. Samples for method A are taken every time step, and samples for methods B and C are taken from all node-crossing walkers. After an iteration of projected gradient descent, we propagate  $f$  for an extra time (around 0.1 Ha<sup>-1</sup>) to adjust it to the new  $\Psi_t$  before taking samples for the next gradient. This is necessary to do for method A because after a shift of the node,  $f$  does not go to zero at the nodes fast enough for the expectation of the sum of Eq. (19) to be finite, due to the  $\Psi_t^{-1}$  factor. For method A, requiring  $f$  to go to zero sufficiently fast also requires us to use an accept-reject step.

The gradient has stochastic error that prevents PGD of the parameters  $\theta_i$  from fully settling to an accurate minimum, where the signal to noise ratio of the gradient diverges. Instead,  $\theta_i$  continues to fluctuate around the minimum after some PGD iterations. We argue in Appendix A that to first order in the  $a_i$  of Eq. (1), the variance of fluctuations is proportional to both the variance of the gradient error and  $a_i$ . Thus we can suppress the fluctuations through smaller  $a_i$  and/or by reducing the statistical error of the energy derivative with more samples.

An important choice to make is the value of  $a_i$ . The optimal value depends on many factors, such as the parameters  $\theta_i$ , the form of  $\Psi_t$ , the number of PGD iterations, and the

stochastic gradient error. Generally, decreasing  $a_i$  requires more PGD iterations to reach a minimum, and may exacerbate the problem of getting stuck in local minima, while increasing  $a_i$  increases fluctuations from the gradient error and may cause overshooting of  $\theta_i$  from its optimum value.

The most basic gradient-descent scheme uses a common value of  $a_i$  for all parameters and iterations. We experiment with this and with the two adaptive gradient-descent algorithms Adam [26] and RMSprop [27]. The adaptive algorithms make  $a_i$  inversely proportional to a trailing root mean square of past gradient components, and are described more in Appendix B. The  $a_i$  of all our algorithms have a common factor which we label the learning rate  $\alpha$ , which we determine empirically. We sometimes decrease this value with PGD iterations to suppress fluctuations of  $\theta_i$ . In Appendix C we examine the effective range of  $\alpha$  for different gradient-descent algorithms and systems.

#### IV. RESULTS

We tested our methods for three systems—atomic Be,  $\text{Li}_2$ , and atomic Ne—and we compared our results to those obtained by the DMC calculation of Umrigar *et al.* [23] where the fermion nodes were optimized using VMC.

##### A. Derivative of the energy calculation

We present tests of the energy gradient calculation methods A, B, and C in this subsection. They examine one parameter of a basis we added to the simple Be trial function that is fully described in Ref. [23]. The basis has quantum numbers  $n = 2$ ,  $l = 1$ , and  $m = 0$ , and the parameter is the non-normalized coefficient of the basis. We varied this parameter and obtained its partial first and second energy derivatives by means of a fixed-node DMC calculation, which is shown in the upper panel of Fig. 2. As can be noticed, the energy obtained from the fixed-node DMC is well approximated with a parabola. We wish to compare the derivative of the parabola with the energy derivatives obtained by our three methods which are shown in the center panel of Fig. 2 and are well approximated with lines intersecting the origin.

The energy derivative obtained by method C agrees with that from the parabola to within 1%, while the energy derivative obtained using methods A and B is proportional to that from the parabola but off by a factor of 2, which we attribute to the approximation of Eq. (14). We expect that when each of the components of the energy derivative has the correct sign it results in shifting the parameters towards the direction where the energy is lower. In addition, many gradient-descent algorithms divide each gradient component by a trailing root mean square of past gradient components; thus, the exact value of the slope is not needed by such algorithms. As a result methods A and B may still be useful.

The bottom panel of Fig. 2 shows how the energy derivative obtained with method C depends on the cutoff time  $\tau_c$ . Our choice of  $\tau_c$  when optimizing with method C is based on this dependence.

We compare the efficiency of the three methods in Fig. 3 by plotting the error of the energy derivative versus wall time, which is proportional to the sample size. In this figure we

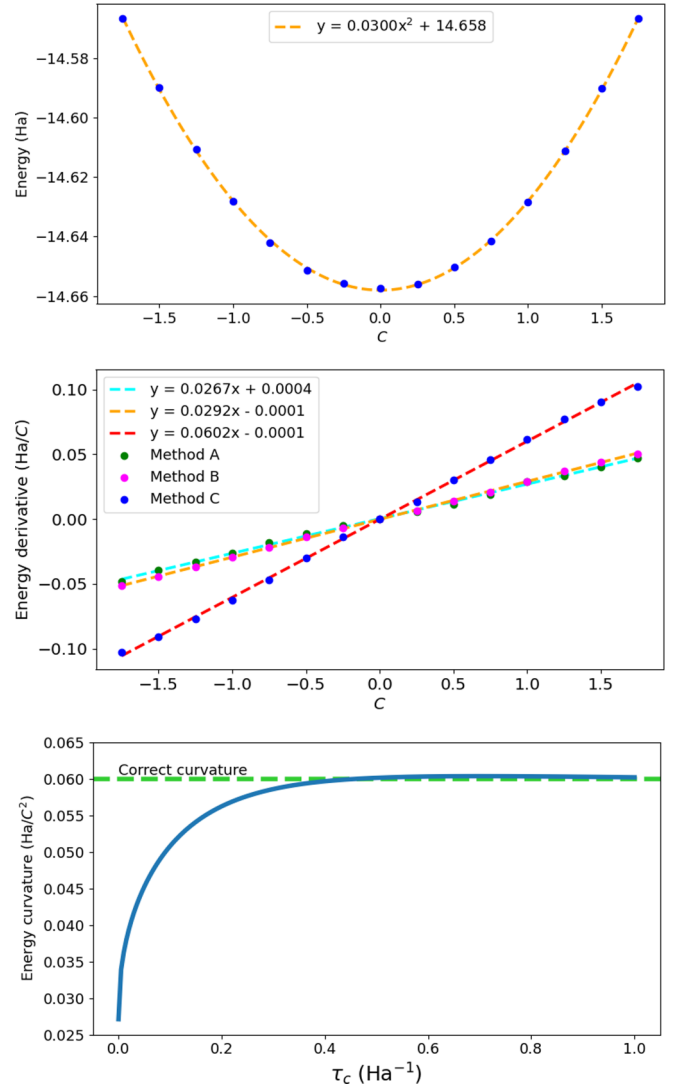


FIG. 2. We compare the DMC energy (top) of a Be trial function for a varied parameter with the energy derivative with respect to said parameter (center) that is calculated with methods A–C. Data points are shown as dots with an associated best fit line or parabola. The parabola curvature is within 1% of the slope of method C, which uses  $\tau_c = 1(\text{Ha}^{-1})$ . We examine how the slope of method C (bottom) depends on  $\tau_c$ .

show the best-fit power laws which, as expected, are approximately proportional to the inverse square root of sample size. Methods B and C seem to require similar wall time to reach a given error level when using  $\tau_c = 0.25 \text{ Ha}^{-1}$ . Method A on the other hand is considerably more expensive, though a fair comparison is difficult partly because method A may not require the same time step, and the error depends both on the associate parameter and on  $\Psi_l$ .

##### B. Projected gradient-descent evolution

In order to demonstrate the utility of the method, we start from random values for the parameters  $C$  and  $\xi$  of all-electron single-determinant wave functions of Be,  $\text{Li}_2$ , and Ne, and we perform PGD iterations with our cusp condition projection.

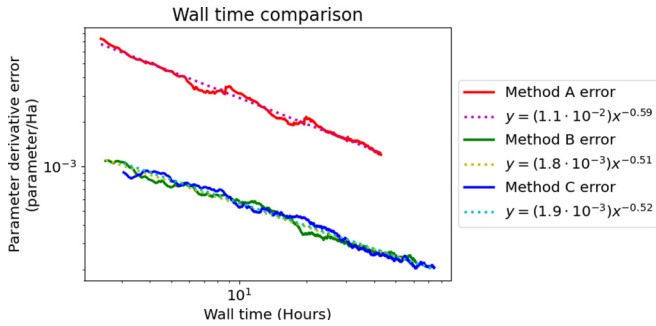


FIG. 3. We compare the efficiency of methods A, B, and C by plotting their energy derivative error for the parameter of Fig. 2 as wall time (sample size) is increased, using a common time step. Best-fit (dotted lines) power laws are roughly proportional to the inverse square root of sample size. Methods B and C require similar wall time to achieve the same error, while method A requires on the order of 50 times more.

We present a history of the energy of their nodes in Fig. 4. One can see the DMC energies rapidly decrease with iteration to the energy of the VMC optimized nodal surface obtained by Umrigar *et al.* [23] (red dotted lines). We continue to fluctuate after initial decrease, which we attribute to statistical error of the energy gradient and an excessively large decent rate [ $a_i$  term of Eq. (1)]. Adam was used for all examples, and we decreased the descent rate [ $\alpha$  parameter of Eq. (B3)] with iteration number. Method C was used for all except Be, where method A was used, demonstrating the utility of the approximation of Eq. (16).

We chose these three systems because a full description of their VMC optimized wave functions was present in the literature [23], allowing us to directly compare our PGD optimization using wave functions of the same form. There are DMC calculations for larger atoms and many atom systems; however, many of the atomic electrons (the inner) are absent in most of these calculations because they are using pseudopotentials and, therefore, do not have to implement the cusp condition. Nevertheless, we have applied PGD to a somewhat larger system of an all electron single determinant of  $F_2$ , shown in Fig. 5. Unfortunately we do not have a wave function of the same form optimized by other means for comparison. However, a DMC calculation of  $F_2$  was done by Giner *et al.* [28], who used a configuration-interaction optimized wave function consisting of  $10^5$  determinants. They achieved an energy of  $-199.2977(1)$  Ha, while our lowest PGD iteration was  $-196.7(3)$  Ha. This discrepancy can partly be attributed to the fact that our wave function had just one determinant (a current limitation of our code), and a single-particle basis of just 20 Slater-type orbitals. The decrease in error of the  $F_2$  energy with iteration number suggests that PGD also indirectly reduces local energy variance.

### C. Energy upper-bound estimation

The set  $\{E_1, E_2, \dots, E_n, \dots\}$  formed by the values of the energy  $E_i$  at the  $i$ th step of the PGD method shown in Fig. 4 is to be taken as a set of energy upper bounds. As already discussed the variations from step to step are due to the fact that the sign of the energy derivative with respect to the

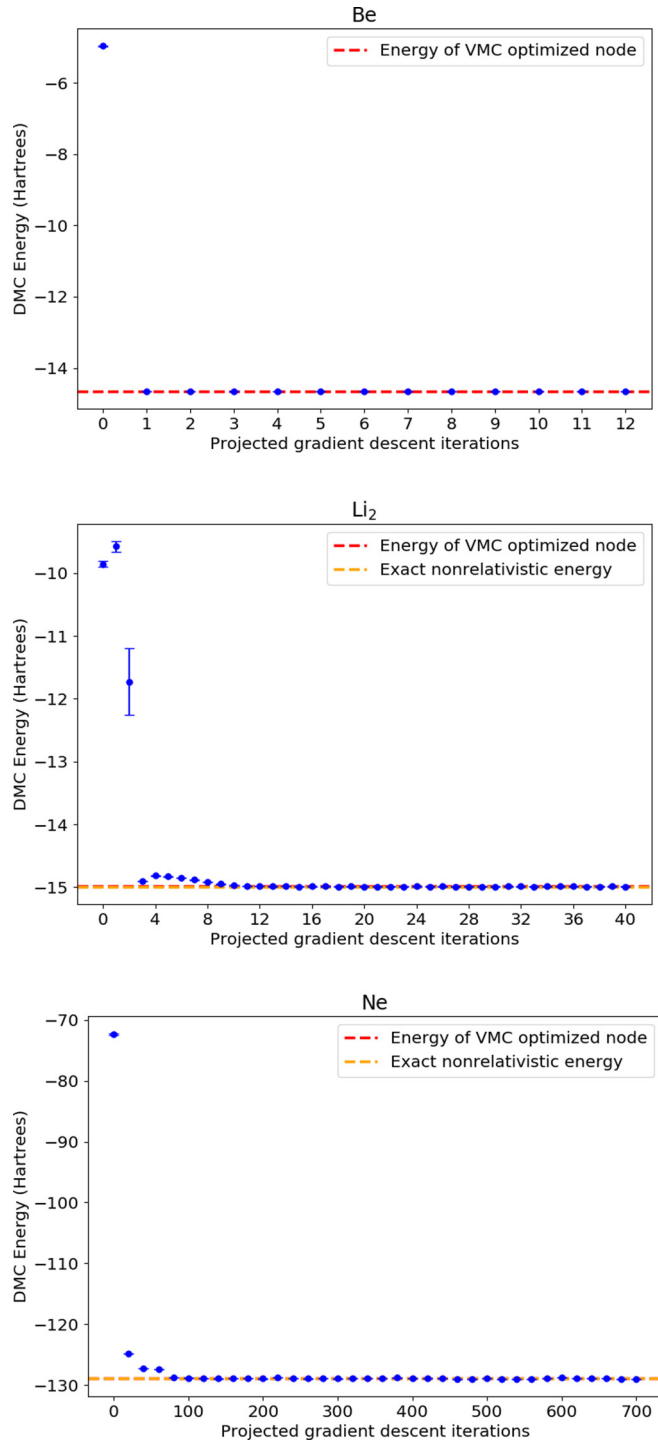


FIG. 4. Projected gradient descent of initially random parameters. See text for detailed explanation.

PGD parameter  $\alpha$  is sampled and, thus, its choice has some random element. For every such choice, the DMC evolution can be carried out with a controlled level of error because the nodal surface is fixed. Therefore, the best energy upper bound  $\mathcal{E}$  corresponds to the lowest value of the energy in the set, i.e.,  $\mathcal{E} = \min\{E_1, E_2, \dots, E_n, \dots\}$ . The last iteration energy obtained for the three systems studied in this paper

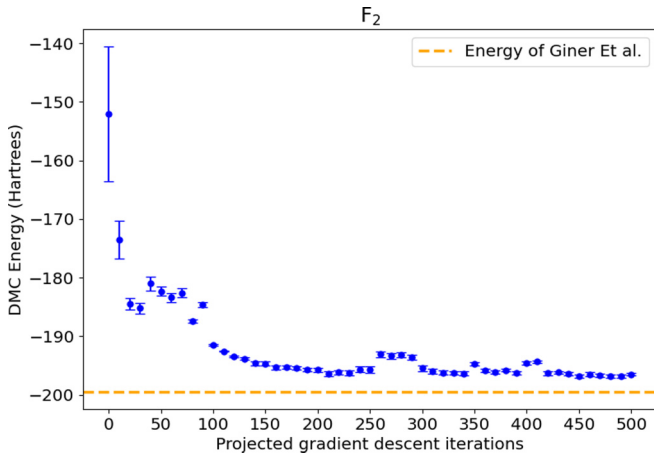


FIG. 5. See text for detailed explanation.

are compared with the best DMC energy obtained with nodes determined by a VMC optimization [23] in Table I.

The present calculation suggests that the nodes obtained with VMC optimization are not far from the optimum at least within the error of the present calculation. Therefore, the small discrepancy with the exact nonrelativistic values could be attributed to the fact that the form of the one-body factor is limited with the space of a single Slater determinant and the lack of spin-spin correlations because the one-body factor is taken to be a product of the form  $D^\uparrow D^\downarrow$ .

## V. DISCUSSION

We are not yet comfortable claiming that one of the methods A, B, or C is ideal, and the best choice may depend on circumstance. For example, although method A was least efficient in our test, it may not require as small of a time step, and we think it would be the simplest to implement into existing code due to its use of a standard walker distribution. Although the efficiency of method C was comparable to that of method B in our test, it requires additional wall time to calculate  $\xi$ , which also introduces the additional parameter  $\tau_c$  for which we do not yet have a good procedure for determining.

Nevertheless, method C performed the best in our tests and also has the advantage of not relying on the quality of  $\Psi_T$ . We recommend pairing method C with Adam. In our tests we were able to optimize the location of the nodes to roughly the same energy as with the VMC optimized nodes using any of the methods and without using Jastrow factors in our guiding function.

TABLE I. DMC energy comparison between the last PGD iteration nodes with VMC optimized nodes. Exact nonrelativistic values are presented.

System	Last PGD iteration	VMC optimized nodes	Exact energy
Be	-14.6567(5)	-14.6571(1)	-14.66736 [29]
Li <sub>2</sub>	-14.989(1)	-14.9898(1)	-14.9954 [29]
Ne	-128.92(1)	-128.919(3)	-128.939 [30]

Without the requirement of an optimized Jastrow function, our method is independent of prerequisite optimization by VMC. Nevertheless, simultaneously optimizing the Jastrow and Slater parameters may be a fruitful avenue, as a quality Jastrow function reduces many sources of error (except for the fixed-node error) and should improve the approximation of methods A and B. This could be done with the standard VMC method, or possibly with DMC and gradient descent, which require a gradient different than one of the DMC energy, perhaps one of the local energy variance.

An important advantage of PGD optimization is that the use of DMC propagation produces correct correlations in the fixed-node wave function, whereas VMC optimizes the nodal surface within a limited Jastrow form, which captures only pairwise correlations. Thus, the exact path to optimization due to the interplay between the adjustment of these correlations as the position of the nodal surface changes and vice versa was not fully accounted for in previous DMC studies.

We would like to give an example of a system which demonstrates the significance of the present paper. Variational calculations of liquid <sup>3</sup>He, where just a Jastrow-Slater wave function is used, yield [31] an unphysical result that the fully polarized liquid is of lower energy at its equilibrium density than the unpolarized liquid at its equilibrium density. Adding three-body correlation factors in the wave function does not change the above conclusion. It is when one includes back-flow correlations, which are state dependent and modify the nodal surface of the variational wave function, that one finds [31] that the unpolarized state of liquid <sup>3</sup>He is energetically favorable.

Therefore, let us pretend that we do not know the fact that the naturally occurring liquid <sup>3</sup>He is spin unpolarized; we might then be naive and use a Slater-Jastrow variational calculation to determine the optimum polarization. The nodes that would be determined by such a VMC optimization would be those that correspond to a polarized liquid. A subsequent fixed-node DMC calculation with these nodes is not going to change this result, and we would reach an unphysical conclusion. The PGD method described in the present paper, however, should find that the unpolarized state is the ground state.

Another advantage of PGD is that it requires only one gradient calculation to shift all the parameters in the correct direction, whereas if one were to vary the parameters and select values with lower calculated DMC energy, the number of required energy calculations would scale with the number of parameters used. We should point out that in our PGD evolution examples, far more CPU time was spent calculating the energy of a single nodal surface than the entire PGD process.

PGD optimization opens up a possible way to optimize atomic positions to minimize DMC energy. If one were to combine the nuclei and electrons into one walker type, then DMC imaginary-time propagation combined with the PGD method should naturally relax the atomic positions. Unfortunately we could not attempt this because the electron-nuclear cusp of the particular type of  $\Psi_T$  we used depends on both the atomic positions and the variational parameters, requiring the atoms be fixed. Using a  $\Psi_T$  that does not have this dependence would be an interesting avenue to explore.



## VI. SUMMARY

We presented a general and self-reliant method using DMC and projected gradient descent that optimizes the nodal surface of a trial function to minimize the fixed-node ground state, i.e., the DMC energy. We derived three methods for calculating the DMC energy gradient from walker distribution (methods A, B, and C). Methods B and C required comparable CPU time in our test, while method A was more expensive.

We combined the three methods with several gradient-descent algorithms and a projection operation that maintains the cusp condition. We benchmarked this projected gradient-descent method to trial functions with randomized parameters of Be, Li<sub>2</sub>, and Ne. Their energies were lowered to the same level as VMC optimized parameters, without the use of a Jastrow factor.

Our paper is a proof of concept using simple systems, but we see no reason it cannot be applied to larger and periodic systems, or be implemented into existing DMC code for various types of node determining wave functions. However, before attempting these goals, our paper may be further improved with better and more adaptive choices for various parameters, by simultaneously optimizing Jastrow parameters during PGD, by better suppressing PGD fluctuations, and by establishing criteria for convergence. In addition, we see no reason why the method could not be extended to DMC on a discrete lattice or to finite temperature path-integral Monte Carlo.

## ACKNOWLEDGMENTS

We would like to thank the Florida State University high performance computing center for the allocation of CPU time to carry out the calculations presented in this paper.

## APPENDIX A: PARAMETER FLUCTUATIONS

After a number of PGD iterations, the parameters will continue to fluctuate around a local minimum due to stochastic error of the energy gradient. Here, we investigate the variance of these fluctuations. First let us perform a coordinate shift so that the parameters are zero at the local minima, then let us rotate them so that the Hessian  $\frac{\partial^2 E}{\partial \theta_i \partial \theta_j}$  is diagonal. Then, if we expand  $E$  to second order in  $\theta_i$ , one iteration of gradient descent using Eq. (1) will result in new parameters  $\theta'_i$  given by

$$\theta'_i = \theta_i - a_i(K_i\theta_i + \sigma_i\eta), \quad (\text{A1})$$

where  $K_i \equiv \frac{\partial^2 E}{\partial \theta_i^2}$ ,  $\sigma_i$  is the standard deviation of stochastic error, and  $\eta$  is a random number with  $\langle \eta \rangle = 0$  and  $\langle \eta^2 \rangle = 1$ .

Let us take the variance of both sides of Eq. (A1):

$$\langle \theta_i'^2 \rangle = \langle \theta_i^2 \rangle + a_i^2(K_i^2\langle \theta_i^2 \rangle + \sigma_i^2) - 2a_iK_i\langle \theta_i^2 \rangle. \quad (\text{A2})$$

After enough PGD iterations, the variance of fluctuations will become constant, i.e.,

$$\langle \theta_i'^2 \rangle = \langle \theta_i^2 \rangle. \quad (\text{A3})$$

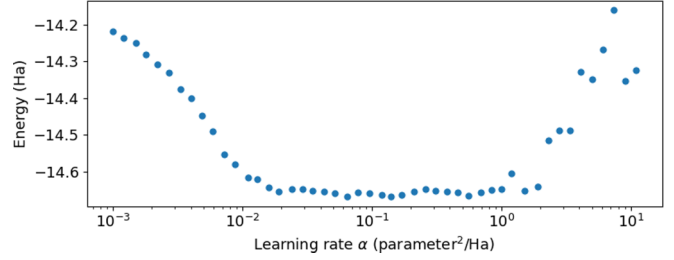


FIG. 6. An example of how the energy, calculated after 100 PGD iterations, depends on different learning rates  $\alpha$ . These data were used for the sixth row of Table II.

If this is the case, we can use Eq. (A1) to evaluate the variance of latter PGD iterations, namely,

$$\langle \theta_i'^2 \rangle = \frac{a_i\sigma_i^2}{2K_i - a_iK_i^2}. \quad (\text{A4})$$

We notice that the variance blows up if any  $a_i$  approach  $2/K_i$ , which is the result of overshooting the minimum by more than double the distance. For  $a_i \ll 2/K_i$  we see that the variance has the property

$$\langle \theta_i'^2 \rangle \propto a_i\sigma_i^2. \quad (\text{A5})$$

## APPENDIX B: GRADIENT-DESCENT ALGORITHMS

We list the gradient-descent algorithms we tested. We use  $g_i \equiv \frac{\partial E}{\partial \theta_i}$  as the gradient, and determine  $\alpha$  empirically. We sometimes decrease  $\alpha$  with iteration number to suppress parameter fluctuations around the minimum.

Basic gradient descent:

$$\theta_i \rightarrow \theta_i - \alpha g_i. \quad (\text{B1})$$

RMSprop:

$$\begin{aligned} v_i &\rightarrow \beta v_i + (\beta - 1)g_i^2, \\ \theta_i &\rightarrow \theta_i - \frac{\alpha g_i}{\sqrt{v_i}}. \end{aligned} \quad (\text{B2})$$

Adam:

$$\begin{aligned} t &\rightarrow t + 1, \\ m_i &\rightarrow \beta_1 m_i + (1 - \beta_1)g_i, \\ v_i &\rightarrow \beta_2 v_i + (\beta_2 - 1)g_i^2, \\ \theta_i &\rightarrow \theta_i - \frac{\alpha m_i g_i}{(1 - \beta_1^t)\sqrt{v_i/(1 - \beta_2^t)}}. \end{aligned} \quad (\text{B3})$$

## APPENDIX C: EFFECTIVE $\alpha$

In Table II we show the effective ranges of the learning rate  $\alpha$ , which is a hyperparameter proportional to a shift in parameters, as shown in Appendix B. We define this effective range as the values that bring the energy of a randomized  $\Psi_S$  down to the optimum level, and show an example in Fig. 6. We empirically determine this range for 100 PGD iterations for different combinations of gradient-descent algorithms, atomic systems, and methods A, B, and C. We display the end points of the effective range along with the logarithmic distance between them.

TABLE II. The range of the effective learning rate  $\alpha$  for 100 iterations for three gradient-descent algorithms: basic gradient descent (common  $a_i$ ), RMSprop, and Adam (see Fig. 6 for an example of the tenth row). RMSprop used  $\beta = 0.99$ . Adam used  $\beta_1 = 0.9$ ,  $\beta_2 = 0.99$ . Method C used  $\tau_c = 0.3 \text{ Ha}^{-1}$ .

Algorithm	Gradient method	$\Psi_5$	Minimum effective $\alpha$	Maximum effective $\alpha$	$\log_{10}(\frac{\alpha_{\max}}{\alpha_{\min}})$
Basic gradient descent	B	Be	$1.2 \times 10^{-1}$	$1.3 \times 10^0$	1.0
	C	Be	$2.4 \times 10^{-2}$	$6.4 \times 10^{-1}$	1.4
	A	Li <sub>2</sub>	$4.8 \times 10^{-3}$	$2.0 \times 10^{-1}$	1.6
RMSprop	B	Ne	$4.5 \times 10^{-2}$	$1.8 \times 10^0$	1.6
	B	Be	$5.5 \times 10^{-2}$	$1.0 \times 10^0$	1.3
	C	Be	$2.0 \times 10^{-2}$	$1.0 \times 10^{-1}$	1.7
	A	Li <sub>2</sub>	$1.8 \times 10^{-2}$	$2.1 \times 10^3$	5.1
Adam	B	Ne	$5.6 \times 10^{-2}$	$1.5 \times 10^0$	1.4
	B	Be	$3.8 \times 10^{-2}$	$1.0 \times 10^0$	1.4
	C	Be	$1.1 \times 10^{-2}$	$1.0 \times 10^0$	2.0
	A	Li <sub>2</sub>	$1.6 \times 10^{-1}$	$7.0 \times 10^2$	3.6
	B	Ne	$2.4 \times 10^{-1}$	$9.2 \times 10^0$	1.6

Compared to basic gradient descent, the adaptive algorithms RMSprop and Adam have more consistent end points for methods A, B, and C and for different atomic systems, making the choice of  $\alpha$  easier. We attribute the more consistent end points to the division of each gradient component by a trailing root mean square

of its prior values. This suppresses the large variations of gradient magnitudes for different atomic systems. It also suppresses large variations of  $\xi(R, \tau_c)$  of Eq. (35) used for method C, which has a factor of  $\exp(\Delta E \tau_c)$ , with  $\Delta E$  the difference between estimated and correct energy.

- [1] J. B. Anderson, A random-walk simulation of the schrödinger equation: H+3, *J. Chem. Phys.* **63**, 1499 (1975).
- [2] D. M. Ceperley and B. J. Alder, Ground State of the Electron Gas by a Stochastic Method, *Phys. Rev. Lett.* **45**, 566 (1980).
- [3] D. M. Ceperley and M. H. Kalos, Quantum many-body problems, in *Monte Carlo Methods in Statistical Physics*, edited by K. Binder (Springer-Verlag, Berlin, 1986), pp. 145–194.
- [4] E. Manousakis, The spin-1/2 Heisenberg antiferromagnet on a square lattice and its application to the cuprous oxides, *Rev. Mod. Phys.* **63**, 1 (1991).
- [5] J. Carlson, Green’s function Monte Carlo study of light nuclei, *Phys. Rev. C* **36**, 2026 (1987).
- [6] S. Gandolfi, D. Lonardonì, A. Lovato, and M. Piarulli, Atomic nuclei from quantum Monte Carlo calculations with chiral  $\text{e}ft$  interactions, *Front. Phys.* **8**, 117 (2020).
- [7] S. C. Pieper, R. B. Wiringa, and J. Carlson, Quantum Monte Carlo calculations of excited states in  $a = 6 - 8$  nuclei, *Phys. Rev. C* **70**, 054325 (2004).
- [8] M. Boninsegni and E. Manousakis, Green’s-function Monte Carlo study of the t-j model, *Phys. Rev. B* **46**, 560 (1992).
- [9] M. Boninsegni and E. Manousakis, Two-hole d-wave binding in the physical region of the t-j model: A Green’s-function Monte Carlo study, *Phys. Rev. B* **47**, 11897 (1993).
- [10] C. S. Hellberg and E. Manousakis, Phase Separation at all Interaction Strengths in the t-j Model, *Phys. Rev. Lett.* **78**, 4609 (1997).
- [11] S. Zhang, J. Carlson, and J. E. Gubernatis, Constrained Path Quantum Monte Carlo Method for Fermion Ground States, *Phys. Rev. Lett.* **74**, 3652 (1995).
- [12] S. Zhang and H. Krakauer, Quantum Monte Carlo Method Using Phase-Free Random Walks with Slater Determinants, *Phys. Rev. Lett.* **90**, 136401 (2003).
- [13] M. Qin, C.-M. Chung, H. Shi, E. Vitali, C. Hubig, U. Schollwöck, S. R. White, and S. Zhang (Simons Collaboration on the Many-Electron Problem), Absence of Superconductivity in the Pure Two-Dimensional Hubbard Model, *Phys. Rev. X* **10**, 031016 (2020).
- [14] R. Jastrow, Many-body problem with strong forces, *Phys. Rev.* **98**, 1479 (1955).
- [15] J. C. Slater, The theory of complex spectra, *Phys. Rev.* **34**, 1293 (1929).
- [16] K. M. Rasch, S. Hu, and L. Mitás, Communication: Fixed-node errors in quantum Monte Carlo: Interplay of electron density and node nonlinearities, *J. Chem. Phys.* **140**, 041102 (2014).
- [17] P. J. Reynolds, D. M. Ceperley, B. J. Alder, and W. A. Lester, Fixed-node quantum Monte Carlo for molecules, *J. Chem. Phys.* **77**, 5593 (1982).
- [18] C. J. Umrigar, K. G. Wilson, and J. W. Wilkins, Optimized Trial Wave Functions for Quantum Monte Carlo Calculations, *Phys. Rev. Lett.* **60**, 1719 (1988).
- [19] C. J. Umrigar and C. Filippi, Energy and Variance Optimization of Many-Body Wave Functions, *Phys. Rev. Lett.* **94**, 150201 (2005).
- [20] F. A. Reboredo, R. Q. Hood, and P. R. C. Kent, Self-healing diffusion quantum Monte Carlo algorithms: Direct reduction of the fermion sign error in electronic structure calculations, *Phys. Rev. B* **79**, 195117 (2009).

- [21] M. E. Foulaadvand and M. Zarenia, Optimization of quantum Monte Carlo wave function: Steepest decent method, *Int. J. Mod. Phys. C* **21**, 523 (2010).
- [22] D. H. Berman, Boundary effects in quantum mechanics, *Am. J. Phys.* **59**, 937 (1991).
- [23] C. J. Umrigar, M. P. Nightingale, and K. J. Runge, A diffusion Monte Carlo algorithm with very small time-step errors, *J. Chem. Phys.* **99**, 2865 (1993).
- [24] T. Kato, On the eigenfunctions of many-particle systems in quantum mechanics, *Commun. Pure Appl. Math.* **10**, 151 (1957).
- [25] E. Steiner, Charge densities in atoms, *J. Chem. Phys.* **39**, 2365 (1963).
- [26] D. P. Kingma and J. Ba, Adam: A method for stochastic optimization, [arXiv:1412.6980](https://arxiv.org/abs/1412.6980).
- [27] G. Hinton, *Neural Networks for Machine Learning, Lecture* (Department of Computer Science, the University of Toronto, Toronto, Canada, 2012).
- [28] E. Giner, A. Scemama, and M. Caffarel, Fixed-node diffusion Monte Carlo potential energy curve of the fluorine molecule f2 using selected configuration interaction trial wavefunctions, *J. Chem. Phys.* **142**, 044115 (2015).
- [29] D. Bressanini, G. Morosi, and S. Tarasco, An investigation of nodal structures and the construction of trial wave functions, *J. Chem. Phys.* **123**, 204109 (2005).
- [30] E. R. Davidson, S. A. Hagstrom, S. J. Chakravorty, V. M. Umar, and C. F. Fischer, Ground-state correlation energies for two- to ten-electron atomic ions, *Phys. Rev. A* **44**, 7071 (1991).
- [31] E. Manousakis, S. Fantoni, V. R. Pandharipande, and Q. N. Usmani, Microscopic calculations for normal and polarized liquid  $^3\text{He}$ , *Phys. Rev. B* **28**, 3770 (1983).