






This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0157512

1 **DeepQMC: an open-source software suite for variational optimization of**
2 **deep-learning molecular wave functions.**

3 Z. Schätzle ^{1, a)} P. B. Szabó ^{1, a)} M. Mezera ¹ J. Hermann ² and F. Noé
4 ^{1, 2, 3, 4}

5 ¹⁾*FU Berlin, Department of Mathematics and Computer Science, Arnimallee 6,*
6 *14195 Berlin, Germany*

7 ²⁾*Microsoft Research AI4Science, Karl-Liebknecht Str. 32, 10178 Berlin,*
8 *Germany*

9 ³⁾*FU Berlin, Department of Physics, Arnimallee 14, 14195 Berlin,*
10 *Germany*

11 ⁴⁾*Rice University, Department of Chemistry, Houston, TX 77005,*
12 *USA*

13 (Dated: 17 July 2023)

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0157512

14 Computing accurate yet efficient approximations to the solutions of the electronic
15 Schrödinger equation has been a paramount challenge of computational chemistry for
16 decades. Quantum Monte Carlo methods are a promising avenue of development as their
17 core algorithm exhibits a number of favorable properties: it is highly parallel, and scales fa-
18 vorably with the considered system size, with an accuracy that is limited only by the choice
19 of the wave function ansatz. The recently introduced machine-learned parametrizations of
20 quantum Monte Carlo ansatzes rely on the efficiency of neural networks as universal func-
21 tion approximators to achieve state of the art accuracy on a variety of molecular systems.
22 With interest in the field growing rapidly, there is a clear need for easy to use, modular, and
23 extendable software libraries facilitating the development and adoption of this new class
24 of methods. In this contribution, the DEEPQMC program package is introduced, in an
25 attempt to provide a common framework for future investigations by unifying many of the
26 currently available deep-learning quantum Monte Carlo architectures. Furthermore, the
27 manuscript provides a brief introduction to the methodology of variational quantum Monte
28 Carlo in real space, highlights some technical challenges of optimizing neural network
29 wave functions, and presents example black-box applications of the program package. We
30 thereby intend to make this novel field accessible to a broader class of practitioners both
31 from the quantum chemistry as well as the machine learning communities.

^{a)}Z. Schätzle and P. B. Szabó contributed equally to this work.

32 I. INTRODUCTION

33 Recently, the application of machine learning to a wide range of problems from the natural sci-
34 ences has proven to be highly successful. Computational chemistry is a field of particular activity:
35 machine learning force fields model complicated quantum mechanical effects at the resolution of
36 atoms, while machine learned functionals elevate density functional theory to unprecedented accu-
37 racy¹⁻³. These approaches utilize supervised training to learn from accurate quantum mechanical
38 reference calculations, and make predictions for unseen configurations. While this results in fast
39 yet accurate approximations to the quantum many-body problem, it inherently depends on high
40 quality training data, which represents a major bottleneck of these methods.

41 An orthogonal way to incorporate machine learning into computational chemistry is its appli-
42 cation to improve ab-initio calculations. Notably, over the course of the last years a new family of
43 deep-learning quantum Monte Carlo (deep QMC) methods has developed, incorporating advance-
44 ments from the field of machine learning⁴. Common to the deep QMC methods is the utilization
45 of neural networks to parametrize highly expressive ansatzes, efficiently approximating the so-
46 lutions of the time-independent electronic Schrödinger equation, thereby providing a complete
47 description of the system's electronic properties. Originating from spin lattices⁵, deep-learning
48 ansatzes were soon applied to molecules in real-space⁶. With the development of PauliNet⁷ and
49 FermiNet⁸, the accuracy of neural-network wave functions became the state of the art within vari-
50 ational Monte Carlo. Subsequent works have further increased the accuracy of these ansatzes⁹⁻¹²,
51 extended them to the simulation of excited states¹³ as well as periodic systems^{14,15}, combined
52 them with pseudo-potentials¹⁶, used them in the calculation of interatomic forces¹⁷, utilized them
53 in diffusion Monte Carlo simulations^{18,19}, and extended them to share parameters across multiple
54 molecular geometries²⁰⁻²² or distinct molecules^{23,24}.

55 Although the method of optimizing deep-learning wave function ansatzes using variational
56 quantum Monte Carlo was developed only a few years ago, it already competes with some of the
57 most accurate traditional quantum chemistry methods on molecules with up to ~ 100 electrons.
58 Exhibiting competitive scaling with the number of electrons, it has the potential to be extended to
59 larger systems in the near future. Achieving this will no doubt require further method development
60 as well as efficient implementations of the core algorithms, creating the need for open source
61 libraries that facilitate experimentation and contribution from the community.

62 Accompanying the above summarized research, various software libraries for variational op-

63 timization of deep-learning wave functions have been released^{25–28}. While NETKET²⁵ provides
 64 a general implementation of variational optimization of machine learning wave functions mainly
 65 for lattice systems with recent extensions to continuous space, the research for molecular machine
 66 learning wave functions was carried out across various repositories and is lacking a unified frame-
 67 work. The presented DEEPQMC program package aims to provide a unified implementation of
 68 the developments in the field of deep-learning molecular wave functions. It intends to be easy to
 69 use out of the box, while allowing full control and flexibility over the ansatz and variational train-
 70 ing for advanced users. The library is designed to be modular, facilitating the rapid development
 71 and testing of individual components, and easing the implementation of new features. It makes use
 72 of the composable function transformations and just-in-time compilation provided by the JAX li-
 73 brary²⁹ to express performant GPU accelerated algorithms using concise Python³⁰ syntax. Neural
 74 network models are encapsulated in higher-level objects, using the haiku deep-learning library³¹.
 75 The project is open-source and distributed online under the MIT license²⁶.

76 II. THEORY

77 A. The electronic structure problem

78 In computational chemistry, molecular systems are often described by the non-relativistic
 79 molecular Hamiltonian using the Born–Oppenheimer approximation:

$$\hat{H} = \sum_{i=1}^N \left(-\frac{1}{2} \Delta_{\mathbf{r}_i} - \sum_{I=1}^M \frac{Z_I}{|\mathbf{r}_i - \mathbf{R}_I|} + \sum_{j=1}^{i-1} \frac{1}{|\mathbf{r}_i - \mathbf{r}_j|} \right), \quad (1)$$

80 where \mathbf{r}_i denotes the position of the i th electron, while Z_I and \mathbf{R}_I are the charge and position of the
 81 I th nucleus, respectively. To determine the electronic structure of these molecular systems, one
 82 must solve the associated time independent Schrödinger equation

$$\hat{H}\psi(\mathbf{x}_1, \dots, \mathbf{x}_N) = E\psi(\mathbf{x}_1, \dots, \mathbf{x}_N), \quad (2)$$

83 where $\mathbf{x}_i = (\mathbf{r}_i, \sigma_i)$ comprise the positions of the electrons and their spin. A solution is an eigen-
 84 function of the Hamiltonian, the electronic wave function ψ , and its corresponding energy eigen-
 85 value E . With the exact wave function at hand, any observable electronic property of the system
 86 can in principle be computed, as the wave function gives a complete description of the system's
 87 electronic state. Since electrons have half-integer spin, their wave functions must be antisymmetric

88 with respect to electron exchanges

$$\begin{aligned} \psi(\dots, \mathbf{x}_i, \dots, \mathbf{x}_j, \dots) = \\ -\psi(\dots, \mathbf{x}_j, \dots, \mathbf{x}_i, \dots). \end{aligned} \quad (3)$$

89 While general wave functions are complex-valued, the solutions of the time independent Schrödinger
90 equation can be chosen as real without loss of generality, due to the hermiticity of the molecular
91 Hamiltonian. Therefore in all of the following discussions, as well as in the DEEPQMC library,
92 only real valued wave functions are considered.

93 B. Variational optimization

94 Even with the aforementioned approximations, the electronic Schrödinger equation involving
95 the molecular Hamiltonian can only be solved analytically for hydrogenic atoms – the special
96 case of the two-body problem. This makes computational quantum chemistry a mainly numerical
97 field, where different methods yield approximate solutions at varying trade-offs of accuracy and
98 computational cost. The class of variational quantum chemistry methods phrases the solution of
99 the Schrödinger equation as a minimization problem. The ground state of the Hamiltonian is
100 approximated by optimizing the parameters θ of a trial wave function (ansatz) ψ_θ , to minimize
101 the expectation value of the Hamiltonian

$$\theta' = \underset{\theta}{\operatorname{argmin}} \langle \hat{H} \rangle_{\psi_\theta}. \quad (4)$$

102 This objective is rooted in the variational principle of quantum mechanics, which states that the
103 ground state energy of the Hamiltonian is a lower bound for the energy expectation value of any
104 wave function from the associated antisymmetric Hilbert space H^-

$$E_0 \leq \min_{\psi} \langle \hat{H} \rangle_{\psi} \quad \psi \in H^-. \quad (5)$$

105 The variational methods can be categorized based on the means of calculating the expectation
106 value $\langle \cdot \rangle$, and choice of ansatz ψ_θ .

107 The DEEPQMC program package implements VMC in real space (first quantization) with
108 neural network wave functions. In VMC, expectation values are estimated through a stochastic

109 sampling of electron configurations

$$\begin{aligned}
 \langle \hat{H} \rangle_{\psi_{\theta}} &= \frac{\langle \psi_{\theta} | \hat{H} | \psi_{\theta} \rangle}{\langle \psi_{\theta} | \psi_{\theta} \rangle} \\
 &= \frac{\int d\mathbf{r}_1, \dots, d\mathbf{r}_N \psi_{\theta}^*(\mathbf{r}_1, \dots, \mathbf{r}_N) \hat{H} \psi_{\theta}(\mathbf{r}_1, \dots, \mathbf{r}_N)}{\int d\mathbf{r}_1, \dots, d\mathbf{r}_N |\psi_{\theta}(\mathbf{r}_1, \dots, \mathbf{r}_N)|^2} \\
 &= \frac{\int d\mathbf{r}_1, \dots, d\mathbf{r}_N |\psi_{\theta}(\mathbf{r}_1, \dots, \mathbf{r}_N)|^2 E_{\text{loc}}[\psi_{\theta}](\mathbf{r})}{\int d\mathbf{r}_1, \dots, d\mathbf{r}_N |\psi_{\theta}(\mathbf{r}_1, \dots, \mathbf{r}_N)|^2} \\
 &= \mathbb{E}_{\mathbf{r} \sim |\psi_{\theta}|^2} [E_{\text{loc}}[\psi_{\theta}](\mathbf{r})] \\
 &\approx \frac{1}{n} \sum_{\mathbf{r} \sim |\psi_{\theta}|^2}^n E_{\text{loc}}[\psi_{\theta}](\mathbf{r}).
 \end{aligned} \tag{6}$$

110 Because the molecular Hamiltonian does not depend on the spin, it is possible to compute the
 111 energy using the spatial wave function $\psi(\mathbf{r}_1, \dots, \mathbf{r}_N)$, where fixed spins are assigned to the electrons
 112 and spin-up and spin-down electrons are treated as distinguishable³². The convention is to sort the
 113 electrons by spin and consider the first N^{\uparrow} electrons to have spin-up and the remaining $N^{\downarrow} = N - N^{\uparrow}$
 114 electrons to possess spin-down.

115 In practice, a VMC simulation then consists of choosing an ansatz (see Section III), and opti-
 116 mizing it in an alternating scheme of sampling and parameter updates. The expectation value in
 117 (6) is approximated by sampling the probability density given by the square of the wave function
 118 (see Section V), followed by a parameter update using the gradient of the expectation value (see
 119 Section IV).

120 C. Neural network wave functions

121 Being exact in principle, the choice of the wave function ansatz is crucial for the efficiency
 122 of a VMC simulation. Recently, neural network parametrizations of real-space molecular wave
 123 functions were introduced by PauliNet⁷, and FermiNet⁸. They both rely on generalized Slater
 124 determinants, that augment the single particle orbitals of conventional Slater determinants with
 125 many-body correlation,

$$126 \quad \psi_{\theta}(\mathbf{r}_1, \dots, \mathbf{r}_N) = \sum_p c_p \det[\mathbf{A}^p(\mathbf{r})], \tag{7}$$

$$127 \quad A_{ik}^p = \phi_k^p(\mathbf{r}_i, \{\mathbf{r}^{\uparrow}\}, \{\mathbf{r}^{\downarrow}\}) \times \varphi_k^p(\mathbf{r}_i). \tag{8}$$

128 Here, ϕ_k^p are many-body orbitals, and φ_k^p are single particle envelopes that ensure the correct
 129 asymptotic behavior of the wave function with increasing distance from the nuclei. The set no-
 129 tation $\{\cdot\}$ is to be understood as a permutation invariant dependence on the spin-up electrons \mathbf{r}^{\uparrow} ,

130 and spin-down electrons \mathbf{r}^\downarrow , respectively. The ansatz may be a linear combination of multiple
 131 generalized Slater determinants, which are distinguished with the p index. The form of ϕ_k^p in
 132 (8) is closely related to the backflow transformation³³, which introduces quasi-particles to obtain
 133 many-body orbital functions. The key observation motivating this augmentation is that the anti-
 134 symmetry of Slater determinants constructed from many-body orbitals is preserved as long as the
 135 orbital functions are equivariant with the exchange of electrons,

$$P_{ij}^{\parallel} \phi_k(\mathbf{r}_i, \{\mathbf{r}^\uparrow\}, \{\mathbf{r}^\downarrow\}) = \phi_k(\mathbf{r}_j, \{\mathbf{r}^\uparrow\}, \{\mathbf{r}^\downarrow\}), \quad (9)$$

136 where P_{ij}^{\parallel} is the operator exchanging same-spin electrons i and j .

137 Most of the currently used deep-learning molecular wave functions^{7,8,11} share the functional
 138 form of (7), and (8), and differ only in the parametrization of the many-body orbitals ϕ_k^p and
 139 single-particle envelopes φ_k^p . DEEPQMC aims to provide a general framework for variational
 140 optimization of deep-learning molecular wave functions, facilitating the investigation of the design
 141 space spanned by the PauliNet, FermiNet, and DeepErwin neural network ansatzes.

142 D. Pseudopotentials

143 Despite the favorable asymptotic scaling of VMC with the number of electrons, systems con-
 144 taining heavy nuclei remain challenging due to a variety of reasons. The high energy of electrons
 145 near these nuclei complicates simulations by spoiling the optimization and reducing the effective-
 146 ness of MCMC sampling. Furthermore, the kinetic energy of these electrons reaches the relativis-
 147 tic regime, requiring the treatment of relativistic effects that are not accounted for in the standard
 148 non-relativistic molecular Hamiltonian of (1). On the other hand, while the core regions of heav-
 149 ily charged nuclei contribute dominantly to the total energy, they are typically unchanged during
 150 chemical processes and thus have little effect on the chemically relevant relative energies. There-
 151 fore most quantum chemistry methods targeted at computing relative energies reduce the above
 152 outlined difficulties, by treating the outer (valence) electrons separately from the inner (core) elec-
 153 trons.

154 The approach most suited for implementation in the context of variational optimization of deep-
 155 learning wave functions is the use of pseudopotentials, which has been previously demonstrated
 156 by Li et al.¹⁶. In this method, the core electrons are excluded from the explicit calculation and
 157 replaced by additional terms in the Hamiltonian, to simulate their influence on the remaining N_v

158 valence electrons. The modified Hamiltonian reads as

$$\hat{H} = \sum_{i=1}^{N_v} \left(-\frac{1}{2} \Delta_{\mathbf{r}_i} + \sum_{j=1}^{i-1} \frac{1}{|\mathbf{r}_i - \mathbf{r}_j|} \right) + \hat{V}_{\text{PP}}. \quad (10)$$

159 The \hat{V}_{PP} pseudopotential term is in turn decomposed to local and non-local parts

$$\hat{V}_{\text{PP}} = \sum_{I=1}^M \sum_{i=1}^{N_v} \left(V^I(r_{iI}) + \sum_{l=0}^{l_{\text{max}}} W_l^I(r_{iI}) \hat{P}_l^{iI} \right), \quad (11)$$

160 where $r_{iI} = |\mathbf{r}_i - \mathbf{R}_I|$, V^I and W_l^I are sets of scalar functions describing the local and non-local
161 pseudopotential contributions, while $\hat{P}_l^{iI} = \sum_{m=-l}^{m=l} |lm\rangle_{iI} \langle lm|_{iI}$ is a projection operator of the i -th
162 electron on spherical harmonics centered on the I -th nucleus. To evaluate the contribution of the
163 non-local part of the pseudopotential (second term of (11)) one considers integrals of the form

$$\frac{\langle \mathbf{r} | W_l^I \hat{P}_l^{iI} | \psi \rangle}{\langle \mathbf{r} | \psi \rangle} = W_l^I(r_{iI}) \sum_{m=-l}^l Y_{lm}(\Omega_{iI}) \times \int Y_{lm}(\Omega'_{iI})^* \frac{\Psi(\mathbf{r}_1, \dots, (r_{iI}, \Omega'_{iI}), \dots, \mathbf{r}_N)}{\Psi(\mathbf{r}_1, \dots, (r_{iI}, \Omega_{iI}), \dots, \mathbf{r}_N)} d\Omega'_{iI} \quad (12)$$

164 where Y_{lm} is a spherical harmonic and (r_{iI}, Ω_{iI}) denotes the position vector of the i -th electron \mathbf{r}_i ,
165 expressed in spherical coordinates centered on nucleus I . Following the first implementation of
166 pseudopotentials for deep-learning molecular wave functions by Li et al.¹⁶, the above integral is
167 approximated using an icosahedral quadrature of 12-points.

168 The scalar functions V^I and W_l^I are typically pre-computed by expanding them in a Gaussian
169 basis, and fitting the expansion parameters directly to a database of reference energies. The DEEP-
170 QMC program package currently includes the widely used BFD³⁴, and the most recent ccECP³⁵
171 pseudopotentials, with an application of the latter presented in Section VIII B.

172 III. WAVE FUNCTION DESIGN SPACE

173 DEEPQMC implements a variety of options to obtain the equivariant many-body orbitals ϕ_k^p
174 and the accompanying envelopes ϕ_k^p , covering PauliNet, FermiNet, DeepErwin and their deriva-
175 tives. In the following, the main architectural concepts of these wave function ansatzes will be
176 described in more detail. For ease of use DEEPQMC provides predefined configuration files to
177 obtain the above mentioned ansatzes, while allowing their interpolation through a manual choice
178 of hyperparameters.

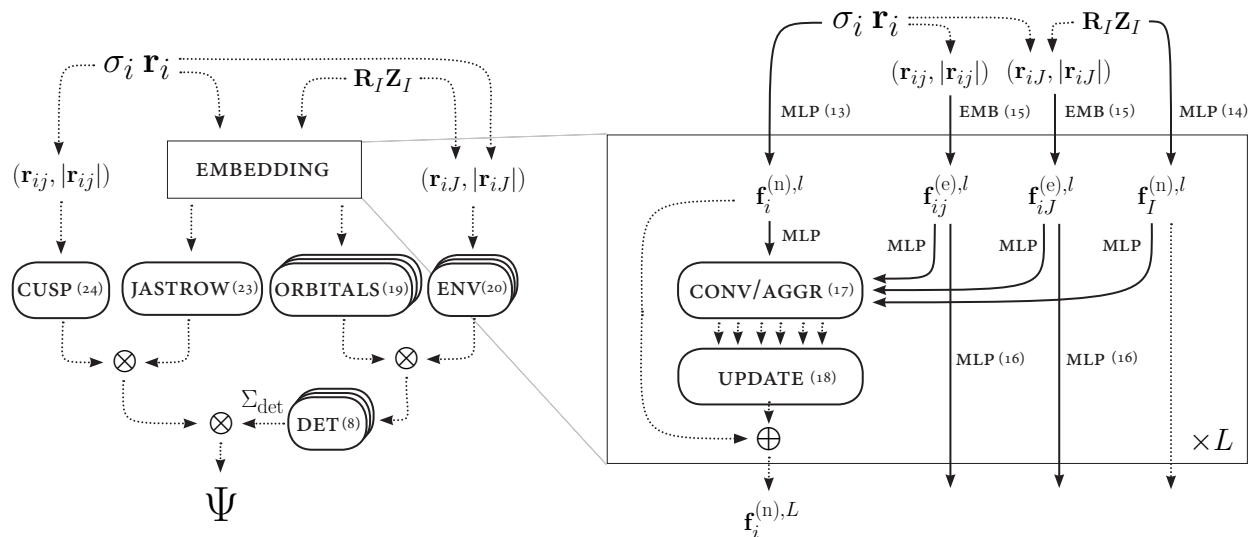


FIG. 1. **Sketch of a general neural network wave function ansatz and its graph neural network architecture.** This sketch comprises the implemented design space for the neural network wave function (left) and the GNN architecture (right). Solid lines can carry MLPs and dotted lines correspond to forwarding without further change. Numbers in parentheses refer to the corresponding equations of the main text. The choice among the drawn connections, the depths of the associated MLPs as well as the aggregation and update rules distinguish the previously published works PauliNet, FermiNet, and DeepErwin.

179 A. Graph Neural Networks

180 Central to the neural network wave function ansatzes is the computation of equivariant many-
 181 body embedding vectors for the electrons, which are used downstream to obtain the entries of the
 182 generalized Slater determinant. Many strategies of obtaining these embeddings can be unified in
 183 the framework of graph neural networks (GNNs). GNNs are well suited to model functions on
 184 graphs that exhibit certain permutational symmetries, and can be adapted to describe electrons of
 185 molecules in real-space.

186 An electronic configuration of a molecule can be encoded as a graph, where the nodes are
 187 electrons and nuclei, and the connecting edges carry pairwise features, e.g. difference vectors.
 188 GNNs are functions of these graphs, yielding high-dimensional latent space embeddings of the
 189 nodes. The electronic node embeddings are initialized with single-electron features and iteratively
 190 updated to incorporate many-body information of the electronic environment. Using graph con-
 191 volutions, the updates are invariant under the exchange of electrons in the environment, and the

192 conditions of (9) are fulfilled.

193 The most relevant aspects of the GNN architecture implemented in DEEPQMC are sketched
194 on the right pane of Figure 1, and are discussed in detail in the following. Electron positions
195 (spins) are denoted with \mathbf{r} (σ), while \mathbf{R} and Z indicate nuclei positions and charges. Node and
196 edge quantities are denoted with the superscripts (n) and (e), respectively. Furthermore, l indexes
197 the GNN interaction layers, θ denotes functions parametrized by MLPs, and t runs over the dif-
198 ferent edge types (those between electrons of identical or opposite spins, or between electrons and
199 nuclei), node types (electron or nuclei nodes) or message types. Lastly, vertical brackets indicate
200 the different options implemented in DEEPQMC for the computation of various quantities.

201 *The graph representation:*

202 A graph is a natural way to denote the electronic configuration of a molecule in real-space.
203 The nodes of the graph represent particles (electrons and nuclei), carrying information such as
204 spin or nuclear charge. The edges support the difference vectors between the particles, resulting
205 in a representation invariant under global translation. Note that using internal coordinates that
206 are invariant under global rotation may not be sufficient to represent all wave functions (simple
207 counterexamples are atomic wave functions with P symmetry), and can only be employed with a
208 careful treatment of the spatial symmetries.

209 To implement a variety of wave function ansatzes, DEEPQMC provides configuration options
210 to define the specifics of the graph construction outlined above. Most importantly, the nodes
211 corresponding to nuclei, and their respective nuclei-electron edges can optionally be excluded
212 from the graph. In this case, electron-nuclei information can still be introduced to the GNN, by
213 initializing the electron embeddings using a concatenation of the difference vectors between the
214 positions of the given electron and all nuclei (see the second case of (13)).

215 *Node features:*

216 The output of DEEPQMC GNNs are electron node embeddings $\mathbf{f}_i^{(n)}$ which are subsequently
217 used to generate the many-body orbitals that constitute the entries of the generalized Slater deter-
218 minants and an optional trainable Jastrow factor. To enforce the equivariance of these quantities
219 with respect to the exchange of electrons, the initialization of the electron embeddings has to be

220 chosen appropriately. In DEEPQMC one can either use identical embeddings for all electrons
221 of the same spin (invariant under permutation of same-spin electrons), or a concatenation of the
222 electron nuclei difference vectors (equivariant under electron permutations)

$$\mathbf{f}_i^{(n),0,\text{el}} = \begin{cases} \mathcal{E}_\theta^{(n),\text{el}}(\sigma_i) \\ \mathcal{E}_\theta^{(n),\text{el}}(\mathbf{r}_i - \mathbf{R}_1, \dots, \mathbf{r}_i - \mathbf{R}_N), \end{cases} \quad (13)$$

223 where $\mathcal{E}_\theta^{(n),\text{el}}$ are parameterized node embedding functions implemented through MLPs.

224 If the GNN is chosen to explicitly consider electron-nuclei interactions, the embeddings associ-
225 ated with the nuclear nodes have to be initialized besides the electronic embeddings. DEEPQMC
226 implements fixed nuclear node embeddings $\mathbf{f}_I^{(n)}$, that either distinguish all nuclei or depend on the
227 respective atom type:

$$\mathbf{f}_I^{(n),0,\text{nuc}} = \begin{cases} \mathcal{E}_\theta^{(n),\text{nuc}}(I) \\ \mathcal{E}_\theta^{(n),\text{nuc}}(Z_I). \end{cases} \quad (14)$$

228 **Edge features :**

229 The edges of the graph hold the pairwise differences of node positions (\mathbf{r}_{ij}) and their embed-
230 dings are consequently initialized as

$$\mathbf{f}_{ij}^{(e),0,t^{(e)}} = \mathcal{E}^{(e),t^{(e)}}(\mathbf{r}_{ij}), \quad (15)$$

231 where $\mathcal{E}^{(e),t^{(e)}}$ is an edge type dependent featurization based on the pairwise differences. This may
232 correspond to directly feeding the difference vectors, using the pairwise distances, expanding in
233 a basis of Gaussians or working with rescaled difference vectors amongst other options. In later
234 interaction layers, the original edge embeddings are either reused or iteratively updated,

$$\mathbf{f}_{ij}^{(e),l,t^{(e)}} = \begin{cases} \mathbf{f}_{ij}^{(e),0,t^{(e)}} \\ u_\theta^{l,t^{(e)}}(\mathbf{f}_{ij}^{(e),l-1,t^{(e)}}), \end{cases} \quad (16)$$

235 with the latter option making use of a parametrized update function $u_\theta^{l,t^{(e)}}$, thus increasing the
236 effective depth of the architecture at the cost of additional MLPs. The parameters of the update
237 function $u_\theta^{l,t^{(e)}}$ may be shared across different edge types.

238 **Message generation:**

239 The electron embeddings are updated in each interaction layer by aggregating messages passed
240 along the graph edges. These messages are constructed via an elementwise product between func-
241 tions of the sending node and edge embeddings (graph convolution),

$$\mathbf{m}_i^{l,t^{(m)}} = \sum_{j \in t^{(n)}} w_{\theta}^{l,t^{(e)}}(\mathbf{f}_{ij}^{(e),l,t^{(e)}}) \cdot h_{\theta}^{l,t^{(n)}}(\mathbf{f}_j^{(n),l,t^{(n)}}). \quad (17)$$

242 The superscript $t^{(n)}$ on the node features specifies the subset of sending nodes and the superscript
243 $t^{(e)}$ on the edge features depends on the type of the sending and the receiving nodes respectively.
244 The choice of how to distinguish electron-electron messages based on their spin (relative spin
245 of sending and receiving electron, spin of sending electron, or no distinction between messages
246 from spin-up and spin-down electrons) is another hyperparameter of the GNN. Optionally, the
247 above sum over the edges can be normalized by dividing it with the number of edges. Note that
248 messages depending only on node (edge) information can be obtained by setting the function $w_{\theta}^{l,t^{(e)}}$
249 ($h_{\theta}^{l,t^{(n)}}$) to return identity. The superscript $t^{(m)}$ runs over all the constructed messages, which may
250 include different choices of $w_{\theta}^{l,t^{(e)}}$ and $h_{\theta}^{l,t^{(n)}}$.

251 **Electron embedding update:**

252 To obtain updated electron embeddings, messages from various edge types are combined and
253 added to a residual connection. DEEPQMC implements a few protocols for the combination of
254 messages, that can be summarized as follows

$$\mathbf{f}_i^{(n),l+1,\text{el}} = \mathbf{f}_i^{(n),l,\text{el}} + \begin{cases} \sum_{t^{(m)}} g_{\theta}^{l,t^{(m)}}(\mathbf{m}_i^{l,t^{(m)}}) \\ g_{\theta}^l(\sum_{t^{(m)}} \mathbf{m}_i^{l,t^{(m)}}) \\ g_{\theta}^l(\oplus_{t^{(m)}} \mathbf{m}_i^{l,t^{(m)}}), \end{cases} \quad (18)$$

255 where \oplus refers to the concatenation of the messages. Additionally to the messages constructed
256 according to equation (17) the message types $t^{(m)}$ can include a residual connection $\mathbf{f}_i^{(n),l}$ such that
257 the trainable self-interaction of FermiNet and DeepErwin can be reproduced.

258 In the above outlined general GNN framework a wide variety of ansatzes can be obtained.
259 Furthermore, the implementation of DEEPQMC and its GNN framework focus on facilitating
260 rapid extensions with new ansatz variants either by exploration within the existing hyperparameter
261 space or by extending it with new features.

262 B. Orbital construction

263 The entries of the generalized Slater determinant in (8) are obtained as products of many-
264 body orbitals ϕ_k^p and envelopes φ_k^p . The many-body orbitals are functions of the final equivariant
265 electron embeddings

$$\phi_k^p(\mathbf{r}_i, \{\mathbf{r}^\uparrow\}, \{\mathbf{r}^\downarrow\}) = \kappa_\theta(\mathbf{f}_i^{(n),L}), \quad (19)$$

266 where κ_θ is an MLP applied electronwise, projecting the embedding dimension to the required
267 number of orbitals. For the φ_k^p envelopes, DEEPQMC implements linear combinations of expo-
268 nentials centered on the nuclei

$$\varphi_k^p(\mathbf{r}_i) = \sum_I \sum_{\beta_I} \omega_{k\beta_I}^p \exp\left(-\|\Sigma_{k\beta_I}^p(\mathbf{r}_i - \mathbf{R}_I)\|^\alpha\right), \quad (20)$$

269 where $\omega_{k\beta_I}^p$ and $\Sigma_{k\beta_I}^p$ are trainable parameters and β_I indexes the basis function centered on atom I .
270 The hyperparameter $\alpha \in (1, 2)$ represents the choice of Slater type orbitals with $\alpha = 1$, and Gaus-
271 sian type orbitals (GTOs) with $\alpha = 2$. DEEPQMC provides an option to restrict the envelopes to
272 be isotropic ($\Sigma_{k\beta_I}^p := \sigma_{k\beta_I}^p \cdot \mathbf{I}$). The GTOs can be initialized from the molecular orbital coefficients
273 of reference solutions with standard quantum chemistry basis sets obtained in PySCF³⁶.

274 C. Determinant construction

275 Because the antisymmetry of the wave function is required with respect to the exchange of
276 same-spin particles only, Slater determinants in VMC are typically considered block diagonal and
277 are factored into a spin-up and a spin-down component

$$\psi_\theta = \sum_p c_p \det[\mathbf{A}^{\uparrow,p}(\mathbf{r})] \det[\mathbf{A}^{\downarrow,p}(\mathbf{r})]. \quad (21)$$

278 Additionally, DEEPQMC implements the full determinant option explored by Lin et al.¹⁰, which
279 constructs a single determinant using both spin-up and spin-down electrons

$$\psi_\theta = \sum_p c_p \det[\mathbf{A}^{1,p}(\mathbf{r})]. \quad (22)$$

280 It's noted that since the many-body orbitals are not equivariant under the exchange of opposite
281 spin electrons, the full determinant ansatz is still not antisymmetric with respect to these per-
282 mutations. Instead, a full determinant can practically be understood as an expansion in multiple

283 spin-factorized determinants, e.g. by relying on the generalized Laplace expansion of determinants
284 to expand $\det[\mathbf{A}^{1l,p}(\mathbf{r})]$ according to the rows corresponding to spin-up electrons

$$\det[\mathbf{A}^{1l}(\mathbf{r})] = \sum_S \varepsilon_S \det[\mathbf{A}^{\uparrow,S}(\mathbf{r})] \det[\mathbf{A}^{\downarrow,\bar{S}}(\mathbf{r})]. \quad (23)$$

285 Here, S runs over all subsets of the orbitals that contain as many elements as the number of spin-up
286 electrons, \bar{S} stands for the complement subset of S , $\mathbf{A}^{\uparrow,S}(\mathbf{r})$ denotes the submatrix of $\mathbf{A}^{1l}(\mathbf{r})$ formed
287 from the orbitals in S and the spin-up electrons, while ε_S is the sign of the permutation defined by
288 the subset S . For the block diagonal matrices of (21), the determinants for all subsets of spin-up
289 orbitals containing off-diagonal elements evaluate to zero and the sum in (23) reduces to a single
290 product of a spin-up and spin-down determinant. Note that since the many-body orbitals defined in
291 (19) are not equivariant with respect to exchanges of electrons with opposite spins, the terms on the
292 right hand side of (23) with different S s will in general be unrelated. In practice, it is conceivable
293 that due to the concrete form of parametrization of the many-body orbitals, there remains some
294 structure in the set of factorized determinants, that makes the full determinant more effective than
295 using an equivalent number of spin-factorized determinants formed from independent orbitals.

296 D. Jastrow factor and cusp correction

297 The antisymmetry of the wave function is retained when multiplying it with a global correc-
298 tion term symmetric under the exchange of the same spin particles. This symmetric correction,
299 traditionally called a Jastrow factor, is well suited to introduce known asymptotics to the ansatz.
300 DEEPQMC implements a learnable Jastrow factor e^J , where J is computed from the permutation
301 invariant sum of many-body electron embeddings

$$J = \eta_{\theta} \left(\sum_i \mathbf{f}_i^{(n),L} \right), \quad (24)$$

302 with η_{θ} again being implemented by an MLP. Furthermore, DEEPQMC provides a fixed Jastrow
303 factor that implements the known asymptotic behavior³⁷ when two electrons approach

$$\gamma(\mathbf{r}) = \sum_{i<j} -\frac{\alpha c_{ij}}{1 + \alpha |\mathbf{r}_i - \mathbf{r}_j|}, \quad (25)$$

304 where c_{ij} is $\frac{1}{4}$ if the electrons i and j are of the same spin and $\frac{1}{2}$ if the electrons possess opposite
305 spins and the hyperparameter alpha scales the width of the correction term. If cusplless Gaussian

306 envelopes are used, a similar cusp correction can be employed for the nuclei

$$\gamma(\mathbf{r}, \mathbf{R}) = \sum_{i,I} \frac{\alpha Z_I}{1 + \alpha |\mathbf{r}_i - \mathbf{R}_I|}, \quad (26)$$

307 serving as a simple replacement for the more involved technique utilized by Hermann et al.⁷.

308 E. Log-representation of the wave function

309 The output of the (unnormalized) neural network wave functions typically spans many orders
310 of magnitude, potentially resulting in instabilities due to finite floating-point precision. In order to
311 improve numerical stability, DEEPQMC represents wave functions in the log-domain

$$\psi = (\text{sign}(\psi), \log(|\psi|)). \quad (27)$$

312 We mitigate over- and underflow problems during the computation of the determinant by perform-
313 ing it directly in the log-domain using the appropriate `slogdet` function provided by JAX. In order
314 to perform the summation over multiple determinants ϕ^p we apply the log-sum-exp trick

$$\log\left(\left|\sum_p \phi^p\right|\right) = \max\{\log(|\phi^p|)\} + \log\left(\left|\sum_p \text{sign}(\phi^p) \exp(\log(|\phi^p|) - \max\{\log(|\phi^p|)\})\right|\right). \quad (28)$$

315 Note that for the variational principle to remain valid, it is sufficient to ensure the antisymmetry of
316 the trial wave function, and its explicit sign is not needed for the evaluation of any of the quantities
317 involved in the optimization (30).

318 IV. TRAINING

319 In this section, some technical aspects of the variational optimization of deep-learning trial
320 wave functions are discussed. While these ansatzes are trained within the standard VMC frame-
321 work, the characteristics of their optimization differ markedly from other VMC ansatzes, mainly
322 due to the greatly increased parameter count introduced by neural networks. On the other hand,
323 it is also distinct from most other deep-learning settings owing to the unusual complexity of the
324 loss function and the self-supervised setting, where training data is generated in parallel to the
325 optimization.

326 A. Loss function and gradient trick

327 As discussed in Section II B, VMC relies on the variational principle by optimizing the wave
328 function ansatz to minimize the expectation value of the local energies. From a machine learning
329 perspective, this translates to considering the loss function

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{r} \sim |\psi_{\boldsymbol{\theta}}|^2} [E_{\text{loc}}[\psi_{\boldsymbol{\theta}}](\mathbf{r})]. \quad (29)$$

330 Naively computing the gradient of this loss would involve taking derivatives of the local energies
331 $E_{\text{loc}}[\psi_{\boldsymbol{\theta}}]$ with respect to the ansatz parameters $\boldsymbol{\theta}$. However, evaluating the local energy already
332 involves second derivatives of the wave function with respect to the electron coordinates due to the
333 Laplacian in (1). Consequently, this naive gradient computation would necessitate taking mixed
334 third derivatives of the ansatz.

335 To reduce the computational costs and numerical instabilities associated with higher order
336 derivatives, a different unbiased estimator of the loss gradient is utilized,

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = 2\mathbb{E}_{\mathbf{r} \sim |\psi_{\boldsymbol{\theta}}|^2} [(E_{\text{loc}}[\psi_{\boldsymbol{\theta}}](\mathbf{r}) - \mathcal{L}(\boldsymbol{\theta})) \nabla_{\boldsymbol{\theta}} \log |\psi_{\boldsymbol{\theta}}|]. \quad (30)$$

337 The derivation of this estimator exploits the hermiticity of the Hamiltonian and is given in full
338 detail by Inui et al.³⁸. It replaces the derivatives of the local energy with a simple gradient of the
339 wave function, therefore it is expected to be more efficient and numerically stable to evaluate than
340 the direct gradient computation.

341 B. Local energy evaluation

342 The evaluation of the local energies of the wave function ansatz is by far the most computation-
343 ally demanding part of the training (and evaluation)

$$\begin{aligned} E_{\text{loc}}[\psi_{\boldsymbol{\theta}}](\mathbf{r}) &= -\frac{1}{2} \sum_i \left(\frac{\Delta_{\mathbf{r}_i} \psi_{\boldsymbol{\theta}}(\mathbf{r})}{\psi_{\boldsymbol{\theta}}(\mathbf{r})} \right) + V(\mathbf{r}) \\ &= -\frac{1}{2} \sum_i \left(\Delta_{\mathbf{r}_i} \log |\psi_{\boldsymbol{\theta}}(\mathbf{r})| + (\nabla_{\mathbf{r}_i} \log |\psi_{\boldsymbol{\theta}}(\mathbf{r})|)^2 \right) + V(\mathbf{r}). \end{aligned} \quad (31)$$

344 While the potential energy term is very cheap to evaluate, the cost of the Laplacian within the
345 kinetic energy term scales steeply with the number of electrons. In this step, one obtains sec-
346 ond derivatives of the wave function with respect to the electron coordinates. We obtain these
347 derivatives of the wave function by applying automatic differentiation in backward-forward mode,

348 which we confirmed to be a good choice in the setting of molecular wave functions. Further dis-
349 cussions regarding the memory bottleneck associated with the Laplacian and details regarding the
350 implementation choices are presented in the Appendix A.

351 C. Pretraining

352 Choosing initial values for ansatz parameters is a non-trivial question common to many compu-
353 tational chemistry methods. One need only think of the sensitivity of the self-consistent iterations
354 to the initial guess in Hartree–Fock (HF) and related methods^{39–41}. The case of deep-learning
355 VMC ansatzes is no different – a random initialization of the neural network parameters according
356 to some of the widely adopted schemes of the machine learning community can lead to the opti-
357 mization diverging or converging to local minima. This problem becomes increasingly severe with
358 growing system size, presumably due to the higher-dimensional, more complex wave functions of
359 larger molecules and their intricate nodal structure.

360 A practical solution to this issue is the initialization of the wave function based on a cheap
361 reference solution. To that end DEEPQMC interfaces with PYSCF³⁶, enabling the initialization
362 of wave functions from the coefficients of a preceding HF or multi-configurational self consistent
363 field (MCSCF) calculation. While this allows the direct initialization of the neural network wave
364 function ansatz as introduced by Hermann et al.⁷, subsequent work suggested that explicitly incor-
365 porating an approximate reference wave function in the model can deteriorate performance¹¹. In-
366 stead a short, supervised pretraining with respect to a reference solution before the self-supervised
367 variational optimization is recommended. In this step, the many-body orbitals of the ansatz are
368 trained to match the reference by minimizing the pretraining loss

$$\mathcal{L}_p(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{r} \sim |\psi_{\theta}|^2} \left[\sum_{ki} (\varphi_k^{\text{ref}}(\mathbf{r}_i) - \phi_k(\mathbf{r}_i, \{\mathbf{r}^{\uparrow}\}, \{\mathbf{r}^{\downarrow}\})) \times \varphi_k(\mathbf{r}_i) \right]^2, \quad (32)$$

369 where φ_k^{ref} are the occupied orbitals of the HF/MCSCF wave function. Unlike the variational loss
370 of (29), computing \mathcal{L}_p does not involve evaluating the Laplacian of the ansatz, which means that
371 pretraining requires significantly less computational resources than variational training. Initializa-
372 tion with pretrained orbitals, as introduced by Pfau et al.⁸, improves the convergence properties
373 of the variational training and, if well balanced with the subsequent variational optimization, can
374 even slightly boost the final accuracy, as Gerard et al. recently demonstrated¹¹.

375 D. Gradient clipping

376 Despite the utilization of sophisticated gradient estimators and pretraining, the convergence
377 of the variational optimization is still often hindered by outliers in the training batches of local
378 energies. The existence of these outliers is not surprising, considering that the electrostatic energy
379 is singular when two particles coincide, while the kinetic energy is singular at the nodes of the wave
380 function – energy contributions that the shape of the wave function precisely levels out in later
381 stages of the training. While the outliers are valid contributions to the energy expectation value,
382 their presence can inject a lot of noise into the gradient estimates. To reduce their contribution to
383 the parameter update, the loss and its gradient (given in (29) and (30)) are evaluated using clipped
384 local energies $\hat{E}_{\text{loc}}^{\mu,\sigma}$, where μ is the center and σ is the half-width of the clipping window.

385 Regarding concrete choices for μ and σ , some empirical findings have been reported in the
386 related literature. Investigating transition metal atoms using pseudopotentials, Li and coworkers
387 report¹⁶, that choosing $\sigma = 10 \times \text{std}(E_{\text{loc}})$ significantly outperforms all other options they've con-
388 sidered. More recently, von Glehn et al.¹² found that centering the clipping window at the median
389 of local energies, and using the mean absolute deviation from the median to determine σ , im-
390 proves the training of multiple deep-learning ansatzes. Considering the practical importance of
391 the clipping mechanism, DEEPQMC implements the algorithm of von Glehn et al.¹², along with
392 an analogous logarithmically scaling “soft” clipping scheme introduced by Hermann et al.⁷, but
393 also offers full flexibility to the user in specifying custom clipping functions.

394 Finally, it should be highlighted that the local energies are only to be clipped for computing the
395 gradient of the loss during optimization. Since clipping can introduce a bias to the estimate of the
396 energy expectation value, variational energy estimates can only be obtained from unclipped local
397 energies.

398 E. Optimizer

399 Utilizing natural gradient descent optimization⁴² or Kronecker-factored approximations thereof⁴³
400 has proven to be a crucial ingredient to the success of variational optimization of deep-learning
401 wave functions on molecular systems^{8,11,12,44,45}. Consequently, DEEPQMC makes use of the
402 Kronecker-Factored Approximate Curvature (KFAC) optimizer implementation of A. Botev and
403 J. Martens⁴⁶. To showcase the importance of the choice of the optimizer, the performance of

404 KFAC is compared to the commonly employed first-order optimizer AdamW⁴⁷, on variational
 405 trainings on six test systems. The obtained training energy curves are plotted in Figure 2. To ac-
 406 count for the 10–25% longer per iteration run time of KFAC compared to AdamW, the wall clock
 407 time of the training (instead of the usual iteration count) is shown on the horizontal axes. The
 408 results show that the slightly increased per-iteration cost is offset by the significantly improved
 409 per-iteration convergence speed of the KFAC optimizer. Furthermore, it is found that the increase
 410 in the relative cost of KFAC over AdamW optimization is smaller for systems with larger numbers
 411 of electrons. In practice, this means that the last percents of correlation energy can be recovered
 412 much more efficiently with KFAC, resulting in improved final energies for a given computational
 413 budget.

414 The effectiveness of natural gradient descent in this setting can be rationalized through its con-
 415 nection to the stochastic reconfiguration method^{8,48}, known from traditional variational quantum
 416 Monte Carlo optimization^{49,50}.

417 These higher order methods utilize the Fisher information of the unnormalized density associ-
 418 ated with the wave function as a preconditioner to the gradients. KFAC extends the application
 419 range of natural gradient descent by low-rank approximating the Fisher information, facilitating
 420 its computation for neural network wave functions with large numbers of model parameters.

421 Instead of following the steepest descent in parameter space, an optimization step with the
 422 preconditioned gradient is in the direction of steepest descent in distribution space, with distance
 423 defined by the Kullback–Leibler (KL) divergence⁵¹. Considering that in VMC the predicted quan-
 424 tity $\psi_{\theta}(\mathbf{r})$ directly defines the distribution $p(\mathbf{r}|\theta) \propto |\psi_{\theta}(\mathbf{r})|^2$, one concludes that a natural gradient
 425 step is in the direction of maximal loss decrease for a given KL divergence between ψ_{θ} and $\psi_{\theta+d\theta}$.
 426 This is an advantageous property, as relying on the KL divergence results in updates that are inde-
 427 pendent of the way ψ_{θ} is parameterized, as opposed to the steepest descent where the Euclidean
 428 metric introduces strong dependence.

430 V. SAMPLING

431 An important characteristic of VMC is that the data (electron positions) used to fit the model
 432 is generated in tandem with the optimization, by sampling the probability distribution of the elec-
 433 tronic degrees of freedom defined by the square of the wave function. This sampling task comes
 434 with its own challenges, due to its tight coupling with the training. For the variational principle to

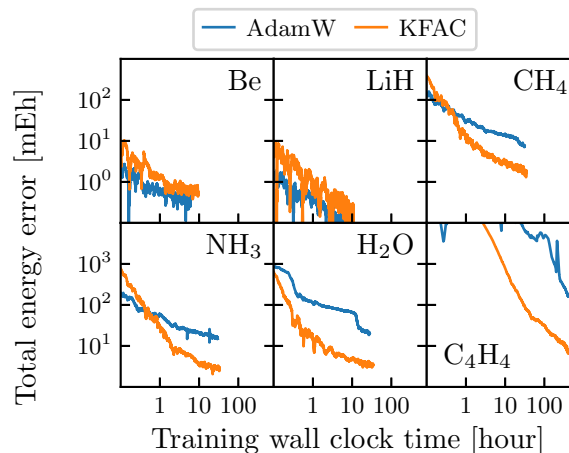


FIG. 2. **Comparing the performance of the AdamW and KFAC optimizers.** Total energy errors during the training process are shown for beryllium, lithium hydride, methane, ammonia, water, and cyclobutadiene. The horizontal axes show the wall clock time of the training, measured on a single Nvidia GTX 1080 Ti GPU. To obtain smooth training curves, the exponential moving average of the training energy is plotted. While on smaller systems (Be and LiH) AdamW converges slightly faster, due to its lower per-iteration cost, on the larger systems the benefit of using KFAC is clear.

435 remain valid, the samples used to evaluate (6) must be equilibrated according to the distribution
 436 $\mathbf{r} \sim |\psi_{\theta}(\mathbf{r})|^2$. Furthermore, since ψ_{θ} is updated in every training iteration, the sampling must ac-
 437 count for the corresponding changes in the distribution of the electron positions. To carry out this
 438 demanding sampling task in a computationally efficient manner, the DEEPQMC program package
 439 offers two optimized Markov chain Monte Carlo (MCMC) algorithms. Along with the Random
 440 Walk MCMC algorithm^{52,53}, referred to as Metropolis sampler, the Metropolis-Adjusted Langevin
 441 Algorithm⁵⁴ (MALA), referred to as Langevin sampler, is also implemented, that proposes walker
 442 updates using overdamped Langevin dynamics. The implemented MALA includes the correction
 443 proposed by Hermann et al.⁷, which scales the electron step sizes around the nuclei to avoid "over-
 444 shooting" the latter. Additionally, changes of the wave function during training can be accounted
 445 for by re-equilibration after each gradient step or using a batch reweighting scheme. In the fol-
 446 lowing sections, these MCMC samplers along with the above described sampling difficulties are
 447 characterized in more detail.

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0157512

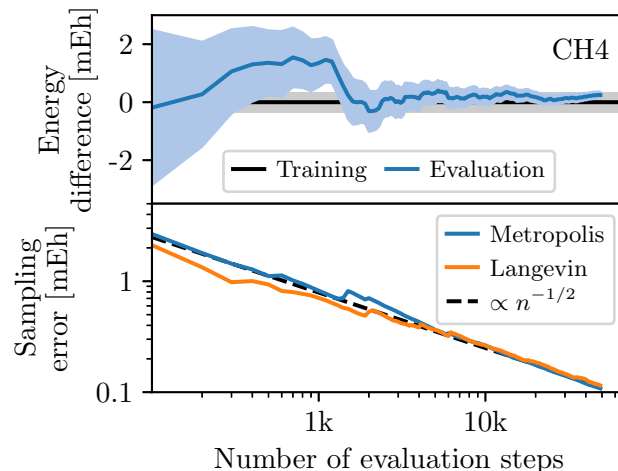


FIG. 3. **Top: Evaluation of a trained wave function ansatz, bottom: sampling error for the evaluation of a trained ansatz.** The convergence of the energy expectation value is depicted during the evaluation of an optimized CH4 ansatz using two thousand MCMC walkers. In the top pane, sampling is performed with the Random Walk MCMC algorithm. The evaluated energy is compared to the final training energy, with shaded areas showing error estimates. In the bottom pane, the convergence of the sampling errors of evaluations with the Metropolis sampler and the Langevin sampler are compared. The sampling errors converge as $n^{-1/2}$ with the number of samples n , as expected from the central limit theorem.

448 A. Energy convergence

449 First, the convergence of the energy expectation value estimate is investigated, when sampling
 450 an unchanging, previously trained ansatz. In order to draw $n = n_b \times n_s$ electron samples $\{\mathbf{r}\}_{ij}$,
 451 distributed according to $|\psi_\theta(\mathbf{r})|^2$, a batch of n_b many walkers is propagated for n_s MCMC steps.
 452 Based on the electron samples the energy expectation value is estimated as

$$\langle E \rangle = \frac{1}{n} \sum_{i=1}^{n_b} \sum_{j=1}^{n_s} E_{\text{loc}}[\psi_\theta](\{\mathbf{r}\}_{ij}). \quad (33)$$

453 Following the central limit theorem^{32,55}, the sampling error of such estimates decays proportional
 454 to $n^{-1/2}$. To approximate the sampling error, we utilize the nonoverlapping batch means estimator,
 455 as reviewed by Flegal et al.⁵⁵. We first obtain independent estimates of the energy by averaging
 456 the local energies over the walker trajectories (batches)

$$\langle E \rangle^i = \frac{1}{n_s} \sum_{j=1}^{n_s} E_{\text{loc}}[\psi_\theta](\{\mathbf{r}\}_{ij}). \quad (34)$$

457 Considering these batches, the sampling error is then estimated as

$$\langle \sigma_E \rangle = \sqrt{\frac{\sum_{i=1}^{n_b} (\langle E \rangle^i - \langle E \rangle)^2}{n_b(n_b - 1)}}. \quad (35)$$

458 The convergence of the energy estimate and its error bar throughout the evaluation of an ansatz
459 trained on the CH₄ molecule is plotted in Figure 3. In the top pane, the final value of the exponen-
460 tial walking mean of the training energies, and its estimated error are also shown with a horizontal
461 line and shaded area. It can be seen from this plot that the energy estimate of the evaluation con-
462 verges gradually towards the final training energy as expected, while its sampling error converges
463 towards zero. Note that due to the parameter updates during the optimization, the energy estimate
464 from the training is an unreliable estimate and a thorough evaluation of the energy expectation
465 value requires sampling the ansatz with fixed parameters.

466 On the bottom pane of Figure 3, the convergence of the estimated sampling error is compared
467 between the Metropolis sampler and the Langevin sampler. Importantly, the expected $n^{-1/2}$ con-
468 vergence behavior is observed for both methods. Comparing the two algorithms, it can be seen
469 that the error of MALA converges slightly faster than that of Random Walk MCMC, indicating
470 a lower degree of correlation between the subsequent positions of the walkers of the Langevin
471 sampler.

472 B. Decorrelated sampling

473 To characterize the phenomenon of correlated samples hinted at in Section V A, autocorrelation
474 functions of the local energy samples are investigated. The autocorrelation function of the local
475 energies sampled by a single MCMC chain is defined as:

$$\rho_{E_{\text{loc}}}(t) = \int_{-\infty}^{\infty} \left(E_{\text{loc}}^{t'+t} - \mu_{E_{\text{loc}}} \right) \left(E_{\text{loc}}^{t'} - \mu_{E_{\text{loc}}} \right) dt', \quad (36)$$

476 where E_{loc}^t denotes the local energy sampled at step t , and $\mu_{E_{\text{loc}}}$ is the mean of the local ener-
477 gies over the entire trajectory. The autocorrelation time of the local energy is then computed as
478 $\tau' = 2 \int_0^{\infty} \rho_{E_{\text{loc}}}(t) dt$. Finally, τ is obtained by taking the mean of τ' s over all propagated MCMC
479 chains, providing a simple measure of local energy autocorrelation.

480 The autocorrelation times for five atoms of increasing size and the cyclobutadiene molecule
481 are plotted on the upper pane of Figure 4, for both the Metropolis and the Langevin sampler. The

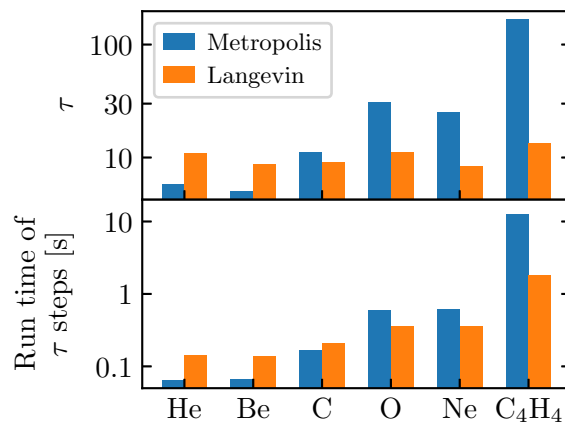


FIG. 4. **Autocorrelation times of the local energy samples.** The top pane shows the MCMC sampling autocorrelation time τ , as defined in Section V B, for a sequence of atoms and the ground state of cyclobutadiene. The bottom pane shows the run time of performing τ sampling steps for the same systems. The electrons are sampled using either the Metropolis sampler or the Langevin sampler. Following the suggestion of Sokal⁵⁶, the autocorrelation times are estimated using MCMC chains of length $\approx 5\tau$.

482 general trend of longer autocorrelation times for larger systems can be observed for the Metropo-
 483 lis sampler. One of the main causes of this trend is the increasing nuclear charge, which induces
 484 higher and higher peaks in the distribution of the electrons near the nuclei. These pronounced
 485 peaks necessitate shorter update proposal radii, ultimately resulting in a higher correlation between
 486 subsequent samples. Furthermore, the autocorrelation time is expected to grow with the increas-
 487 ing complexity of the wave functions and their nodal surfaces. On the other hand, the Langevin
 488 sampler seems less affected by this trend, delivering largely constant autocorrelation times for all
 489 systems. It is reasonable to assume that by explicitly making use of information about the gradient
 490 of the wave function, the MALA update proposal retains better decorrelation efficiency than Ran-
 491 dom Walk MCMC, when considering more and more complicated wave functions. Finally, the
 492 showcased autocorrelation times are in reasonably good agreement with the fact that the default
 493 number of decorrelating steps performed between parameter updates is chosen between 10–30 in
 494 the currently used neural wave function program packages.

495 The experiments depicted in Figure 4 also demonstrate a slightly smaller correlation between
 496 subsequent samples of the Langevin sampler in comparison to those of the Metropolis sampler,
 497 for all but the smallest of systems. On the bottom pane of Figure 4, the wall clock run time of
 498 performing τ sampling steps are shown for each system, to account for the slightly increased com-

499 putational cost of the MALA update proposal. Considering wall clock run times, the Metropolis
500 sampler is more efficient on atoms up to carbon, while the Langevin sampler performs slightly
501 to considerably better on the larger atoms and cyclobutadiene. While we find MALA to be more
502 efficient than Random Walk Metropolis, we observe that for larger systems with heavier nuclei
503 it could result in less stable optimization. To improve the black-boxed nature of the method, we
504 applied Random Walk MCMC in all subsequent experiments.

505 VI. SCALING

506 Understanding the scaling of a method's computational cost with the considered system size
507 is of utmost importance in the field of quantum chemistry, where a pervasive caveat of the most
508 accurate approaches is their unfavorable scaling behavior. Given its high accuracy, the asymptotic
509 scaling of VMC (typically listed with N^4)³² is considered favorable. Although this general scaling
510 is indeed much better than e.g. the N^7 scaling of the gold-standard CCSD(T) method, and on par
511 with the scaling of hybrid density functionals (such as DM21³), deep QMC calculations incur a
512 larger prefactor, resulting in much higher practical costs on systems of intermediate size. While
513 reducing this prefactor is an important long term goal of method developers in the field, investigat-
514 ing the method's scaling is also of interest, to estimate the prospect of system sizes feasible with
515 further improvement and serve as a baseline for future developments. In this section, the scaling
516 of the computational cost of the variational training of deep-learning ansatzes is investigated using
517 the DEEPQMC program package. Further scaling aspects of the pseudopotential implementation,
518 and design choices regarding the major computational bottlenecks of the algorithm are discussed
519 in Appendix A.

520 The theoretical scaling of VMC (N^4) is obtained when combining the N^3 cost of the determi-
521 nant evaluation with an additional factor of N from the Laplacian required in the computation of
522 the kinetic energy. In practice, however, for simulations with the currently feasible system sizes,
523 the determinant evaluation makes up only a fraction of the computational cost, which is instead
524 dominated by the evaluation of the neural networks of the ansatz. To investigate the practical
525 scaling of a variational training step in DEEPQMC, single iteration run times are compared across
526 atoms with increasing atomic numbers, as well as across chains containing an increasing number
527 of hydrogen atoms (Figure 5). Although the systems contain different numbers of particles, due to
528 the parametrization with GNNs the total parameter count of the wave function ansatz changes only

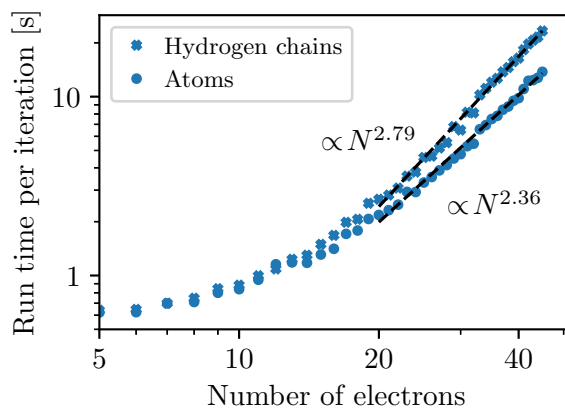


FIG. 5. **Scaling of the computational cost with system size.** The figure depicts the time in seconds per variational optimization step for systems with up to 45 electrons. The timings were obtained for training steps with a batch size of 2000 run, on a single A100 GPU. An exponential fit gives the scaling exponent of 2.79 for the hydrogen chains and 2.37 for the atoms respectively.

marginally between systems. On the other hand, owing to the varying numbers of nuclei, isoelectronic species can have slightly different computational requirements. The system classes of atoms and hydrogen chains were chosen, as they represent the lower and upper bounds respectively, on the number of nuclei a neutral system with a fixed number of electrons can contain. Consequently, the scaling of the run time with the number of electrons is also expected to be bounded by these system classes. With the tight empirical bounds of $N^{2.36-2.79}$ depicted in Figure 5, the observed scaling of DEEPQMC is still far below the theoretical estimate of N^4 , highlighting the potential for extension to larger systems.

VII. ANSATZ VALIDATION

Relying on the general framework introduced above, the DEEPQMC software suite enables the use of many of the previously published deep-learning ansatzes by providing configuration files to reproduce PauliNet⁷, FermiNet⁸ and DeepErwin¹¹. To validate our implementation of these ansatzes, the hyperparameters of sampling, optimization, and GNN architecture are compared in depth to those of the respective reference implementations. Additionally, it is verified that when using the same parameters, the DEEPQMC implementations predict the same wave function value and local energy as their reference counterparts for a given configuration of electrons and nuclei.

545 Note that we have refrained from exactly matching the cusp-corrected GTOs of PauliNet, because
546 subsequent work has demonstrated that explicitly including a reference solution is limiting the ac-
547 curacy of the ansatz. However, by combining Gaussian envelopes initialized from the coefficients
548 of a reference calculation with a nuclear cusp correction in the Jastrow factor (26) it is possible
549 to obtain a variant of PauliNet within DEEPQMC, that matches the characteristics of the original
550 ansatz.

551 In Figure 6 the empirical performance of the various ansatzes is checked against results pub-
552 lished in the literature for a small set of molecules. It can be seen that our DEEPQMC implemen-
553 tation of PauliNet, FermiNet, and DeepErwin matches the reference energies well. The remaining
554 discrepancies of FermiNet result from slightly different experimental setups, such as an increased
555 number of reference optimization steps (200 000 compared to 50 000 used here) and batch size
556 (4096 compared to 2048 used here), or an older TensorFlow-based implementation being used
557 in case of N_2 . The impact of these changes on the deviations of the model accuracy highlights
558 the importance and difficulty of comparing ansatzes implemented in different codebases under the
559 same experimental conditions.

560 As a further contribution, we introduce and analyze the performance of a new default ansatz
561 for the DEEPQMC program package, which we refer to as PauliNet2. This exemplary ansatz was
562 optimized to have a good trade-off between accuracy and trainable model parameters. Despite
563 achieving a similar accuracy as FermiNet and DeepErwin for the small systems under investiga-
564 tion (see Figures 6 and 7), the PauliNet2 ansatz has only about a third of the model parameters of
565 FermiNet and a quarter of DeepErwin (i.e. for the CO molecule 239 829, 766 944, and 998 816 pa-
566 rameters respectively). The ansatz combines the SchNet-like graph convolutions of PauliNet (17)
567 with the iterative update of the edge embeddings of FermiNet (18). Edge features are constructed
568 from difference vectors between the electrons and isotropic exponentials are used as envelopes.
569 Furthermore, the ansatz comprises a trainable Jastrow factor (24) as well as the fixed electronic
570 cusp correction (25). While these hyperparameters are found to be suitable for the presented ex-
571 periments, it is conceivable that an extended hyperparameter search targeting specific applications
572 could further improve its performance. The detailed settings of the discussed ansatzes can be
573 found in the respective configuration files shipped with the DEEPQMC package.

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0157512

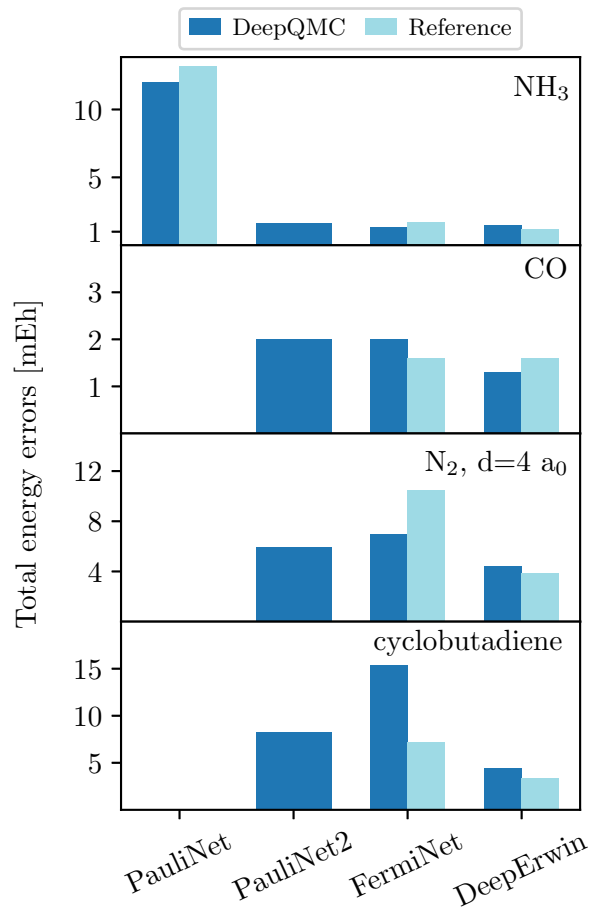


FIG. 6. Validating the DEEPQMC implementation of various ansatzes by comparing their accuracy to published results obtained with their respective reference implementations. Note that results obtained with the DEEPQMC or DeepErwin codebases were computed using 50 000 variational optimization steps and a batch size of 2048, while the FermiNet reference results used 200 000 training steps and 4096 samples in a batch. Results computed with the reference implementations are taken from the works of Pfau⁸, Spencer⁴⁴, and Gerard¹¹.

575 VIII. APPLICATION EXAMPLES

576 In this section, the ease of applying the DEEPQMC program package as a black-box method
 577 to obtain electronic energies is demonstrated on benchmark datasets. Two widely different exam-
 578 ple problems are chosen in order to showcase the general applicability of the presented method.
 579 The experiments are performed using DEEPQMC command line interface, which exposes all con-
 580 figuration options of the software suite while also allowing for effortless submission of simple
 581 calculations. Short usage examples of the DEEPQMC command line interface are provided in

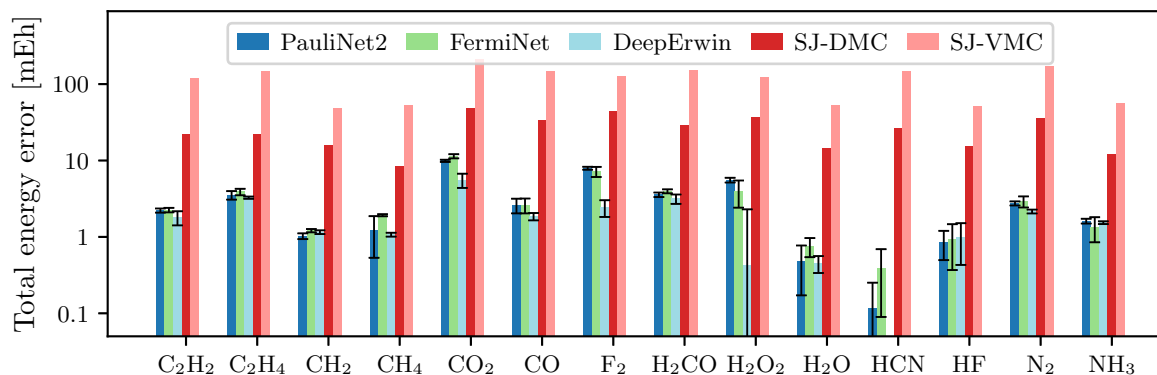


FIG. 7. **Total energy deviations for small molecules of H, C, O, N, and F atoms.** The presented molecules participate in the reactions investigated by Nemeč⁵⁷. Reference total and HF energies are taken from the review of O’Neill⁶⁰, while SJ VMC and SJ DMC results are taken from the work of Nemeč⁵⁷. Results for the hydrogen molecule are omitted, as it is described nearly exactly by all depicted methods. Error bars denote the sampling error as estimated according to (35).

582 Appendix B.

583 A. Small molecule reactions

584 The electronic contributions to the reaction energies of 12 reactions involving small inorganic
 585 molecules and hydrocarbons are investigated. These reactions were used by Nemeč⁵⁷ to bench-
 586 mark the accuracy of Slater–Jastrow (SJ) type trial wave functions, constructed following Drum-
 587 mond et al.⁵⁸ using electron–electron, electron–nucleus, and electron–electron–nucleus terms in the
 588 Jastrow factor. The 14 participating molecules are built from H, C, O, N, and F atoms, containing
 589 at least 2 and at most 22 electrons. To facilitate the comparison with the DMC results of Nemeč⁵⁷,
 590 the same molecular geometries are considered, obtained from the work of Feller⁵⁹. Reference
 591 energies are taken from the review of O’Neill⁶⁰. All electron, complete basis set extrapolated
 592 CCSD(T) energies are computed in house, using the PySCF program package³⁶.

593 First, single-point electronic energies obtained for the participating molecules are compared
 594 in Figure 7. On the vertical axis, the error of the recovered total energy is plotted, for VMC
 595 and DMC calculations utilizing SJ type trial wave functions, and for VMC with deep-learning
 596 ansatzes. Looking at Figure 7, one can observe that the total energy errors of SJ-VMC ansatzes
 597 are consistently above 47 mEh (with a mean of 114 mEh), while the associated DMC errors are

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0157512

598 in the range of 8 - 50 mEh (26 mEh on average). In comparison, deep-learning ansatzes exhibit
599 at maximum only 11 mEh total energy error, with a mean deviation of 2.6 mEh. While the main
600 goal of quantum chemistry methods is to accurately model energy differences, rather than recover
601 exact total energies, it is encouraging to see that DEEPQMC and deep-learning ansatzes in general
602 are very competitive in this area.

603 The accuracy of the energy differences obtained with SJ-DMC, CCSD(T), and deep-learning
604 QMC methods are compared in Figure 8. Note that energy differences obtained with SJ-VMC are
605 not shown in this figure, as they are an order of magnitude less accurate than the depicted ap-
606 proaches. Comparing SJ-DMC results with those obtained from DEEPQMC, one concludes that
607 combining the VMC method with expressive deep-learning ansatzes greatly increases its accuracy,
608 surpassing SJ-DMC on eleven out of twelve reactions. The accuracy advantage of DEEPQMC's
609 PauliNet2, FermiNet, and DeepErwin ansatzes is similarly clear when comparing their respective
610 reaction energy mean absolute deviations (MADs) of 2.4 mEh and 2.3 mEh, and 1.5 mEh to the 7.6
611 mEh of SJ-DMC. As a final comparison, Figure 8 also shows the reaction energy differences ob-
612 tained from a complete basis set extrapolated, all-electron CCSD(T). Not surprisingly, CCSD(T)
613 performs outstandingly on these small, single reference systems in equilibrium geometry, achiev-
614 ing a MAD of 3.4 mEh, and chemical accuracy (less than 1 kcal/mol or 1.6 mEh error) on four
615 reactions. In comparison, the MAD value of PauliNet2 for this exemplary study with DEEPQMC
616 is found to be below that of CCSD(T), and chemical accuracy is achieved on seven out of twelve
618 reactions.

619 **B. Transition metal oxides**

620 The effects of utilizing pseudopotentials in variational optimization of deep-learning molec-
621 ular wave function are evaluated on a series of four first-row transition metal oxides. The bond
622 lengths of the ScO, TiO, VO, and CrO molecules are taken from the experimental results of Annab-
623 erdieyev et al.⁶¹. The latest ccECP pseudopotentials⁶¹ are applied to the transition metal atoms
624 only, replacing neon-like cores of 10 electrons. Although replacing argon cores (18 electrons)
625 with pseudopotentials would result in even larger computational savings, this is avoided as the
626 third shell electrons are known to play a non-negligible role in the bond formation of transition
627 metal atoms⁶². Apart from the introduction of pseudopotentials, the ansatz employed on small
628 molecule reactions (Section VIII A) is utilized here without further modifications.

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0157512

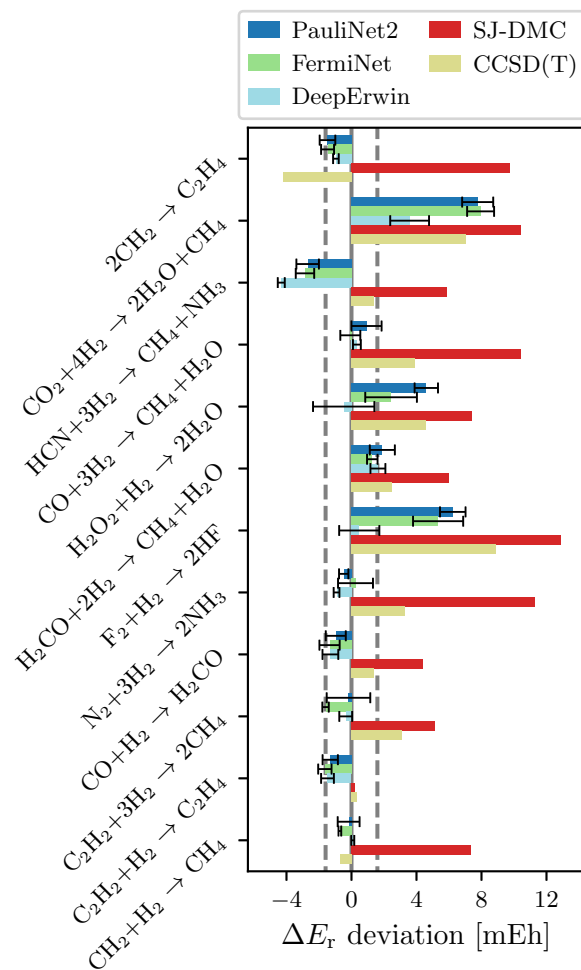


FIG. 8. **Reaction energy deviations for the reactions involving small molecules of H, C, O, N, and F atoms.** Reference reaction energies are computed from the electronic energies reviewed by O'Neill⁶⁰, SJ-DMC results are obtained from Nemeč⁵⁷, while complete basis set extrapolated CCSD(T) values are calculated in house using PySCF³⁶. The energies for PauliNet2, FermiNet and DeepErwin are obtained with the DEEPQMC program package. The range of ± 1 kcal/mol deviation (often referred to as chemical accuracy) is highlighted with dashed lines.

629 Comparing the technical details of pseudopotential calculations to all-electron ones, the advan-
 630 tage of the former is clear. Due to the exclusion of the fastest moving core electrons, the length
 631 of the electron position updates is sixfold increased, and higher accuracy is achieved in a given
 632 number of training steps, at about half of the computational cost. Next, the dissociation energies
 633 of the four transition metal oxides are estimated. The dissociation energy for transition metal X is

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0157512

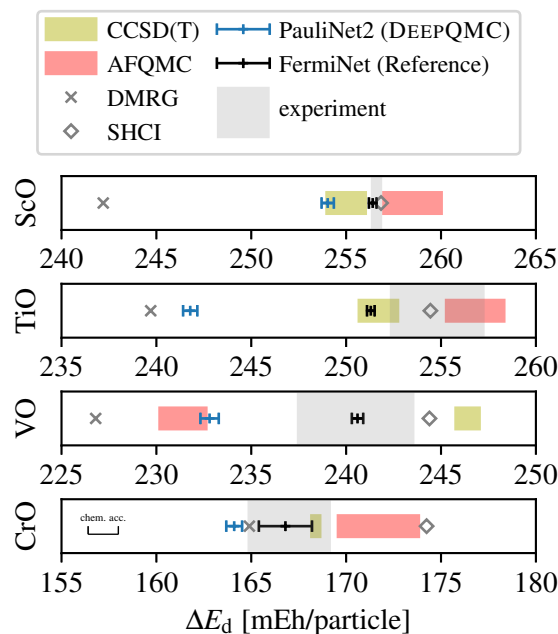


FIG. 9. **Dissociation energy of transition metal oxides calculated using different methods.** DEEP-QMC+ECP result was obtained using 55 000 training steps and 5 000 evaluation steps. The results for FermiNet+ECP are taken from Li et al.¹⁶, where they used 10 times more steps and consequently achieve higher accuracy. Other results are from the references^{62,64}.

634 defined as

$$\Delta E_d^{XO} = E^X + E^O - E^{XO}, \quad (37)$$

635 where $E^O = -75.0631(1)$ Ha is the result of an all-electron calculation with the same hyper-
 636 parameters. Figure 9 compares the obtained dissociation energies to experimental values⁶³, and
 637 some other accurate computational methods like CCSD(T)⁶¹, FermiNet¹⁶, auxiliary field quantum
 638 Monte Carlo (AFQMC), semi-stochastic heat batch configuration interaction (SHCI), and density
 639 matrix renormalization group (DMRG)⁶⁴. Apart from the TiO case, the accuracy of DEEPQMC
 640 with pseudopotentials is comparable to other theoretical methods, such as CCSD(T) or AFQMC.
 641 The fact that the dissociation energy estimates with DEEPQMC are systematically lower than the
 642 experimental results, indicates that the single atoms are described more accurately than the oxide
 643 molecules. This can be counteracted by increasing the expressiveness of the ansatz and investing
 644 more compute. Note that results obtained with FermiNet¹⁶ utilized a larger ansatz and trained for
 645 about ten times more training iterations than done in this study.

646 IX. SUMMARY AND CONCLUSIONS

647 We have presented the DEEPQMC program package – a general variational quantum Monte
648 Carlo framework for optimizing deep-learning molecular wave functions. The implementation fo-
649 cuses on modularity, facilitating rapid development of new ansatzes, and provides maximal free-
650 dom in choosing the specifics of the variational training setup. The ansatz shipped with DEEP-
651 QMC attempts to unify most of the currently existing deep-learning molecular wave functions,
652 while remaining easy to extend as new models emerge. To reduce the computational complexity
653 associated with heavy nuclei, some popular precomputed pseudopotentials are also implemented.

654 Using the framework provided by DEEPQMC, the most important practical aspects of varia-
655 tional optimization of deep-learning molecular wave functions are discussed. The importance of a
656 proper gradient estimator along with robust gradient clipping is highlighted. For consistent ansatz
657 initialization, supervised pretraining to HF wave functions is suggested. The advantages of using
658 the second-order KFAC optimizer are demonstrated, along with a rationalization of its effective-
659 ness. The theoretical convergence of the Markov Chain Monte Carlo sampling error is verified,
660 and MALA is shown to be more effective in obtaining decorrelated samples than the widely uti-
661 lized Random Walk MCMC algorithm. The empirical scaling of the method's computational cost
662 is found to be more favorable than the most popular post-HF approaches, while its large prefactor
663 is identified as an obstacle on the path to wider adoption.

664 The black-box application of the program package is demonstrated in two significantly different
665 settings. The electronic reaction energies of 12 small molecule reactions are computed with a
666 mean absolute deviation of 1.5 mEh, and compared to the 3.4 mEh achieved by CCSD(T) and 7.6
667 mEh achieved by DMC with SJ ansatzes. Using the same ansatz hyperparameters, dissociation
668 energies are computed for a series of transition metal monoxides, utilizing the latest ccECP⁶¹
669 pseudopotential. Improved training characteristics compared to all-electron calculations highlight
670 the benefit of employing pseudopotentials. The accuracies of the predicted dissociation energies
671 are on par with or exceed those of some other recently popularized methods such as auxiliary field
672 quantum Monte Carlo, or density matrix renormalization group.

673 To conclude, the presented method shows great promise to become an easy-to-use, general,
674 black-box method accurately describing the molecular electronic structure. Especially encourag-
675 ing is the favorable scaling of computational requirements with increasing system size. It is easy
676 to envision that after further development reducing the large prefactor of the computational costs,

677 the DEEPQMC package will prove valuable to the wider quantum chemistry community.

678 DATA AVAILABILITY

679 The data that support the findings of this study will be made openly available upon publication
680 of the manuscript.

681 ACKNOWLEDGMENTS

682 We would like to thank Jannes Nys for helpful discussions regarding the nuclear cusp correc-
683 tion, and Leon Gerard for providing reference data for the ansatz comparison. The work of ZS
684 and PBS is funded and supported by the Berlin Mathematics Research Center MATH+ (Project
685 AA2-22). The work of MM has received funding from the European Research Council (ERC)
686 under the European Union's Horizon 2020 research and innovation programme (Grant agreement
687 No. 772230). Computing time granted by the HLRN and provided on the supercomputer Lise
688 compute at NHR@ZIB is gratefully acknowledged.

689 REFERENCES

- 690 ¹F. Noé, A. Tkatchenko, K.-R. Müller, and C. Clementi, Annual review of physical chemistry
691 **71**, 361 (2020).
- 692 ²H. Kulik, T. Hammerschmidt, J. Schmidt, S. Botti, M. Marques, M. Boley, M. Scheffler,
693 M. Todorović, P. Rinke, C. Oses, *et al.*, Electronic Structure **4**, 023004 (2022).
- 694 ³J. Kirkpatrick, B. McMorro, D. H. P. Turban, A. L. Gaunt, J. S. Spencer, A. G. D. G. Matthews,
695 A. Obika, L. Thiry, M. Fortunato, D. Pfau, L. R. Castellanos, S. Petersen, A. W. R. Nelson,
696 P. Kohli, P. Mori-Sánchez, D. Hassabis, and A. J. Cohen, Science **374**, 1385 (2021).
- 697 ⁴J. Hermann, J. Spencer, K. Choo, A. Mezzacapo, W. M. C. Foulkes, D. Pfau, G. Carleo,
698 and F. Noé, “Ab-initio quantum chemistry with neural-network wave functions,” (2022),
699 arXiv:2208.12590 [physics, stat].
- 700 ⁵G. Carleo and M. Troyer, Science **355**, 602 (2017).
- 701 ⁶J. Han, L. Zhang, and W. E, Journal of Computational Physics **399**, 108929 (2019).
- 702 ⁷J. Hermann, Z. Schätzle, and F. Noé, Nature Chemistry **12**, 891 (2020), arXiv:1909.08423
703 [physics, stat].

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0157512

- 704 ⁸D. Pfau, J. S. Spencer, A. G. D. G. Matthews, and W. M. C. Foulkes, *Physical Review Research*
705 **2**, 033429 (2020), publisher: American Physical Society.
- 706 ⁹Z. Schätzle, J. Hermann, and F. Noé, *The Journal of Chemical Physics* **154**, 124108 (2021).
- 707 ¹⁰J. Lin, G. Goldshlager, and L. Lin, *Journal of Computational Physics* **474**, 111765 (2023).
- 708 ¹¹L. Gerard, M. Scherbela, P. Marquetand, and P. Grohs, “Gold-standard solutions to the
709 Schrödinger equation using deep learning: How much physics do we need?” (2022),
710 arXiv:2205.09438 [physics].
- 711 ¹²I. von Glehn, J. S. Spencer, and D. Pfau, “A Self-Attention Ansatz for Ab-initio Quantum
712 Chemistry,” (2022), arXiv:2211.13672 [physics].
- 713 ¹³M. T. Entwistle, Z. Schätzle, P. A. Erdman, J. Hermann, and F. Noé, *Nature Communications*
714 **14**, 274 (2023), number: 1 Publisher: Nature Publishing Group.
- 715 ¹⁴G. Pescia, J. Han, A. Lovato, J. Lu, and G. Carleo, *Physical Review Research* **4**, 023138 (2022),
716 publisher: American Physical Society.
- 717 ¹⁵M. Wilson, S. Moroni, M. Holzmann, N. Gao, F. Wudarski, T. Vegge, and A. Bhowmik, *Physical*
718 *Review B* **107**, 235139 (2023).
- 719 ¹⁶X. Li, C. Fan, W. Ren, and J. Chen, *Physical Review Research* **4**, 013021 (2022),
720 arXiv:2108.11661 [physics.chem-ph].
- 721 ¹⁷Y. Qian, W. Fu, W. Ren, and J. Chen, *The Journal of Chemical Physics* **157**, 164104 (2022).
- 722 ¹⁸M. Wilson, N. Gao, F. Wudarski, E. Rieffel, and N. M. Tubman, arXiv preprint
723 arXiv:2103.12570 (2021).
- 724 ¹⁹W. Ren, W. Fu, X. Wu, and J. Chen, *Nature Communications* **14**, 1860 (2023).
- 725 ²⁰M. Scherbela, R. Reisenhofer, L. Gerard, P. Marquetand, and P. Grohs, “Solving the elec-
726 tronic Schrödinger equation for multiple nuclear geometries with weight-sharing deep neural
727 networks,” (2021), arXiv:2105.08351 [physics].
- 728 ²¹N. Gao and S. Günemann, “Ab-Initio Potential Energy Surfaces by Pairing GNNs with Neural
729 Wave Functions,” (2022), arXiv:2110.05064 [physics].
- 730 ²²N. Gao and S. Günemann, “Sampling-free Inference for Ab-Initio Potential Energy Surface
731 Networks,” (2023), arXiv:2205.14962 [physics].
- 732 ²³M. Scherbela, L. Gerard, and P. Grohs, “Towards a Foundation Model for Neural Network
733 Wavefunctions,” (2023), arXiv:2303.09949 [physics].
- 734 ²⁴N. Gao and S. Günemann, “Generalizing Neural Wave Functions,” (2023), arXiv:2302.04168
735 [physics, physics:quant-ph].

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0157512

- 736 ²⁵F. Vicentini, D. Hofmann, A. Szabó, D. Wu, C. Roth, C. Giuliani, G. Pescia, J. Nys, V. Vargas-
737 Calderón, N. Astrakhantsev, and G. Carleo, *SciPost Phys. Codebases*, 7 (2022).
- 738 ²⁶J. Hermann, Z. Schätzle, P. B. Szabó, M. Mezera, and DeepQMC Contributors, “DEEPQMC,”
739 (2023).
- 740 ²⁷J. S. Spencer, D. Pfau, and FermiNet Contributors, “FERMINET,” (2020).
- 741 ²⁸M. Scherbela, L. Gerard, and R. Reisenhofer, “DEEPERWIN,” (2020).
- 742 ²⁹J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula,
743 A. Paszke, J. VanderPlas, S. Wanderman-Milne, and Q. Zhang, “JAX: composable transfor-
744 mations of Python+NumPy programs,” (2018).
- 745 ³⁰G. Van Rossum and F. L. Drake, *Python 3 Reference Manual* (CreateSpace, Scotts Valley, CA,
746 2009).
- 747 ³¹T. Hennigan, T. Cai, T. Norman, and I. Babuschkin, “Haiku: Sonnet for JAX,” (2020).
- 748 ³²W. M. C. Foulkes, L. Mitas, R. J. Needs, and G. Rajagopal, *Reviews of Modern Physics* **73**, 33
749 (2001).
- 750 ³³R. P. Feynman and M. Cohen, *Physical Review* **102**, 1189 (1956), publisher: American Physical
751 Society.
- 752 ³⁴M. Burkatzki, C. Filippi, and M. Dolg, *Journal of Chemical Physics* **126**, 234105 (2007).
- 753 ³⁵M. C. Bennett, C. A. Melton, A. Annaberdiyev, G. Wang, L. Shulenburger, and L. Mitas,
754 *Journal of Chemical Physics* **147**, 224106 (2017).
- 755 ³⁶Q. Sun, X. Zhang, S. Banerjee, P. Bao, M. Barbry, N. S. Blunt, N. A. Bogdanov, G. H. Booth,
756 J. Chen, Z.-H. Cui, J. J. Eriksen, Y. Gao, S. Guo, J. Hermann, M. R. Hermes, K. Koh, P. Koval,
757 S. Lehtola, Z. Li, J. Liu, N. Mardirossian, J. D. McClain, M. Motta, B. Mussard, H. Q. Pham,
758 A. Pulkin, W. Purwanto, P. J. Robinson, E. Ronca, E. R. Sayfutyarova, M. Scheurer, H. F.
759 Schurkus, J. E. T. Smith, C. Sun, S.-N. Sun, S. Upadhyay, L. K. Wagner, X. Wang, A. White,
760 J. D. Whitfield, M. J. Williamson, S. Wouters, J. Yang, J. M. Yu, T. Zhu, T. C. Berkelbach,
761 S. Sharma, A. Y. Sokolov, and G. K.-L. Chan, *The Journal of Chemical Physics* **153**, 024109
762 (2020), publisher: American Institute of Physics.
- 763 ³⁷T. Kato, *Communications on Pure and Applied Mathematics* **10**, 151 (1957).
- 764 ³⁸K. Inui, Y. Kato, and Y. Motome, *Physical Review Research* **3**, 043126 (2021), publisher: Amer-
765 ican Physical Society.
- 766 ³⁹S. Sæther, T. Kjærgaard, H. Koch, and I.-M. Høyvik, *Journal of Chemical Theory and Compu-
767 tation* **13**, 5282 (2017), publisher: American Chemical Society.

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0157512

- 768 ⁴⁰X. He and K. M. J. Merz, *Journal of Chemical Theory and Computation* **6**, 405 (2010), publisher:
769 American Chemical Society.
- 770 ⁴¹D. L. Yeager and P. Jørgensen, *The Journal of Chemical Physics* **71**, 755 (1979), publisher:
771 American Institute of Physics.
- 772 ⁴²S.-i. Amari, in *Advances in Neural Information Processing Systems*, Vol. 9, edited by M. C.
773 Mozer, M. Jordan, and T. Petsche (MIT Press, 1996).
- 774 ⁴³J. Martens and R. Grosse, in *Proceedings of the 32nd International Conference on Machine*
775 *Learning* (PMLR, 2015) pp. 2408–2417, iSSN: 1938-7228.
- 776 ⁴⁴J. S. Spencer, D. Pfau, A. Botev, and W. M. C. Foulkes, “Better, Faster Fermionic Neural
777 Networks,” (2020), arXiv:2011.07125 [physics].
- 778 ⁴⁵G. Pescia, J. Nys, J. Kim, A. Lovato, and G. Carleo, “Message-Passing Neural Quantum States
779 for the Homogeneous Electron Gas,” (2023).
- 780 ⁴⁶A. Botev and J. Martens, “KFAC-JAX,” (2022).
- 781 ⁴⁷I. Loshchilov and F. Hutter, “Decoupled Weight Decay Regularization,” (2019),
782 arXiv:1711.05101 [cs, math].
- 783 ⁴⁸Y. Nomura, A. S. Darmawan, Y. Yamaji, and M. Imada, *Physical Review B* **96**, 205152 (2017),
784 publisher: American Physical Society.
- 785 ⁴⁹S. Sorella, *Physical Review Letters* **80**, 4558 (1998).
- 786 ⁵⁰S. Sorella, *Physical Review B* **64**, 024512 (2001).
- 787 ⁵¹J. Martens, *The Journal of Machine Learning Research* **21**, 146:5776 (2020).
- 788 ⁵²N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller, *The Journal of*
789 *Chemical Physics* **21**, 1087 (1953), publisher: American Institute of Physics.
- 790 ⁵³W. K. Hastings, *Biometrika* **57**, 97 (1970).
- 791 ⁵⁴J. Besag, *Journal of the Royal Statistical Society: Series B (Methodological)* **56**, 581 (1994).
- 792 ⁵⁵J. M. Flegal, M. Haran, and G. L. Jones, *Statistical Science* **23**, 250 (2008), publisher: Institute
793 of Mathematical Statistics.
- 794 ⁵⁶A. Sokal (Springer, 1996).
- 795 ⁵⁷N. Nemeč, M. D. Towler, and R. J. Needs, *The Journal of Chemical Physics* **132**, 034111 (2010),
796 publisher: American Institute of Physics.
- 797 ⁵⁸N. D. Drummond, M. D. Towler, and R. J. Needs, *Physical Review B* **70**, 235119 (2004).
- 798 ⁵⁹D. Feller, K. A. Peterson, and D. A. Dixon, *The Journal of Chemical Physics* **129**, 204105
799 (2008), publisher: American Institute of Physics.

800 ⁶⁰D. P. O'Neill and P. M. W. Gill, *Molecular Physics* **103**, 763 (2005).

801 ⁶¹A. Annaberdiyev, G. Wang, C. A. Melton, M. C. Bennett, L. Shulenburger, and L. Mitas, *The*
802 *Journal of Chemical Physics* **149**, 134108 (2018), arXiv:1805.06436 [cond-mat.mtrl-sci].

803 ⁶²M. Dolg, U. Wedig, H. Stoll, and H. Preuss, *Journal of Chemical Physics* **86**, 866 (1987).

804 ⁶³Y. A. Aoto, A. P. de Lima Batista, A. Köhn, and A. G. S. de Oliveira-Filho, *Journal of Chemical*
805 *Theory and Computation* **13**, 5291 (2017).

806 ⁶⁴K. T. Williams, Y. Yao, J. Li, L. Chen, H. Shi, M. Motta, C. Niu, U. Ray, S. Guo, R. J. Anderson,
807 J. Li, L. N. Tran, C.-N. Yeh, B. Mussard, S. Sharma, F. Bruneval, M. van Schilfhaarde, G. H.
808 Booth, G. K.-L. Chan, S. Zhang, E. Gull, D. Zgid, A. Millis, C. J. Umrigar, and L. K. Wagner
809 (Simons Collaboration on the Many-Electron Problem), *Phys. Rev. X* **10**, 011041 (2020).

810 **Appendix A: Additional scaling experiments**

811 **1. Pseudopotentials**

812 Figure 10 shows the scaling of the run time of the non-local pseudopotential evaluation (second
813 term in (12)) on the third and fourth-row atoms. This term dominates the total computational
814 overhead of using pseudopotentials overwhelmingly. From the 5 nested summations of (12), only
815 the sum over the valence electrons scales with the number of electrons, hinting at an approximate
816 linear scaling with system size. Looking at Figure 10, the obtained empirical scaling of $N^{1.19}$ is in
817 good agreement with expectations. The sudden jump in run time from 20 to 21 electrons is caused
818 by the reduction of valence electrons, as the utilized ccECP pseudopotentials use a larger core for
819 4p elements than for 3d ones.

820 **2. Memory requirement**

821 For all investigated applications, the memory requirement bottleneck is presented by the com-
822 putation of the Laplacian of the wave function, $\sum_i^{3N} \frac{\partial^2}{\partial^2 r_i} \psi_\theta(\mathbf{r})$. In this step, one obtains second
823 derivatives of the wave function with respect to the electron coordinates. We obtain the derivatives
824 of the wave function by applying automatic differentiation in backward-forward mode. While the
825 gradient $\nabla_{\mathbf{r}} \psi_\theta$ is obtained in one backward pass for all coordinates, the diagonal of the Hessian
826 $(\frac{\partial^2 \psi_\theta}{\partial^2 r_i})$ requires an additional $3N$ forward mode differentiations to compute, one for each electron
827 coordinate. Due to the flexible function transformations of JAX, both the serial and parallel exe-

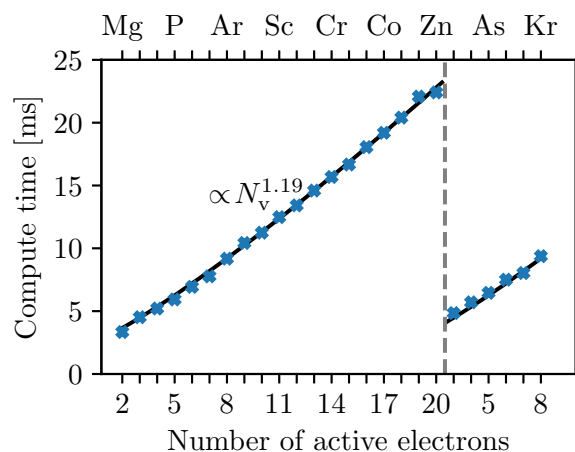


FIG. 10. **Scaling of the compute time of the non-local part of the ccECP pseudopotential.** The evaluation of the second term of Eq. (11) on a single NVIDIA GeForce RTX 3090 GPU. An exponential fit gives a scaling exponent of 1.19 ± 0.04 with the number of valence (explicitly treated) electrons. The pseudopotential uses a Neon core for elements up to Zn and a larger core for heavier atoms.

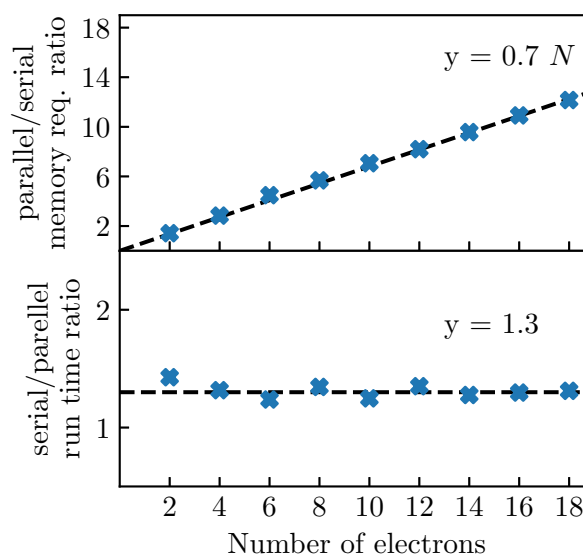


FIG. 11. **Comparing the memory requirements and run times of the serial and parallel Laplacian computations.** Data is obtained by evaluating the Laplacian of an untrained DEEPQMC ansatz with respect to all electron coordinates. The parallel implementation computes all diagonal elements of the Hessian at once, while the serial version computes one entry at a time. The number of diagonal Hessian entries scales linearly with the number of electrons. Run times measured on a single Nvidia GTX 1080 Ti GPU.

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0157512

828 cution of the $3N$ forward mode differentiation passes can be implemented in a few lines of code,
829 with the two implementations presumably differing in how they trade computational efficiency for
830 memory requirement.

831 To decide between the serial and parallel approach to the Laplacian computation, benchmark
832 calculations on a series of atoms with increasing nuclear charges are performed. The obtained
833 relative memory requirements of the parallel and serial computations are presented on the left
834 vertical axis of Figure 11. The observed linear scaling of the relative memory requirement between
835 parallel and serial evaluations can be understood by considering that the parallel implementation
836 holds data for all $3N$ backward passes in memory, while the serial approach stores data for a single
837 pass at a time. However, the prefactor of the scaling curve is significantly less than three, which
838 indicates that JAX performs some optimizations on the parallel code that reduce the naive $3N$
839 memory requirement multiplier. Considering run times of the two versions (lower panel of 11)
840 it is found that the relative timings of the serial implementation over the parallel implementation
841 do not scale with the system size. In fact, the ratio of run times between the serial and parallel
842 implementations appears to converge around 1.5. Taking the above observations into account, the
843 serial implementation of the Laplacian evaluation is chosen, due to its favorable scaling memory
844 requirements which outweigh the slight, non-scaling run time edge of the parallel implementation.

845 **Appendix B: Usage of DEEPQMC**

846 In this section, we provide a few minimal examples of the usage of the DEEPQMC command
847 line interface. The interface is based on HYDRA, which provides a modular way to configure and
848 execute complex jobs. DEEPQMC implements a wide variety of configuration options for the
849 wave function ansatz as well as the hyperparameters of training and evaluation. For ease of use,
850 the package includes predefined configuration files, which can be augmented using the command
851 line or extended with custom configuration files. For a thorough tutorial and API documentation
852 the reader is referred to the DEEPQMC documentation. For examples of typical DEEPQMC
853 commands see Figure 12.

This is the author's peer reviewed, accepted manuscript. However, the online version of record will be different from this version once it has been copyedited and typeset.

PLEASE CITE THIS ARTICLE AS DOI: 10.1063/5.0157512

```
# Optimize the default ansatz for H2O using decorrelated Langevin Sampling and  
# a reduced KFAC norm constraint  
deepqmc hamil/mol=H2O task/sampler=decorr_langevin task.opt.norm_constraint=0.0005  
  
# Now use a Ferminet ansatz along with its default hyperparameters  
deepqmc ansatz=ferminet task=train_ferminet hamil/mol=H2O hydra.run.dir=fn  
  
# Evaluate the previously trained ansatz using 10k inference steps  
deepqmc task=evaluate task.restdir=fn +task.steps=10000
```

FIG. 12. Example usage of the DEEPQMC program package through its command line interface.