

Ab-Initio Solution of the Many-Electron Schrödinger Equation with Deep Neural Networks

David Pfau, James S. Spencer, and Alexander G. de G. Matthews
DeepMind, 6 Pancras Square, London N1C 4AG

W. M. C. Foulkes

Department of Physics, Imperial College London, South Kensington Campus, London SW7 2AZ

(Dated: September 6, 2019)

Given access to accurate solutions of the many-electron Schrödinger equation, nearly all chemistry could be derived from first principles. Exact wavefunctions of interesting chemical systems are out of reach because they are NP-hard to compute in general,¹ but approximations can be found using polynomially-scaling algorithms. The key challenge for many of these algorithms is the choice of wavefunction approximation, or Ansatz, which must trade off between efficiency and accuracy. Neural networks have shown impressive power as accurate practical function approximators^{2,3} and promise as a compact wavefunction Ansatz for spin systems,⁴ but problems in electronic structure require wavefunctions that obey Fermi-Dirac statistics. Here we introduce a novel deep learning architecture, the Fermionic Neural Network, as a powerful wavefunction Ansatz for many-electron systems. The Fermionic Neural Network is able to achieve accuracy beyond other variational Monte Carlo Ansätze on a variety of atoms and small molecules. Using no data other than atomic positions and charges, we predict the dissociation curves of the nitrogen molecule and hydrogen chain, two challenging strongly-correlated systems, to significantly higher accuracy than the coupled cluster method,⁵ widely considered the gold standard for quantum chemistry. This demonstrates that deep neural networks can outperform existing ab-initio quantum chemistry methods, opening the possibility of accurate direct optimisation of wavefunctions for previously intractable molecules and solids.

I. INTRODUCTION

The ground state wavefunction $\psi(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$ and energy E of a chemical system with n electrons may be found by solving the time-independent Schrödinger equation,

$$\begin{aligned} \hat{H}\psi(\mathbf{x}_1, \dots, \mathbf{x}_n) &= E\psi(\mathbf{x}_1, \dots, \mathbf{x}_n) \\ \hat{H} &= -\frac{1}{2} \sum_i \nabla_i^2 + \sum_{i>j} \frac{1}{|\mathbf{r}_i - \mathbf{r}_j|} \\ &\quad - \sum_{I'} \frac{Z_{I'}}{|\mathbf{r}_i - \mathbf{R}_{I'}|} + \sum_{I>J} \frac{Z_I Z_J}{|\mathbf{R}_I - \mathbf{R}_J|} \end{aligned} \quad (1)$$

where $\mathbf{x}_i = \{\mathbf{r}_i, \sigma_i\}$ are the coordinates of electron i , with $\mathbf{r}_i \in \mathbb{R}^3$ the position and $\sigma_i \in \{\uparrow, \downarrow\}$ the spin, ∇_i^2

is the Laplacian with respect to \mathbf{r}_i , and \mathbf{R}_I and Z_I are the position and atomic number of nucleus I . We work in the Born-Oppenheimer approximation,⁶ where the nuclear positions are fixed input parameters, and Hartree atomic units are used throughout. The Schrödinger differential operator is spin independent but the electron spin matters because the wavefunction must obey Fermi-Dirac statistics — it must be antisymmetric under the simultaneous exchange of the position and spin coordinates of any two electrons: $\psi(\dots, \mathbf{x}_i, \dots, \mathbf{x}_j, \dots) = -\psi(\dots, \mathbf{x}_j, \dots, \mathbf{x}_i, \dots)$.

Many approaches in quantum chemistry start from a finite set of one-electron orbitals ϕ_1, \dots, ϕ_N and expand the many-electron wavefunction as a linear combination of antisymmetrised tensor products (Slater determinants) of those functions:

$$\sum_{\mathcal{P}} \text{sign}(\mathcal{P}) \prod_i \phi_i^k(\mathbf{x}_{\mathcal{P}_i}) = \begin{vmatrix} \phi_1^k(\mathbf{x}_1) & \dots & \phi_1^k(\mathbf{x}_n) \\ \vdots & & \vdots \\ \phi_n^k(\mathbf{x}_1) & \dots & \phi_n^k(\mathbf{x}_n) \end{vmatrix} = \det[\phi_i^k(\mathbf{x}_j)] = \det[\Phi^k], \quad (2)$$

$$\psi(\mathbf{x}_1, \dots, \mathbf{x}_n) = \sum_k \omega_k \det[\Phi^k], \quad (3)$$

where $\{\phi_1^k, \dots, \phi_n^k\}$ is a subset of n of the N orbitals, the sum in Eqn. 2 is taken over all permutations \mathcal{P} of

the electron indices, and the sum in Eqn. 3 is over all subsets of n orbitals. The difficulty is that the num-

ber of possible Slater determinants rises exponentially with the system size, restricting this “full configuration-interaction” (FCI) approach to tiny molecules, even with recent advances.⁷

To address problems of practical interest, a more compact representation of the wavefunction is needed. The choice of function class used to approximate the wavefunction is known as the wavefunction Ansatz. For most applications of quantum Monte Carlo (QMC) methods to quantum chemistry, the default choice is the Slater-Jastrow Ansatz,⁸ which takes a truncated linear combination of Slater determinants and adds a multiplicative term — the Jastrow factor — to capture close-range correlations. The Jastrow factor is normally a product of functions of the distances between pairs or triplets of particles. There are many alternative Ansätze,^{9,10} but for continuous-space many-electron problems in three dimensions the Slater-Jastrow Ansatz remains the default.

Here, we greatly improve the accuracy of the variational quantum Monte Carlo (VMC) method by using a neural network we dub the Fermionic Neural Network, or Fermi Net, as a more flexible Ansatz. This avoids the use of finite basis set, a significant source of error for other Ansätze, and models higher-order electron-electron interactions compactly. The use of neural networks as a compact wavefunction Ansatz has been studied before for lattice spin systems^{4,11,12} and small systems of bosons in continuous space.¹³ Applications of neural networks to continuous fermionic systems have been limited to date, presumably due to the complexity of Fermi-Dirac statistics. Existing work has been restricted to very small numbers of electrons,¹⁴ or has been of very low accuracy.¹⁵ Unlike these other approaches, we use the Slater determinant as the starting point for our Ansatz, and then extend it by generalising the single-electron orbitals to include generic exchangeable nonlinear interactions of *all* electrons.

The Fermi Net is not only an improvement over existing Ansätze for VMC, but is competitive with and in some cases superior to more sophisticated quantum chemistry algorithms. Projector methods such as diffusion quantum Monte Carlo (DMC)⁸ and auxiliary field quantum Monte Carlo (AFQMC)¹⁶ generate stochastic trajectories that sample the ground state wavefunction without the need for an explicit representation, although accurate explicit trial wavefunctions are still required for good performance and numerical stability. We find the Fermi Net is competitive with projector methods on all systems investigated, in contrast with the conventional wisdom that VMC is less accurate. Coupled cluster methods⁵ use an Ansatz that multiplies a reference wavefunction by an exponential of a truncated sum of creation and annihilation operators. This proves remarkably accurate for equilibrium geometries, but conventional reference wavefunctions are insufficient for systems with many low-lying excited states. We evaluate the Fermi Net on a variety of stretched systems and find that it outperforms coupled cluster in all cases.

II. FERMIONIC NEURAL NETWORKS

A. Fermionic Neural Network Architecture

To construct an expressive neural network Ansatz, we note that nothing requires the orbitals in the matrix in Eqn. 2 to be functions of the coordinates of a single electron. The only requirement for the determinant of a matrix-valued function of $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ to be antisymmetric is that exchanging any two input variables, \mathbf{x}_i and \mathbf{x}_j , exchanges two rows or columns of the output matrix, leaving the rest invariant. This observation allows us to replace the single-electron orbitals $\phi_i(\mathbf{x}_j)$ in the determinant by multi-electron functions $\phi_i(\mathbf{x}_j; \mathbf{x}_1, \dots, \mathbf{x}_{j-1}, \mathbf{x}_{j+1}, \dots, \mathbf{x}_n) = \phi_i(\mathbf{x}_j; \{\mathbf{x}_{/j}\})$, where $\{\mathbf{x}_{/j}\}$ denotes the set of all electron states except \mathbf{x}_j , so long as these functions are invariant to any change in the order of the arguments after \mathbf{x}_j . The construction of a set of these permutation-equivariant functions with a neural network is the main innovation of the Fermi Net. A diagram of the full Fermi Net architecture is given in Figure 1.

The Fermionic Neural Network takes features of single electrons and pairs of electrons as input. As input to the single-electron stream of the network, we include both the difference in position between each electron and nucleus $\mathbf{r}_i - \mathbf{R}_I$ and the distance $|\mathbf{r}_i - \mathbf{R}_I|$. The input to the two-electron stream is similarly the differences $\mathbf{r}_i - \mathbf{r}_j$ and distances $|\mathbf{r}_i - \mathbf{r}_j|$. Adding the absolute distances between particles directly as input removes the need to include a separate Jastrow factor after a determinant. As the distance is a non-smooth function at zero, the neural network is capable of expressing the non-smooth behavior of the wavefunction when two particles coincide — the wavefunction cusps. Accurately modeling the cusps is critical for correctly estimating the energy and other properties of the system. We denote the concatenation of all features for one electron \mathbf{h}_i^0 , or $\mathbf{h}_i^{0\alpha}$ if we explicitly index its spin $\alpha \in \{\uparrow, \downarrow\}$; the features of two electrons are denoted \mathbf{h}_{ij}^0 or $\mathbf{h}_{ij}^{0\alpha\beta}$.

Intermediate layers of the Fermionic Neural Network mix information together in a permutation-equivariant way, by taking the mean of activations from different streams of the network. These mean activations are then concatenated together and appended to the single-electron streams of the network. For a single layer this is a purely linear operation, but when combined with a nonlinear activation function after each layer it becomes a flexible architecture for building permutation-equivariant functions¹⁷. Information from both the other one-electron streams and the two-electron streams are fed into the one-electron streams. However, to reduce the computational overhead, no information is transferred between two-electron streams — these are multilayer perceptrons running in parallel. If the outputs of the one-electron stream at layer ℓ with spin α are denoted $\mathbf{h}_i^{\ell\alpha}$ and outputs of the two-electron stream are $\mathbf{h}_{ij}^{\ell\alpha\beta}$, then

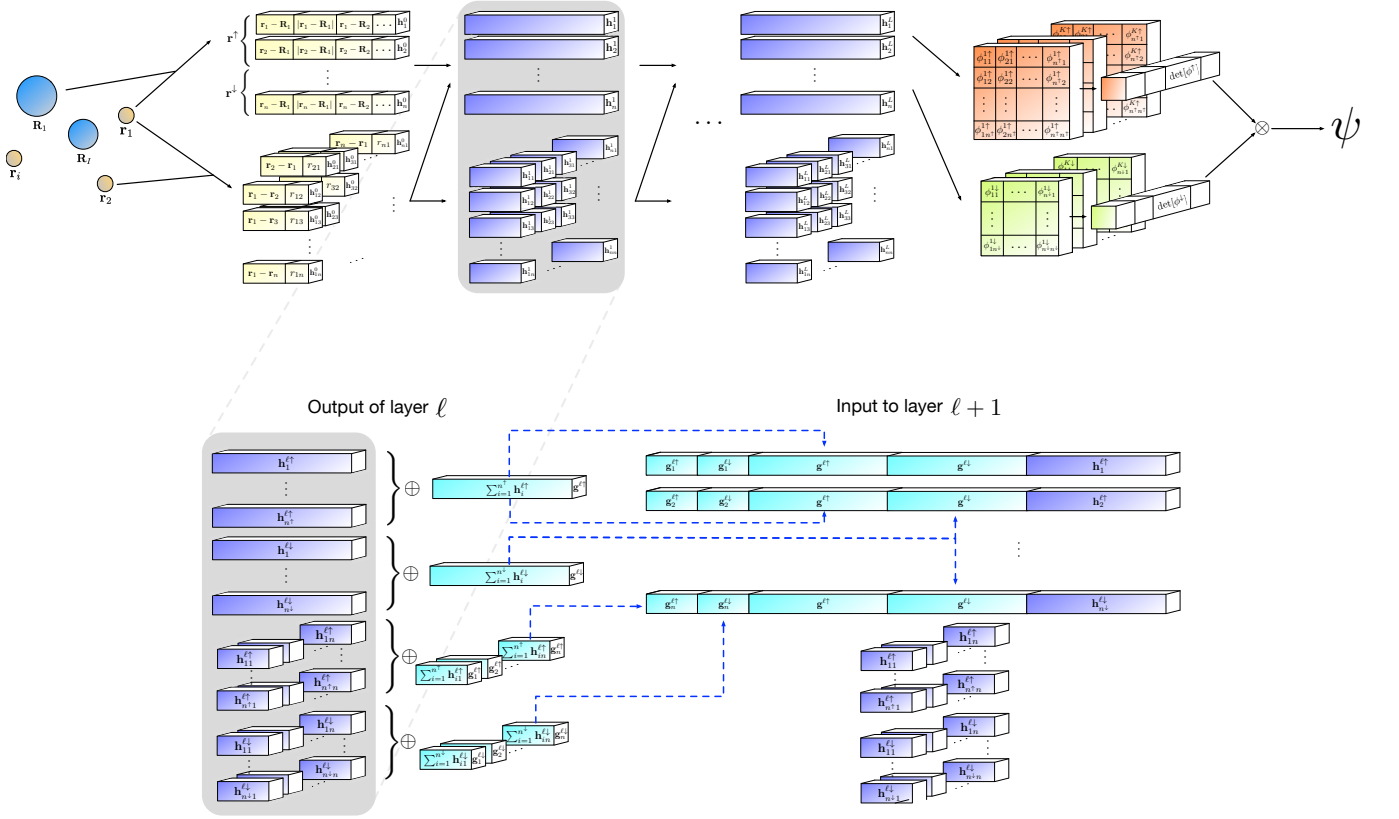


FIG. 1: The Fermionic Neural Network (Fermi Net). Top: Global architecture. Features of one or two electron positions are inputs to different streams of the network. These features are transformed through several layers, a determinant is applied, and the wavefunction at that position is given as output. Bottom: Detail of a single layer. The network averages features of electrons with the same spin together, then concatenates these features to construct an equivariant function of electron position at each layer.

the input to the one-electron stream for electron i with spin α at layer $\ell+1$ is

$$\left(\frac{1}{n^\uparrow} \sum_{j=1}^{n^\uparrow} \mathbf{h}_j^{\ell, \uparrow}, \frac{1}{n^\downarrow} \sum_{j=1}^{n^\downarrow} \mathbf{h}_j^{\ell, \downarrow}, \frac{1}{n^\uparrow} \sum_{j=1}^{n^\uparrow} \mathbf{h}_{ij}^{\ell, \alpha \uparrow}, \frac{1}{n^\downarrow} \sum_{j=1}^{n^\downarrow} \mathbf{h}_{ij}^{\ell, \alpha \downarrow}, \mathbf{h}_i^{\ell, \alpha} \right) = \left(\mathbf{g}^{\ell, \uparrow}, \mathbf{g}^{\ell, \downarrow}, \mathbf{g}_i^{\ell, \alpha \uparrow}, \mathbf{g}_i^{\ell, \alpha \downarrow}, \mathbf{h}_i^{\ell, \alpha} \right), \quad (4)$$

which is the concatenation of the mean activation for spin up and down parts of the one and two electron streams, respectively. This concatenated vector is then input into a linear layer followed by a tanh nonlinearity. A residual connection is also added between layers of the same shape, for both one and two electron streams.

After the last intermediate layer of the network, a final spin-dependent linear transformation is applied to the activations, and the output is multiplied by a weighted sum of exponentially-decaying envelopes, which enforces the boundary condition that the wavefunction goes to

zero far away from the nuclei:

$$\phi_i^{k\alpha}(\mathbf{r}_j^\alpha; \{\mathbf{r}_{/j}^\alpha\}; \{\mathbf{r}^{\bar{\alpha}}\}) = (\mathbf{w}_i^{k\alpha} \cdot \mathbf{h}_j^{L\alpha} + g_i^{k\alpha}) \sum_m \pi_{im}^{k\alpha} \exp(-|\boldsymbol{\Sigma}_{im}^{k\alpha}(\mathbf{r}_j^\alpha - \mathbf{R}_m)|), \quad (5)$$

where $\bar{\alpha}$ is \downarrow if α is \uparrow or vice versa, $\mathbf{h}_j^{L\alpha}$ is an output from the L -th (final) layer of the intermediate single-electron features network for electrons of spin α , and $\mathbf{w}_i^{k\alpha}$ ($g_i^{k\alpha}$) are the weights (biases) of the final linear transformation for determinant k . The learned parameters $\pi_{im}^{k\alpha}$ and $\boldsymbol{\Sigma}_{im}^{k\alpha} \in \mathbb{R}^{3 \times 3}$ control the anisotropic decay to zero far from each nucleus. The functions $\{\phi_i^{k\alpha}(\mathbf{r}_j^\alpha)\}$ are then used as the input to multiple determinants, and the full wavefunction is taken as a weighted sum of these determinants:

$$\psi(\mathbf{r}_1^\uparrow, \dots, \mathbf{r}_n^\downarrow) = \sum_k \omega_k \left(\det \left[\phi_i^{k\uparrow}(\mathbf{r}_j^\uparrow; \{\mathbf{r}_{/j}^\uparrow\}; \{\mathbf{r}^\downarrow\}) \right] \det \left[\phi_i^{k\downarrow}(\mathbf{r}_j^\downarrow; \{\mathbf{r}_{/j}^\downarrow\}; \{\mathbf{r}^\uparrow\}) \right] \right). \quad (6)$$

Eq. 6 uses the fact that the full determinant $\det[\Phi] = \det[\phi_i(\mathbf{x}_j; \{\mathbf{x}_{/j}\})]$ may be replaced by a product of spin-

up and spin-down terms since the spin does not appear explicitly in the Schrödinger equation⁸:

$$\det [\phi_i^\uparrow(\mathbf{r}_j^\uparrow; \{\mathbf{r}_{/j}^\uparrow\}; \{\mathbf{r}^\downarrow\})] \det [\phi_i^\downarrow(\mathbf{r}_j^\downarrow; \{\mathbf{r}_{/j}^\downarrow\}; \{\mathbf{r}^\uparrow\})] = \det [\Phi^\uparrow] \det [\Phi^\downarrow]. \quad (7)$$

The new wavefunction is only antisymmetric under exchange of electrons of the same spin, $\{\mathbf{r}^\uparrow\}$ or $\{\mathbf{r}^\downarrow\}$, but nevertheless yields correct expectation values of spin-independent observables and the fully antisymmetric wavefunction can be reconstructed if required. This factorisation allows spin-dependence to be handled explicitly rather than as input to the network.

B. Wavefunction Optimisation

As in the standard setting for wavefunction optimisation for variational Monte Carlo, we sought to minimise the energy expectation value of the wavefunction Ansatz:

$$\mathcal{L}(\theta) = \frac{\langle \psi_\theta | \hat{H} | \psi_\theta \rangle}{\langle \psi_\theta | \psi_\theta \rangle} = \frac{\int d\mathbf{X} \psi_\theta^*(\mathbf{X}) \hat{H} \psi_\theta(\mathbf{X})}{\int d\mathbf{X} \psi_\theta^*(\mathbf{X}) \psi_\theta(\mathbf{X})},$$

where θ are the parameters of the Ansatz, \hat{H} is the Hamiltonian of the system as given in Eqn. I, and $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ denotes the full state of all electrons. As \hat{H} is time-reversal invariant and Hermitian, its eigenfunctions and eigenvalues are real. If the minimisation is taken over all real normalisable functions, the minimum of the energy occurs when $\psi_\theta(\mathbf{X})$ is the ground-state eigenfunction of \hat{H} ; for a more restricted Ansatz, the minimum lies above the ground-state eigenvalue. When samples from the probability distribution defined by the wavefunction Ansatz $p(\mathbf{X}) \propto \psi_\theta^2(\mathbf{X})$ are available, unbiased estimates of the gradient of the energy with respect to θ can be computed as follows:

$$E_L(\mathbf{X}) = \psi^{-1}(\mathbf{X}) \hat{H} \psi(\mathbf{X}),$$

$$\nabla_\theta \mathcal{L} = \mathbb{E}_{p(\mathbf{X})} [(E_L - \mathbb{E}_{p(\mathbf{X})}[E_L]) \nabla_\theta \log|\psi|], \quad (8)$$

where E_L is the *local energy* and we have dropped the dependence of ψ on θ for clarity.

For all wavefunction Ansätze used in this paper, the determinants were computed in the log domain, and the final network output gave the log of the absolute value of the wavefunction, along with its sign. The local energy was computed directly in the log domain using the formula:

$$E_L(\mathbf{X}) = \psi^{-1}(\mathbf{X}) \hat{H} \psi(\mathbf{X})$$

$$= -\frac{1}{2} \sum_i \left[\left. \frac{\partial^2 \log|\psi|}{\partial r_i^2} \right|_{\mathbf{X}} + \left(\left. \frac{\partial \log|\psi|}{\partial r_i} \right|_{\mathbf{X}} \right)^2 \right] + V(\mathbf{X}),$$

where $V(\mathbf{X})$ is the potential energy of the state \mathbf{X} and the index i runs over all $3N$ dimensions of the electron position vector.

To optimise the wavefunction, we used a modified version of Kronecker Factorised Approximate Curvature (KFAC),¹⁸ an approximation to natural gradient

descent¹⁹ appropriate for neural networks. KFAC maintains the advantage over first order methods that other second order and hybrid methods have for wavefunction optimisation²⁰ (see Figure 2), but unlike exact second order methods, KFAC scales linearly in the number of layers of a neural network. Natural gradient descent updates for optimising \mathcal{L} with respect to parameters θ have the form $\delta\theta \propto \mathcal{F}^{-1} \nabla_\theta \mathcal{L}(\theta)$, where \mathcal{F} is the Fisher Information Matrix (FIM):

$$\mathcal{F}_{ij} = \mathbb{E}_{p(\mathbf{X})} \left[\frac{\partial \log p(\mathbf{X})}{\partial \theta_i} \frac{\partial \log p(\mathbf{X})}{\partial \theta_j} \right].$$

This is equivalent to stochastic reconfiguration²¹ when the probability density is unnormalised (see Supplementary Information) and closely related to the linear method of Toulouse and Umrigar.²²

For large neural networks with thousands to millions of parameters, solving the linear system $\mathcal{F} \delta\theta = \nabla_\theta \mathcal{L}$ becomes impractical. KFAC ameliorates this with two approximations. First, any terms \mathcal{F}_{ij} are assumed to be zero when θ_i and θ_j are in different layers of the network. This makes the FIM block-diagonal and significantly more efficient to invert. The second approximation is based on the structure of the gradient for a linear layer in a neural network. If W_ℓ is the weight matrix for layer ℓ of a network, then the block of the FIM for that weight is, in vectorised form:

$$\mathbb{E}_{p(\mathbf{X})} \left[\frac{\partial \log p(\mathbf{X})}{\partial \text{vec}(\mathbf{W}_\ell)} \frac{\partial \log p(\mathbf{X})}{\partial \text{vec}(\mathbf{W}_\ell)}^T \right] = \mathbb{E}_{p(\mathbf{X})} \left[(\mathbf{a}_\ell \otimes \mathbf{e}_\ell) (\mathbf{a}_\ell \otimes \mathbf{e}_\ell)^T \right] \quad (9)$$

where \mathbf{a}_ℓ are the forward activations and \mathbf{e}_ℓ are the backward sensitivities for that layer. KFAC approximates the inverse of this block as the Kronecker product of the inverse second moments:

$$\mathbb{E}_{p(\mathbf{X})} \left[(\mathbf{a}_\ell \otimes \mathbf{e}_\ell) (\mathbf{a}_\ell \otimes \mathbf{e}_\ell)^T \right]^{-1} \approx \mathbb{E}_{p(\mathbf{X})} [\mathbf{a}_\ell \mathbf{a}_\ell^T]^{-1} \otimes \mathbb{E}_{p(\mathbf{X})} [\mathbf{e}_\ell \mathbf{e}_\ell^T]^{-1} \quad (10)$$

Further details can be found in Martens and Grosse (2015).¹⁸

The original KFAC derivation assumed the density to be estimated was normalised, but we wish to extend it to stochastic reconfiguration for unnormalised wavefunctions. In the supplementary information, we show that if we only have access to an unnormalised wavefunction, terms in the FIM can be expressed as:

$$\mathcal{F}_{ij} \propto \mathbb{E}_{p(\mathbf{X})} [(\mathcal{O}_i - \mathbb{E}_{p(\mathbf{X})}[\mathcal{O}_i])(\mathcal{O}_j - \mathbb{E}_{p(\mathbf{X})}[\mathcal{O}_j])]$$

where $\mathcal{O}_i = \frac{\partial \log|\psi|}{\partial x_i}$. The terms in the FIM for the weights of a linear neural network layer would then be:

$$\begin{aligned}\mathbb{E}_{p(\mathbf{X})} \left[\frac{\partial \log p(\mathbf{X})}{\partial \text{vec}(\mathbf{W}_\ell)} \frac{\partial \log p(\mathbf{X})}{\partial \text{vec}(\mathbf{W}_\ell)}^T \right] &\propto \mathbb{E}_{p(\mathbf{X})} \left[(\mathbf{a}_\ell \otimes \mathbf{e}_\ell - \mathbb{E}_{p(\mathbf{X})}[\mathbf{a}_\ell \otimes \mathbf{e}_\ell]) (\mathbf{a}_\ell \otimes \mathbf{e}_\ell - \mathbb{E}_{p(\mathbf{X})}[\mathbf{a}_\ell \otimes \mathbf{e}_\ell])^T \right] \\ &= \mathbb{E}_{p(\mathbf{X})} \left[(\mathbf{a}_\ell \otimes \mathbf{e}_\ell) (\mathbf{a}_\ell \otimes \mathbf{e}_\ell)^T \right] - \mathbb{E}_{p(\mathbf{X})}[\mathbf{a}_\ell \otimes \mathbf{e}_\ell] \mathbb{E}_{p(\mathbf{X})}[\mathbf{a}_\ell \otimes \mathbf{e}_\ell]^T\end{aligned}$$

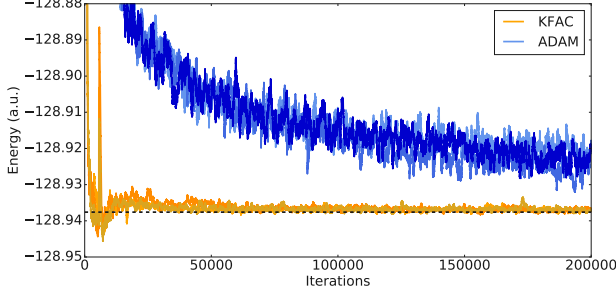


FIG. 2: Optimisation progress for neon atom with KFAC (orange) vs. ADAM (blue), with exact numbers in black. The qualitative advantage of KFAC is obvious. For clarity, a moving average of the energy is given over the last 1000 iterations. Note that the small dip early in optimisation with KFAC is due to the slow equilibration of the MCMC chain and goes away with a larger Metropolis-Hastings proposal step size.

We use a similar approximation for the inverse to that of conventional KFAC, replacing the uncentered second moments with mean-centered covariances:

$$\mathbb{E}_{p(\mathbf{X})} \left[\frac{\partial \log p(\mathbf{X})}{\partial \text{vec}(\mathbf{W}_\ell)} \frac{\partial \log p(\mathbf{X})}{\partial \text{vec}(\mathbf{W}_\ell)}^T \right] \approx \mathbb{E}_{p(\mathbf{X})} [\hat{\mathbf{a}}_\ell \hat{\mathbf{a}}_\ell^T]^{-1} \otimes \mathbb{E}_{p(\mathbf{X})} [\hat{\mathbf{e}}_\ell \hat{\mathbf{e}}_\ell^T]^{-1}, \quad (11)$$

where

$$\begin{aligned}\hat{\mathbf{a}}_\ell &= \mathbf{a}_\ell - \mathbb{E}_{p(\mathbf{X})}[\mathbf{a}_\ell], \\ \hat{\mathbf{e}}_\ell &= \mathbf{e}_\ell - \mathbb{E}_{p(\mathbf{X})}[\mathbf{e}_\ell].\end{aligned}$$

C. Experimental Setup

For all experiments, a Fermionic Neural Network with four layers was used, not counting the final linear layer that outputs the orbitals. Each layer had 256 hidden units for the one-electron stream and 32 hidden units for the two electron stream. A tanh nonlinearity was used for all layers, as a smooth function is needed to guarantee that the Laplacian is well defined and nonzero everywhere. 16 determinants were used where not otherwise specified. For comparison, the conventional VMC results

in Table I from Seth *et al.* (2011)²³ use 50 configuration state functions (CSF). While the exact number of determinants in a CSF will depend on the system, generally this will be on the order of hundreds to thousands of determinants. With this configuration of the Fermi Net there were approximately 700,000 parameters in the network, although the exact number depends on the number of atoms in the system due to the way we construct the input features and exponentially-decaying envelope.

For the baseline Slater-Jastrow network, multilayer perceptrons with 3 hidden layers of 128 units were used for the orbitals. The same multiplicative envelope employed in the Fermionic Neural Network was included; this can be seen as an extension to the electron-nuclear Jastrow factor. The Jastrow factor is of the standard form:²⁴

$$\begin{aligned}J(\{\mathbf{r}^\uparrow\}, \{\mathbf{r}^\downarrow\}, \{\mathbf{R}\}) &= \\ &\sum_{\alpha \in \{\uparrow, \downarrow\}} \sum_{i=1}^{n^\alpha} \sum_I \chi_j(|\mathbf{r}_i^\alpha - \mathbf{R}_I|) \\ &+ \sum_{\substack{\alpha \in \{\uparrow, \downarrow\} \\ \beta \in \{\uparrow, \downarrow\}}} \sum_{i=1}^{n^\alpha} \sum_{j=1}^{n^\beta} u^{\alpha\beta}(|\mathbf{r}_i^\alpha - \mathbf{r}_j^\beta|) \\ &+ \sum_{\substack{\alpha \in \{\uparrow, \downarrow\} \\ \beta \in \{\uparrow, \downarrow\}}} \sum_{i=1}^{n^\alpha} \sum_{j=1}^{n^\beta} \sum_I f_k^{\alpha\beta}(|\mathbf{r}_i^\alpha - \mathbf{r}_j^\beta|, |\mathbf{r}_i^\alpha - \mathbf{R}_I|, |\mathbf{r}_j^\beta - \mathbf{R}_I|),\end{aligned} \quad (12)$$

where χ_j , $u^{\alpha\beta}$ and $f_k^{\alpha\beta}$ are all 3-layer perceptrons with 64 hidden units.

To sample from $\psi^2(\mathbf{X})$ we used the standard Metropolis-Hastings algorithm.⁸ The proposed moves were Gaussian distributed with a fixed, isotropic covariance. All electron positions were updated simultaneously. Due to slow equilibration of the Markov Chain Monte Carlo (MCMC) sampling, the computed energy sometimes overshoot the true value, but always reequilibrated after a few thousand iterations. We experimented with Hamiltonian Monte Carlo to give faster mixing and lower bias in the gradients, but found this led to significantly higher variance in the local energy and lower overall performance.

Before using the local energy as an optimization objective we pretrained the network to match Hartree-Fock (HF) orbitals computed using PySCF²⁵. There were two

reasons for this. First, we found that the numerical stability of the subsequent local energy optimization was improved. On large systems, the determinants in the Fermionic Neural Network would often numerically underflow if no pretraining was used, causing the optimization to fail. Pretraining with HF orbitals as a guide meant that the main optimization started in a region of relatively low variance, with comparatively stable determinant evaluations and electron walkers in representative configurations. Second, we found that time was saved by not optimizing the local energy through a region that we knew to be physically uninteresting, given that it had an energy higher than that of a straightforward mean field approximation. The pretraining did not seem to strand the neural network in a poor local optimum, as the energy minimisation always gave consistent results capturing roughly 99% of the correlation energy. This is consistent with the conventional wisdom in the machine learning community that issues with local minima are less severe in wider, deeper neural networks. Further, stochasticity in the optimization procedure helps break symmetry and escape bad minima.

The pretraining loss is:

$$\mathcal{L}^{\text{pre}}(\theta) = \int \left[\sum_{\alpha \in \{\uparrow, \downarrow\}} \sum_{ijk} \left(\phi_i^{k\alpha}(\mathbf{r}_j^\alpha; \{\mathbf{r}_{/j}^\alpha\}; \{\mathbf{r}^{\bar{\alpha}}\}) - \phi_{i\alpha}^{\text{HF}}(\mathbf{r}_j^\alpha) \right)^2 \right] p^{\text{pre}}(\mathbf{X}) d\mathbf{X},$$

where $\phi_{i\alpha}^{\text{HF}}(\mathbf{r}_j^\alpha)$ denotes the value of the i -th Hartree-Fock orbital for spin α at the position of electron j , $\bar{\alpha}$ is \downarrow if α is \uparrow or vice versa, and $\phi_i^{k\alpha}(\mathbf{r}_j^\alpha; \{\mathbf{r}_{/j}^\alpha\}; \{\mathbf{r}^{\bar{\alpha}}\})$ is the corresponding entry in the input to the k -th determinant of the Fermionic Neural Network. We use a minimal (STO-3G) basis set for the Hartree-Fock computation as we require only a stable initialisation in the rough vicinity of the mean field solution, not an accurate mean field result. The probability distribution $p^{\text{pre}}(\mathbf{X})$ is an equal mixture of the product of Hartree-Fock orbitals and the output of the Fermionic Neural Network:

$$p^{\text{pre}}(\mathbf{X}) = \frac{1}{2} \left(\prod_{\alpha \in \{\uparrow, \downarrow\}} \prod_{ij} (\phi_{i\alpha}^{\text{HF}}(\mathbf{r}_j^\alpha))^2 + \psi^2(\mathbf{X}) \right).$$

We chose not to use the distribution from the Hartree-Fock determinant because we wanted sample coverage at every point where the orbitals were large, but in practice the difference to using the anti-symmetrized distribution was marginal. The inclusion of the neural network density helps to increase the sampling probability in areas where the neural network wavefunction is spuriously high. We approximate the expectation for the loss by using MCMC to draw half the samples in the batch from ψ^2 and half from the product of Hartree-Fock orbitals using MCMC.

Initial MCMC configurations were drawn from Gaussian distributions centred on each atom in the molecule. Electrons were assigned to atoms according to the nuclear charge and spin polarisation of the ground state of the isolated atom, with the atomic spins orientated such that the total spin projection of the molecule was correct, which was possible for systems studied here. We used ADAM with default parameters as the optimiser. After pretraining, we reinitialized the electron walker positions and then had a burn in MCMC period with target distribution ψ^2 before we began local energy minimization.

For the Fermi Net, all code was implemented in TensorFlow, and each experiment was run in parallel on 8 V100 GPUs. With a smaller batch size we were able to train on a single GPU but convergence was significantly slower. For instance, ethene converged after just 2 days of training with 8 GPUs, while several weeks were required on a single GPU. We expect further engineering improvements will reduce this number. 10 Metropolis-Hastings steps were taken between every parameter update, and it typically required $O(10^5 - 10^6)$ parameter updates to reach convergence (results in the paper used 2×10^5 parameter updates). Conventional VMC wavefunction optimisation will perform $O(10^1 - 10^2)$ parameter updates and $O(10^4 - 10^6)$ MCMC steps between updates, so we require roughly the same number of wavefunction evaluations as conventional VMC.

Accurate and stable convergence was highly dependent on the hyperparameters used; the default values for all experiments are included in Table III. These hyperparameters do seem to be generalisable — we have observed good performance on every system for which we used this configuration (all systems except bicyclobutane, due to memory limitations). For some larger systems, stability was improved by using more pretraining iterations. Getting good performance from KFAC requires careful tuning, and we found that the damping and norm constraint parameters critically affect the asymptotic performance. If the damping is too high, KFAC behaves like gradient descent near a local minimum and converges too slowly. If the damping is reduced, training quickly becomes unstable unless the norm constraint (a generalisation of gradient clipping) is lowered in tandem. Surprisingly, we found little advantage to using momentum, and sometimes it even seemed to reduce training performance, so we set it to zero for all experiments.

To reduce the variance in the parameter updates, we clipped the local energy when computing the gradients but not when evaluating the total energy of the system. This is a commonly used strategy to improve the accuracy of QMC²⁶. We computed the total variation of each batch, $\frac{1}{N} \sum_i |E_L(\mathbf{X}_i) - \tilde{E}_L|$, where \tilde{E}_L is the median local energy of that batch. This is to the ℓ_1 norm what the standard deviation is to the ℓ_2 norm, and we prefer it to the standard deviation as it is more robust to outliers. We clip any local energies more than 5 times further from the median than this total variation and compute the gradient in Eqn. 8 with the clipped energy in place

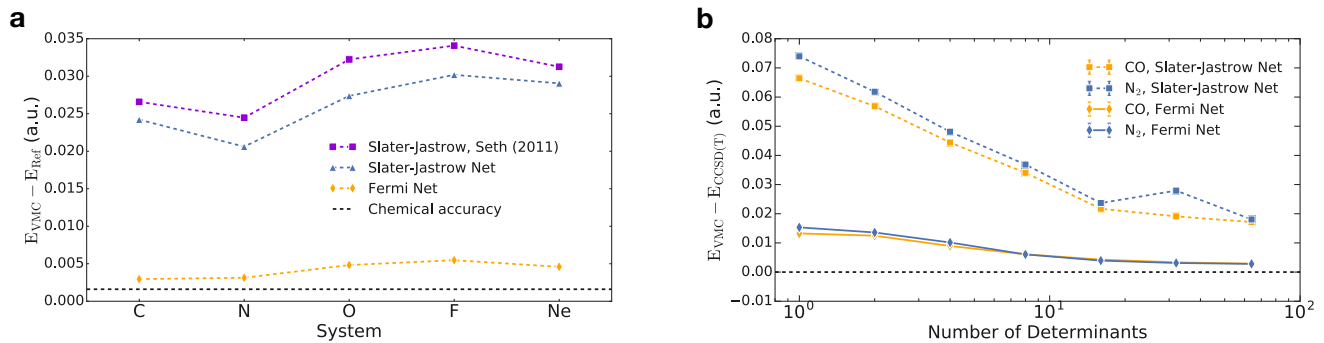


FIG. 3: Comparison of the Fermi Net against the Slater-Jastrow Ansatz. **a:** First-row atoms with a single determinant. Baseline numbers are from Chakravorty *et al.*³². The Slater-Jastrow neural network yields slightly lower energies than VMC with a conventional Slater-Jastrow Ansatz, while the Fermi Net is substantially more accurate. **b:** The CO and N₂ molecules (bond lengths 2.17328 a_0 and 2.13534 a_0 respectively) with increasing numbers of determinants. All-electron CCSD(T)/CBS results are used as the baseline. No matter how many determinants are used, the Fermi Net far exceeds the accuracy of the Slater-Jastrow net.

of E_L . The aforementioned KFAC norm constraint enforces gradient clipping in a manner which respects the information geometry of the model.

We used PySCF²⁵ to perform all-electron CCSD(T) calculations on atoms and dimers (Table I). PSI4²⁷ was used to perform all-electron CCSD(T) calculations on methane, ethene, ammonia and bicyclobutane, and CCSD(T) and FCI calculations on H₄. Cholesky decomposition²⁸ was used to reduce the memory requirements for bicyclobutane, which we verified introduces an error in the total energies of $\mathcal{O}(10^{-5})$ hartrees with the aug-cc-pCVTZ basis set. The H₄ calculations used a cc-pVXZ (X=T, Q, 5) basis set. All other CCSD(T) calculations used aug-cc-pCVXZ (X=T, Q, 5) basis sets. An unrestricted Hartree-Fock reference was used for atoms and dimers, with restricted Hartree-Fock used otherwise. We extrapolated energies to the CBS limit using standard methods^{29,30}. CBS Hartree-Fock energies for Li, Be and Li₂ were taken from aug-cc-pCV5Z calculations, in which the basis set error was below 10^{-4} hartrees. CBS Hartree-Fock energies for other systems were obtained by fitting the function $E_{\text{HF}}(X) = E_{\text{HF}}(\text{CBS}) + ae^{-bX}$, where X is the cardinality of the basis; CCSD, CCSD(T) and FCI correlation energies were extrapolated to the CBS by fitting the energies from quadruple- and quintuple-zeta basis sets (triple- and quadruple-zeta for bicyclobutane) to the function $E_c(X) = E_c(\text{HF}) + aX^{-3}$. The total energy is given by the sum of the Hartree-Fock energy and correlation energy. To compare the dissociation potential of N₂ against experiment, we used the MLR₄(6, 8) potential from Le Roy *et al.* (2006),³¹ which is based on fitting 19 lines of the N₂ vibrational spectrum.

III. RESULTS

A. Comparison of Slater-Jastrow and Fermi Net Ansätze

To demonstrate the expressive power of the Fermi Net, we investigated its performance relative to a standard Slater-Jastrow Ansatz with varying numbers of determinants. Rather than using Hartree-Fock orbitals and a closed-form Jastrow factor with few free parameters, our Slater-Jastrow Ansatz used a neural network to represent the one-electron orbitals and Jastrow factor, making it much more flexible. To fairly compare our calculations against previous work, we first looked at single-determinant Ansätze for first-row atoms. Figure 3a compares the Fermi Net results with numbers already available in the literature.²³ The neural network Slater-Jastrow Ansatz already outperforms the numbers from the literature by a few milli-Hartrees (mE_h). This could be due to the lack of basis set approximation error when using a neural network to represent the orbitals and Jastrow factor. The Fermi Net cuts the error relative to both the Slater-Jastrow neural network and the baseline by almost an order of magnitude. Just a single Fermi Net determinant is sufficient to come within a few mE_h of chemical accuracy, defined as 1 kcal/mol (1.594 mE_h), which is the typical standard for a quantum chemical calculation to be considered “correct.”

Not only is the Fermi Net a significant improvement over the Slater-Jastrow Ansatz with one determinant, but only a few Fermi Net determinants are necessary to saturate performance. Figure 3b shows the Slater-Jastrow network and Fermi Net energies of the nitrogen and carbon monoxide molecules as functions of the number of determinants. As FCI calculations are impractical for these systems, we compare against the unrestricted coupled cluster singles, doubles, and pertur-

Atom	Fermi Net	VMC ²³	Ground state energy (E_h)					Ionisation potential (mE_h)			Electron affinity (mE_h)		
			DMC ²³	CCSD(T)/CBS	HF/CBS	Exact ³²	% corr	Fermi Net	Expt. ³⁵	ΔE	Fermi Net	Expt. ³⁵	ΔE
Li	-7.47798(1)	-7.478034(8)	-7.478067(5)	-7.478157	-7.432747	-7.47806032	99.82(3)	198.10(4)	198.147	0.04(4)	21.82(20)	22.716	0.89(20)
Be	-14.66733(3)	-14.66719(1)	-14.667306(7)	-14.66737	-14.57301	-14.66736	99.97(3)	342.77(18)	342.593	-0.17(18)	-	-	-
B	-24.65370(3)	-24.65337(4)	-24.65379(3)	-24.65373	-24.53316	-24.65391	99.83(3)	304.86(4)	304.979	0.12(4)	9.03(11)	10.336	1.31(11)
C	-37.84471(5)	-37.84377(7)	-37.84446(6)	-37.8448	-37.6938	-37.8450	99.81(3)	413.98(8)	414.014	0.03(8)	46.18(9)	46.610	0.43(9)
N	-54.58882(6)	-54.5873(1)	-54.58867(8)	-54.5894	-54.4047	-54.5892	99.80(3)	534.80(12)	534.777	-0.03(12)	-	-	-
O	-75.06655(7)	-75.0632(2)	-75.0654(1)	-75.0678	-74.8192	-75.0673	99.70(3)	500.29(26)	500.453	0.17(26)	53.55(19)	53.993	0.44(19)
F	-99.7329(1)	-99.7287(2)	-99.7318(1)	-99.7348	-99.4168	-99.7339	99.69(3)	640.86(41)	640.949	0.09(41)	125.71(26)	125.959	0.25(26)
Ne	-128.9366(1)	-128.9347(2)	-128.9366(1)	-128.9394	-128.5479	-128.9376	99.74(3)	794.30(52)	794.409	0.11(52)	-	-	-

TABLE I: Ground state energy, ionisation potential and electron affinity for first-row atoms. The QMC method (Fermi Net, conventional VMC or DMC) closest to the exact ground state energy for each atom is in bold. Electron affinities for Be, N and Ne are not computed as their anions are unstable. Experimental ionisation potentials and electron affinities have had estimated relativistic effects³⁵ removed. All ground state energies are within chemical accuracy of the exact numerical solution, and all electron affinities and all ionisation potentials except neon are within chemical accuracy of experimental results. If no citation is provided, the number was from our own calculation.

Molecule	Bond length (a_0)	Fermi Net (E_h)	CCSD(T)/CBS (E_h)	HF/CBS (E_h)	Exact (E_h)	% corr
LiH	3.015	-8.07050(1)	-8.070696	-7.98737	-8.070548 ³⁶	99.94(1)
Li ₂	5.051	-14.99475(1)	-14.99507	-14.87155	-14.9954 ³⁷	99.47(1)
CO	2.173	-113.3218(1)	-113.3255	-112.7871		99.32(3)
N ₂	2.068	-109.5388(1)	-109.5425	-108.9940	-109.5423 ³⁷	99.36(2)
NH ₃		-56.56295(8)	-56.5644	-56.2247		99.57(2)
CH ₄		-40.51400(7)	-40.5150	-40.2171		99.66(3)
Ethene		-78.5824(2)	-78.5888	-78.0705		98.77(4)
Bicyclobutane		-155.9155(8)*	-155.9575	-154.9372		95.88(8)

TABLE II: Ground state energy at equilibrium geometry for diatomics and small molecules. The percentage of correlation energy captured by the Fermi Net relative to the exact energy (where available) or CCSD(T) is given in the rightmost column. If no citation is provided, the number was from our own calculation. *Due to computational limitations, the Fermi Net was run with half as many determinants and half the batch size for bicyclobutane as for other systems. Geometries for larger molecules are given in Supplementary Information.

bative triples method (CCSD(T)) in the complete basis set (CBS) limit, typically considered the gold standard for molecules of this size at equilibrium geometry. As the Slater-Jastrow network optimises all orbitals separately, the results from the Slater-Jastrow network should be a lower bound on the performance of a Slater-Jastrow Ansatz with a given number of determinants. As expected, the Slater-Jastrow network is still far from the accuracy of CCSD(T) at 64 determinants. The 64-determinant Fermi Net, in contrast, comes within a few mE_h of CCSD(T). The Fermi Net energies begin to plateau after only a few tens of determinants, suggesting that large linear combinations of Fermi Net determinants are not required. Despite recent advances in optimal determinant selection,^{33,34} conventional Slater-Jastrow VMC calculations typically require tens of thousands of determinants for systems of this size and rarely match CCSD(T) accuracy even then.

B. Equilibrium Geometries

Tables I and II show that the same 16-determinant Fermi Net with the same training hyperparameters generalises well to a wide variety of atoms and diatomic and small organic molecules. As a baseline, we used a combination of experimental and exact computational results where available,^{32,36,37} and all-electron CBS CCSD(T) otherwise. On all atoms, as well as LiH, Li₂, methane and ammonia, the Fermi Net error was within chemical accuracy. In comparison, energies from VMC using a conventional Slater-Jastrow-backflow Ansatz for first-row atoms²³ are uniformly worse than the Fermi Net, despite using at least an order of magnitude more determinants. The VMC-based Fermi Net energies are more comparable in quality to diffusion Monte Carlo (DMC), which is typically much more accurate than VMC. On the largest molecules investigated, ethene (C₂H₄) and bicyclobutane (C₄H₆), we recovered over 98% and 95% of the correlation energy respectively – remarkably good for a variational calculation. Bicyclobutane is an especially

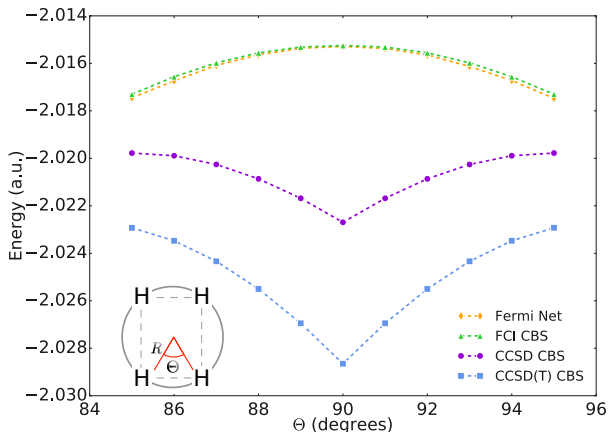


FIG. 4: The H_4 rectangle, $R = 3.2843a_0$. Coupled cluster methods incorrectly predict a cusp and energy minimum at $\Theta = 90^\circ$, while the Fermi Net approach agrees with exact FCI results.

challenging system due to its high ring strain and large number of electrons – too many for exact methods like FCI. We also computed the first ionisation potentials, $E(X^+) - E(X)$ for element X , and electron affinities, $E(X) - E(X^-)$, for first-row atoms (Table I) and compare to experimental data³⁵ with relativistic effects removed. Agreement with experiment is excellent (mean absolute error of 0.09 mE_h for ionisation potentials and 0.66 mE_h for electron affinities), demonstrating that the Fermi Net Ansatz is capable of representing charged and neutral species with similar accuracy.

C. The H_4 Rectangle

While CCSD(T) is considered the gold-standard for equilibrium geometries, it often fails for molecules with low-lying excited states or stretched, twisted or otherwise out-of-equilibrium geometries. Understanding these systems is critical for predicting many chemical properties. A model system small enough to be solved exactly by FCI for which coupled cluster fails is the rectangle of four hydrogen atoms, parametrised by the distance R of the atoms from the centre and the angle θ between neighbouring atoms.³⁸ FCI shows that the energy varies smoothly with θ and is maximised when the atoms are at the corners of a square ($\theta = 90^\circ$). The coupled cluster results are qualitatively incorrect, predicting an energy *minimum* with a non-analytic downward-facing cusp at 90° , caused by a crossing of two Hartree–Fock states with different symmetries.³⁹ Figure 4 shows that the Fermi Net does not suffer from the same errors as coupled cluster and is in essentially perfect agreement with FCI. We

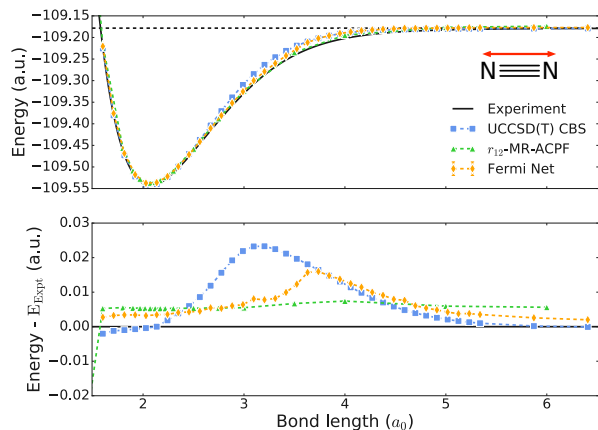


FIG. 5: The dissociation curve for the nitrogen triple-bond. The difference from experimental data³¹ is given in the bottom panel. In the region of largest UCCSD(T) error, the Fermi Net prediction is comparable to highly-accurate r_{12} -MR-ACPF results.⁴⁰

attribute the small discrepancy between the Fermi Net and FCI energies to errors arising from the basis set extrapolation used for the FCI energies.

D. The Nitrogen Molecule

A problem more relevant to real chemistry that troubles coupled cluster methods is the dissociation of the nitrogen molecule. The triple bond is challenging to describe accurately and the stretched N_2 molecule has several low-lying excited states, leading to errors when using single-reference coupled cluster methods.⁴¹ Experimental values for the dissociation potential have been reconstructed from spectroscopic measurements using the Morse/Long-range potential.³¹ These closely match calculations using r_{12} -MR-ACPF,⁴⁰ which is highly accurate but scales factorially. A comparison between unrestricted CCSD(T), the Fermi Net, and these high-accuracy methods is given in Figure 5. The total Fermi Net error is significantly smaller than UCCSD(T), and in the region of largest UCCSD(T) error the Fermi Net reaches accuracy comparable to r_{12} -MR-ACPF, but scales much more favourably with system size. While coupled cluster could in theory be made more accurate by extending to full triples or quadruples, or using multireference methods, CCSD(T) is generally considered the largest coupled cluster approximation that can reasonably scale beyond small molecules. This shows that, without any specific tuning to the system of interest, the Fermi Net is a clear improvement over coupled cluster for modelling a strongly correlated real-world chemical system.

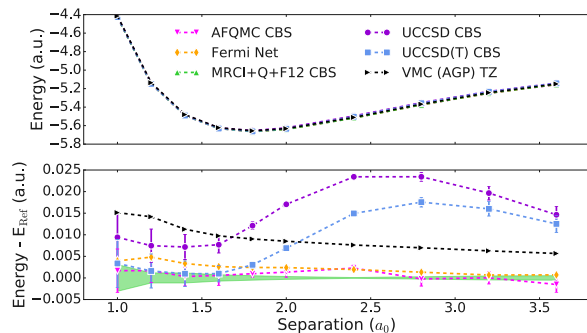


FIG. 6: The H_{10} chain. All energies except the Fermi Net are taken from Motta *et al.* (2017)⁴². Highly accurate MRCI+Q+F12 results are used as reference energies in the bottom panel, where the shaded region indicates an estimate of the basis-set extrapolation error. The errors in the coupled cluster and conventional VMC energies are large at medium atomic separations but the Fermi Net remains comparable to AFQMC.

E. The Hydrogen Chain

Finally, we investigated the performance of the Fermi Net on the evenly-spaced linear hydrogen chain. The hydrogen chain is of great interest as a system that bridges model Hamiltonians and real material systems and may undergo an insulator-to-metal transition as the separation of the atoms is decreased. Consequently, results obtained using a wide range of many-electron methods have been rigorously evaluated and compared.⁴² We compare the performance of the Fermi Net against many of these methods in Figure 6. Of the two projector QMC methods studied by Motta *et al.*, AFQMC gave slightly better results than lattice regularized DMC and so we omit the latter for clarity. Without changing the network architecture or hyperparameters, we are again able to outperform coupled cluster methods and conventional VMC and obtain results competitive with state-of-the-art approaches.

IV. DISCUSSION

We have shown that antisymmetric neural networks can be constructed and optimised to enable high-accuracy quantum chemistry calculations of challenging systems. The Fermionic Neural Network Ansatz makes the simple and straightforward VMC method competitive with DMC, AFQMC and CCSD(T) methods for equilibrium geometries and better than CCSD(T) for many out-of-equilibrium geometries. Importantly, one network architecture with one set of training parameters has been able to attain high accuracy on every system examined. The use of neural networks means that we do not have to choose a basis set or worry about basis-set extrapolation, a common source of error in computational quantum chemistry. There are many possible applications of the Fermi Net beyond VMC, for instance as a trial wavefunction for projector QMC methods. We expect further work investigating the tradeoffs of different antisymmetric neural networks and optimisation algorithms to lead to greater computational efficiency, higher representational capacity, and improved accuracy on larger systems. This has the potential to bring to quantum chemistry the same rapid progress that deep learning has enabled in numerous fields of artificial intelligence.

ACKNOWLEDGMENTS

We would like to thank J. Jumper, J. Kirkpatrick, T. Green, S. Mohamed and A. Cohen for helpful discussions, B. McMorrow for providing data, J. Martens and P. Buchlovsky for assistance with code, and A. Obika, S. Nelson, C. Meyer, T. Back, P. Kohli, K. Kavukcuoglu and D. Hassabis for support and guidance. Additional thanks to the rest of the DeepMind team for support, ideas and encouragement.

Correspondence and requests for materials should be addressed to D.P. (pfau@google.com).

¹ M. Troyer and U. Wiese, Phys. Rev. Lett. **94**, 170201 (2005).

² A. Krizhevsky, I. Sutskever, and G. E. Hinton, in *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1*, NIPS'12 (2012) pp. 1097–1105.

³ D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, *et al.*, Nature **529**, 485 (2016).

⁴ G. Carleo and M. Troyer, Science **356**, 602 (2017).

⁵ R. J. Bartlett and M. Musiał, Rev. Modern Phys. **79**, 291 (2007).

⁶ M. Born and R. Oppenheimer, Annalen der Physik **389**, 457 (1927).

⁷ G. H. Booth and A. Alavi, J. Chem. Phys. **132**, 174104 (2010).

⁸ W. M. C. Foulkes, L. Mitás, R. J. Needs, and G. Rajagopal, Rev. Modern Phys. **73**, 33 (2001).

⁹ M. Bajdich, L. Mitás, G. Drobný, L. Wagner, and K. Schmidt, Phys. Rev. Lett. **96**, 130201 (2006).

¹⁰ R. Orús, Ann. Phys. **349**, 117 (2014).

- ¹¹ K. Choo, G. Carleo, N. Regnault, and T. Neupert, Phys. Rev. Lett. **121**, 167204 (2018).
- ¹² A. Nagy and V. Savona, Phys. Rev. Lett. **122**, 250501 (2019).
- ¹³ H. Saito, J. Phys. Soc. Japan **87**, 074002 (2018).
- ¹⁴ J. Kessler, C. Calcavechia, and T. D. Kühne, arXiv preprint arXiv:1904.10251 (2019).
- ¹⁵ J. Han, L. Zhang, and W. E, arXiv preprint arXiv:1807.07014 (2018).
- ¹⁶ S. Zhang, “Ab initio electronic structure calculations by auxiliary-field quantum monte carlo,” in *Handbook of Materials Modeling : Methods: Theory and Modeling*, edited by W. Andreoni and S. Yip (2018) pp. 1–27.
- ¹⁷ J. Shawe-Taylor, in *ICANN* (1989) pp. 158–162.
- ¹⁸ J. Martens and R. Grosse, in *Proceedings of the 32nd International Conference on Machine Learning*, Proceedings of Machine Learning Research, Vol. 37 (2015) pp. 2408–2417.
- ¹⁹ S. Amari, Neural Comput. **10**, 251 (1998).
- ²⁰ L. Otis and E. Neuscamman, Phys. Chem. Chem. Phys. **21**, 14491 (2019).
- ²¹ S. Sorella, Phys. Rev. Lett. **80**, 4558 (1998).
- ²² J. Toulouse and C. Umrigar, J. Chem. Phys. **126**, 084102 (2007).
- ²³ P. Seth, P. L. Ríos, and R. J. Needs, J. Chem. Phys. **134**, 084105 (2011).
- ²⁴ R. Needs, M. Towler, N. Drummond, , and P. L. Rios, “CASINO: User’s Guide Version 2.13,” (2015).
- ²⁵ Q. Sun, T. C. Berkelbach, N. S. Blunt, G. H. Booth, S. Guo, Z. Li, J. Liu, J. D. McClain, E. R. Sayfutyarova, S. Sharma, S. Wouters, and G. K.-L. Chan, Wiley Interdisciplinary Reviews: Computational Molecular Science **8**, e1340 (2018).
- ²⁶ C. Umrigar, M. Nightingale, and K. Runge, J. Chem. Phys. **99**, 2865 (1993).
- ²⁷ R. M. Parrish, L. A. Burns, D. G. A. Smith, A. C. Simmonett, A. E. DePrince, E. G. Hohenstein, U. Bozkaya, A. Y. Sokolov, R. Di Remigio, R. M. Richard, J. F. Gonthier, A. M. James, H. R. McAlexander, A. Kumar, M. Saitow, X. Wang, B. P. Pritchard, P. Verma, H. F. Schaefer, K. Patkowski, R. A. King, E. F. Valeev, F. A. Evangelista, J. M. Turney, T. D. Crawford, and C. D. Sherrill, J. Chem. Theory Comput. **13**, 3185–3197 (2017).
- ²⁸ A. DePrince and C. Sherrill, J. Chem. Theory Comput. **9**, 2687 (2013).
- ²⁹ D. Feller, J. Chem. Phys. **96**, 6104 (1992).
- ³⁰ T. Helgaker and K. W., J. Chem. Phys. **106**, 9639 (1997).
- ³¹ R. J. Le Roy, Y. Huang, and C. Jary, J. Chem. Phys. **125**, 164310 (2006).
- ³² S. J. Chakravorty, S. R. Gwaltney, E. R. Davidson, F. A. Parpia, and C. F. Fischer, Phys. Rev. A **47**, 3649 (1993).
- ³³ E. Giner, R. Assaraf, and J. Toulouse, Mol. Phys. **114**, 910 (2016).
- ³⁴ M. Dash, S. Moroni, A. Scemama, and C. Filippi, J. Chem. Theory Comput. **14**, 4176 (2018).
- ³⁵ W. Klopper, B. R. A., T. D. P., , and H. C., Phys. Rev. A **81**, 022503 (2010).
- ³⁶ W. Cencek and J. Rychlewski, Chem. Phys. Lett. **320**, 549 (2000).
- ³⁷ C. Filippi and C. Umrigar, J. Chem. Phys. **105**, 213 (1996).
- ³⁸ T. Van Voorhis and M. Head-Gordon, J. Chem. Phys. **113**, 8873 (2000).
- ³⁹ H. Burton and A. Thom, J. Chem. Theory Comput. **12**, 67 (2016).
- ⁴⁰ R. Gdanitz, Chem. Phys. Lett. **283**, 253 (1998).
- ⁴¹ D. Lyakh, M. Musiał, V. Lotrich, and R. Bartlett, Chem. Rev. **112**, 182 (2011).
- ⁴² M. Motta, D. M. Ceperley, G. K.-L. Chan, J. A. Gomez, E. Gull, S. Guo, C. A. Jiménez-Hoyos, T. N. Lan, J. Li, F. Ma, A. J. Millis, N. V. Prokof’ev, U. Ray, G. E. Scuseria, S. Sorella, E. M. Stoudenmire, Q. Sun, I. S. Tupitsyn, S. R. White, D. Zgid, and S. Zhang, Phys. Rev. X **7**, 031059 (2017).
- ⁴³ D. Petz and C. Sudár, J. Math. Phys. **37**, 2662 (1996).
- ⁴⁴ M. B. Giles, in *Advances in Automatic Differentiation*, edited by C. H. Bischof, H. M. Bücker, P. Hovland, U. Naumann, and J. Utke (2008) pp. 35–44.
- ⁴⁵ L. A. Curtiss, K. Raghavachari, P. C. Redfern, V. Rassolov, and J. Pople, J. Chem. Phys. **109**, 7764 (1998).

SUPPLEMENTARY INFORMATION

A. Scaling and Computation Time

One of the main claims of the Fermi Net is that it scales favorably compared to other ab-initio quantum chemistry methods. The ability to run at all on systems the size of bicyclobutane proves the Fermi Net scales more favorably than exact methods like FCI, but the scaling relative to other approximate methods is a more subtle question. Both the size of the Fermi Net (number of hidden units, number of layers, number of determinants) as well as the number of training iterations required to reach a certain level of accuracy is unknown, and likely depends on the system being studied. What can be quantified is the computational complexity of a single iteration of training, which can be seen as a lower bound on the computational complexity of training the Fermi Net to a certain level of accuracy.

For a system with N electrons and a Fermi Net with L hidden layers, H hidden units per layer and K determinants, evaluating the one-electron stream of the network scales as $\mathcal{O}(NLH^2)$, evaluating the two-electron stream scales as $\mathcal{O}(N^2LH^2)$ and evaluating the determinants scales as $\mathcal{O}(N^3K)$, so the determinant calculation should dominate for large systems. While evaluating the gradient of a function has the same asymptotic complexity as evaluating the function, evaluating the local energy scales as N times the complexity of evaluating the function, as computing the Laplacian has the same complexity as computing the Hessian with respect to the inputs, giving an asymptotic complexity of $\mathcal{O}(N^4K)$ as N grows. The KFAC update takes an additional $\mathcal{O}(K^3 + LH^3)$ computation, so it does not contribute anything to the scaling with N , giving an overall quartic asymptotic scaling with system size for the optimisation update.

We give an empirical analysis of the scaling in Figure 7 on atoms from lithium to zinc, using the default training configuration with 8 GPUs. For larger atoms, we were not able to run optimisation to convergence, but we were able to execute enough updates to get an accurate estimate of the timing for a single iteration. Fitting polynomials of different order to the curve, we find a cubic fit is

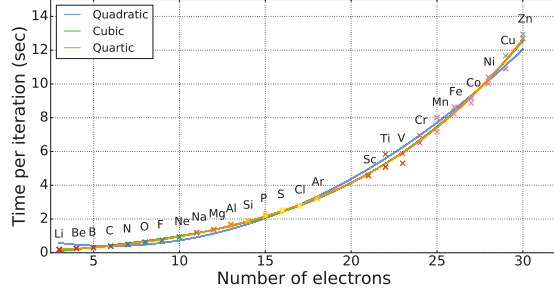


FIG. 7: Comparison of the runtime for one optimisation iteration on atoms up to zinc. Polynomial regressions up to fourth order are fit to the data. The small difference between the cubic and quartic fit suggests that the determinant computation is not the dominant factor at this scale.

able to accurately match the scaling, suggesting that for systems of this size the computation is dominated by the $O(N^2)$ evaluation of the two-electron stream of the Fermi Net, while the determinant only becomes dominant for much larger systems.

B. Equivalence of Natural Gradient Descent and Stochastic Reconfiguration

Here we provide a derivation illustrating that stochastic reconfiguration is equivalent to natural gradient descent for unnormalised distributions. Though many authors have investigated extensions of the Fisher information metric to quantum systems,⁴³ this particular connection between methods in machine learning and quantum chemistry seems not to be well appreciated by either community.

We denote the density proportional to $\psi^2(\mathbf{X})$ by $p(\mathbf{X})$, and the normalizing factor by $Z(\theta)$. In addition, let $\tilde{p}(\mathbf{X}) = \psi^2(\mathbf{X})$ denote the unnormalized density. In stochastic reconfiguration, the entries of the preconditioner matrix \mathcal{M} have the form

$$\mathcal{M}_{ij} = \mathbb{E}_{p(\mathbf{X})} [(\mathcal{O}_i - \mathbb{E}_{p(\mathbf{X})}[\mathcal{O}_i]) (\mathcal{O}_j - \mathbb{E}_{p(\mathbf{X})}[\mathcal{O}_j])],$$

where

$$\mathcal{O}_i(\mathbf{X}) = \psi(\mathbf{X})^{-1} \frac{\partial \psi(\mathbf{X})}{\partial \theta_i} = \frac{\partial \log |\psi(\mathbf{X})|}{\partial \theta_i} = \frac{1}{2} \frac{\partial \log \tilde{p}(\mathbf{X})}{\partial \theta_i}.$$

The term $\mathbb{E}_{p(\mathbf{X})}[\mathcal{O}_i]$ can be expressed in terms of the

normalizing factor:

$$\begin{aligned} \mathbb{E}_{p(\mathbf{X})}[\mathcal{O}_i] &= \frac{1}{2} \mathbb{E}_{p(\mathbf{X})} \left[\frac{\partial \log \tilde{p}(\mathbf{X})}{\partial \theta_i} \right] \\ &= \frac{1}{2} \int \frac{\partial \log \tilde{p}(\mathbf{X})}{\partial \theta_i} p(\mathbf{X}) d\mathbf{X} \\ &= \frac{1}{2} \int \frac{\partial \log \tilde{p}(\mathbf{X})}{\partial \theta_i} \frac{\tilde{p}(\mathbf{X})}{Z(\theta)} d\mathbf{X} \\ &= \frac{1}{2} \int \frac{1}{\tilde{p}(\mathbf{X})} \frac{\partial \tilde{p}(\mathbf{X})}{\partial \theta_i} \frac{\tilde{p}(\mathbf{X})}{Z(\theta)} d\mathbf{X} \\ &= \frac{1}{2Z(\theta)} \int \frac{\partial \tilde{p}(\mathbf{X})}{\partial \theta_i} d\mathbf{X} \\ &= \frac{1}{2Z(\theta)} \frac{\partial}{\partial \theta_i} \int \tilde{p}(\mathbf{X}) d\mathbf{X} \\ &= \frac{1}{2Z(\theta)} \frac{\partial Z(\theta)}{\partial \theta_i} \\ &= \frac{1}{2} \frac{\partial \log Z(\theta)}{\partial \theta_i}. \end{aligned}$$

Plugging this into the expression for \mathcal{M}_{ij} yields

$$\begin{aligned} \mathcal{M}_{ij} &= \mathbb{E}_{p(\mathbf{X})} [(\mathcal{O}_i - \mathbb{E}_{p(\mathbf{X})}[\mathcal{O}_i]) (\mathcal{O}_j - \mathbb{E}_{p(\mathbf{X})}[\mathcal{O}_j])] \\ &= \frac{1}{4} \mathbb{E}_{p(\mathbf{X})} \left[\left(\frac{\partial \log \tilde{p}(\mathbf{X})}{\partial \theta_i} - \frac{\partial \log Z(\theta)}{\partial \theta_i} \right) \left(\frac{\partial \log \tilde{p}(\mathbf{X})}{\partial \theta_j} - \frac{\partial \log Z(\theta)}{\partial \theta_j} \right) \right] \\ &= \frac{1}{4} \mathbb{E}_{p(\mathbf{X})} \left[\frac{\partial \log p(\mathbf{X})}{\partial \theta_i} \frac{\partial \log p(\mathbf{X})}{\partial \theta_j} \right], \end{aligned}$$

which, up to a constant, is the Fisher information metric for $p(\mathbf{X})$.

C. Numerically Stable Computation of the Log Determinant and Derivatives

For numerical stability, the Fermionic Neural Network outputs the *logarithm* of the absolute value of the wavefunction (along with its sign), and we compute log determinants rather than determinants. Even if some of the matrices are singular, this is not an issue for numerical stability on the forward pass, because these matrices will have zero contribution to the overall sum of determinants the network outputs:

$$\log |\psi(\mathbf{r}_1^\uparrow, \dots, \mathbf{r}_{n^\downarrow}^\downarrow)| = \log \left| \sum_k \omega_k \det [\Phi^{k\uparrow}] \det [\Phi^{k\downarrow}] \right|.$$

We use the “exp-normalise trick” to compute the sum — that is, we subtract off the largest log determinant before exponentiating and computing the weighted sum, and add it back in after the logarithm at the end. This avoids numerical underflow if the log determinants are not well scaled.

Naively applying automatic differentiation frameworks to compute the gradient and Laplacian of the log wavefunction will not work if one of the matrices is singular. However, the first and second derivatives are still well defined, and we show how to express these derivatives in closed form appropriate for reverse-mode automatic differentiation. Several of the results used here, as well as the notation, are based on the collected matrix derivative results of Giles (2008)⁴⁴.

From Jacobi's formula, the gradient of the determinant of a matrix is given by

$$\frac{\partial \det(\mathbf{A})}{\partial \mathbf{A}} = \det(\mathbf{A}) \mathbf{A}^{-T} = \text{Adj}(\mathbf{A})^T = \text{Cof}(\mathbf{A}),$$

where $\text{Cof}(\mathbf{A})$ is the cofactor matrix of \mathbf{A} . Let $\mathbf{C} = \text{Cof}(\mathbf{A})$. Then, by the product rule, we can express the reverse-mode gradient of $\text{Cof}(\mathbf{A})$ as

$$\bar{\mathbf{A}} = \mathbf{A}^{-T} [\text{Tr}(\bar{\mathbf{C}}^T \text{Cof}(\mathbf{A})) \mathbf{I} - \bar{\mathbf{C}}^T \text{Cof}(\mathbf{A})],$$

where $\bar{\mathbf{C}}$ is the reverse-mode sensitivity. Unfortunately, this expression becomes undefined if the matrix \mathbf{A} is singular. Even so, both the cofactor matrix and its derivative are still well defined. To see this, we express the cofactor in terms of the singular value decomposition of \mathbf{A} . Let $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ be the singular value decomposition of \mathbf{A} , then

$$\begin{aligned} \text{Cof}(\mathbf{A}) &= \det(\mathbf{A}) \mathbf{A}^{-T} \\ &= \det(\mathbf{U}) \det(\mathbf{\Sigma}) \det(\mathbf{V}) \mathbf{U} \mathbf{\Sigma}^{-1} \mathbf{V}^T. \end{aligned}$$

Since \mathbf{U} and \mathbf{V} are orthonormal matrices, their determinant is just the sign of their determinant. To avoid clutter, we drop the $\det(\mathbf{U})$ and $\det(\mathbf{V})$ terms until the very end. Let σ_i be the i th diagonal element of $\mathbf{\Sigma}$, then we have $\det(\mathbf{\Sigma}) = \prod_i \sigma_i$, and cancelling terms in the expression, we get (up to a sign factor)

$$\text{Cof}(\mathbf{A}) = \mathbf{U} \mathbf{\Gamma} \mathbf{V}^T,$$

where $\mathbf{\Gamma}$ is a diagonal matrix with elements γ_i defined as

$$\gamma_i = \prod_{j \neq i} \sigma_j$$

because the σ_i^{-1} term in $\mathbf{\Sigma}^{-1}$ cancels out one term in $\det(\mathbf{\Sigma})$.

The gradient of the cofactor is more complicated, but once again terms cancel. Again neglecting a sign factor, the reverse-mode gradient can be expanded in terms of the singular vectors as:

$$\begin{aligned} \bar{\mathbf{A}} &= \mathbf{A}^{-T} [\text{Tr}(\bar{\mathbf{C}}^T \text{Cof}(\mathbf{A})) \mathbf{I} - \bar{\mathbf{C}}^T \text{Cof}(\mathbf{A})] \\ &= \mathbf{U} \mathbf{\Sigma}^{-1} \mathbf{V}^T [\text{Tr}(\bar{\mathbf{C}}^T \mathbf{U} \mathbf{\Gamma} \mathbf{V}^T) \mathbf{I} - \bar{\mathbf{C}}^T \mathbf{U} \mathbf{\Gamma} \mathbf{V}^T] \\ &= \mathbf{U} [\text{Tr}(\bar{\mathbf{C}}^T \mathbf{U} \mathbf{\Gamma} \mathbf{V}^T) \mathbf{\Sigma}^{-1} - \mathbf{\Sigma}^{-1} \mathbf{V}^T \bar{\mathbf{C}}^T \mathbf{U} \mathbf{\Gamma}] \mathbf{V}^T \\ &= \mathbf{U} [\text{Tr}(\mathbf{M} \mathbf{\Gamma}) \mathbf{\Sigma}^{-1} - \mathbf{\Sigma}^{-1} \mathbf{M} \mathbf{\Gamma}] \mathbf{V}^T, \end{aligned}$$

where $\mathbf{M} = \mathbf{V}^T \bar{\mathbf{C}}^T \mathbf{U}$, and we have taken advantage of the invariance of the trace of matrix products to cyclic permutation in the last line.

Now, in the expression inside the square brackets in the last line, terms conveniently cancel that prevent the expression from becoming undefined should $\sigma_i = 0$ for some singular value. Denote this term Ξ , the off-diagonal terms of Ξ only depend on the second term $\mathbf{\Sigma}^{-1} \mathbf{M} \mathbf{\Gamma}$:

$$\begin{aligned} \Xi_{ij} &= -M_{ij} \sigma_i^{-1} \gamma_j \\ &= -M_{ij} \sigma_i^{-1} \prod_{k \neq j} \sigma_k \\ &= -M_{ij} \prod_{k \neq i, j} \sigma_k, \end{aligned}$$

and the diagonal terms have the form

$$\begin{aligned} \Xi_{ii} &= \sigma_i^{-1} \sum_j M_{jj} \gamma_j - M_{ii} \sigma_i^{-1} \gamma_i \\ &= \sum_{j \neq i} M_{jj} \sigma_i^{-1} \gamma_j \\ &= \sum_{j \neq i} M_{jj} \prod_{k \neq i, j} \sigma_k. \end{aligned}$$

Putting this all together, we get

$$\bar{\mathbf{A}} = \text{Sgn}(\det(\mathbf{U})) \text{Sgn}(\det(\mathbf{V})) \mathbf{U} \Xi \mathbf{V}^T,$$

with

$$\Xi_{ij} = \begin{cases} \sum_{j \neq i} M_{jj} \rho_{ij}, & \text{if } i = j, \\ -M_{ij} \rho_{ij}, & \text{otherwise,} \end{cases}$$

$$\rho_{ij} = \prod_{k \neq i, j} \sigma_k,$$

$$\mathbf{M} = \mathbf{V}^T \bar{\mathbf{C}} \mathbf{U}.$$

This allows us to compute second derivatives of the matrix determinant even for singular matrices. To handle degenerate matrices gracefully, we fuse everything from the computation of the log determinant to the final network output into a single TensorFlow op, with a custom gradient and gradient-of-gradient that includes the expression above.

D. Molecular structures

Molecular structures were taken from the G3 database⁴⁵ where available. We reproduce the atomic positions for all molecules studied in Tables IV-VII.

	Parameter	Value
	Batch size	4096
	Training iterations	2e5
	Pretraining iterations	1e3
	Learning rate	$(1e4 + t)^{-1}$
	Local energy clipping	5.0
KFAC	Momentum	0
KFAC	Covariance moving average decay	0.95
KFAC	Norm constraint	1e-3
KFAC	Damping	1e-3
MCMC	Proposal std dev (per dimension)	0.02
MCMC	Steps between parameter updates	10

TABLE III: Default hyperparameters for all experiments in the paper. For bicyclobutane, the batch size was halved and the pretraining iterations were increased by an order of magnitude.

Atom	Position (a_0)
C1	(0.0, 2.13792, 0.58661)
C2	(0.0, -2.13792, 0.58661)
C3	(1.41342, 0.0, -0.58924)
C4	(-1.41342, 0.0, -0.58924)
H1	(0.0, 2.33765, 2.64110)
H2	(0.0, 3.92566, -0.43023)
H3	(0.0, -2.33765, 2.64110)
H4	(0.0, -3.92566, -0.43023)
H5	(2.67285, 0.0, -2.19514)
H6	(-2.67285, 0.0, -2.19514)

TABLE VII: Atomic positions for bicyclobutane (C_4H_6).

Atom	Position (a_0)
N	(0.0, 0.0, 0.22013)
H1	(0.0, 1.77583, -0.51364)
H2	(1.53791, -0.88791, -0.51364)
H3	(-1.53791, -0.88791, -0.51364)

TABLE IV: Atomic positions for NH_3 .

Atom	Position (a_0)
C	(0.0, 0.0, 0.0)
H1	(1.18886, 1.18886, 1.18886)
H2	(-1.18886, -1.18886, 1.18886)
H3	(1.18886, -1.18886, -1.18886)
H4	(-1.18886, 1.18886, -1.18886)

TABLE V: Atomic positions for CH_4 .

Atom	Position (a_0)
C1	(0.0, 0.0, 1.26135)
C2	(0.0, 0.0, -1.26135)
H1	(0.0, 1.74390, 2.33889)
H2	(0.0, -1.74390, 2.33889)
H3	(0.0, 1.74390, -2.33889)
H4	(0.0, -1.74390, -2.33889)

TABLE VI: Atomic positions for ethene (C_2H_4).