Is attention all you need to solve the correlated electron problem?



Khachatur Nazaryan



Timothy Zaklama



Liang Fu

Max Geier arXiv:2502.05383 subMIT, 07/01/2025



Philosophy of neural network wavefunctions

Many-parameter ansatz for unbiased approximation of *any* wavefunction, **but** network structure should **reflect basic structure** of ground states for efficiency and accuracy.

→ Generative Al model: The networks learns to represent the electron correlations.
 → Physics-informed: Pauli principle enforced.



Minimal human bias

Electron correlations from generalized orbitals

Describe how the state of each electron is affected by the state of all other electrons

 $\phi_j(\boldsymbol{r}_i) \rightarrow \phi_j(\boldsymbol{r}_i; \{\boldsymbol{r}_{/i}\})$

Slater determinant of *correlated* orbitals describes many-body wavefunction $\Psi(\boldsymbol{r}_1, \dots, \boldsymbol{r}_N) = \det_{i,j} \phi_j(\boldsymbol{r}_i; \{\boldsymbol{r}_{/i}\})$





A particle moving in an electron liquid distorts the motion of all other particles along its way.

Is attention all you need to solve the correlated electron problem?

Employ attention to learn correlations between electrons.



How to construct a neural network variational wavefunction

r х О у О

feature





feature

embed: $\mathbf{h}^0 = W^0 \mathbf{r}$





Universal approximation theorem

The family of neural networks lies dense in the function space.

Hornik, Stinchcombe, and White, *Neural networks* **2** (5) 359-366 (1989)

Neural network Hartree-Fock

Neural network Hartree-Fock



Universal approximation theorem

Because single-particle orbitals are densely approximated, so are Slater determinants.

Capturing correlations

$$\Psi(\mathbf{r}_1, \dots, \mathbf{r}_N) = \sum_k \omega_k \det\left(\Phi_j^k(\mathbf{r}_i; \{\mathbf{r}_{/i}\})\right)$$

Slater determinants of generalized orbitals $\Phi_j(\mathbf{r}_i; \{\mathbf{r}_{/i}\})$ with **permutation equivariant** dependence on remaining electrons $\Phi_j(\mathbf{r}_i; \{\dots, \mathbf{r}_l, \dots, \mathbf{r}_k, \dots\} = \Phi_j(\mathbf{r}_i; \{\dots, \mathbf{r}_k, \dots, \mathbf{r}_l, \dots\})$

Any antisymmetric function can be written as a Slater determinant of generalized orbitals*

Pfau et al., Phys Rev. Research 2, 033429 (2020)

Permutation equivariant functions

Extend on Hartree-Fock: Which functions $\mathbf{f}: \mathbb{R}^{d_i \times N} \to \mathbb{R}^{d_i}$ are permutation equivariant in N - 1 arguments $\{\mathbf{h}_{/i}\}$, nonlinear, and depend explicitly on \mathbf{h}_i ? $\mathbf{f}(\mathbf{h}_i; ..., \mathbf{h}_l, ..., \mathbf{h}_k, ..., \mathbf{h}_N) = \mathbf{f}(\mathbf{h}_i; ..., \mathbf{h}_k, ..., \mathbf{h}_l, ..., \mathbf{h}_N)$



Permutation equivariant functions

Extend on Hartree-Fock: Which functions $\mathbf{f}: \mathbb{R}^{d_i \times N} \to \mathbb{R}^{d_i}$ are permutation equivariant in N - 1 arguments $\{\mathbf{h}_{/i}\}$, nonlinear, and depend explicitly on \mathbf{h}_i ? $\mathbf{f}(\mathbf{h}_i; ..., \mathbf{h}_l, ..., \mathbf{h}_k, ..., \mathbf{h}_N) = \mathbf{f}(\mathbf{h}_i; ..., \mathbf{h}_k, ..., \mathbf{h}_l, ..., \mathbf{h}_N)$

$$\mathbf{f}_{i} \mathbf{W}^{o} \mathbf{h}_{i}^{o} \mathbf{h}_{i}^{o$$

Capturing correlations with self-attention



von Glehn et al., arXiv:2211.13672 (2022)

Self-attention wavefunction ansatz



Optimizing the variational wavefunction with variational Monte Carlo

Variational Monte Carlo



How to sample an unnormalized probability density $P(\mathbf{R})$

Metropolis algorithm.

(i) Start walker at **R**.

(ii) Propose move to \mathbf{R}' with probability $T(\mathbf{R} \leftarrow \mathbf{R}')$.

(iii) Accept with probability

 $A(\mathbf{R} \leftarrow \mathbf{R}') = \operatorname{Min}\left(1, \frac{T(\mathbf{R} \leftarrow \mathbf{R}')P(\mathbf{R}')}{T(\mathbf{R}' \leftarrow \mathbf{R})P(\mathbf{R})}\right)$

Record \mathbf{R}' if accepted, else record \mathbf{R} . (iv) Repeat at (ii).

=> Walker density $n(\mathbf{R})$ proportional to $P(\mathbf{R})$.



Estimating energy

For local Hamiltonian
$$\widehat{H}_{\mathbf{R}\mathbf{R}'} = \widehat{H}_{\mathbf{R}}\delta_{\mathbf{R}\mathbf{R}'}$$

 $\langle \widehat{H} \rangle = \frac{1}{\int d\mathbf{R} |\Psi(\mathbf{R})|^2} \int d\mathbf{R} \Psi^*(\mathbf{R}) \widehat{H}_{\mathbf{R}} \Psi(\mathbf{R})$
 $= \frac{1}{\int d\mathbf{R} |\Psi(\mathbf{R})|^2} \int d\mathbf{R} |\Psi(\mathbf{R})|^2 \left[\Psi^{-1}(\mathbf{R}) \widehat{H}_{\mathbf{R}} \Psi(\mathbf{R}) \right]$

$$\langle \hat{H} \rangle \approx \frac{1}{M} \sum_{\mathbf{R}} E_L(\mathbf{R}) \text{ with } E_L(\mathbf{R}) = \Psi^{-1}(\mathbf{R}) \hat{H}_{\mathbf{R}} \Psi(\mathbf{R})$$

Properties:

(i) $P(E_L(\mathbf{R}))$ approaches normal distribution (ii) In eigenstate $\widehat{H}\Psi = E_0\Psi$: $E_L(\mathbf{R}) = E_0$, $var(E_L) = 0$ (iii) $\lim_{M\to\infty} \frac{1}{M} \sum_{\mathbf{R}} E_L(\mathbf{R})$ larger than ground state energy



Natural gradient descent

The steepest descent direction $d\mathbf{w}$ of cost function $\mathcal{L}(\mathbf{w})$ minimizes $\mathcal{L}(\mathbf{w} + d\mathbf{w})$ for fixed infinitesimal distance $||d\mathbf{w}||^2 = \epsilon$ where ||...|| is a distance measure of distributions $\Psi_{\mathbf{w}}$.

Distance measure on Hilbert space:

$$\|d\mathbf{w}\|^2 = 1 - |\langle \Psi_{\mathbf{w}+d\mathbf{w}}|\Psi_{\mathbf{w}}\rangle|^2 = \sum_{ij} g_{ij}(\mathbf{w})dw_i dw_j + \mathcal{O}(d\mathbf{w}^4)$$

Steepest descent direction:

$$d\mathbf{w} = \frac{1}{2\lambda} g^{-1}(\mathbf{w}) \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w})$$

Equivalent to *imaginary time-evolution*.

 Ψ_{w_1}

In practice: (i) Approximate curvature $g^{-1}(\mathbf{w})$ (ii) Coarse-grain curvature $g(\mathbf{w}) \rightarrow g(\mathbf{w}) + \delta \mathbb{I}$ (iii) Use only $|\Psi_{\mathbf{w}}|$ for curvature estimate (iv) Stochastic gradient estimates

Results

Two-dimensional Coulomb gas with periodic potential

$$\widehat{H} = -\frac{1}{2} \sum_{i} V_{\mathbf{r}_{i}}^{2} + \sum_{i < j} \frac{1}{|\mathbf{r}_{i} - \mathbf{r}_{j}|}$$
$$+ V_{0} \sum_{j} \cos \mathbf{G}_{j} \mathbf{r}_{i}$$

Filling fraction:
$$v = \frac{2}{3}$$

Spin-polarized

Geier, Nazaryan, Zaklama, and Fu (2025)

Generalized Wigner crystal in WS₂/WSe₂ bilayers

Li et al., Nature 597, 650-654 (2021)

Convergence



Fermi liquid to generalized Wigner crystal transition



Benchmark with band-projected exact diagonalization



27 lattice sites

# sites ϵ		sε	self-attention NN	Hartree-Fock	BP-ED
	27	10	-32.070(7)	-31.35(2)	-31.32443
-	27	5	-59.127(9)	-58.01(3)	-57.80848

Scaling law

● 6e ● 8e

 $N_{\rm par}^{10^5}$

– 18e

106

Number of parameters required until convergence



Is attention all you need to solve the correlated electron problem?

Many-parameter neural network provides an ansatz without human bias, where correlations are described by self-attention, achieving quantitative accuracy.

Spontaneous symmetry breaking [1]

Topological order [2]

[1] M. Geier, K. Nazaryan, T. Zaklama, and L. Fu, arXiv:2502.05383 (2025)
[2] Y. Teng, D. Dai, and L. Fu, arXiv:2412.00618 (2024)

