

# Teaching Quantum Chemistry to a Deep Learning Model

Transformers for the many body Schrodinger Equation

Jorge Munoz Laredo, Angel Flores

January 7, 2026

- 1 The Schrödinger Wave Function and the physical laws that rule
  - Schrödinger Equation
  - Physical laws and conditions
  - Optimizing an Ansatz
- 2 Transformers
  - Attention Mechanism
- 3 Psiformer
  - Fermi Net
  - Psi Former
  - Practical Implementation Details

# The Schrödinger equation

On 1926 Schrodinger derived his equatin:

$$\hat{H} \Psi = E \Psi \quad (1)$$

- $\Psi$  is a complex value function called **wave function**.
- $\hat{H}$  is called the **Hamiltonian Operator**.

## Hamiltonian

$$\hat{H} = \frac{\hat{\vec{P}}^2}{2m} + \hat{V} = -\frac{\hbar^2}{2m} \nabla^2 + \hat{V} \quad (2)$$

- Find the electrostatic potential  $V$  of the system.

# Wave Function as Probability Density

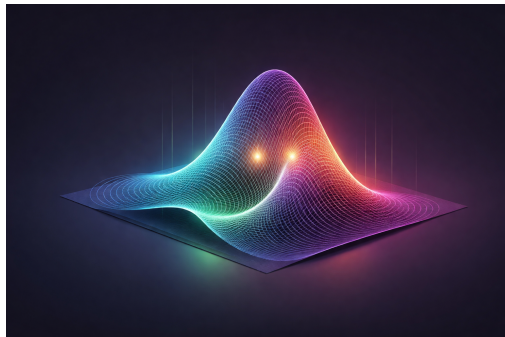


Figure 1:  $|\Psi(\mathbf{R})|^2$  represent the probability to find a particle near the position  $\mathbf{R}$ .

# Many-Body System

When considering  $n$  bodies, we have:

$$\hat{H}\psi(\mathbf{x}_1, \dots, \mathbf{x}_n) = E\psi(\mathbf{x}_1, \dots, \mathbf{x}_n) \quad (3)$$

With  $\mathbf{x}_i = \{\mathbf{r}_i, \sigma\}$ , where  $\mathbf{r}_i \in \mathbb{R}^3$  is the position of each particle and  $\sigma \in \{\uparrow, \downarrow\}$  is the so called spin.

## Considerations

- Each particle interact with all the another particles.
- For atoms, consider all the protons, electrons and neutrons.
- Solution obey physical laws.

# Setting up the Hamiltonian

The first step is obtain a practical form of the **Hamiltonian**.

- Kinetic energy:  $T = -\frac{1}{2} \sum_{i=1}^N \nabla_i^2$ .
- Electron-electron repulsion:  $V_{ee} = \sum_{i < j} \frac{1}{r_{ij}}$ .

$$\hat{H} = - \sum_{i=1}^N \frac{1}{2} \nabla_i^2 - \sum_{i=1}^N \frac{Z_I}{|\mathbf{r}_i - \mathbf{R}_I|} + \sum_{1 \leq i < j \leq N} \frac{1}{|\mathbf{r}_i - \mathbf{r}_j|} \quad (4)$$

# Fermi-Dirac statistics

All the fermions follow the Fermi-Dirac Statistics, this is.

- Exchanging two electrons flips the wavefunction's sign:  
 $\Psi(\dots i, j \dots) = -\Psi(\dots j, i \dots).$

## Slater Determinant

Enforce it using a determinant.

$$\psi = \begin{vmatrix} \phi_1^k(\mathbf{x}_1) & \dots & \phi_1^k(\mathbf{x}_n) \\ \vdots & & \vdots \\ \phi_n^k(\mathbf{x}_1) & \dots & \phi_n^k(\mathbf{x}_n) \end{vmatrix} \quad (5)$$

Where  $\phi$  are called spin orbitals

# Kato cusp conditions, Jastrow Factor

The potential are:

$$\sum \frac{1}{|\mathbf{r}_i - \mathbf{r}_j|}$$

- Coulomb potentials cause a sharp cusp in  $\Psi$  when particles overlaps.

## Jastrow Factor $\exp(\mathcal{J})$

In this work we are going to use this specific form:

$$\mathcal{J}_\theta(x) = \sum_{i < j; \sigma_i = \sigma_j} -\frac{1}{4} \frac{\alpha_{par}^2}{\alpha_{par} + |\mathbf{r}_i - \mathbf{r}_j|} + \sum_{i, j; \sigma_i \neq \sigma_j} -\frac{1}{2} \frac{\alpha_{anti}^2}{\alpha_{anti} + |\mathbf{r}_i - \mathbf{r}_j|} \quad (6)$$



# Loss: Variational Principle

Variational principle states:

$$E[\Psi] = \frac{\langle \Psi | H | \Psi \rangle}{\langle \Psi | \Psi \rangle} \geq E_0$$

Minimizing  $E[\Psi]$  drives the ansatz toward the ground state.

$$E[\Psi] = \mathcal{L}(\Psi_\theta) = \frac{\langle \Psi_\theta | \hat{H} | \Psi_\theta \rangle}{\langle \Psi_\theta | \Psi_\theta \rangle} = \frac{\int d\mathbf{R} \Psi^*(\mathbf{R}) \hat{H} \Psi(\mathbf{R})}{\int d\mathbf{R} \Psi^*(\mathbf{R}) \Psi(\mathbf{R})}$$

Define:

$$p_\theta(\mathbf{R}) = |\Psi_\theta(\mathbf{R})|^2 \frac{1}{\int d\mathbf{R}' \Psi_\theta^2(\mathbf{R}')} \wedge E_L(\mathbf{R}) = \frac{\hat{H} \Psi_\theta(\mathbf{R})}{\Psi_\theta(\mathbf{R})}$$

Then:

$$\mathcal{L}_\theta = \mathbb{E}_{\mathbf{R} \sim |\Psi_\theta|^2} [E_L(\mathbf{R})] \quad (7)$$

## Quantum Monte Carlo

With the samples  $\mathbf{R}_1, \dots, \mathbf{R}_M \sim |\Psi|_{\theta}^2(\mathbf{R})$  we can make:

$$\mathcal{L}_{\theta} = \mathbb{E}_{\mathbf{R} \sim \Psi_{\theta}^2}[E_L(\mathbf{R})] \approx \frac{1}{M} \sum_{i=1}^M E_L(\mathbf{R}_k) \quad (8)$$

With:

$$E_L(\mathbf{R}_k) = \frac{\hat{H}\psi(\mathbf{R}_k)}{\psi(\mathbf{R}_k)} = -\frac{1}{2} \frac{\nabla^2 \psi(\mathbf{R}_k)}{\psi(\mathbf{R}_k)} + V(\mathbf{R}_k)$$

Obtain  $\mathbf{R}_k \rightarrow$  Metropolis-Hastings Algorithm

# Metropolis-Hastings Algorithm

**Goal: Generate many samples  $\mathbf{R} \sim \rho$ , Requirement:  $C\rho$**

1.  $\mathbf{X}_0 \in E$  arbitrary:
2. Propose  $\mathbf{X}' = \mathbf{X}_0 + \eta$ , where  $\eta \sim q(\eta)$ , (Normal Gaussian)
3. Compute the quantity:

$$A(\mathbf{X}_0, \mathbf{X}') = \min \left( 1, \frac{\rho(\mathbf{X}')}{\rho(\mathbf{X}_0)} \right)$$

4. Generate a uniform number  $U \in [0, 1]$ . If:  $U < A(\mathbf{X}_0, \mathbf{X}')$  then  $\mathbf{X}_1 = \mathbf{X}'$ , otherwise try another  $\mathbf{X}'$ . Accept or decline.

# Metropolis-Hastings Algorithm

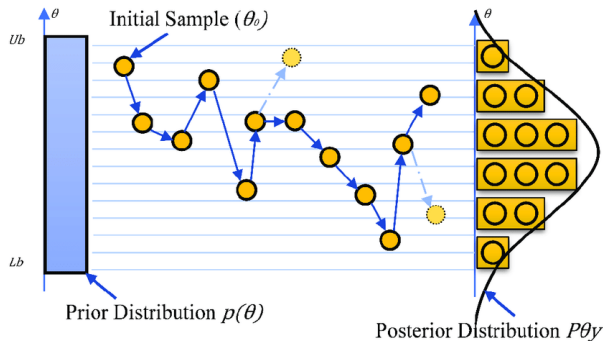


Figure 2: Metropolis-Hastings Walkers

- Obtain its own data (no dataset  $\mathcal{D}$ ),  $\rightarrow \mathcal{L}_\theta = \mathbb{E}_{\mathbf{R} \sim |\Psi_\theta|^2} [E_L(\mathbf{R})]$ .
- Third order derivatives
- $|\Psi_\theta|^2$  changes over time, you are just optimizing just the **energy**, not the wave function itself.

# Solution: Log Derivative Trick

$$\begin{aligned}\mathcal{L}_\theta &= \mathbb{E}_{\mathbf{R} \sim |\Psi_\theta|^2} [E_L(\mathbf{R})] \\ &\quad \downarrow \\ \nabla_\theta \mathcal{L} &= 2 \mathbb{E}_{\mathbf{R} \sim \Psi^2} [(E_L(\mathbf{R}) - \mathbb{E}_p[E_L]) \nabla_\theta \log \psi] \\ &\quad \uparrow \\ \mathcal{L}(\theta) &= 2 \mathbb{E}_{\mathbf{R} \sim \Psi^2} [\underbrace{(E_L(\mathbf{R}) - \mathbb{E}_p[E_L])}_{\text{detach}} \log \psi]\end{aligned}$$

REINFORCE

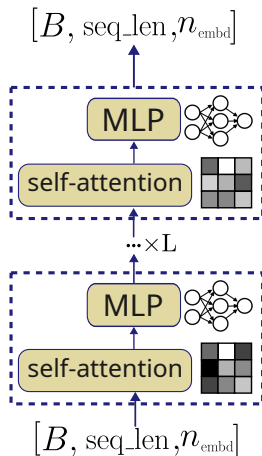
$$J(\theta) = \mathbb{E}_\tau [R(\tau)]$$

$\downarrow$

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi} [(R(\tau) - b) \nabla_\theta \log \pi_\theta(\tau)]$$

- 1 The Schrödinger Wave Function and the physical laws that rule
  - Schrödinger Equation
  - Physical laws and conditions
  - Optimizing an Ansatz
- 2 Transformers
  - Attention Mechanism
- 3 Psiformer
  - Fermi Net
  - Psi Former
  - Practical Implementation Details

# Transformer Architecture



## Multi Head Attention $\rightarrow$ Self Attention

- $n_{\text{embd}}$  the embedding dimension
- $n_h$  the number of attention heads
- $d_h$  the dimension per head
- $\mathbf{h}_t \in \mathbb{R}^{n_{\text{embd}}}$  the hidden dimension.

Figure 3: Tranformer backbone

# Attention on the room

The learnable matrices are:

$$W^Q, W^K, W^V \in \mathbb{R}^{n_{\text{embd}} \times n_{\text{embd}}}$$

$$\mathbf{k}_i = W^K \mathbf{h}_i, \mathbf{q}_i = W^Q \mathbf{h}_i, \mathbf{v}_i = W^V \mathbf{h}_i$$

$$[\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_{n_h}] = \mathbf{q}$$

$$[\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_{n_h}] = \mathbf{k}$$

$$[\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{n_h}] = \mathbf{v}$$

In the  $i$ -th head:

$$\mathbf{o}_{t,i} = \sum_{j=1}^t \text{Softmax} \left( \frac{\mathbf{q}_{t,i}^T \mathbf{k}_{j,i}}{\sqrt{d_h}} \right) \mathbf{v}_{j,i} \quad (9)$$

$W^O$  the output projection matrix.

$$\mathbf{u}_t = W^O[\mathbf{o}_{t,1}; \mathbf{o}_{t,2}; \dots; \mathbf{o}_{t,n_h}]$$



- 1 The Schrödinger Wave Function and the physical laws that rule
  - Schrödinger Equation
  - Physical laws and conditions
  - Optimizing an Ansatz
- 2 Transformers
  - Attention Mechanism
- 3 Psiformer
  - Fermi Net
  - Psi Former
  - Practical Implementation Details

**Ansatz:** Proposal model that you propose guided by intuition and that you optimize.

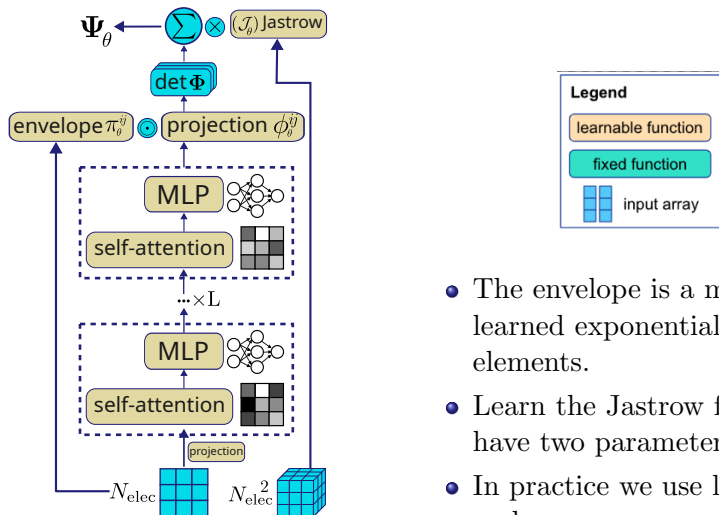
- be *antisymmetric* under particle exchange,
- capture strong *electron–electron correlations*.

## Proposed ansatz

Motivated by these constraints, we propose a Slater–Jastrow form:

$$\Psi_{\theta}(\mathbf{R}) = \underbrace{\exp(\mathcal{J}_{\theta}(\mathbf{R}))}_{\text{Coulomb correlations}} \times \underbrace{\sum \omega_k \det[\phi_{\theta}^k(\mathbf{R})]}_{\text{antisymmetry}}$$

# Psi Former Architecture



- The envelope is a matrix with learned exponential decay as elements.
- Learn the Jastrow factor only have two parameters.
- In practice we use logarithm scale.

Figure 4: Psi Former Architecture

# Psiformer Shapes

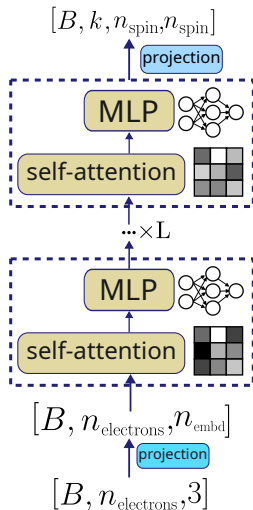


Figure 5: Psiformer shapes handling

Hyperparameter	Small	Large
Layers $L$	2	4
Heads $H$	4	8
Model Dim $d$	256	512
MLP Dim $d_{\text{ff}}$	1024	2048
Determinants $K$	1	2
MCMC walkers $N_w$	1024	2048
MCMC steps / iter	10	10
Learning rate	$2 \times 10^{-4}$	$1 \times 10^{-4}$
Total parameters	50441	3M

Figure 6: Psiformer Torch Small and Large

**Goal:** Obtain accurate ground state-energies using an Ansatz created with **Torch** library 

## Training Loop:

Set the spins and electron number

1. Sample configurations  $\mathbf{R}_k$
2. Estimate local energy  $E_L$
3. Backpropagation using **A.D**
4. Update parameters  $\theta$  using **AdamGrad**.
5. Track metrics with **Wandb**.

# Implementation : Laplacian Computation

The kinetic energy requires the Laplacian:  $\nabla^2\Psi(\mathbf{R}) = \sum_i \frac{\partial^2\Psi}{\partial R_i^2}$

```
R = R_o.requires_grad_(True)          # particle coordinates
psi = model(R)                         # neural wavefunction
# first derivative: gradient
grad_psi = torch.autograd.grad(
    psi, R,
    create_graph=True,
    retain_graph=True
)[0]
# second derivative: Laplacian
laplacian = 0.0
for i in range(R.shape[-1]):
    laplacian += torch.autograd.grad(
        grad_psi[..., i], R,
        create_graph=True,
        retain_graph=True
    )[0][..., i]
```

# Determinant Stability in Psiformer

$$\Psi_{\theta}(\mathbf{R}) \propto \det[\Phi_{\theta}(\mathbf{R})]$$

**Derivative of a determinant.**

$$\frac{\partial \det(\mathbf{A})}{\partial \mathbf{A}} = \det(\mathbf{A}) \mathbf{A}^{-T}$$

**Key instability.** If  $\Phi$  becomes singular or nearly singular:

$$\det \mathbf{A} \rightarrow 0 \quad \Rightarrow \quad [a_{ii}^{-1}] \rightarrow \infty$$

so the gradient can explode even when the wavefunction itself is small.

# Fix: Custom Operation

**Idea:** Use SVD ( $A = U\Sigma V^T$ ) to compute  $\nabla \log \det A = A^{-T}$  without explicit inversion (Appendix D).

```
import torch

class StableLogDet(torch.autograd.Function):
    @staticmethod
    def forward(ctx, A):
        # Decompose A: S are singular values
        U, S, Vh = torch.linalg.svd(A)
        ctx.save_for_backward(U, S, Vh)
        return S.log().sum()

    @staticmethod
    def backward(ctx, g):
        U, S, Vh = ctx.saved_tensors
        # Reconstruct  $A^{-T} = U * \text{diag}(1/S) * Vh$ 
        inv_S = torch.diag_embed(1.0 / S)
        grad_A = U @ inv_S @ Vh
        return g * grad_A
```



# Optimizer: Adam vs. AdamW

**Core Difference:** AdamW *decouples* weight decay from the gradient update to fix regularization on adaptive optimizers.

## 1. Adam (Entangled L2 Regularization)

- Decay is added to the gradient, so it gets scaled by the adaptive variance.

$$g_t = \nabla \mathcal{L} + \lambda \theta_t$$
$$\theta_{t+1} = \theta_t - \text{AdamStep}(g_t)$$

## 2. AdamW (Decoupled Weight Decay)

- Decay is applied directly, bypassing the adaptive scaling mechanism.

$$g_t = \nabla \mathcal{L}$$
$$\theta_{t+1} = \theta_t - \text{AdamStep}(g_t) - \eta \lambda \theta_t$$

# Keeping the GPU Busy: Batched Energy Evaluation

**Naive training issue.** Initial training evaluated energies step-by-step over MCMC samples.

GPU utilization  $\approx 30\%$

**Solution: batched evaluation.** MCMC samples are reshaped and flattened:

$$(\text{mc\_steps}, B, n_e, 3) \rightarrow (\text{mc\_steps} \times B, n_e, 3)$$

**Result.**

GPU utilization  $\approx 99\%$

# Results: Convergence Curve

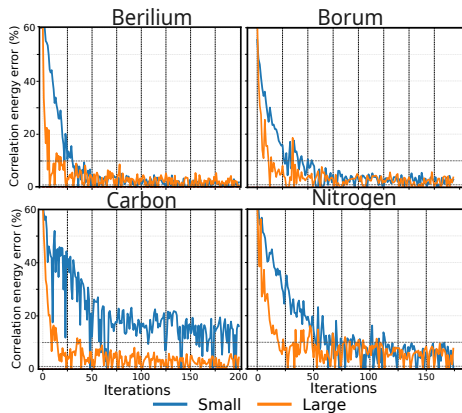


Figure 7: Convergence Curve

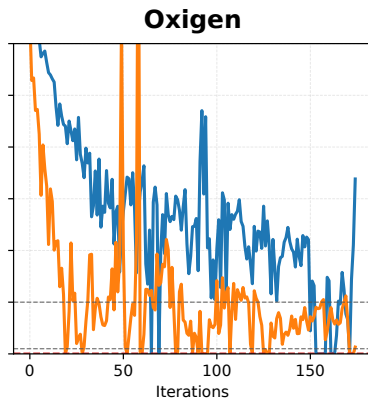


Figure 8: Oxygen Convergence

# Results: Energy Estimates

Atom	$E_b$	$E_s$	$E_l$	$\Delta_s$	$\Delta_l$	$\Delta_{l-s}$
H	-0.500	-0.492	<b>-0.498</b>	0.008	0.002	<b>-0.006</b>
He	-2.903	-2.801	<b>-2.893</b>	0.102	0.010	<b>-0.092</b>
Li	-7.478	-7.097	<b>-7.243</b>	0.381	0.235	<b>-0.146</b>
Be	-14.667	-13.901	<b>-14.237</b>	0.766	0.430	<b>-0.336</b>
B	-24.653	-24.042	<b>-24.567</b>	0.611	0.086	<b>-0.525</b>
C	-37.845	-35.492	<b>-36.457</b>	2.353	1.388	<b>-0.965</b>
N	-54.589	-50.492	<b>-51.700</b>	4.097	2.889	<b>-1.208</b>
O	-75.067	-63.492	<b>-72.139</b>	11.575	2.928	<b>-8.647</b>

Figure 9: Ground state energies (Ha) for H-O, baseline vs Psiformer

- Pretraining using external data.
- Laplacian Bottleneck
- KFCA Optimizer (Natural Gradient Descent)
- Flash Attention
- Learning transferability
- Scaling Laws



F. Hermann, Z. Schätzle, and F. Noé,  
*A Self Attention Ansatz for Ab Initio Quantum Chemistry*,  
arXiv:2309.12345, 2023.



J. Pfau, J. S. Spencer, A. G. D. G. Matthews, and W. M. C. Foulkes,  
*Ab Initio Solution of the Many-Electron Schrödinger Equation with Deep Neural Networks*,  
Physical Review Research, 2, 033429, 2020.



A. Vaswani et al.,  
*Attention Is All You Need*,  
Advances in Neural Information Processing Systems (NeurIPS), 2017.

# Thanks!

I thank the Computer Science Faculty for providing access to GPU resources, in particular NVIDIA GeForce RTX 4080 Super which enabled the training and evaluation of Psiformer models reported in this work.

