# L12 - Statistical Mechanics 4

# Neural networks

# What is a Neural Network?

# Plain vanilla
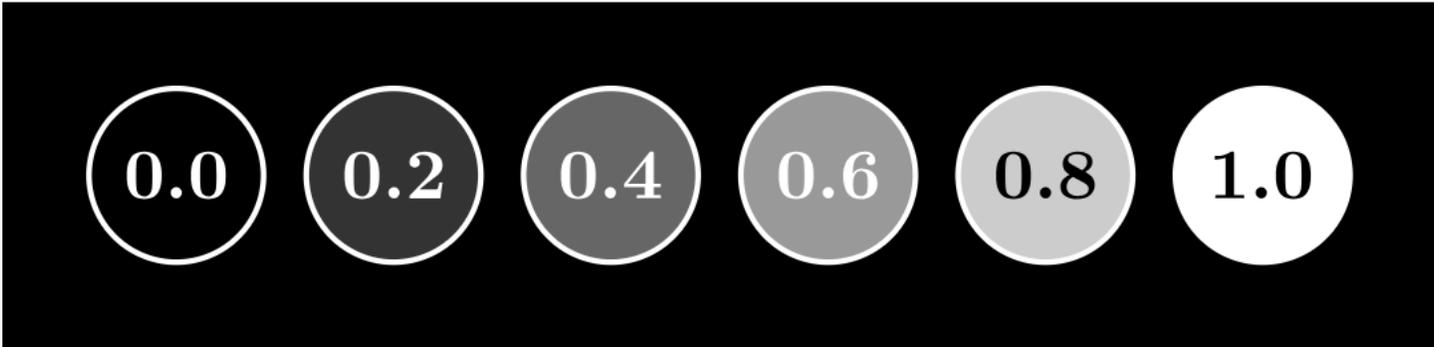## (aka "multilayer perceptron")

$$28 \times 28 = 784$$

$784 \{$

"Hidden layers"

784

$$\sigma\left(w_1 a_1 + w_2 a_2 + w_3 a_3 + \cdots + w_n a_n \boxed{-10}\right)$$

"bias"

Only activate meaningfully
when weighted sum > 10

$784 \times 16 + 16 \times 16 + 16 \times 10$
weights

$16 + 16 + 10$
biases

**13,002**

Superscript corresponds to the layer

$$a_0^{(1)} = \sigma\left(w_{0,0} a_0^{(0)} + w_{0,1} a_1^{(0)} + \cdots + w_{0,n} a_n^{(0)} + b_0\right)$$

Subscript corresponds to a neuron in the layer

Sigmoid

$$a_0^{(1)} = \sigma\left( w_{0,0}\ a_0^{(0)} + w_{0,1}\ a_1^{(0)} + \cdots + w_{0,n}\ a_n^{(0)} + b_0 \right)$$

Bias

Sigmoid

$$a_0^{(1)} = \sigma\left( w_{0,0}\, a_0^{(0)} + w_{0,1}\, a_1^{(0)} + \cdots + w_{0,n}\, a_n^{(0)} + b_0 \right)$$

Bias

$$\begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix}$$

Sigmoid

$$a_0^{(1)} = \sigma\left( w_{0,0}\, a_0^{(0)} + w_{0,1}\, a_1^{(0)} + \cdots + w_{0,n}\, a_n^{(0)} + b_0 \right)$$

Bias

$$\begin{bmatrix} w_{0,0} & w_{0,1} & \cdots & w_{0,n} \\ w_{1,0} & w_{1,1} & \cdots & w_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k,0} & w_{k,1} & \cdots & w_{k,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix}$$

$$\sigma\left(\begin{bmatrix} w_{0,0} & w_{0,1} & \cdots & w_{0,n} \\ w_{1,0} & w_{1,1} & \cdots & w_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k,0} & w_{k,1} & \cdots & w_{k,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_n \end{bmatrix}\right)$$

# Neural network function



784 inputs

13,002 weights
and biases

10 outputs

A NN with $L$ hidden layers takes an input vector **x** and returns an output vector **y** though the forward equations chain

$$\mathbf{a}^{(1)} = \sigma\left(\mathbf{w}^{(0)}\mathbf{x} + \mathbf{b}^{(0)}\right)$$

$$\mathbf{a}^{(2)} = \sigma\left(\mathbf{w}^{(1)}\mathbf{a}^{(1)} + \mathbf{b}^{(1)}\right)$$

$$\vdots$$

$$\mathbf{a}^{(L)} = \sigma\left(\mathbf{w}^{(L-1)}\mathbf{a}^{(L-1)} + \mathbf{b}^{(L-1)}\right)$$

$$\mathbf{y} = \sigma\left(\mathbf{w}^{(L)}\mathbf{a}^{(L)} + \mathbf{b}^{(L)}\right)$$

The NN is a function of **x** with a parametric dependence on **w** and **b**

$$\mathbf{y} = f\left(\mathbf{x}; \mathbf{w}, \mathbf{b}\right)$$

# The cost of learning

Adapted from:
www.3blue1brown.com/lessons/gradient-descent

(0, 0) (6, 6) (3, 3) (6, 6) (7, 7) (8, 8) (0, 0) (9, 9)

(5, 5) (4, 4) (3, 3) (6, 6) (5, 5) (8, 8) (9, 9) (5, 5)

(4, 4) (4, 4) (7, 7) (2, 2) (0, 0) (3, 3) (2, 2) (8, 8)

(9, 9) (1, 1) (9, 9) (2, 2) (2, 2) (7, 7) (9, 9) (4, 4)

(8, 8) (7, 7) (4, 4) (1, 1) (3, 3) (1, 1) (5, 5) (3, 3)

(2, 2) (3, 3) (9, 9) (0, 0) (9, 9) (9, 9) (1, 1) (5, 5)

(8, 8) (4, 4) (7, 7) (7, 7) (4, 4) (4, 4) (4, 4) (2, 2)

(0, 0) (7, 7) (2, 2) (4, 4) (8, 8) (2, 2) (6, 6) (9, 9)

(9, 9) (2, 2) (8, 8) (7, 7) (6, 6) (1, 1) (1, 1) (2, 2)

(3, 3) (9, 9) (1, 1) (6, 6) (5, 5) (1, 1) (1, 1) (0, 0)

Train on these

5 → 5
0 → 0
4 → 4
1 → 1
9 → 9
2 → 2
1 → 1
3 → 3
1 → 1
4 → 4

Test on these

3 → 3
5 → 5
3 → 3
6 → 6
1 → 1
7 → 7
2 → 2
8 → 8
6 → 6
9 → 9

Random initialization
of **w** and **b**

784

0
1
2
3
4
5
6
7
8
9

Utter trash

What's the "loss" of this difference?

784

Utter trash

What's the "loss"
of this difference?



Utter trash

Loss of 

$$\begin{cases} (0.43 - 0.00)^2+ \\ (0.28 - 0.00)^2+ \\ (0.19 - 0.00)^2+ \\ (0.88 - 1.00)^2+ \\ (0.72 - 0.00)^2+ \\ (0.01 - 0.00)^2+ \\ (0.64 - 0.00)^2+ \\ (0.86 - 0.00)^2+ \\ (0.99 - 0.00)^2+ \\ (0.63 - 0.00)^2 \end{cases}$$

# Cost function



13,002 weights
and biases

$(\boxed{9}, 9)(\boxed{0}, 0)(\boxed{2}, 2)(\boxed{6}, 6)$
$(\boxed{0}, 0)(\boxed{4}, 4)(\boxed{6}, 6)(\boxed{7}, 7)$
$(\boxed{7}, 7)(\boxed{8}, 8)(\boxed{3}, 3)(\boxed{1}, 1)$
$(\boxed{1}, 1)(\boxed{1}, 1)(\boxed{6}, 6)(\boxed{3}, 3)$
$(\boxed{1}, 1)(\boxed{1}, 1)(\boxed{0}, 0)(\boxed{4}, 4)$

Lots of training data

3.37

One number

For each training data $k$, with true value $\hat{\mathbf{y}}_k$
get the **loss function** over the output neurons

$$\mathcal{L}_k\left(\mathbf{w},\mathbf{b}\right) = \frac{1}{N_{out}} \sum_{i=1}^{N_{out}} \left( NN\left(x_{k,i};\mathbf{w},\mathbf{b}\right) - \hat{y}_{k,i} \right)^2$$

Considering all training data, get the **cost function**

$$C\left(\mathbf{w},\mathbf{b}\right) = \frac{1}{N_{train}} \sum_{k=1}^{N_{train}} \mathcal{L}_k\left(\mathbf{w},\mathbf{b}\right)$$

The optimized NN is the one with **w** and **b** that minimizes $C(\mathbf{w}, \mathbf{b})$

$$\min_{\mathbf{w},\mathbf{b}} C\left(\mathbf{w},\mathbf{b}\right) \Rightarrow \nabla_{\mathbf{w},\mathbf{b}} C\left(\mathbf{w},\mathbf{b}\right) = 0$$

# How to Train Your Neural Network

Adapted from:

$$\frac{dC}{dw}(w) = 0$$

$C(w)$

$w$

$w_{\min}$

$C(w)$

$w$

$w$

$C(w)$

$w$

$w$

Input space

$y$

4

3

2

Which direction decreases
$C(x, y)$ most quickly?

$x$

-3    -2    -1         1    2    3

-1

-2

-3

Compute $\nabla C$
Small step in $-\nabla C$ direction
Repeat.

13,002 weights and biases

How to nudge all
weights and biases

$$
\vec{\mathbf{W}} = \begin{bmatrix} 2.25 \\ -1.57 \\ 1.98 \\ \vdots \\ -1.16 \\ 3.82 \\ 1.21 \end{bmatrix}
$$

$$
-\nabla C(\vec{\mathbf{W}}) = \begin{bmatrix} 0.18 \\ 0.45 \\ -0.51 \\ \vdots \\ 0.40 \\ -0.32 \\ 0.82 \end{bmatrix}
$$

$$-\nabla C(\vec{\mathbf{W}}) = \begin{bmatrix} 0.31 \\ 0.03 \\ -1.25 \\ \vdots \\ 0.78 \\ -0.37 \\ 0.16 \end{bmatrix}$$

$w_0$  should increase somewhat

$w_1$  should increase a little

$w_2$  should decrease a lot

$w_{13,000}$  should increase a lot

$w_{13,001}$  should decrease somewhat

$w_{13,002}$  should increase a little

784

This weight matters a lot

0
1
2
3
4
5
6
7
8
9

Who even cares about this weight?

The minimization is done with **gradient descent**:

$$\mathbf{w}_{m+1} = \mathbf{w}_m - \gamma_m \frac{\partial C\left(\mathbf{w}_m, \mathbf{b}_m\right)}{\partial \mathbf{w}_m} \qquad \left(\gamma_m - \text{learning rate}\right)$$

$$\mathbf{b}_{m+1} = \mathbf{b}_m - \gamma_m \frac{\partial C\left(\mathbf{w}_m, \mathbf{b}_m\right)}{\partial \mathbf{b}_m}$$

Usually, the gradient is computed for a sub-set of **w** and **b** components chosen at random. It is called **stochastic gradient descent**

en.wikipedia.org/wiki/Gradient_descent

Now that we have optimal **w** and **b**, we can use the neural network to identify images that were not in the training set.

# This is the basics. But there is so much more...

- **Physics-informed neural networks (PINN)**
  Integrates partial differential equations expressing physical laws into the NN cost function
  youtu.be/G_hIppUWcsc

- **Graph neural networks (GNN)**
  NN for processing data that can be represented as vertices connected by edges
  distill.pub/2021/gnn-intro/

- **Convolutional neural networks (CNN)**
  NN with convolutional layers that capture local patterns and global features in the input data
  tinyurl.com/convnnet

- **Generative adversarial networks (GAN)**
  Two NN – a generator and a discriminator – compete to create realistic-looking outputs
  tinyurl.com/ganetintro

- **Transformer neural networks**
  NN architecture for encoding words, word position, and word's contextual relation with others in the sentence (self-attention).
  https://youtu.be/zxQyTK8quyY

# Machine learning in practice

# Machine learning development network

```python
import torch
import torch.nn as nn

class SimpleNeuralNetwork(nn.Module):
    def __init__(self):
        super(SimpleNeuralNetwork, self).__init__()
        self.layer1 = nn.Linear(784, 16)
        self.layer2 = nn.Linear(16, 16)
        self.output_layer = nn.Linear(16, 10)

    def forward(self, x):
        x = torch.sigmoid(self.layer1(x))
        x = torch.sigmoid(self.layer2(x))
        x = self.output_layer(x)
        return x

# Instantiate the model
model = SimpleNeuralNetwork()
```

**Deep Learning With PyTorch**

Codemy.com

19 videos  19,968 views  Last updated on Oct 16, 2023

▶ Play all          ⤬ Shuffle

In this playlist we'll learn A.I. Deep Learning with Pytorch and Python.  If you're new to Artificial Intelligence, this is for you!

tinyurl.com/pytorchlearn

# ML for Modeling Materials

AI for theoretical chemistry has been used to

- Search the chemical space of compounds

- Perform dimensionality reduction, clustering, and pattern recognition

- Improve or accelerate quantum chemical methods

- Predict properties as a surrogate approach

Search the chemical space of compounds

# GNN-based discovery of new materials

# GNN-based discovery added 381,000 new stable materials to the database

Perform dimensionality reduction, clustering, and pattern recognition

Cytosine

Hierarchical protocol for the automatic analysis of the ring deformation in surface hopping

Based on
- dimensionality reduction (PCA)
- clustering (DBSCAN + agglomerative clustering)

Zhu *et al.* PCCP **2022**, *24*, 24362

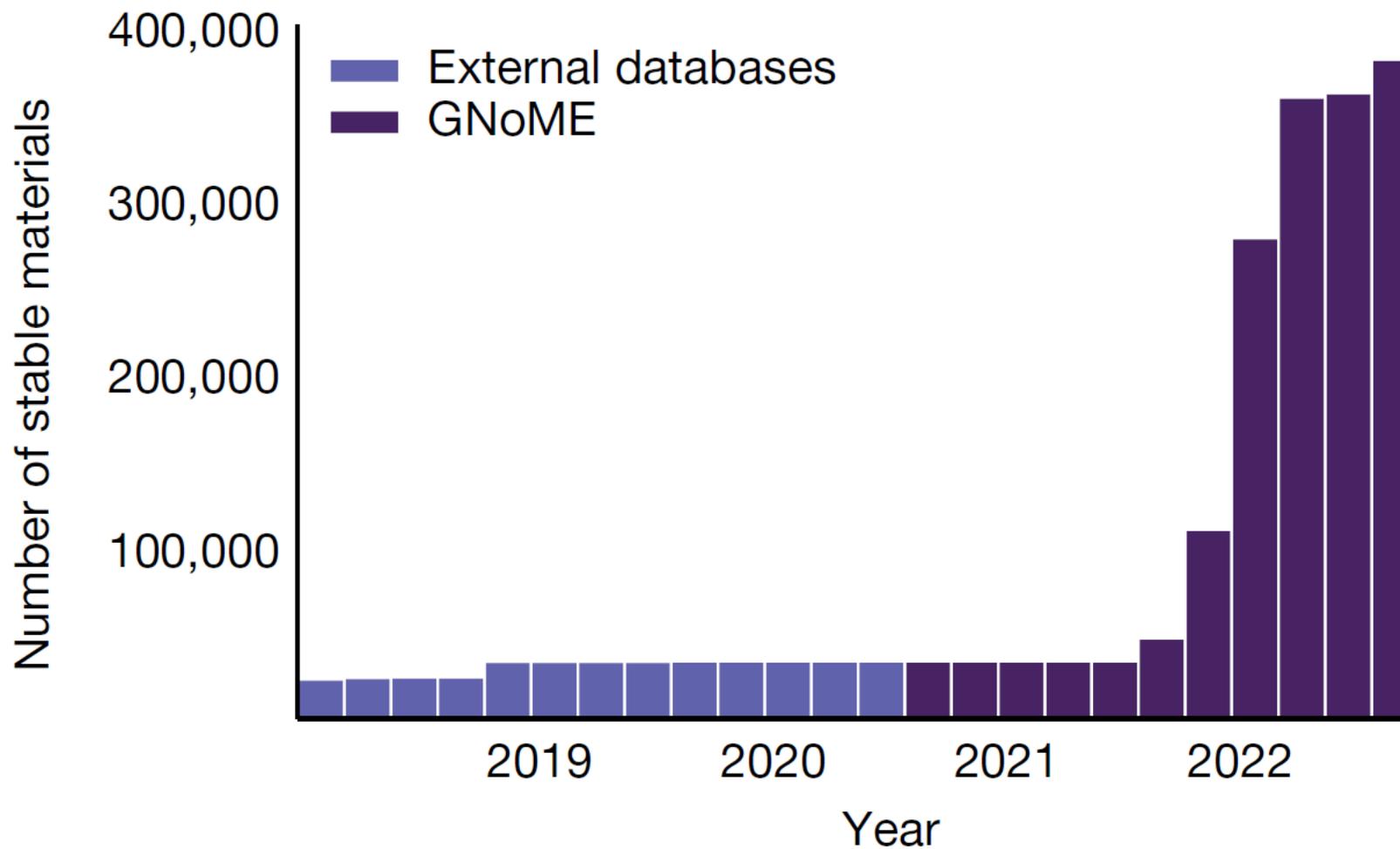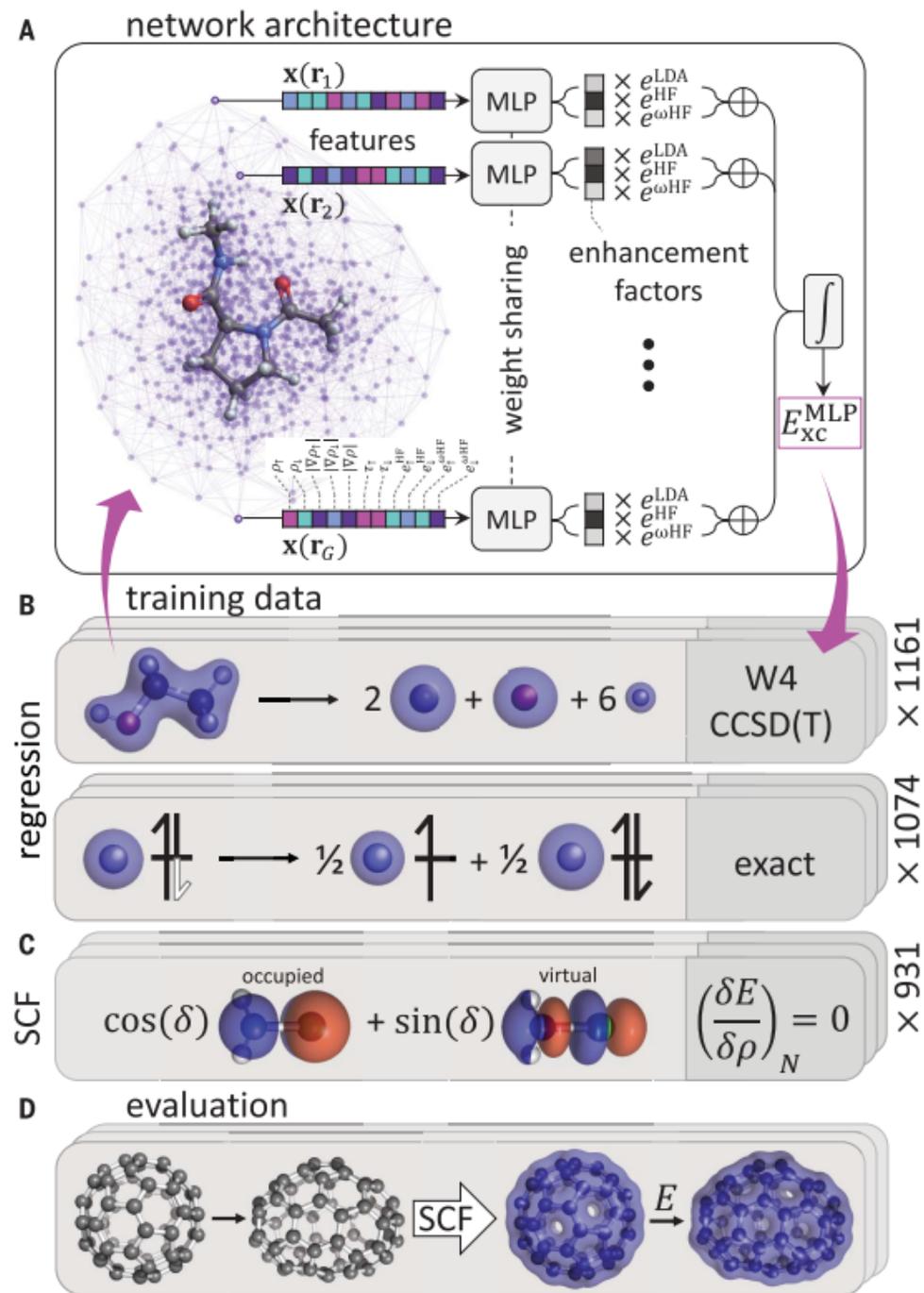| Channel | Important motion |
|---|---|
| **Cluster A1B1** 10.2% | C1-puckering C10-puckering C=O out-of-plane motion[−][a] |
| **Cluster A1B2C1** 56.2% | C1-puckering $NH_2$ out-of-plane motion[−][a] |
| **Cluster A1B2C2** 12.1% | C1-puckering |
| **Cluster A1B2C3** 6.3% | C1-puckering $NH_2$ out-of-plane motion[+][a] C1–N6 bond stretching |
| **Cluster A1B3** 7.7% | C1-puckering C=O out-of-plane motion[+][a] |
| **Cluster A2** 7.5% | C=O stretching motion |

# Improve or accelerate quantum chemical methods

# Density functional from an NN

Input features:
- charge density r,
- norm of charge density
- electron kinetic energy density
- local HF exchange energy densities

"The resulting functional, DM21 (DeepMind 21), correctly describes typical examples of artificial charge delocalization and strong correlation and performs better than traditional functionals on thorough benchmarks for main-group atoms and molecules. DM21 accurately models complex systems such as hydrogen chains, charged DNA base pairs, and diradical transition states."

Kirkpatrick *et al. Science* **2021**, *374*, 1385
Quanta magazine: tinyurl.com/qmdm21

# Predict properties as a surrogate approach:
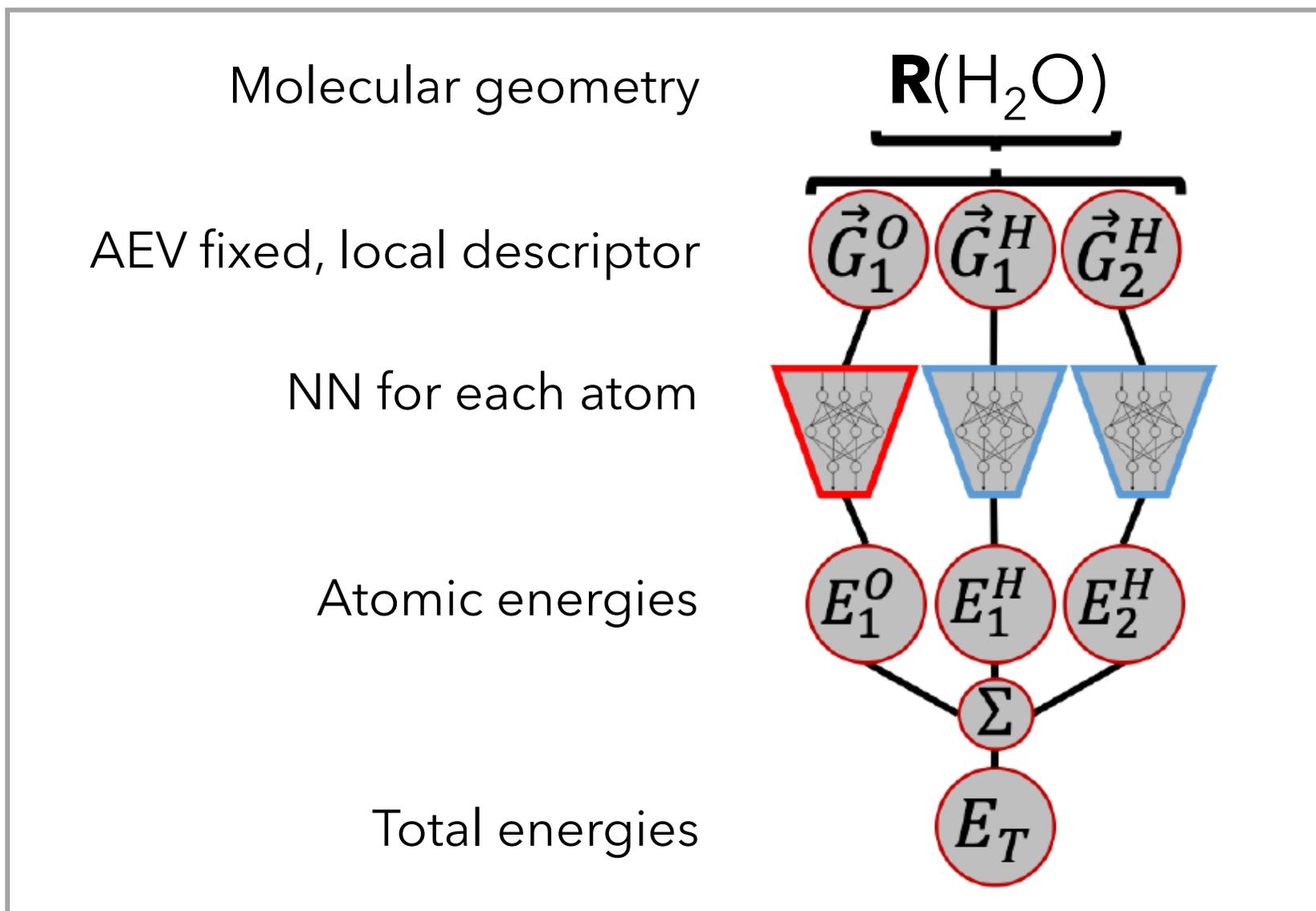# ML Potentials

$$\mathbf{R} \rightarrow \quad \Big(T_{elec}(\mathbf{r}) + V(\mathbf{r}, \mathbf{R})\Big)\varphi(\mathbf{r}; \mathbf{R}) = E(\mathbf{R})\varphi(\mathbf{r}; \mathbf{R}) \quad \rightarrow E(\mathbf{R})$$

$\mathbf{R} \rightarrow \mathbf{G}$    ML Algorithm    $\rightarrow E(\mathbf{R})$

Descriptor

Pinheiro Jr *et al.* Chem Sci **2021**, *12*, 14396

# Example: ANI ML Potential



Molecular geometry     $\mathbf{R}(H_2O)$

AEV fixed, local descriptor     $\vec{G}_1^O$   $\vec{G}_1^H$   $\vec{G}_2^H$

NN for each atom

Atomic energies     $E_1^O$   $E_1^H$   $E_2^H$

    $\Sigma$

Total energies     $E_T$

Gao *et al. J Chem Inf Model* **2020**, *60*, 3408

Legend: sGDML, GAP-SOAP, ANI, DPMD, PhysNet

y-axis: RMSE in energies [$kcal\ mol^{-1}$]

x-axis categories: benzene, ethanol, malonaldehyde, uracil, toluene, salicylic, paracetamol, naphthalene, azobenzene, aspirin

- MD17 Database
- Energy + Force
- $N_{train}$ = 1k; $N_{model}$ = 20; $N_{test}$ = 20k

Pinheiro Jr *et al*. *Chem Sci* **2021**, *12*, 14396

Legend: GAP-SOAP, KRR-CM, KREG, ANI, DPMD, PhysNet

Y-axis: RMSE on energy [$kcal\ mol^{-1}$]

X-axis: Training set size (1000, 10000, 50000)

- MD17 Database
- Energy only
- $N_{test}$ = 20k

Pinheiro Jr *et al*. *Chem Sci* **2021**, *12*, 14396

Pinheiro Jr *et al*. *Chem Sci* **2021**, *12*, 14396

## To know more:

3Blue1Brown Course on NN
- [www.3blue1brown.com/topics/neural-networks](www.3blue1brown.com/topics/neural-networks)

Kernel Methods
- Pinheiro Jr; Dral, In *Quantum chemistry in the age of machine learning,* **2023**; *pp 205*

ML Potentials
- Pinheiro Jr *et al*. *Chem Sci* **2021**, *12*, 14396

Deep Learning Applied to Computational Mechanics
- Vu-Quoc; Humer. *Comput Model Eng Sci* **2023**, *137*, 1069

# Computational modeling of nanosystems