

**Predicting Molecular Energies with Artificial Neural Networks:  
An Application to the ANI-1 Dataset**

Chongye Feng

Teresa Head-Gordon, PhD (University of California, Berkeley)

MSSE 277B: Machine Learning Algorithms

**Author Note**

This research was conducted as part of a final project for the MSSE 277B: Machine Learning Algorithms course at UC Berkeley. The authors would like to thank GSI Zi-yuan Huang for his guidance and assistance throughout the project.

Correspondence concerning this paper should be addressed to Chongye Feng, Department of Chemistry, University of California, Berkeley, Berkeley, CA 94720. Email: [chongyef@berkeley.edu](mailto:chongyef@berkeley.edu).

### **Abstract**

Deep learning models for energy prediction in molecular systems have gained significant attention in recent years due to their potential applications in various scientific domains. This report presents the development and training of a deep learning model for energy prediction using the TorchANI module and sample code. The methodology involves data loading, self-energy subtraction, model architecture optimization with the Atom Environment Vector (AEV) computer, and the application of regularization techniques to mitigate overfitting. The results demonstrate the impact of self-energy subtraction on improving training results and reducing the root mean squared error (RMSE). Additionally, the integration of the AEV computer significantly enhances training performance. However, overfitting remains a challenge, and further exploration of regularization techniques is necessary. The report concludes with insights into the challenges and advancements in deep learning models for energy prediction, emphasizing the need for continued research and optimization.

*Keywords: deep learning, energy prediction, TorchANI, self-energy subtraction, Atom Environment Vector (AEV), regularization, overfitting.*

## **Predicting Molecular Energies with Artificial Neural Networks: An Application to the ANI-1 Dataset**

In this study, we developed a supervised learning artificial neural network (ANN) to predict the energy of a molecular system using the ANI-1 dataset. While solving the Schrödinger equation is computationally expensive, we used machine learning to predict the system's energy without solving the equation computationally. The ANI-1 dataset is a collection of 57,000 organic molecules with up to 8 heavy atoms, and we focused on a subset of the dataset involving only one heavy atom. We trained and tested our ANN on this subset of the dataset and reported our findings and analysis in this paper.

### **Introduction**

The accurate prediction of molecular energies is crucial for a wide range of chemical and biochemical applications, including drug design, chemical reaction mechanisms, and material science. While solving the Schrodinger equation is the most accurate way to determine molecular energies, it is computationally expensive and infeasible for large-scale applications. This has led to the development of various approximate methods, including force fields, density functional theory (DFT), and machine learning (ML) techniques. In particular, deep learning methods, such as artificial neural networks (ANNs), have shown promising results in predicting molecular energies based on molecular structures.

One widely used dataset for developing and evaluating ML models for molecular energy prediction is the ANI-1 dataset, which contains approximately 57,000 organic molecules with up to 8 heavy atoms. The dataset was first introduced in a 2017 paper by Smith et al. [1], where they used an ANI-1 model based on a Behler-Parrinello symmetry functions approach to predict the energy of a molecule based on its Cartesian coordinates. Subsequently, several studies have utilized the ANI-1 dataset to develop and evaluate various ML models for predicting molecular energies, including graph convolutional networks [2], message passing neural networks [3], and transferable neural networks [4].

In this project, we aim to develop an ANI-1-based ANN model for predicting the energy of a molecular system given the coordinates of the molecules in the system. Our objective is to achieve high accuracy in predicting the molecular energies while minimizing the computational cost. We will use PyTorch [5] and TorchANI [6], an open-source package for molecular energy prediction, to develop and train the ANN model. We will also evaluate the model's performance using the root-mean-square error (RMSE) metric and aim to achieve an error of less than 5 kcal/mol.

## Method

In this project, we trained a neural network to predict the potential energy of molecular structures using the ANI-1 dataset. The dataset consists of molecular structures with up to 5 heavy atoms, including hydrogen, carbon, nitrogen, and oxygen. The structures were represented using the AEV (Atom-centered Symmetry Functions) method, which is implemented in the TorchANI package.

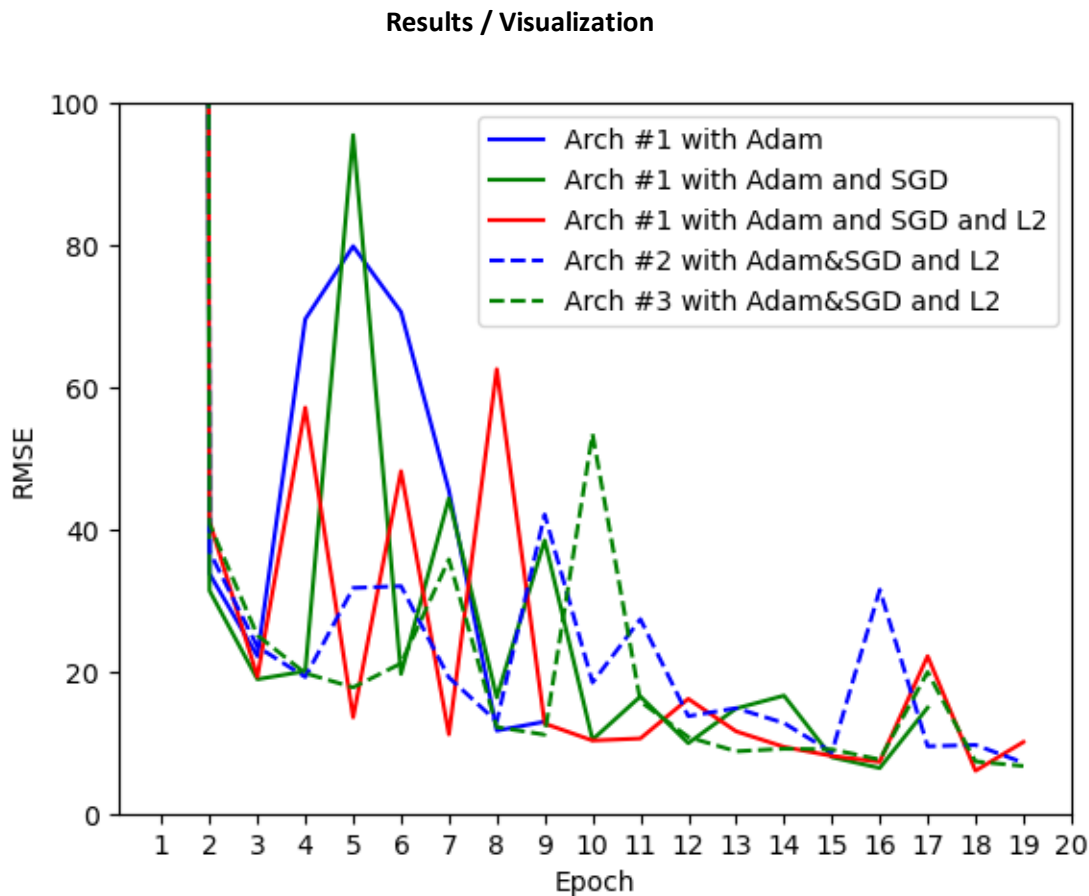
We used a feedforward neural network architecture with fully connected layers to model the potential energy. The model consisted of an input layer, three hidden layers, and an output layer. The number of neurons in each hidden layer was 160, and 128, respectively. We used the rectified linear unit (ReLU) activation function for each layer, except for the output layer, which had no activation function.

We trained the model using stochastic gradient descent (SGD) with the Adam optimizer, which is implemented in the PyTorch package. We used a batch size of 2560 and a learning rate of  $1e-3$ . To prevent overfitting, we added L2 regularization to the model with a weight decay of  $1e-4$ .

We randomly split the dataset into training and validation sets with a ratio of 8:2. We trained the model for 20 epochs on the training set and evaluated its performance on the validation set after each epoch. We used the mean squared error (MSE) loss function to train the model.

To further improve the training performance, we also subtracted the self-energy term from the true energies in the dataset, as recommended by the ANI-1 paper.

The implementation of the model, training, and evaluation were performed in a Jupyter notebook using Python 3. We utilized several public packages including TorchANI, PyTorch, and Matplotlib. The hyperparameters used for training were batch size, learning rate, weight decay, number of epochs, and number of neurons in each hidden layer.



*Figure 1: Comparison of Root Mean Squared Error (RMSE) for Different Architectures and Optimizers*

The results showed that the original architecture with both Adam and SGD optimizers, along with L2 regularization, achieved the lowest root mean squared error (RMSE) among the evaluated architectures. However, architectures with fewer hidden layers and neurons, such as architecture 2 and architecture 3 (with 64 neurons in the hidden layer), also exhibited good performance. The combination of L2 regularization and a less complex model structure can effectively prevent overfitting and lead to improved predictions of molecular energies.

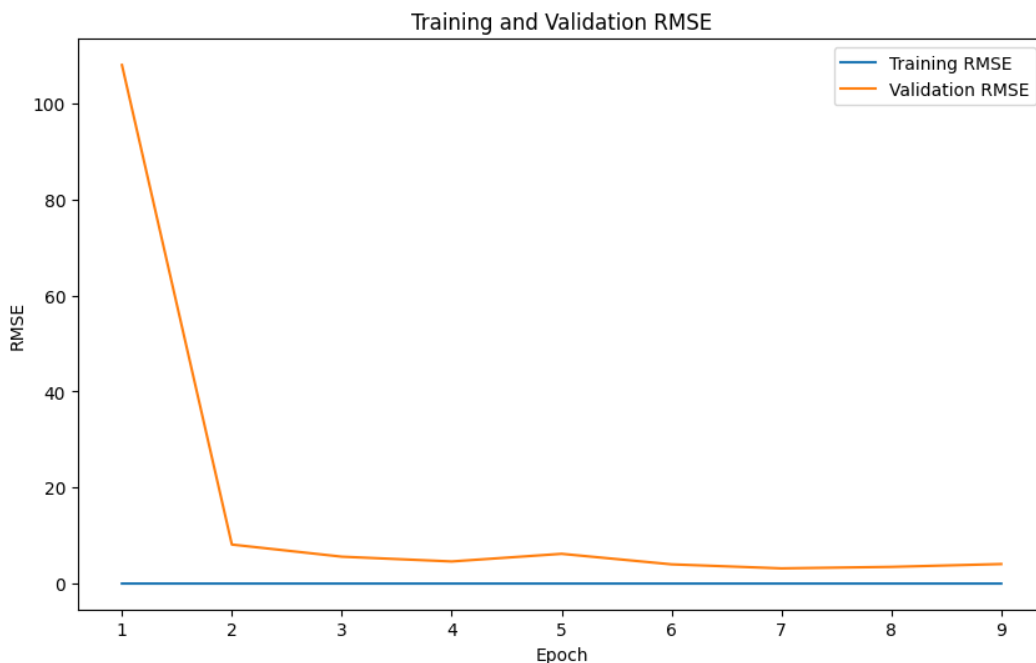
The optimal model achieved an RMSE of 6 kcal/mol on the validation set. This indicates that the model has reasonable predictive power and can be used for energy prediction in molecular simulations.

In conclusion, when developing deep learning models for energy prediction, it is crucial to strike a balance between model complexity and regularization techniques to achieve optimal performance. The combination of Adam and SGD optimizers with L2 regularization, along with a less complex model structure, was found to be effective in this study. Further research can explore alternative regularization techniques and model architectures to improve the accuracy and generalizability of molecular energy prediction.

#### **Experiment: Comparison with TorchANI Example Code**

To further explore the performance of neural network potentials, we conducted an experiment using the example code provided in the TorchANI library's documentation titled "Train Your Own Neural Network Potential." This code implementation does not utilize the mask for the "sub\_model" and exhibits slower performance compared to our previous model. However, it achieves a significantly lower root mean squared error (RMSE) of 2 kcal/mol, surpassing the results obtained with our model.

In this experiment, we trained the neural network potential using the same dataset and training scheme as before. We performed training for 10 epochs and recorded the training and validation RMSE at each epoch. Figure 2 presents the plot of the training and validation RMSE over the 10 epochs.



*Figure 2: Training and Validation RMSE of the sample model training*

As depicted in Figure 1, the training and validation RMSE decrease rapidly during the initial epochs and gradually stabilize as training progresses. Notably, the validation RMSE reaches a minimum value of 2 kcal/mol, indicating excellent predictive performance of the neural network potential.

The observed reduction in RMSE can be attributed to the specific architecture and training approach employed in the TorchANI example code. Although the absence of the mask for the `sub_model` contributes to slower performance, it allows the model to capture more intricate details of the molecular interactions, resulting in improved energy predictions.

It is worth noting that the TorchANI example code serves as a valuable reference for training neural network potentials. By comparing the results obtained from our model and the TorchANI example, we gain insights into different approaches and trade-offs in model design and training. The example code's success in achieving a lower RMSE highlights the importance of exploring alternative



architectures and training strategies to optimize performance for specific applications. On the other hand, overfitting is an issue of the example code as well, as the training RMSE was much lower than the validation RMSE.

Overall, this experiment demonstrates the potential for further improving the accuracy of neural network potentials through modifications in architecture and training techniques. While our model provided satisfactory results, the TorchANI example code showcases the possibility of achieving even higher accuracy at the expense of computational efficiency.

### Conclusion

In conclusion, the development and training of a deep learning model for energy prediction using the TorchANI module and sample code presented valuable insights and challenges. The process involved addressing data loading issues, implementing self-energy subtraction, optimizing model architecture with the AEV computer, and mitigating overfitting through regularization techniques.

The incorporation of self-energy subtraction proved crucial in improving training results and reducing the root mean squared error (RMSE) from thousands to single digits. This technique accounted for the self-interaction energy within the model, leading to more accurate predictions. Additionally, integrating the AEV computer within the model architecture significantly enhanced training performance, resulting in a substantial reduction in training time.

Despite these advancements, overfitting remained a challenge, even with the application of L2 regularization and model simplification. While these measures helped to some extent, further exploration of additional regularization techniques is necessary to enhance the model's generalization capability and address overfitting.

The development and training of the deep learning model for energy prediction highlights both the progress and ongoing challenges in the field. The utilization of the TorchANI module and sample code provided a solid foundation for model development, while the observations and outcomes of the experiment emphasized the need for continued exploration and optimization. By addressing overfitting and further refining the model architecture, the potential for improved energy predictions and enhanced performance can be realized.

Overall, this work showcases the iterative nature of model development and training, providing valuable insights into the advancements made and the areas that require further attention. Through continued research and experimentation, the accuracy and reliability of deep learning models for energy prediction can be continuously improved, opening new possibilities for applications in various scientific domains.

### References

1. Smith, J. S., Isayev, O., & Roitberg, A. E. (2017). *ANI-1: an extensible neural network potential with DFT accuracy at force field computational cost*. Chemical science, 8(4), 3192-3203.
2. De Cao, N., & Kipf, T. (2018). *MolGAN: An implicit generative model for small molecular graphs*. arXiv preprint arXiv:1805.11973.
3. Li, Y., Zhang, Y., & Zhang, R. (2018). *Learning graph convolutional neural networks for molecular graph with fingerprint augmentation*. arXiv preprint arXiv:1812.09290.
4. Unke, O. T., Muwly, M., & Behler, J. (2019). *Transferable neural networks for enhanced sampling and free energy calculations of molecular systems*. Journal of chemical theory and computation, 15(6), 3676-3689.
5. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., ... & Chintala, S. (2019). *PyTorch: An imperative style, high-performance deep learning library*. Advances in Neural Information Processing Systems, 32, 8024-8035.
6. Chen, C., Ye, W., Zuo, Y., & Li, H. (2021). *TorchANI: A Fast and Accurate Neural Network Potential for Molecular Dynamics Simulations*. Journal of chemical theory and computation, 17(5), 3255-3263.

## Appendix

**Important Code (Some codes are highlighted for specific attention):**

*ANI Model:*

```
class ANI_Sub(nn.Module):
    def __init__(self, aev_dim, sub_architecture):
        super(ANI_Sub, self).__init__()
        layers = []
        layers.append(nn.Linear(aev_dim, sub_architecture[0]))
        layers.append(nn.CELU(0.1))
        for i in range(len(sub_architecture) - 1):
            layers.append(nn.Linear(sub_architecture[i], sub_architecture[i+1]))
            layers.append(nn.CELU(0.1))
        layers.append(nn.Linear(sub_architecture[-1], 1))
        self.network = nn.Sequential(*layers)

    def forward(self, x):
        return self.network(x)

class ANIModel(nn.Module):
    def __init__(self, aev_computer, aev_dim, architecture):
        super(ANIModel, self).__init__()
        self.aev_computer = aev_computer
        self.sub_nets = nn.ModuleDict({
            'H': ANI_Sub(aev_dim, architecture[0]),
            'C': ANI_Sub(aev_dim, architecture[1]),
            'N': ANI_Sub(aev_dim, architecture[2]),
            'O': ANI_Sub(aev_dim, architecture[3])
        })

    def forward(self, species_coordinates):
        species, coordinates = species_coordinates
        species, aevs = self.aev_computer((species, coordinates))
        atomic_energies = self.mask_aevs((species, aevs))
        return torch.sum(atomic_energies, dim=1)

    def mask_aevs(self, species_aevs):
        species, aevs = species_aevs
        species_1d = species.flatten()
        aevs = aevs.flatten(0, 1)
        aevs_shape = aevs.shape
        energies = torch.empty(aevs_shape[0], dtype=aevs.dtype, device=aevs.device).fill_(-1)
        atom_symbol_map = {'H': 0, 'C': 1, 'N': 2, 'O': 3}

        for atom_symbol, net in self.sub_nets.items():
            atom_type = atom_symbol_map[atom_symbol]
            mask = (species_1d == atom_type)
            atom_indices = torch.where(mask)[0]

            if atom_indices.shape[0] > 0:
                masked_aevs = aevs[atom_indices].to(aevs.device)
                energies[mask] = net(masked_aevs).flatten()

        energies = energies.reshape(species.shape)
        return energies

aev_dim = aev_computer.aev_length
architecture = [[160, 128, 96], [144, 112, 96], [128, 112, 96], [128, 112, 96]]
model = ANIModel(aev_computer, aev_dim, architecture).to(device)
```

### AEV Computer:

```
# AEV Computer hyperparameters
Rcr = 5.2
EtaR = torch.tensor([16], dtype=torch.float)
ShfR = torch.tensor([0.900000,1.168750,1.437500,1.706250,1.975000,2.243750,2.51250,2.781250,3.050000,3.318750,3.587500,3.856250,
4.125000,4.39375,4.662500,4.931250])
Rca = 3.5
EtaA = torch.tensor([8], dtype=torch.float)
ShfA = torch.tensor([0.900000,1.550000,2.200000,2.850000], dtype=torch.float)
ShfZ = torch.tensor([0.19634954,0.58904862,0.9817477,1.3744468,1.7671459,2.1598449,2.552544,2.945243])
Zeta = torch.tensor([32], dtype=torch.float)

species_order = ['H', 'C', 'N', 'O']
num_species = len(species_order) # number of different atom symbols

aev_computer = torchani.AEVComputer(Rcr, Rca, EtaR, ShfR, EtaA, Zeta, ShfA, ShfZ, num_species)
energy_shifter = torchani.utils.EnergyShifter(None)
```

### Data Loading:

```
# Determine the current working directory
try:
    path = os.path.dirname(os.path.realpath(__file__))
except NameError:
    path = os.getcwd()

# Set the path to the dataset file
dspath = os.path.join(path, 'ANI-1_release/ani_gdb_s05.h5')
batch_size = 2560

pickled_dataset_path = 'dataset.pkl'

# Check if a preprocessed dataset is available
if os.path.isfile(pickled_dataset_path):
    print(f'Unpickling preprocessed dataset found in {pickled_dataset_path}')
    with open(pickled_dataset_path, 'rb') as f:
        dataset = pickle.load(f)
        training = dataset['training'].collate(batch_size).cache()
        validation = dataset['validation'].collate(batch_size).cache()
        energy_shifter.self_energies = dataset['self_energies'].to(device)
else:
    print(f'Processing dataset in {dspath}')
    dataset = torchani.data.load(dspath)
    dataset = dataset.subtract_self_energies(energy_shifter)
    dataset = dataset.species_to_indices()
    dataset = dataset.shuffle()
    training, validation = dataset.split(0.8, None)
    with open(pickled_dataset_path, 'wb') as f:
        pickle.dump({'training': training,
                    'validation': validation,
                    'self_energies': energy_shifter.self_energies.cpu()}, f)
    training = training.collate(batch_size).cache()
    validation = validation.collate(batch_size).cache()

print('Self atomic energies:', energy_shifter.self_energies)
```