

Guides for Writing your report

For the full document:

- you can use a html output but you need to keep the sectioning
- You are writing a comprehensive report
- The goal is to communicate to others
- Need to focus on communicating
- There should be no code or text not in the appropriate headings below
- Delete/comment out extra notes (like these) or others so your document is as clear as it can be
- Reference your figures, tables, equations in the document.
- Captions are needed on all figures, tables, and equations.
- Think about your audience and clearly communicating.

For the updates:

- please include the previous update for each Submitted Update Document
 - Each update is important to keep for grading
- For the final report remove updates and give a full comprehensive report
- Delete or comment out my notes that you are not using when you submit your documents
 - The extra clutter makes communication difficult

Update 2

- Please put a bulleted list of things you have accomplished since the last update
 - Include things that didn't work but you tried
 - Things you are planning on doing
 - Questions that you might have on your project.
- Reference the sections and figures you are dicussing here

Update 1

- Please put a bulleted list of things you have accomplished since the last update
 - Include things that didn't work but you tried
 - Things you are planning on doing
 - Questions that you might have on your project.
- Reference the sections and figures you are dicussing here

Excuetive Summary

- Summarize the key (This could be a bulleted list)
 - information about your data set
 - major data cleaning
 - findings from EDA

- Model output
- Overall conclusions

Abstract

The accurate prediction of molecular properties from structural information is essential for accelerating discovery in chemistry and materials science. While Density Functional Theory (DFT) provides reliable quantum mechanical predictions, its high computational cost limits its applicability in large-scale screening. In this study, we develop a neural network-based regression model to predict molecular properties—specifically total energy—directly from three-dimensional atomic coordinates. Using the **DFT_all.npz** dataset available from [Zenodo](#), which contains a variety of DFT-computed properties for small organic molecules, we train the model in a supervised manner to learn the structure–property relationship. Our results demonstrate that neural networks can effectively approximate DFT-level accuracy while significantly reducing computation time. This work highlights the potential of machine learning as a scalable alternative to traditional quantum chemical simulations, enabling faster exploration of chemical space for materials and drug design.

Introduction

Predicting molecular properties directly from structural information is a fundamental task in computational chemistry and materials science. Traditionally, this is achieved through quantum mechanical methods such as Density Functional Theory (DFT), which provide accurate predictions but are computationally expensive and limited in scalability. As the demand grows for rapid property evaluation in high-throughput screening and molecular design, data-driven alternatives have gained significant attention.

Recent advances in machine learning, particularly neural networks, have opened new pathways for modeling the complex relationship between a molecule’s structure and its physicochemical properties. These models can learn from large datasets of precomputed molecular structures and properties to make fast, accurate predictions without relying on costly simulations.

In this project, we focus on **predicting molecular properties from 3D molecular structures** using supervised learning with neural networks. We use the **DFT_all.npz** dataset, derived from DFT calculations and available through [Zenodo](#), which contains atomic coordinates and quantum-level properties for a variety of small organic molecules.

Our goal is to **train a neural network to accurately predict key molecular properties—such as total energy—from 3D atomic coordinates**, thereby capturing the structure–property relationship encoded in quantum mechanical simulations. This approach aims to demonstrate how machine learning models can serve as efficient surrogates for DFT, accelerating materials discovery and molecular design through predictive modeling.

Data Science Methods

We decided to realize a Property prediction of the elements from their 3D molecular structure. We will use Supervising training on Neural Network. To do so we will use the 3D properties of the molecules as training

inputs, then we will train the network with the use of the desired Property (Atomization Energy for example) as a label.

The input data should need no pretreatment. The output label should be a continuous value defining the property of the element.

The dataset is composed of 784875 element which is a quite huge amount of data (to compare, MNIST which is a basic digit recognition dataset contains 70000 elements) so the split between training subset and test subset should be relevant.

The main limit of this method is for each property we would like to predict, we would have to entirely redo the training with a different label.

Exploratory Data Analysis

Explanation of your data set

- How many variables?
 - 784875 data elements described by 26
- What are the data classes?
 - compounds : dtype = array
 - atoms : dtype = array
 - freqs: dtype = array
 - vibmodes: dtype = array
 - zpves: dtype = float64
 - U0: dtype = float64
 - U298: dtype = float64
 - H: dtype = float64
 - S: dtype = float64
 - G: dtype = float64
 - Cv: dtype = float64
 - Cp: dtype = float64
 - coordinates: dtype = array
 - Vesp: dtype = array
 - Qmulliken: dtype = array
 - dipole: dtype = array
 - quadrupole: dtype = array
 - octupole: dtype = array
 - hexadecapole: dtype = array
 - rots: dtype = array
 - gap: dtype = float64
 - Eee: dtype = float64
 - Exc: dtype = float64
 - Edisp: dtype = float64
 - Etot: dtype = float64
 - Eatomization: dtype = float64
- How many levels of factors for factor variables?

- atoms \backslash to\$ level \$10\$.
- compounds \backslash to\$ level \$784837\$.
- Is your data suitable for a project analysis?
 - Yes, we think. Sufficient variables are included in this dataset.
- Write your databook, defining variables, units and structures

◦	variables	units	discreption
	compounds		Stoichiometric formulas of the molecules
	atoms		Atomic numbers in the molecule
	freqs	$\text{\text{cm}}^{-1}$	Vibrational frequencies obtained from harmonic frequency analysis.
	vibmodes	$r\{A\}$	Normal modes of vibration represented as displacement vectors.
	U0	Ha	Internal energy at 0 K
	U298	Ha	Internal energy at 298 K
	H	Ha	Enthalpy
	S		Entropy
	G	Ha	Gibbs free energy
	Cv		Heat capacity at constant volume
	Cp		Heat capacity at constant pressure
	coordinates		coordinates (XYZ) of atoms in the molecule.
	Vesp		Electrostatic potential
	Qmulliken		Mulliken atomic charges
	dipole	a.u.	Dipole moment
	quadrupole	a.u.	Quadrupole moment
	octupole	a.u.	Octupole moment
	hexadecapole	a.u.	Hexadecapole moment
	rots	MHz	Rotational constants of the molecule.
	gaps	Ha	HOMO-LUMO energy gap
	Eee	Ha	Electron-electron repulsion energy
	Exc	Ha	Exchange-correlation energy
	Edisp	Ha	Dispersion correction energy
	Etot	Ha	Total electronic energy
	Eatomization	Ha	Atomization energy

Data Cleaning

- We performed data cleaning as follows:

```
# Filter fixed-length molecules (same n_atoms)
n_atoms_arr = np.array([len(a) for a in atoms])
most_common_n = Counter(n_atoms_arr).most_common(1)[0][0]
filtered_data = [
    (coord, atom) for coord, atom in zip(coordinates, atoms)
    if len(atom) == most_common_n
]
filtered_coords = [c for c, _ in filtered_data]
filtered_atoms = [a for _, a in filtered_data]
# Construct X as 3D input: [x, y, z, Z]
X = np.array([
    np.hstack([coord, atom.reshape(-1, 1)]) # shape: (n_atoms, 4)
    for coord, atom in zip(filtered_coords, filtered_atoms)
])
# Select targets
targets = ['U0', 'H', 'gap']
Y_all = np.column_stack([data[prop] for prop in targets])
Y = np.array([
    y for y, atom in zip(Y_all, atoms) if len(atom) == most_common_n
])
```

- Discription of this code
 - 1 Count the number of atoms in each molecule.
 - `python n_atoms_arr = np.array([len(a) for a in atoms])`
 - 2 Find the most common atom count among all molecules.
 - `python most_common_n = Counter(n_atoms_arr).most_common(1)[0][0]`
 - 3 Filter the data to keep only the molecules with that common atom count (when using a CNN, the input data must have consistent dimensions).

- ```
filtered_data = [
 (coord, atom) for coord, atom in zip(coordinates, atoms)
 if len(atom) == most_common_n
]
filtered_coords = [c for c, _ in filtered_data]
filtered_atoms = [a for _, a in filtered_data]
```

- 4 Build the input features X by stacking atomic coordinates `python[x, y, z]` with atomic numbers Z.

- ```
X = np.array([
    np.hstack([coord, atom.reshape(-1, 1)]) # shape: (n_atoms, 4)
    for coord, atom in zip(filtered_coords, filtered_atoms)
])
```

- 5 Define target values Y.

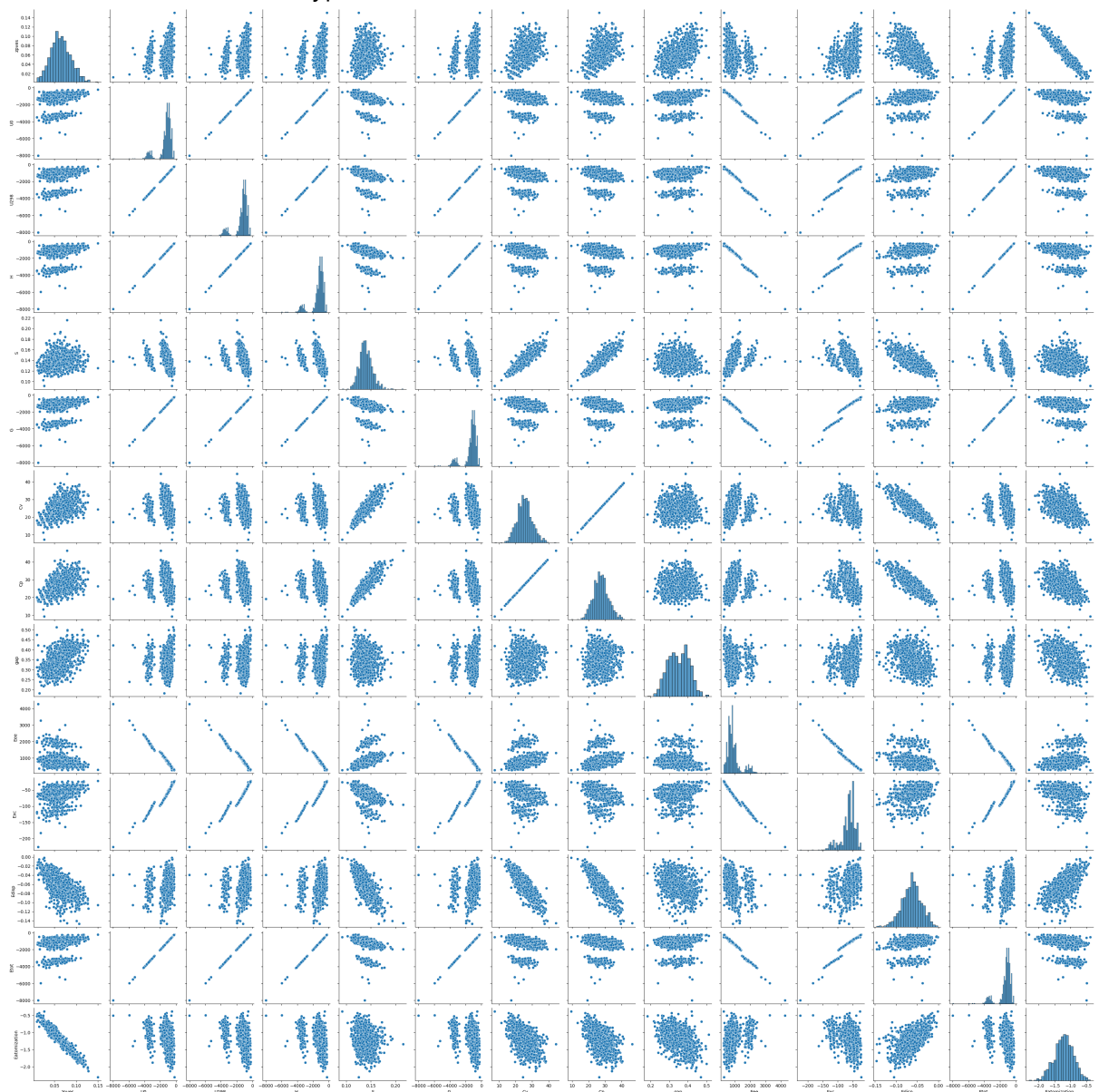
```

■ targets = ['U0', 'H', 'gap']
  Y_all = np.column_stack([data[prop] for prop in targets])
  Y = np.array([
    y for y, atom in zip(Y_all, atoms) if len(atom) == most_common_n
  ])

```

Data Visualizations

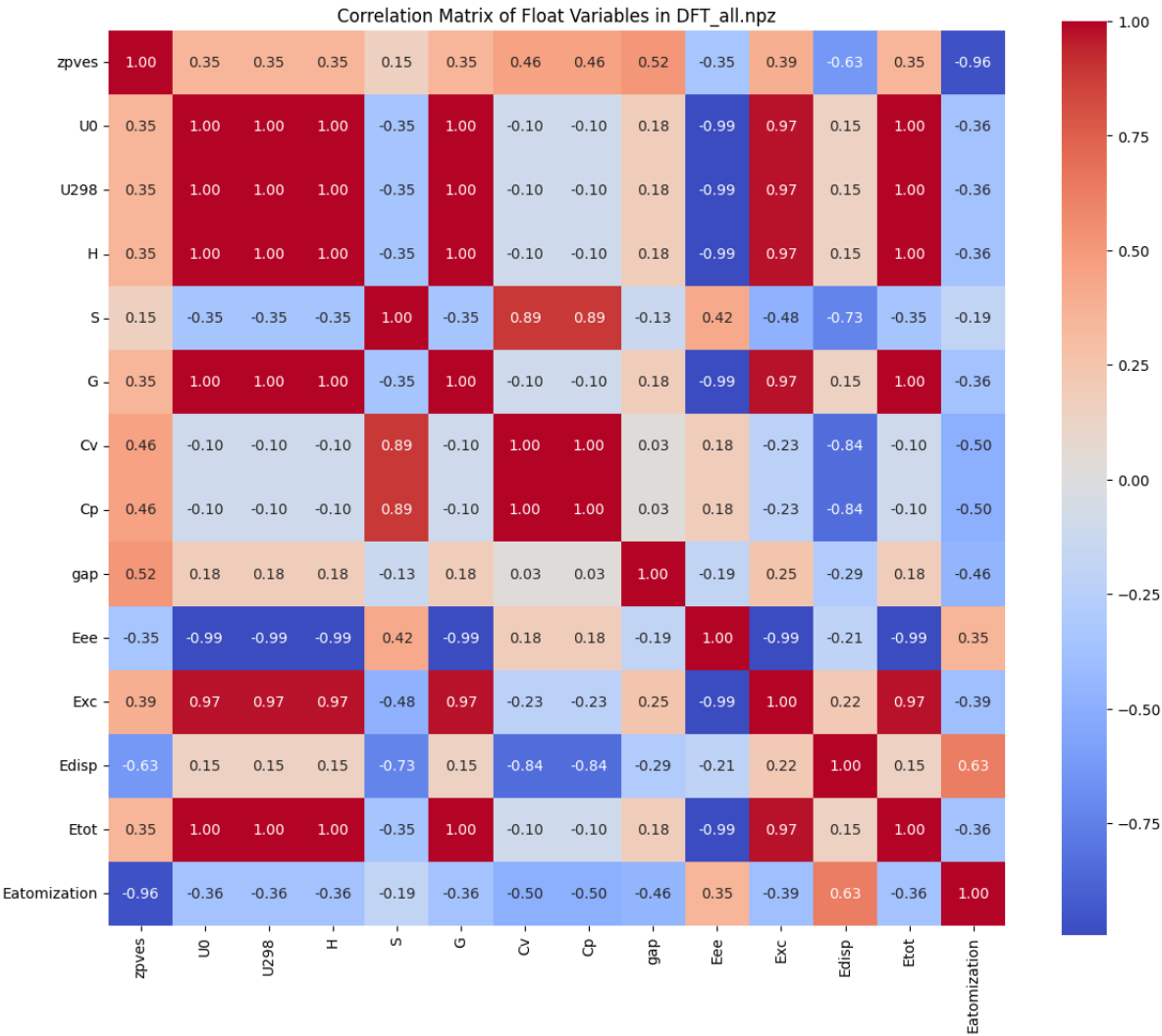
- Here is data visualizations (DFT_all.npz):
 - variable correlations (float type variables)



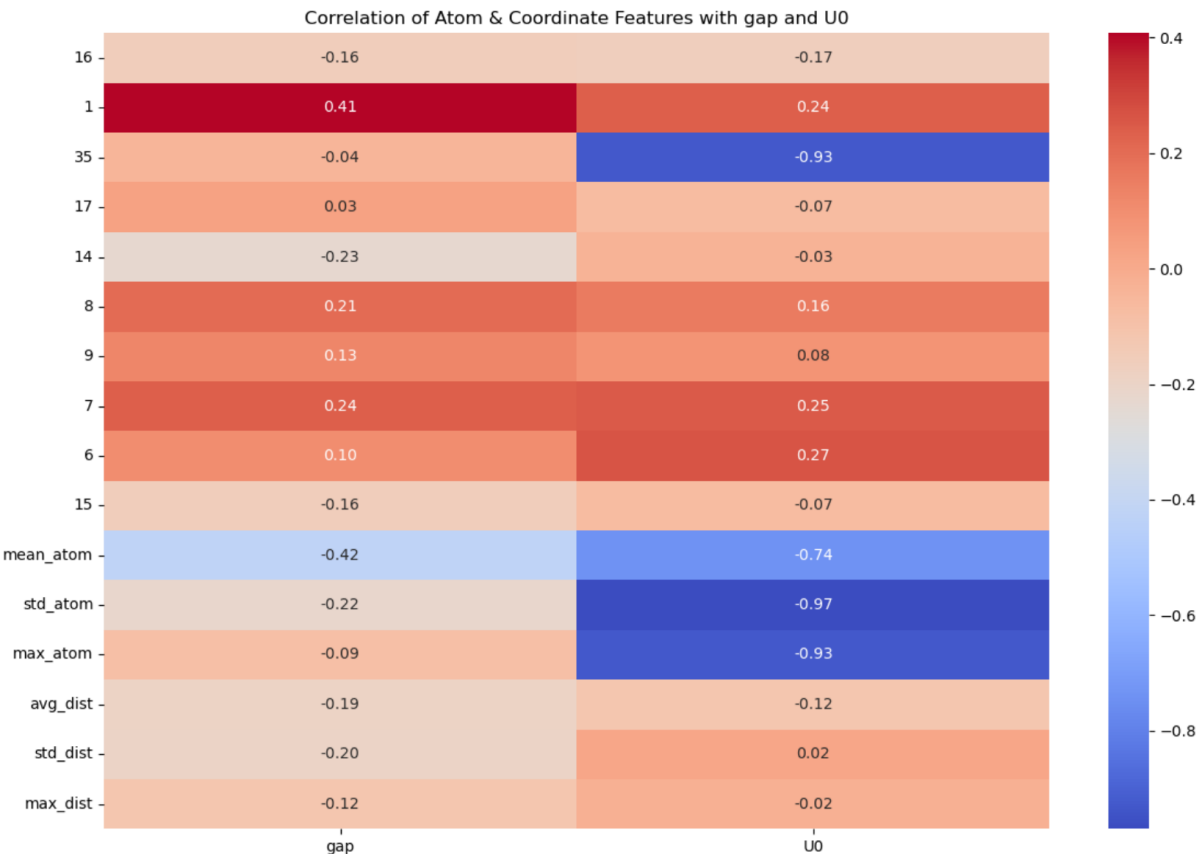
Variable Correlations

- Here is variable correlations (DFT_all.npz):

◦ variable correlations (float type variables)



◦ variable correlations (gap and U0)



- We can see that atoms can not explain about UO sufficiently. This suggests that the prediction accuracy for UO may not be very high.
- The model indicates that the number of hydrogen in the atoms leads to improved band gap prediction accuracy.

Statistical Learning: Modeling & Prediction

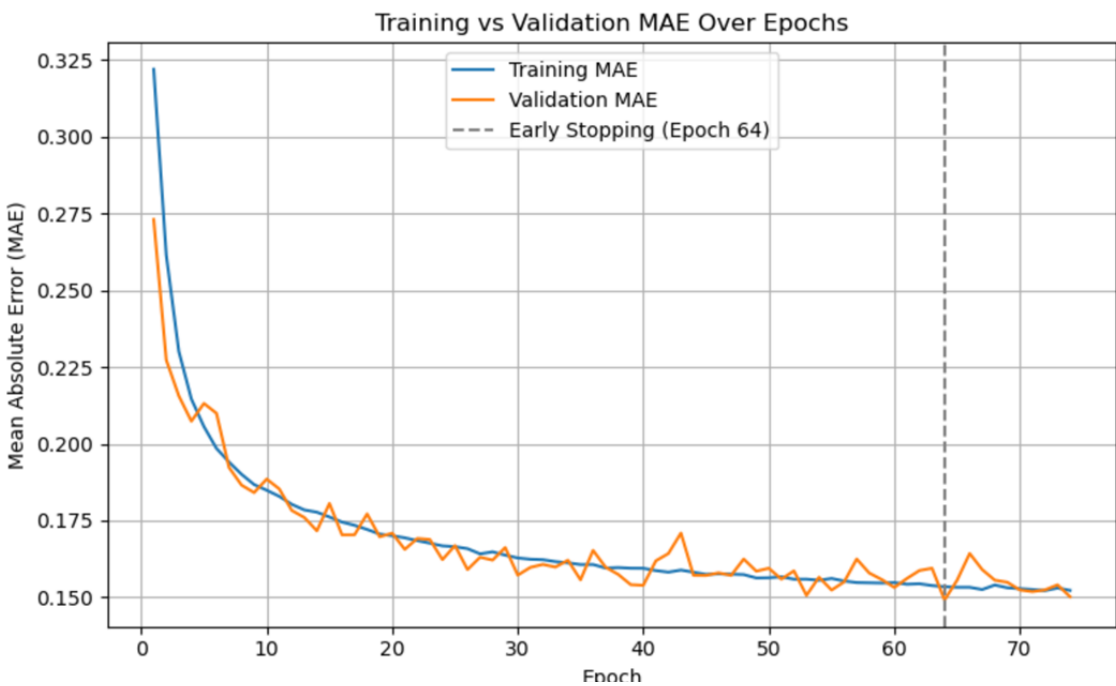
- Case 1 (CNN)

- Results

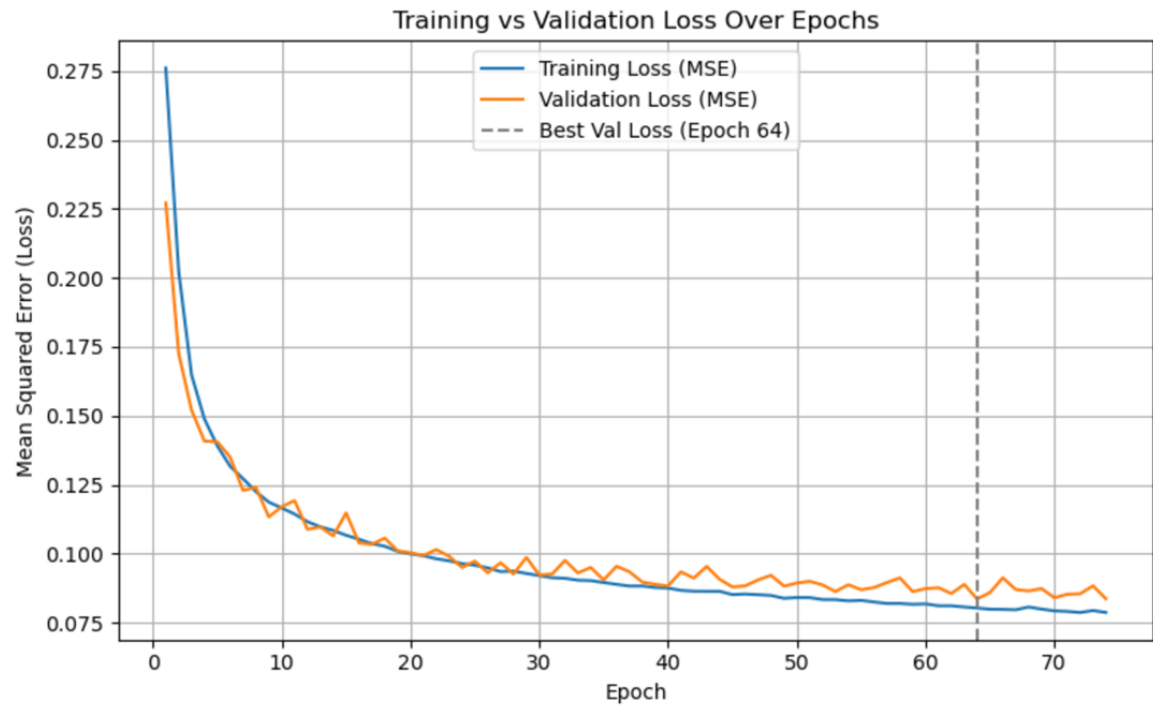
- Validation metric

Metric	Values
Test MAE	0.0852
Test Loss	0.1501
Best Validation MAE	0.1491 at Epoch 64
Best Validation Loss(MSE)	0.0836 at Epoch 64
Early Stopping Epoch	Triggered at 64

- Observation (Trainig vs Validation MAE)



- There is a steady decrease in both MAE.
- **Validatn MAE** stabilizes around **0.15**, which is very low.
- Both MAE curve indicate that there is no **overfitting**.
- The Validation curve flattens out after Epoch 64.



- Observation (Training vs Validation Loss)
 - There is a steady decrease in both Loss.
 - Both cyrves flatten around Epoch 60.
 - No visible sign of **overfitting**, because the validation loss closely follows the training loss.
 - Final MSE values shows a very good score.
- Summary

Insight	Implication
CNN achieved low MAE and MSE	The model captured property–structure relationships well
Properties are continuous, smooth	Ideal for regression with neural nets
Training & validation closely matched	Good generalization, minimal overfitting
Test MAE and loss match validation	Model performance is robust and consistent

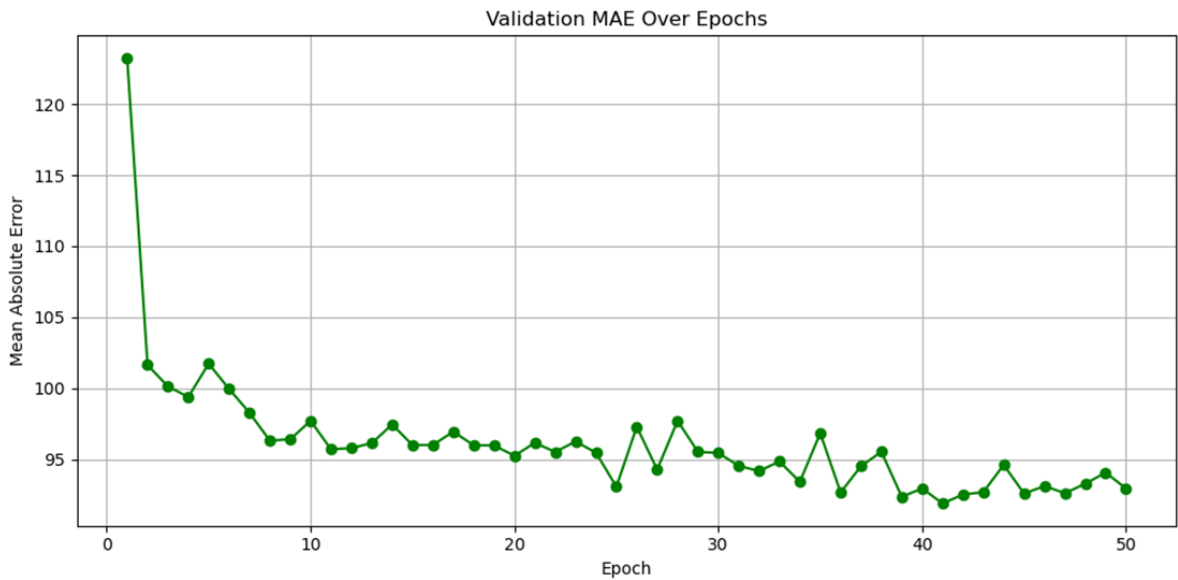
- Case 2 (GNN)

- Results
 - Validation metric

Metric	Values
Test MAE	92.0787

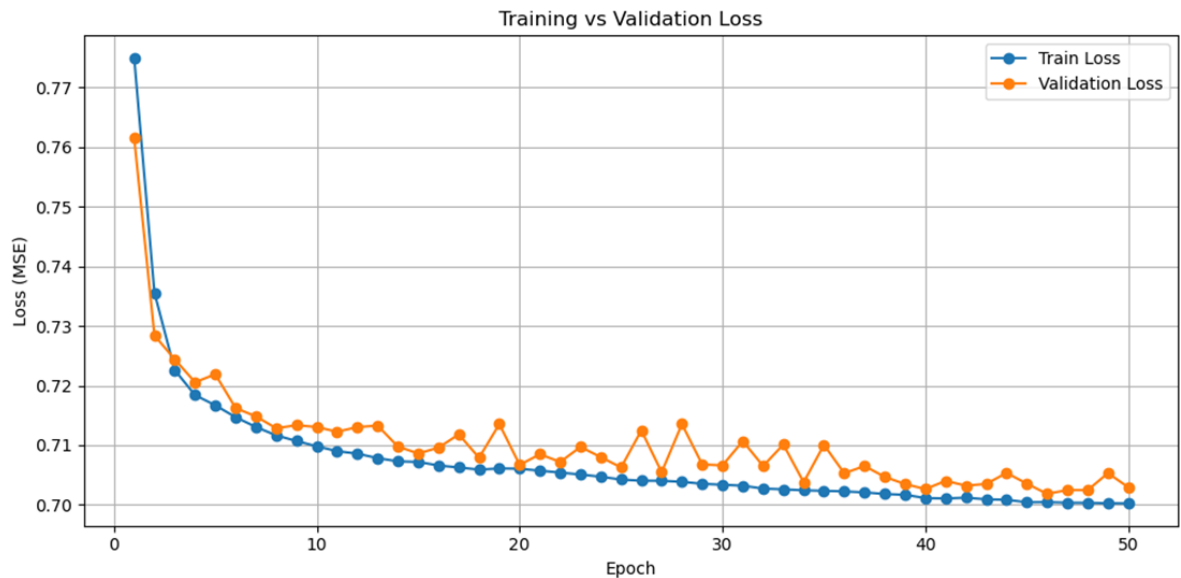
Metric	Values
Test Loss(MSE)	0.718738125
Validation MAE	92.9349441528
validation Loss(MSE)	0.7029326735
Metric	Values
U0	277.143677
gap	0.0369988717
H	277.143005
dip_x	1.40786136
dip_y	1.07787681
dip_z	0.806750417

○ Observation (Validation MAE over Epochs)



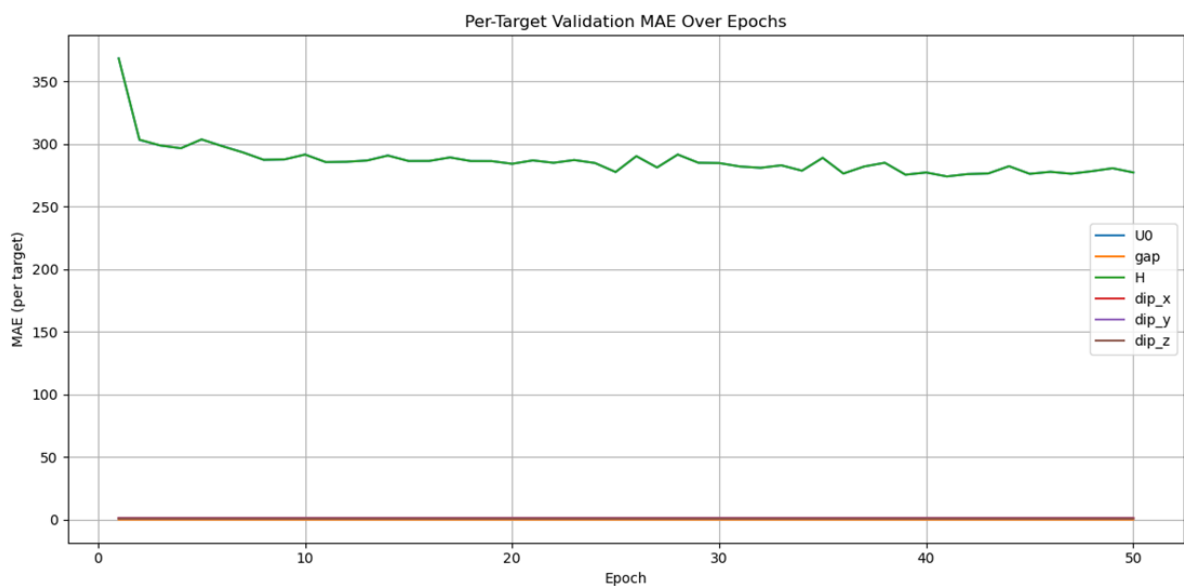
- By observing Epoch 1 and 2, we can see that this model learned **quickly** in training.
- The validation MAE stabilizes around epoch 10. This indicates that this model is **converging**.
- It can be seen that this model has **no overfitting**.
- The final MAE is around 93--94. This shows that this model achieved a relatively low and stable error across the validation set.

◦ Observation (Training vs Validation Loss)



- There is a steady decrease in both curve.
- The training loss remains **sufficiently stable** after epoch 20, whereas the validation loss **fluctuates** around epoch 30.
- **No indications of overfitting** are observed.
- The final MSE is higher than **Case 1**.

◦ Observation (Per-Target Validation MAE)



- The validation MAE for H is relatively high. This indicates that this model is **unable to make accurate predictions for H**.
- It can be seen that the predictions for dip_x, dip_y, and dip_z **does not work**.
- Gap has an order ten time shorter than its classical values, so the gap prediction seems to work.

Summary

Property	MAE Behavior	Implication
H (Enthalpy)	High MAE (~370 → ~275), dominates scale	Likely on a larger physical scale or has greater variability; normalization or separate treatment may help.
U0 (Internal Energy @ 0 K)	Very low MAE (not visible due to scale)	Well-learned; low-error, likely predictable based on molecular structure.
gap (HOMO-LUMO Gap)	Extremely low MAE (~0.037)	One of the best-learned properties — ideal for regression using GNNs.
dip_x, dip_y, dip_z (Dipole components)	Flat curves, MAEs ~0.8–1.4	These are vector quantities and the model predicts them accurately; stable across epochs.

Insight	Implication
GNN achieved low MAE/MSE for gap and dipole properties	Captures local and global molecular structure effectively
High MAE for H indicates scale mismatch	Target normalization or weighted loss is needed
Training vs Validation loss curves closely match	No overfitting; strong generalization
Per-property MAE stable after ~20 epochs	Model has converged; training was successful
Test and validation results are consistent	Model is robust and reliable for unseen data

Discussion

Aspect	CNN	GNN
Overall Accuracy	Superior MAE/MSE across metrics	Higher MAE due to unbalanced target scales
Generalization	Very consistent training/validation trends	Also consistent, no overfitting
Target-specific Learning	Likely trained on one property at a time	Some targets (e.g., H) harder to predict
Early Stopping	Used for optimal training	Not used
Interpretability & Chemistry	Less interpretable	Better modeling of molecular structure

Conclusions

Comparison of CNN and GNN

- CNN experience
 - **(Validation-Train)MAE=0.0639.**

- **(Validation-Train)Loss=0.0665.**
- This training model has an outstanding generalization error. Therefore, it is considered capable of handling new data effectively.
- GNN experience
 - **(Validation-Train)MAE=92.9349441528**
 - **(Validation-Train)Loss<2(per target)**
 - This model has a good generalization error. However, each value of MAE and Loss is higher than CNN experiment.

Future work

- The CNN modeling achieved high accuracy in learning target variables.
- GNN modeling successfully achieved high accuracy in learning a specific target variable (e.g. band gap). It can also be observed that the generalization error remains within a small range.
- By combining DFT and QM9 datasets, it may be possible to achieve even better results.
- To get an easier and more accurate precision, we should redo the GNN training with a normalization of the input.

Acknowledgments

We would like to express sincere thanks to the professor Pawan Tripathi for delivering such an engaging, informative lecture, and variable comments and remarks. Our project featured an in-depth and productive discussion.

References

- K. Danish, B. Anouar, K. Scott, G. F. von Rudorff, and A. von Lilienfeld : *Vector-QM24 (VQM24) dataset*, Zendo, <https://zenodo.org/records/15442257> (2025).