# BOOM: Benchmarking Out-Of-distribution Molecular Property Predictions of Machine Learning Models

**Evan R. Antoniuk[1*], Shehtab Zaman[2*], Tal Ben-Nun[1], Peggy Li[1], James Diffenderfer[1], Busra Demirci[2], Obadiah Smolenski[2], Tim Hsu[1], Anna M. Hiszpanski[1], Kenneth Chiu[2], Bhavya Kailkhura[1], and Brian Van Essen[1]**

[*]Co-first authors
[1]Lawrence Livermore National Laboratory
[2]Binghamton University, School of Computing

## ABSTRACT

Advances in deep learning and generative modeling have driven interest in data-driven molecule discovery pipelines, whereby machine learning (ML) models are used to filter and design novel molecules without requiring prohibitively expensive first-principles simulations. Although the discovery of novel molecules that extend the boundaries of known chemistry requires accurate out-of-distribution (OOD) predictions, ML models often struggle to generalize OOD. Furthermore, there are currently no systematic benchmarks for molecular OOD prediction tasks. We present BOOM, **b**enchmarks for **o**ut-**o**f-distribution **m**olecular property predictions—a benchmark study of property-based out-of-distribution models for common molecular property prediction models. We evaluate more than 140 combinations of models and property prediction tasks to benchmark deep learning models on their OOD performance. Overall, we do not find any existing models that achieve strong OOD generalization across all tasks: even the top performing model exhibited an average OOD error 3x larger than in-distribution. We find that deep learning models with high inductive bias can perform well on OOD tasks with simple, specific properties. Although chemical foundation models with transfer and in-context learning offer a promising solution for limited training data scenarios, we find that current foundation models do not show strong OOD extrapolation capabilities. We perform extensive ablation experiments to highlight how OOD performance is impacted by data generation, pre-training, hyperparameter optimization, model architecture, and molecular representation. We propose that developing ML models with strong OOD generalization is a new frontier challenge in chemical ML model development. This open-source benchmark will be made available on Github.

## 1 Introduction

Molecular discovery is the process by which novel molecular structures with desirable application-specific properties are identified. Given the immense total search space of hypothetical small molecules enumerating approximately $10^{60}$ hypothetical molecules and 100 billion enumerated molecules, molecule discovery has increasingly relied upon machine learning models to efficiently navigate this space.[1–3] Typically, molecule discovery is performed by first training a machine learning (ML) model on a molecule property dataset to learn the property-structure relationship. Then, this trained ML model is used to discover new molecules either by screening a list of enumerated molecules or by guiding a generative model towards molecules of interest [4].

Molecule discovery is inherently an out-of-distribution (OOD) prediction problem. For the discovered molecules to constitute an exciting chemical discovery, the molecules need to either i) exhibit properties that extrapolate beyond those of the known molecules in the training dataset, or ii) possess a new chemical substructure that was previously not considered for the application of interest. In either case, the success of the molecule discovery campaign is dependent on the machine learning model's ability to make accurate predictions on samples that do not follow the same distribution as the known molecules (training data).

Despite the importance of OOD performance to the problem of molecule discovery, the OOD performance of commonly used ML models for molecular property prediction has yet to be systematically explored. The majority of the standardized benchmarks used to assess the performance of chemical property prediction models do not include evaluations of model performance in the case where the test set is drawn from a different distribution than the training data. As a result of the lack of OOD chemistry benchmarks, the development of chemistry ML models are currently driven primarily by maximizing in-distribution performance, which may be hurting model generalization. The lack of OOD chemistry benchmarks has also hindered our understanding of how to develop generalizable chemistry foundation models. Currently, there is little empirical knowledge about how choices regarding the pretraining task, model architecture, and/or dataset diversity impact the generalization performance of chemistry foundation models that are expected to generalize across all chemical systems.

In this work, we develop BOOM, **b**enchmarks for **o**ut-**o**f-distribution **m**olecular property predictions, a standardized

benchmark for assessing the OOD generalization performance of molecule property prediction models.

## 1.1 Main Findings
Our work consists of the following main contributions:

- We develop a robust methodology for evaluating the performance of chemical property prediction models to extrapolate to property values beyond their training distribution. Notably, this methodology is developed in a general manner, allowing it to apply to any material property dataset regardless of the specific model architecture, material property, or chemical system.

- We perform the first large-scale benchmarking of the OOD performance of state-of-the-art ML chemical property prediction models. Across 10 diverse OOD tasks and 12 ML models, we do not find any existing models that show strong OOD generalization across all tasks. We therefore put forth BOOM OOD property prediction as a frontier challenge for chemical foundation models.

- Our work highlights insights into how pretraining strategies, model architecture, molecule representation, and data augmentation impact OOD performance. These findings point towards strategies for the chemistry community to achieve chemical foundation models with strong OOD generalization across all chemical systems.

## 2 BOOM Overview

In general, one can define OOD with respect to either the model inputs (holding out a region of chemical space as the OOD test split) or with respect to the model outputs (holding out a range of chemical property values). In this work, we adopt the latter approach of benchmarking the performance of the models to extrapolate to property values not seen in training. Following the OOD definitions outlined by Farquhar et al., we here define OOD as a complement distribution with respect to the targets [5,6]. Specifically, given a molecule property dataset of chemical structures and their numerical property values, we create our OOD test set to consist of numerical values on the tail ends of the numerical property distribution (see Figure 1). In this way, our OOD benchmarking is directly aligned with the molecule discovery task in that it allows us to evaluate the consistency of ML models to discover molecules with state-of-the-art properties that extrapolate beyond the training data.

### 2.1 Datasets
Overall, BOOM consists of 10 unique molecular property datasets. We collect 8 molecular property datasets from the QM9 Dataset: isotropic polarizability ($\alpha$), heat capacity ($C_v$), HOMO energy, LUMO energy, HOMO-LUMO gap, dipole moment ($\mu$) electronic spatial extent ($R^2$) and zero point vibrational energy (zpve).[7] Each property in the QM9 dataset was determined from Density Functional Theory (DFT) calculations on the same set of 133,886 small molecules made up of CHONF atoms. We also benchmark on the 10k Dataset, which consists of two molecular properties calculated by DFT: density and solid heat of formation.[8] The 10k Dataset was sourced from 10,206 experimentally synthesized CHON small molecules from the Cambridge Crystal Structure Dataset.
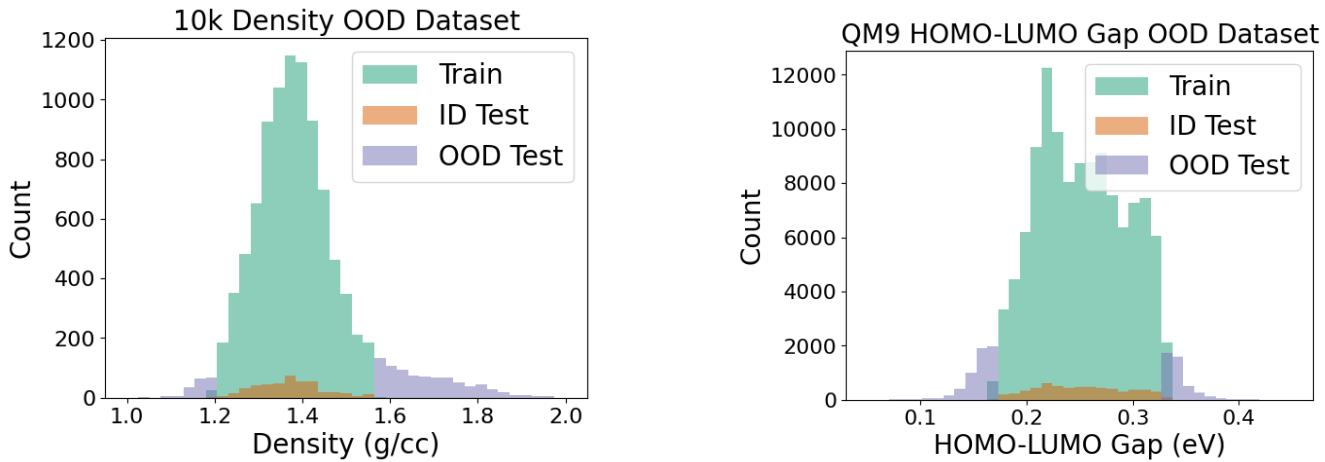
### 2.2 OOD Splitting
We generate our training, ID, and OOD splits based on the property distribution. For each of the 10 molecular properties, we generate OOD splits by first fitting a kernel density estimator (with Gaussian kernel) to the property values and obtain the probability of a molecule given its property. We select the molecules with the lowest probabilities for the OOD split for that property. This results in selecting the molecules at the tail end of the distribution for typical molecular property distributions since these molecules will have low densities in property space. Unlike partitioning by cut-off values, this method of splitting allows us to capture low-probability samples for general distributions that aren't necessarily unimodal.

For QM9 we take the lowest 10% of the probability scores as predicted by the kernel density estimator for the OOD set. We take the lowest 1000 molecules for the 10K dataset. We then randomly sample molecules from the remaining molecules to generate the ID test set. We sample 10% of the molecules in the case of QM9 and 5% of the molecules for ID test split for 10K data. The remaining molecules are used for training and fine-tuning. Two example OOD splits are provided in Figure 1, whereas the remaining datasets can be found in the Appendix 8.1.

### 2.3 Models
In this subsection, we describe the set of models we investigate for our OOD tasks. We use a plethora of traditional ML models, GNNs, and hybrid architectures to compare against large-scale transformer models.

**Figure 1.** Two example OOD datasets included in the BOOM benchmarking. To assess OOD performance, we split each chemical property dataset into an out-of-distribution (OOD) Test Set (blue), a in-distribution (ID) Test Set (orange) and a Train Set (green), as described in Section 2.2. (Left) 10k Dataset density OOD split and (Right) QM9 HOMO-LUMO gap OOD split.

### *Molecule Featurizer-based Models*

Traditional ML models utilize molecular fingerprints or other vector representations of molecules as input to statistical methods. As a baseline, we use RDKit Featurizer[9] coupled with a Random Forest regressor as the baseline structure-to-property model. The RDKit Featurizer, as implemented in the Deepchem package,[10] consists of 125 chemically-informed features (such as molecular weight and number of valence electrons), as well as 86 features describing the fraction of atoms that belong to notable functional groups such as alcohols or amines.

### *Transformers*

We begin our exploration of deep learning with a representative set of transformer models. Transformers, including large language models (LLMs), have revolutionized language modeling and vision tasks. We provide a brief discussion of transformer architectures in the Appendix.

We choose three representative models to cover the major archetypes of transformer models. Molformer [11] is an encoder-decoder model with a T5 [12] backbone originally trained on PubChem. ChemBERTa [13] is an encoder-only model with a BERT [14] backbone trained on PubChem. Finally, we also use Regression Transformer [15], an XLNet-based [16] model that is capable of both masked language modeling as well as autoregressive generation. We also evaluate ModernBERT, a state-of-the-art encoder-only model with architectural improvements such as rotary positional embeddings[17], pre-normalization, and GeGLU activation layers[18]. Along with different architectures, we also investigate the effects of different pre-training and tokenization schemes in our experiments. The training details are presented in Appendix 8.4.

| Model Name | Architecture | Molecule Representation | Symmetry | # of parameters |
|---|---|---|---|---|
| Random Forest | Random Forest | RDKit Molecular Descriptors | N/A | N/A |
| ChemBERTa | Transformer | SMILES | N/A | 83M |
| MolFormer | Transformer | SMILES | N/A | 48M |
| RT | Transformer | SMILES | N/A | 27M |
| ModernBERT | Transformer | SMILES | N/A | 111M |
| Chemprop | GNN | {Atom, Bond} | permutation | 200k |
| TGNN | GNN | {Atom, Bond} | permutation | 200k |
| IGNN | GNN | {Atom, Bond, Pair-wise Distances } | E(3)-invariant + permutation | 217K |
| EGNN | GNN | {Atom, Bond, Atom Positions} | E(3)-equivariant + permutation | 217K |
| MACE | GNN | {Atom, Bond, Pair-wise Distances} | E(3)-equivariant + permutation | 3.9M |
| Graphormer-3D | Hybrid | {Atom, Bond, Pair-wise Distances } | E(3)-invariant + permutation | 47.1M |
| ET | Hybrid | {Atom, Bond, Atom Positions} | E(3)-equivariant + permutation | 6.8M |

**Table 1.** Summary of the model architectures included in the BOOM benchmark, along with their model architecture, molecular representation, model symmetry, and total number of model parameters.

### *Graph Neural Networks*

GNNs are neural networks designed for learning on graph-structured data. Molecules and materials are represented as graphs of atoms and bonds, with 3D Euclidean space providing a natural molecular representation. As a result, message-passing neural networks (MPNNs) serve as the de facto backbone for deep learning-based molecular property prediction[19, 20]. Extensive work compares various GNN algorithms for this task. Instead of focusing on specific GNN variants, we examine the significance of architectural differences in our OOD task, emphasizing the relational inductive bias of molecular graphs and symmetries. (See Appendix 8.3).

3D information and symmetries are fundamental to physical laws governing molecular behavior. Chemprop[21] serves as the baseline for a standard topological (2D) GNN. Additionally, we use three GNNs with topological, E(3) invariant, and E(3) equivariant learned models based on EGNN[22]. MACE is a popular E(3) equivariant GNN, which uses pair-wise distances for message passing and construction[23, 24]. Unlike EGNN, MACE also takes into account higher order interactions, potentially allowing for greater expressivity. To explore the effects of these symmetries, we test these five GNNs for our OOD tasks.

### *Hybrid Architectures*

Recently, we have seen an emergence of hybrid architectures that combine the inductive properties of GNNs and with the flexibility of the attention mechanism in Transformers. Graphormer[25] is a GNN-Transformer model that incorporates a graph-specific encoding mechanism to the input perform attention over structured data rather than sequences. Furthermore, Graphormer adds a bias term to the Query-Key product matrix to bias the attention to include bond information. We evaluate Graphormer-3D, a variant of Graphormer that incorporates inter-atomic distances to introduce 3D information to the attention mechanism. Finally, we also evaluate Equivariant Transformer (ET)[26], a 3D encoder-only transformer model that incorporate E(3) equivariance. Rather than inter-atomic distances, ET operates directly on 3D atomic coordinates.

## 2.4 Metrics

We consider a model to generalize well to the OOD test set if the RMSE prediction performance is comparable to that of the ID test set. However, most of the models tested were unable to reliably predict the exact numerical values of the OOD samples. Short of achieving strong OOD generalization, the next-best case is for the model to achieve a strong correlation on the OOD samples. Even if the model cannot accurately predict the correct numerical property values, achieving a strong correlation between the true and predicted property values is still useful for molecule discovery since the top-performing molecules will still be correctly identified. Altogether, we evaluate the OOD and ID RMSE prediction performance as a primary metric for measuring OOD generalization. We also evaluate the correlation on the OOD samples by calculating a *binned $R^2$* value, which is the average value $R^2$ of the OOD samples in the lower and upper tails of the property distribution.

## 3 Related Work

Modern ML and DL models are designed for little discrepancy between training and test data distribution. OOD predictions present a key challenge for incorporating data-driven models into production pipelines where test time input may significantly shift from training data [27]. A possible solution to the OOD problem is to detect the distribution shift of input data away from the training distribution. OOD detection algorithms involve post-hoc or embedded updates to prediction models to detect when a sample is out of distribution. OOD detection has been approached through the lens of anomaly detection, uncertainty quantification, and open-set detection[6, 28, 29].

Rather than detecting distribution shifts, improving model robustness to unknown distribution shifts also addresses challenges posed by out-of-distribution data. Broadly speaking, the techniques for improving model robustness to distribution shifts can be divided into: (1) improved representation learning of features, (2) improved mapping of features to labels, and (3) improved optimization techniques [30].

In the domain of inorganic crystalline materials, Matbench [31] presents benchmarking on 13 tasks including the optical, thermal, electronic, thermodynamic, tensile, and elastic properties of inorganic materials. The Matbench tasks include 10 regression and 3 classification tasks based on the material composition and crystal structures. Follow-up studies utilized the Matbench datasets to evaluate the OOD performance of graph neural networks in both property and chemical structure space[32]. MatFold provides a convenient and systematic method to generate OOD test splits with respect to the structure of the material.[33] MatFold currently supports the creation of these OOD splits by structure, composition, chemical system, element, periodic group, periodic row, space group, point group, or crystal system. Our work differs from MatFold/Matbench in that i) we focus on OOD generalization in the property (y) space, instead of the input (x) space and ii) BOOM evaluates small molecule properties instead of inorganic crystalline materials.

In the small molecule space, the MoleculeNet is widely used for benchmarking molecular property prediction models. MoleculeNet consists of 17 small molecule prediction tasks, along with four splitting protocols: random, scaffold splitting,

stratified splitting, and time splitting (test set consists of newest data). Segal et. al[34] also explore zero-shot extrapolation of molecular and material property prediction beyond the training data. Similar to our work, they also define OOD samples in the property space. They present bilinear transduction as a technique to enhance extrapolation capabilities. They evaluate on ESOL, Freesolv, Lipophilicity, and BACE binding from MoleculeNet for molecules. Bilinear transduction requires a descriptor-based model for anchor-based regression, while we focus on benchmarking existing deep-learning solutions.

# 4 Results

## 4.1 Survey of Model Architectures

| Model | Split | HoF | Density | HOMO | LUMO | GAP | ZPVE | $R^2$ | $\alpha$ | $\mu$ | $C_v$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Random Forest | ID | 24.3453 | 0.0251 | 0.0061 | 0.007 | 0.0084 | 0.0012 | 52.0985 | 0.9441 | 0.6709 | 0.3785 |
| | OOD | 139.0728 | 0.1813 | 0.0304 | 0.0372 | 0.037 | 0.023 | 364.1557 | 8.4777 | 2.8977 | 3.3779 |
| RT | ID | 28.1941 | 0.0224 | 0.0112 | 0.0135 | 0.0196 | 0.0034 | 95.8337 | 3.2265 | 1.2805 | 0.9291 |
| | OOD | 496.6896 | 19403.8 | 54.6121 | 1.6136 | 0.0367 | 83.0197 | 374.4082 | 183239.6 | 2.9445 | 5.9453 |
| ChemBERTa | ID | 22.86 | 0.0163 | 0.0068 | 0.0088 | 0.0103 | 0.0046 | 50.297 | 1.444 | 0.7134 | 0.4923 |
| | OOD | 99.7253 | 0.1173 | 0.0245 | 0.0267 | 0.0315 | 0.0214 | 306.14 | 6.303 | 2.766 | 3.0175 |
| MolFormer | ID | 9.3608 | 0.0175 | 0.0054 | 0.0054 | 0.0065 | 0.0015 | 40.9528 | 1.0544 | 0.6145 | 0.0826 |
| | OOD | 93.6099 | 0.0924 | 0.0239 | 0.0248 | 0.0279 | 0.0213 | 283.0316 | 7.1359 | 2.2245 | 2.9328 |
| Chemprop | ID | 15.68 | 0.0092 | 0.0041 | 0.0048 | 0.0058 | 0.0018 | 35.68 | 0.8305 | 0.55 | 0.3341 |
| | OOD | 100.6 | 0.0551 | 0.0192 | 0.0187 | 0.0267 | 0.0129 | 234.73 | 4.7729 | 2.3 | 2.149 |
| EGNN | ID | 9.3667 | 0.0076 | 0.0044 | 0.0044 | 0.0062 | 0.0013 | 20.2683 | 0.5655 | 0.4772 | 0.2504 |
| | OOD | 17.7773 | 0.0296 | 0.0203 | 0.0223 | 0.029 | 0.0077 | 184.7985 | 5.4988 | 2.3559 | 2.1071 |
| IGNN | ID | 19.7635 | 0.0086 | 0.0048 | 0.0051 | 0.0069 | 0.0015 | 30.5926 | 0.8374 | 0.5036 | 0.4018 |
| | OOD | 24.5891 | 0.0261 | 0.0202 | 0.0195 | 0.0297 | 0.0075 | 184.1348 | 5.4229 | 2.4993 | 2.2091 |
| TGNN | ID | 14.3601 | 0.0157 | 0.0054 | 0.0076 | 0.0372 | 0.0013 | 211.5326 | 0.7452 | 0.6371 | 0.3449 |
| | OOD | 27.7216 | 0.0145 | 0.0178 | 0.0203 | 0.0129 | 0.0021 | 622.6762 | 2.8217 | 2.5002 | 0.5545 |
| MACE | ID | 5.6737 | 0.0692 | 0.0121 | 0.0118 | 0.0172 | 0.0018 | 9.9994 | 0.3114 | 0.3731 | 0.1147 |
| | OOD | 47.2462 | 0.1115 | 0.0279 | 0.0230 | 0.0392 | 0.0021 | 68.1583 | 1.5223 | 2.0991 | 0.1991 |
| Graphormer (3D) | ID | 11.0218 | 0.008 | 0.0041 | 0.0043 | 0.0056 | 0.0002 | 33.0714 | 0.4326 | 0.6297 | 0.1932 |
| | OOD | 30.6469 | 0.0255 | 0.0214 | 0.0226 | 0.0299 | 0.0141 | 264.556 | 4.4304 | 2.4939 | 1.4963 |
| ET | ID | 50.4611 | 0.0079 | 0.0026 | 0.003 | 0.0042 | 0.0005 | 21.7465 | 0.3234 | 0.369 | 0.1296 |
| | OOD | 108.1276 | 0.0341 | 0.0152 | 0.0137 | 0.0238 | .0031 | 112.7228 | 3.0890 | 2.2832 | 0.9457 |
| ModernBERT | ID | 12.8641 | 0.0113 | 0.0075 | 0.0095 | 0.0112 | 0.001 | 35.6684 | 1.1364 | 0.7342 | 0.5697 |
| | OOD | 38.9671 | 0.0218 | 0.0226 | 0.0267 | 0.0351 | 0.0026 | 237.9504 | 2.6627 | 2.784 | 0.6373 |

**Table 2.** RMSE scores of all models on OOD and ID tasks. Best performing **ID** and **OOD** models are highlighted in **Black** and **Blue** respectively. The worst performing **ID** and **OOD** models are highlighted in **Orange** and **Red** respectively. The graph-based and hybrid models provide the best scores across nearly all tasks for OOD and ID splits. Numerical encoding issues greatly hamper RTs performance and result in large errors.

First, we survey the performance of the most commonly used current machine learning models for predicting molecular properties (Table 2). The goal of this initial survey is to quantify the OOD performance that can be expected if one were to use these models in an 'out-of-the-box' fashion without modifying the model architecture. However, we perform simple optimizations of the training parameters to achieve the highest possible accuracy for each task.
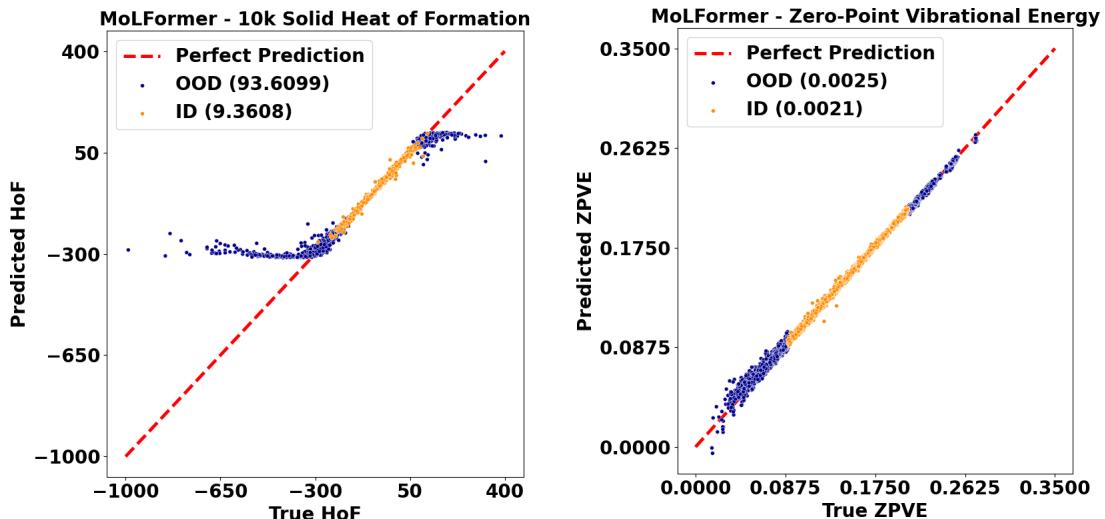
We visualize the relative performance of the models on ID and OOD tasks in Figure 2, whereas Table 2 provides the complete results with all models. We do not find any model that clearly outperforms the others on either ID or OOD performance across all tasks. Among all tested models, the Equivariant Transformer (ET) model achieves the best overall ID performance - achieving the lowest ID RMSE on 4 out of 10 tasks. For OOD prediction, MACE achieves top performance on 5 out of 10 tasks, and ET achieves top performance on 3 out of 10 tasks. We similarly find that ET achieves the best overall ID performance and MACE achieves the best OOD performance when evaluated according to the binned $R^2$ metric (see Table 4). We note that the Regression Transformer model is unique among the tested models in that property prediction is performed via autoregressive numerical token prediction. The large OOD RMSEs noted by Regression Transformer were found to arise from inaccuracies in this autoregressive numerical token generation, for example, predicting '00913', for a true value of '0.913'. Figure 3 shows a common mode of failure for OOD predictions for most models (see parity plots for all models tested in Appendix 9). We find that models performing poorly on OOD splits overwhelmingly produce an S-shaped parity plot. The models are therefore capable of clustering OOD samples together but are unable to extend the prediction region beyond the training data. Such

**Figure 2.** Normalized ID and OOD RMSE values for all models of interest for each property. RMSE values are normalized such that the best performing model has a value of 0 and the worst performing model has a value of 1. We omit the Regression Transformer model to remove the extreme values for better visualization. The models are sorted with respect to average RMSE, with the best model (lowest overall RMSE) appearing at the top of the heatmap and the worst model at the bottom.
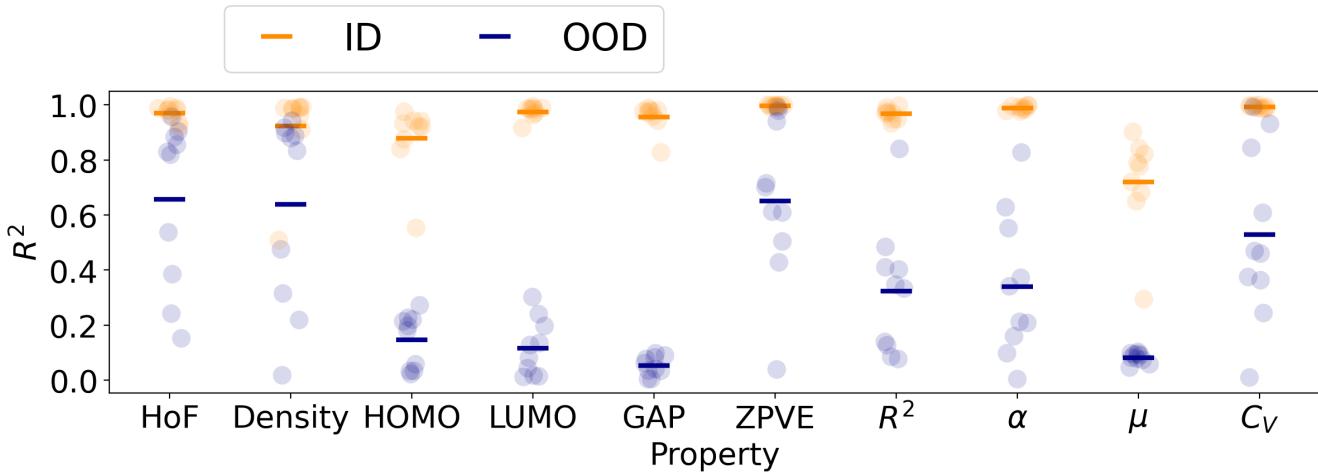
S-shaped behavior is a known failure case that arises when models learn short-cut features that maximize ID performance, but fail to generalize to OOD data [35].

Furthermore, we notice MoLFormer, one of the largest models in our test suite, achieves near state of the art performance on HoF and $C_v$, greatly outperforming other models, including similar Transformer models like ChemBERTa and RT. But ChemBERTa and MolFormer achieve similar results on OOD for both tasks. We notice a similar trend where ID performance is not necessarily correlated with OOD performance. Considering the size of our datasets, we believe large models may be able to overfit to the ID space, while achieving subpar generalization. This suggests the common strategy of pre-training on large datasets and fine-tuning on niche domains may have pitfalls for OOD samples.



**Figure 3.** Representative parity plots illustrating the OOD (blue) and ID (orange) test predictions on the 10K solid heat of formation dataset. (Left) MoLFormer predictions on this task exhibit poor OOD performance with weak correlation on the OOD samples. (Right) MoLFormer displays strong OOD performance on the QM9 zero-point vibrational energy task.

We also provide a task-centric view of our results in Fig. 4. We highlight the mean for each split with a horizontal bar in Fig. 4. The larger the difference between the ID and OOD bars, the higher the discrepancy between ID and OOD performance. As expected, ID performance is better than OOD performance for all model-task pairs. We highlight that good OOD performance is achievable in certain tasks, such as HoF, Density, ZPVE, and $C_v$. Some models achieve near-ID-level performance on these particular tasks. However, Fig. 4 also highlights particular tasks (HOMO, LUMO, Gap, and Dipole Moment($\mu$)) where no models achieve good OOD performance. Since all these properties are related to the electronic structure of molecules, we hypothesize that the inability of any model to generalize well in these tasks is due to the lack of explicit electronic structure information in their molecular representations. It is also important to note that for properties such as $C_v$, although most models achieve similar ID $R^2$ values, there is a large variance in OOD binned $R^2$ values- further highlighting the importance of performing OOD performance evaluations.



**Figure 4.** We present binned $R^2$ scores for OOD and standard $R^2$ scores for ID on each task for all models. The orange and blue bars indicate the performance averaged across all models for ID and OOD, respectively. Nearly all models have significant discrepancies between ID and OOD performance, but some models can reach ID-level accuracy. We observe that OOD performance is highly task-dependent.

## 4.2 Impact of Pre-Training

In this section, we seek to understand to what extent the large scale pretraining of existing chemical foundation models improves their OOD generalization performance. In analogy to the work in language models, these chemical foundation models are trained on vast pretraining datasets of billions of molecules with the hope that this pretraining will provide a general understanding of fundamental chemical principles. We benchmark ChemBERTa and MoLFormer (both MLM pretraining) and Regression Transfomer (PLM pretraining) to understand how the choice of pretraining tasks impacts OOD performance. Notably, the original reports of all three of these foundation models showed that that this large-scale language pretraining strategy can achieve SOTA performance on in-distribution molecular property prediction tasks[11,13], but did not perform evaluations of the OOD performance.
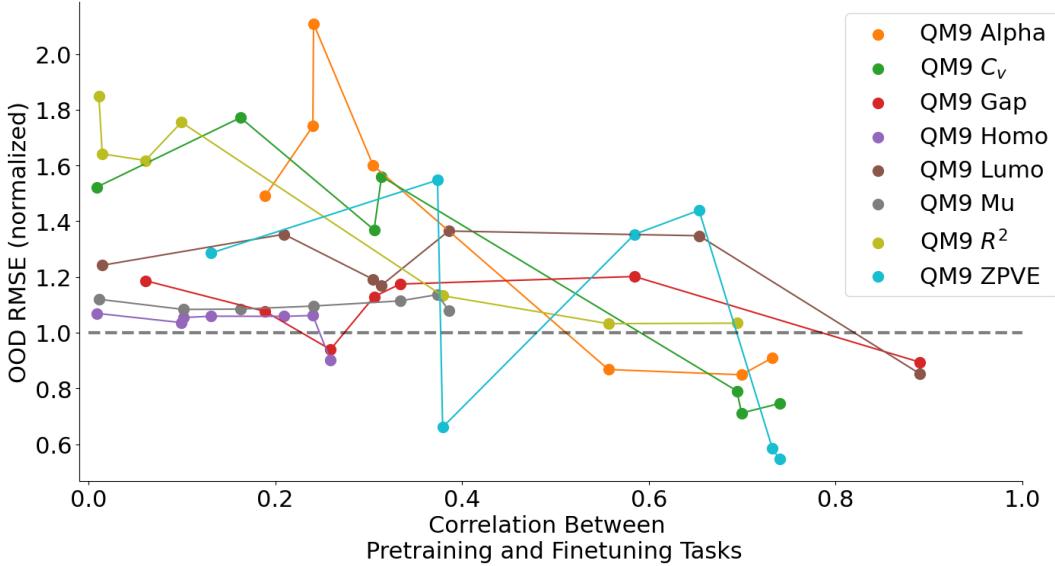
Across all three foundation models, we see a consistent improvement of ID performance across the majority of tasks Figure 11 . Averaged across all 10 tasks, the pretrained models show a sizable improvement in ID RMSE due to pretraining (31% for ChemBERTa, 35% for MoLFormer and 12% for Regression Transformer). We note that these results are consistent with the findings in their original reports. For example, the MoLFormer paper found a 29% reduction in MAE ID performance on the QM9 dataset due to MLM pretraining, whereas Regression Transformer reported up to a 52% reduction in predicting ID QED RMSE from optimizing the pretraining objective. A similar ablation study was not performed in the original ChemBERTa paper.

Surprisingly, we find that all three foundation models do not show any significant improvement in OOD performance due to language modeling pretraining (Fig.5). All three models show a negligible change in average OOD RMSE due to pretraining, and the Binned OOD $R^2$ decreases significantly for both MoLFormer (53%) and ChemBERTa (39%). Although pretraining does provide chemical foundation models with a richer understanding of chemistry, as signified by stronger ID performance, the existing pretraining procedures do not seem to allow for the models to extrapolate well to new chemistries. This result may suggest that the current pretraining tasks used by the foundation models (PLM and MLM) do not convey the relevant chemical information to allow the foundation model to extrapolate well to the downstream OOD property prediction tasks.

**Figure 5.** OOD Performance of chemical foundation models (ChemBERTA, MoLFormer and Regression Transformer) with and without pretraining, averaged across all tasks. We find that current pretraining strategies improve ID performance, but not OOD. The task-specific performances are provided in the Appendix (Figure 11).
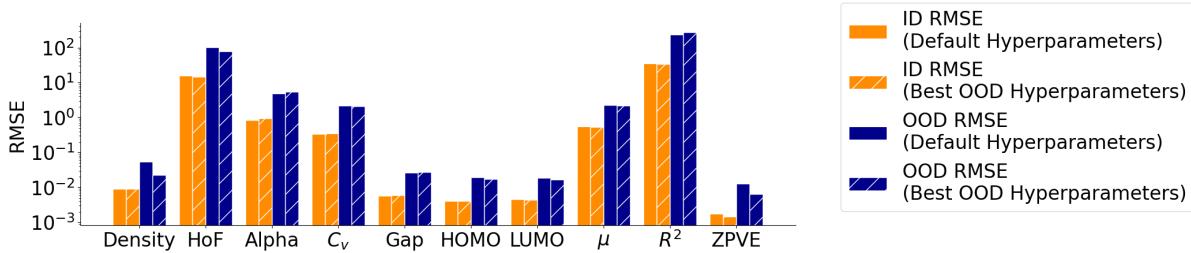
To explore if strong OOD generalization can be achieved through alternative pretraining tasks, we also explore pretraining on a supervised property prediction task. First, we perform supervised pretraining of a Chemprop model on the entirety of one of the eight QM9 property datasets. This pretrained model is then finetuned on only the training set of one of the other seven QM9 property dataset. This isolates only the property to be OOD, as the model has seen all the molecules in another context. We can see in Fig. 6 OOD performance is only improved when the pretraining task is sufficiently related to the downstream finetuning task. Specifically, across all eight datasets, we see a significant reduction in the OOD RMSE (compared against the performance without pretraining) when the Pearson correlation coefficient between the pretraining and finetuning tasks is less than 0.35. This result may explain why the masked language modeling pretraining used in current chemical foundation models resulted in worse OOD performance (Fig. 5).

**Figure 6.** OOD Performance of the Chemprop MPNN model on when pre-trained on different QM9 property datasets. Each line corresponds to the OOD performance on one of the eight QM9 OOD test sets when pre-trained on one of the other seven QM9 properties. The OOD RMSE is plotted against the Pearson correlation coefficient between the pretraining property and the finetuning property in the QM9 dataset. The OOD RMSE is normalized against the Chemprop performance without any pretraining.

### 4.3 Hyperparameter Optimization

The significant gap between the ID and OOD performance in Table 2 may indicate that the models are overfit to the ID molecules, thereby hurting OOD generalization performance. Furthermore, due to the lack of prior OOD benchmarks for molecule property prediction, the default hyperparameters used by these models are also fit to maximize ID performance, which may also negatively impact OOD generalization. In this section, we explore to what extent the OOD performance of models can be improved simply by tuning the model hyperparameters to maximize OOD performance.



**Figure 7.** OOD Performance of the Chemprop MPNN model when using default hyperparameters and the best performing OOD hyperparameters. The best OOD hyperparameters are determined according to the maximum OOD test RMSE for each property. For each property, we randomly sampled 50 values of the following hyperparameters: the number of message-passing steps; the number of layers in the feed-forward network; the dropout probability; the hidden dimension in the message-passing step; and the hidden dimension in the feed-forward network.

As shown in Fig. 7, we compare the OOD performance of Chemprop when using the default hyperparameters and hyperparameters that have been optimized to maximize OOD performance. Overall, we do not find that hyperparameter optimization can provide meaningful improvements to OOD performance. While we see a noticeable reduction in the OOD RMSE in relatively simple properties such as density(-60%), heat of formation(-23%) and ZPVE(-50%) following hyperparameter tuning, the models are still unable to generalize significantly beyond the training regime. It may not solve the problem, but for certain properties, hyperparameter optimization with respect to OOD improved over the default model by 60%, without any significant decrease to the ID performance. The results here highlight that OOD performance should be considered as an important evaluation criterion for future model optimization to ensure that models strike a balance between ID
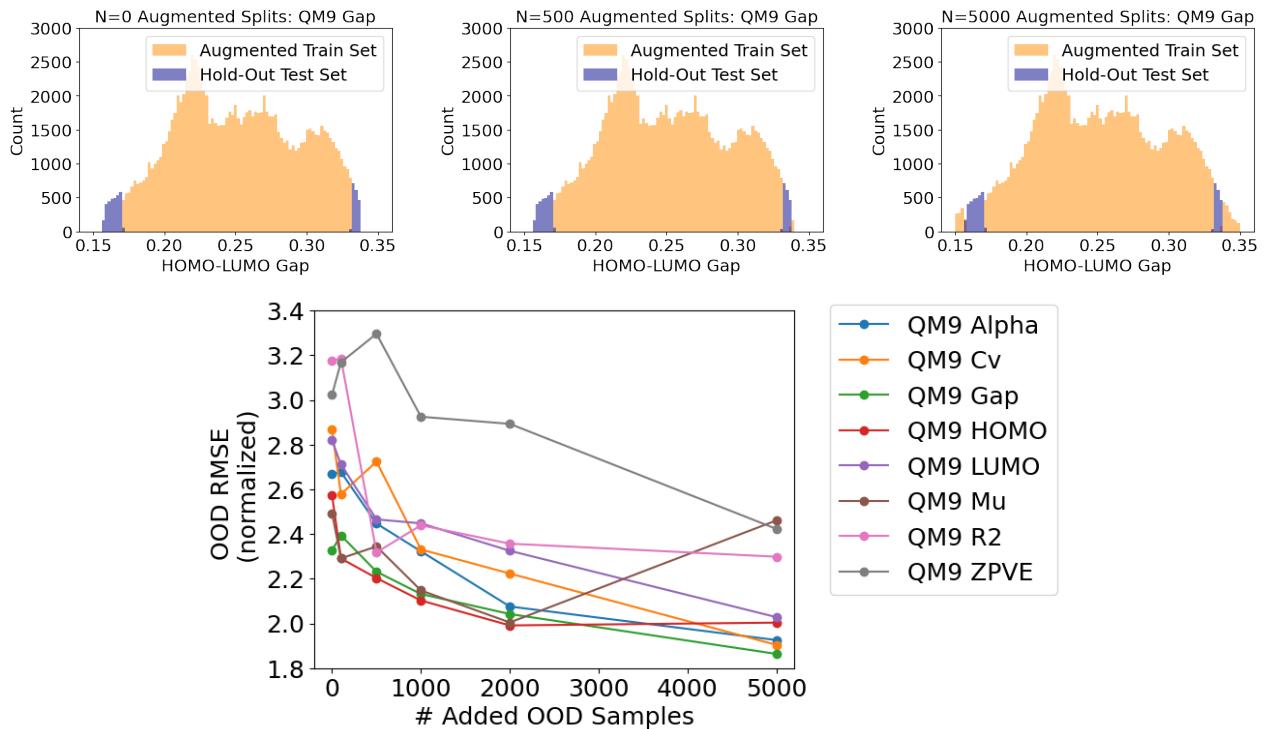
performance and OOD generalization.

## 4.4 Representation

| Representation | Split | HoF | Density | HOMO | LUMO | Gap | ZPVE | $R^2$ | $\alpha$ | $\mu$ | $C_V$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 3D | ID | **11.09** | .0121 | **.0026** | **.0030** | **.0042** | **.0005** | 21.7465 | .3234 | **.3690** | **.1296** |
| | OOD | 21.76 | .0247 | .0152 | .0137 | .0238 | .0031 | 112.7228 | .30890 | 2.2832 | .9457 |
| Explicit Bonds | ID | 15.68 | **.0092** | .0041 | .0048 | .0058 | 0.0014 | 35.68 | .8305 | .55 | .3341 |
| | OOD | 100.6 | .0551 | .0192 | .0187 | .0267 | .0129 | 234.73 | .4772 | 2.3 | 2.149 |
| SMILES | ID | 22.86 | .0163 | .0068 | .0088 | .0103 | .0046 | 50.297 | 1.444 | .7134 | .4923 |
| | OOD | 99.7253 | .1173 | .0245 | .0267 | .0315 | .0214 | 306.14 | 6.303 | 2.766 | 3.0175 |

**Table 3.** RMSE of models on OOD and ID tasks as grouped by input representation. The best performing **ID** and **OOD** models are highlighted in **Black** and **Blue** respectively. The worst performing **ID** and **OOD** models are highlighted in **Orange** and **Red** respectively.

In our study, 3D models with equivariant and invariant symmetries significantly outperform the SMILES-based models in nearly all tasks. Furthermore, the 3D GNN models like EGNN and IGNN are significantly more parameter-efficient. As we can see in Figure 3, the SMILES-based models, namely the transformer models, perform significantly worse than the 3D and explicit structure models in nearly all tasks. SMILES and explicit bonds are interchangeable representations in that SMILES can be converted into a molecular graph and vice-versa. SMILES-based representations present the same atom and topology information present in graph-like explicit-bonds representation but in a sequence format. This suggests the inductive bias present in the graph-based models improve the model performance over attention-based models, especially for OOD splits. Interestingly, the graph-like models also perform comparatively to the transformer-based models if we discount RT. MolFormer, a SMILES based, model has strong ID performance compared to other models as well.

## 4.5 Data Augmentation



**Figure 8.** Performance of the Chemprop MPNN model when various amounts of OOD samples are included in the training. For each property, we separate the 10,000 OOD examples into a hold-out test set (N=5000) and various amounts of the remaining 5000 examples are included during training. The OOD test RMSE is normalized against the validation RMSE of the model trained without any additional OOD examples included during training.
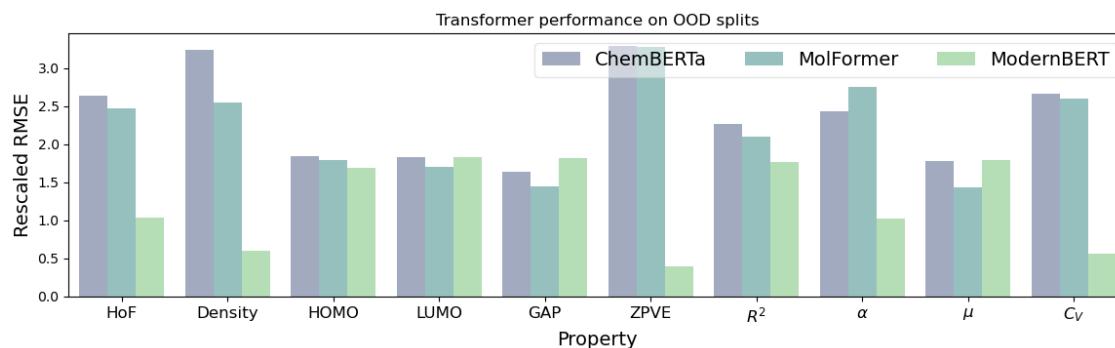
Beyond exploring different model architectures and molecular representations, data generation is a common strategy for improving the generalization capabilities of deep learning models. As a notable example, the authors of the GNoME material dataset showed that the out-of-distribution property prediction performance of a deep learning model can be vastly improved by training on an augmented dataset of millions of generated materials.[36] Similarly, using active learning to iteratively retrain graph neural networks was recently shown to reduce the property prediction error on OOD generated molecules for both density (75% reduction) and solid heat of formation (83% reduction).[8]

Figure 8 investigates improving OOD property prediction by augmenting the QM9 training set ( described in Section 1.3) with extreme-valued molecules from the QM9 OOD test set. We select N =[0, 100, 500, 1000, 2000, 5000] molecules from the original QM9 test set by selecting molecules with properties below and above the 25th and 75th quantiles, respectively.

Across 7 of 8 QM9 tasks, Chemprop's generalization improves with augmented data (Figure 8). The lack of improvement for QM9 dipole moments (Mu) likely stems from Chemprop's graph representation lacking 3D electronic structure. Data augmentation consistently yields sizable generalization improvements, even with a small fraction ( 4%) of augmented data. Further improvements may be achievable with more extensive data generation.

### 4.6 Improving Underperforming Models

Finally, we highlight a significant improvement in OOD performance with ModernBERT among the NLP-style models tested. While all the NLP models tested don't have any chemistry specific design choices, the improvements proposed in ModernBERT translate to the chemistry domain as well. We highlight the task-specific behavior in Figure 9 for OOD performance. ModernBERT performs similar to other transformer models for the difficult tasks highlighted in Figure 4, but improves significantly for the remaining properties. We see improvements in both intensive and extensive properties. ModernBERT decreases OOD HoF and $C_v$ RMSE by more than 58% and 78% respectively over other best performing transformer models. While not in the scope of our current work, understanding the design choices that result in these improvements can inform design choices for future chemistry foundation model design.



**Figure 9.** ModernBERT outperforms transformer models for OOD tasks. ModernBERT bridges the gap between transformer and GNN-based models on OOD splits, especially for HoF, Density, and $C_v$.

## 5 Discussion

In this work, we have developed BOOM, a standardized benchmark for evaluating the OOD performance of chemical property prediction models. Overall, across all 12 tested model architectures we do not find any model that achieves strong performance on all OOD tasks. As a result, we expect that current property prediction models will struggle to consistently discover molecules with properties that extrapolate beyond known molecules. Nevertheless, given the saturation of the most commonly used chemistry benchmarks, we hope that the results presented here inspire the chemistry community to pursue OOD generalization as the next frontier challenge for further developing molecular property prediction models.

Our findings also highlight promising future directions towards improving OOD generalization. Surprisingly, we found that commonly employed molecular pretraining strategies, such as masked language modeling, often result in a decrease in OOD performance. Our experiments show that developing new pretraining tasks whereby the pretraining task and the downstream property prediction tasks are more closely related results in improved OOD generalization. Fig. 6 consistently shows that OOD performance is only improved by pretraining when the chemical information contained in the pretraining task is related to the downstream property prediction task. Randomly sampling model hyperparameters of the Chemprop GNN was found to significantly improve OOD performance, with very little change in ID performance. This result highlights the need to consider OOD generalization when optimizing the model hyperparameters of chemical prediction models.

While high-inductive bias and 3D models perform well on our current tests, scalability remains a significant issue. Small models can provide strong predictive power, but they do not allow for techniques such as in-context learning and test-time compute that may be available to large scale models. Similarly, 3D models are attractive, but high-quality DFT data is not always available. While 3D molecular data is becoming increasingly more available, it dwarfs in comparison to the billions of molecules used in unsupervised molecular pretraining strategies. In general, as our results with ModernBERT show, transformer-based models can potentially catch up to the small models while enabling greater scalability. Numerical encoding is a concern for LLM-like models and was a significant drawback for RT. Improved post-hoc solutions [37] or modern tokenization techniques [38,39] will be key in the development of LLM-based predictive models.

Perhaps unsurprisingly, augmenting the training dataset to minimize the distribution shift between the ID and OOD samples was found to consistently improve model generalization. However, the ease of generating a sufficient quantity of high-quality molecular property data to augment the training set is highly task-dependent, which may limit the applicability of this approach.

Notably, we do not find any strategy that universally improves OOD performance across all property prediction tasks. On relatively simple properties (i.e. density and ZPVE), current SOTA property prediction models exhibit strong generalization with very little difference between ID and OOD performance. However, we note that properties that depend on the electronic structure of the molecules (i.e. Dipole moment and HOMO) showed little improvement in OOD performance across all experimental settings. We anticipate that achieving strong OOD generalization on these properties will require larger datasets, in combination with molecular representations that explicitly capture the molecules' electronic structure.

## 6 Code Availability

The BOOM benchmark datasets will be made available open-source on Github. This repository contains code for generating consistent OOD property prediction evaluations, as well as the benchmarking results presented in this paper.
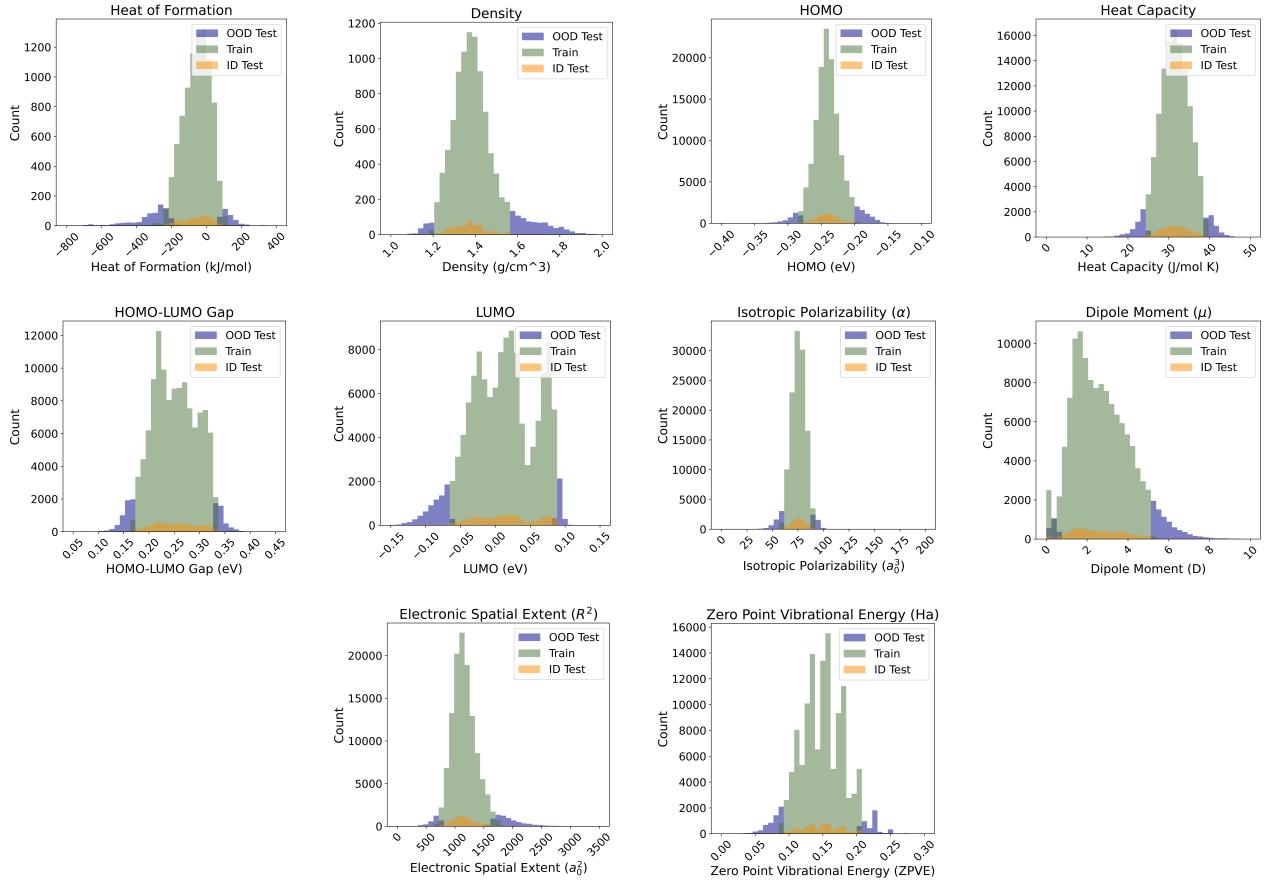
## 7 Acknowledgments

# References

1. Bohacek, R. S., McMartin, C. & Guida, W. C. The art and practice of structure-based drug design: A molecular modeling perspective. *Medicinal Res. Rev.* **16**, 3–50, DOI: 10.1002/(SICI)1098-1128(199601)16: 1⟨3::AID-MED1⟩3.0.CO;2-6 (1996). _eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/%28SICI%291098-1128%28199601%2916%3A1%3C3%3A%3AAID-MED1%3E3.0.CO%3B2-6.

2. Reymond, J.-L. The Chemical Space Project. *Accounts Chem. Res.* **48**, 722–730, DOI: 10.1021/ar500432k (2015). Publisher: American Chemical Society.

3. Kailkhura, B., Gallagher, B., Kim, S., Hiszpanski, A. & Han, T. Y.-J. Reliable and explainable machine-learning methods for accelerated material discovery. *npj Comput. Mater.* **5**, 108 (2019).

4. Wang, H. *et al.* Scientific discovery in the age of artificial intelligence. *Nature* **620**, 47–60 (2023).

5. Farquhar, S. & Gal, Y. What'out-of-distribution'is and is not. In *Neurips ml safety workshop* (2022).

6. Scheirer, W. J., de Rezende Rocha, A., Sapkota, A. & Boult, T. E. Toward open set recognition. *IEEE transactions on pattern analysis machine intelligence* **35**, 1757–1772 (2012).

7. Ramakrishnan, R., Dral, P. O., Rupp, M. & von Lilienfeld, O. A. Quantum chemistry structures and properties of 134 kilo molecules. *Sci. Data* **1**, 140022, DOI: 10.1038/sdata.2014.22 (2014). Publisher: Nature Publishing Group.

8. Antoniuk, E. R. *et al.* Active learning enables extrapolation in molecular generative models. *arXiv preprint arXiv:2501.02059* DOI: 10.48550/arXiv.2501.02059 (2025). ArXiv:2501.02059 [cs].

9. Landrum, G. *et al.* Rdkit: A software suite for cheminformatics, computational chemistry, and predictive modeling. *Greg Landrum* **8**, 5281 (2013).

10. Ramsundar, B. *et al.* *Deep Learning for the Life Sciences* (O'Reilly Media, 2019). https://www.amazon.com/Deep-Learning-Life-Sciences-Microscopy/dp/1492039837.

11. Ross, J. *et al.* Large-scale chemical language representations capture molecular structure and properties. *Nat. Mach. Intell.* **4**, 1256–1264, DOI: 10.1038/s42256-022-00580-7 (2022).

12. Raffel, C. *et al.* Exploring the limits of transfer learning with a unified text-to-text transformer. *J. machine learning research* **21**, 1–67 (2020).

13. Chithrananda, S., Grand, G. & Ramsundar, B. Chemberta: large-scale self-supervised pretraining for molecular property prediction. *arXiv preprint arXiv:2010.09885* (2020).

14. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, 4171–4186 (2019).

15. Born, J. & Manica, M. Regression transformer enables concurrent sequence regression and generation for molecular language modelling. *Nat. Mach. Intell.* **5**, 432–444 (2023).

16. Yang, Z. *et al.* Xlnet: Generalized autoregressive pretraining for language understanding. *Adv. neural information processing systems* **32** (2019).

17. Su, J. *et al.* Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing* **568**, 127063 (2024).

18. Shazeer, N. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202* (2020).

19. Schütt, K. T., Sauceda, H. E., Kindermans, P.-J., Tkatchenko, A. & Müller, K.-R. Schnet–a deep learning architecture for molecules and materials. *The J. Chem. Phys.* **148** (2018).

20. Qiao, Z., Welborn, M., Anandkumar, A., Manby, F. R. & Miller, T. F. Orbnet: Deep learning for quantum chemistry using symmetry-adapted atomic-orbital features. *The J. chemical physics* **153** (2020).

21. Heid, E. *et al.* Chemprop: A machine learning package for chemical property prediction. *J. Chem. Inf. Model.* **64**, 9–17 (2023).

22. Satorras, V. G., Hoogeboom, E. & Welling, M. E (n) equivariant graph neural networks. In *International conference on machine learning*, 9323–9332 (PMLR, 2021).

23. Batatia, I., Kovács, D. P., Simm, G. N. C., Ortner, C. & Csányi, G. MACE: Higher Order Equivariant Message Passing Neural Networks for Fast and Accurate Force Fields. *arXiv preprint arXiv:2206.07697* DOI: 10.48550/arXiv.2206.07697 (2023). ArXiv:2206.07697 [stat].

24. Kovács, D. P., Batatia, I., Arany, E. S. & Csányi, G. Evaluation of the MACE force field architecture: From medicinal chemistry to materials science. *The J. Chem. Phys.* **159**, 044118, DOI: 10.1063/5.0155322 (2023).

25. Ying, C. *et al.* Do transformers really perform badly for graph representation? *Adv. neural information processing systems* **34**, 28877–28888 (2021).

26. Thölke, P. & De Fabritiis, G. Torchmd-net: equivariant transformers for neural network based molecular potentials. *arXiv preprint arXiv:2202.02541* (2022).

27. Yang, J. *et al.* Openood: Benchmarking generalized out-of-distribution detection. *Adv. Neural Inf. Process. Syst.* **35**, 32598–32611 (2022).

28. Bendale, A. & Boult, T. E. Towards open set deep networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 1563–1572 (2016).

29. Bulusu, S., Kailkhura, B., Li, B., Varshney, P. K. & Song, D. Anomalous example detection in deep learning: A survey. *IEEE Access* **8**, 132330–132347 (2020).

30. Liu, J. *et al.* Towards out-of-distribution generalization: A survey. *arXiv preprint arXiv:2108.13624* (2021).

31. Dunn, A., Wang, Q., Ganose, A., Dopp, D. & Jain, A. Benchmarking materials property prediction methods: the matbench test set and automatminer reference algorithm. *npj Comput. Mater.* **6**, 138 (2020).

32. Omee, S. S., Fu, N., Dong, R., Hu, M. & Hu, J. Structure-based out-of-distribution (ood) materials property prediction: a benchmark study. *npj Comput. Mater.* **10**, 144 (2024).

33. Witman, M. D. & Schindler, P. MatFold: systematic insights into materials discovery models' performance through standardized cross-validation protocols. *Digit. Discov.* **4**, 625–635, DOI: 10.1039/D4DD00250D (2025). Publisher: RSC.

34. Segal, N., Netanyahu, A., Greenman, K. P., Agrawal, P. & Gomez-Bombarelli, R. Known unknowns: Out-of-distribution property prediction in materials and molecules. *arXiv preprint arXiv:2502.05970* (2025).

35. Geirhos, R. *et al.* Shortcut learning in deep neural networks. *Nat. Mach. Intell.* **2**, 665–673 (2020).

36. Merchant, A. *et al.* Scaling deep learning for materials discovery. *Nature* **624**, 80–85, DOI: 10.1038/s41586-023-06735-9 (2023). Publisher: Nature Publishing Group.

37. Golkar, S. *et al.* xval: A continuous number encoding for large language models. *arXiv preprint arXiv:2310.02989* (2023).

38. Achiam, J. *et al.* Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).

39. Grattafiori, A. *et al.* The llama 3 herd of models. *arXiv preprint arXiv:2407.21783* (2024).

40. Vaswani, A. *et al.* Attention is all you need. *Adv. neural information processing systems* **30** (2017).

41. Rogers, A., Kovaleva, O. & Rumshisky, A. A primer in bertology: What we know about how bert works. *Transactions association for computational linguistics* **8**, 842–866 (2021).

42. Liu, Y. *et al.* Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).

43. Radford, A., Narasimhan, K., Salimans, T., Sutskever, I. *et al.* Improving language understanding by generative pre-training. *OpenAI blog* (2018).

44. Radford, A. *et al.* Language models are unsupervised multitask learners. *OpenAI blog* **1**, 9 (2019).

45. Brown, T. *et al.* Language models are few-shot learners. *Adv. neural information processing systems* **33**, 1877–1901 (2020).

# 8 Appendix

## 8.1 Datasets



## 8.2 Transformers Overview

Unlike multi-layer perceptrons, CNNs, and RNNs, transformers use learnable attention layers, which capture data-driven associations in input features, offering greater flexibility and expressiveness than convolutional or recurrent layers. While variations exist in tokenization, attention mechanisms, model architecture, and pre-training schemes, the field has converged on a few foundational architectures. Recent advancements focus on improving training methods, tokenization, and scaling.

Transformer architectures fall into three categories: encoder-only, decoder-only, and encoder-decoder models. Encoder-decoder models, such as the original Transformer [40] and T5 [12], use an encoder to process input, with the attention mechanism accessing the entire sequence. Decoder models are autoregressive, with masked attention attending only to previous tokens. They are trained using self-supervised techniques like next-token prediction, masked language modeling, or both. Encoder-only (BERT-like) models [14,41,42] rely solely on the encoder and are primarily trained with masked language modeling. They are fine-tuned for supervised prediction by adding prediction heads to the encoder output. Decoder-only (GPT-like) models [38,43–45] consist only of the transformer's decoder and generate outputs autoregressively, trained with a next-token prediction objective.

## 8.3 GNN Formulation

Symmetries are inherent to the physical laws that dictate molecular properties. Algebraically, they are represented as groups, where each element corresponds to a transformation. For non-chiral molecules, the E(3) group, encompassing rotations, translations, and reflections, is key. Chiral molecules require the SE(3) subgroup, which excludes reflection. Since molecular properties remain invariant under these transformations, learned structure-to-property functions should obey the same symmetries.

GNNs naturally encode these symmetries. MPNNs enforce permutation-invariant message aggregation, making models permutation-invariant. Geometric deep learning models extend this by enabling molecular representations in 3D space, ensuring networks are invariant or equivariant to geometric transformations. Invariance implies properties remain unchanged after

transformation, while equivariance means vector properties transform consistently with applied transformations. Here, we provide rigorous definitions.

For completeness, we reproduce the GNN formulation from[22]. For a given GNN with node features $h_i^{(l)}$ are the features of the $i$-th node for $l$-th layer. $b_{ij}$ are the edge-features between two connected nodes $i$ and $j$ such that $j \in \mathcal{N}_i$. The neighborhood $\mathcal{N}_i$ is the set of nodes connected to node $i$. $W^{(l)}$ is a learnable projection matrix of layer $l$.

*Topological GNN:*

$$h_i^{(l+1)} = h_i^{(l)} W^{(l)} + \sum_{j \in \mathcal{N}_i} \theta(b_{ij}, h_i^{(l)}, h_j^{(l)}) \tag{1}$$

Where $\theta(\cdot)$ is a learnable function of the bond and node features, shared between all node pairs.

*Invariant GNN:*

$$h_i^{(l+1)} = h_i^{(l)} W^{(l)} + \sum_{j \in \mathcal{N}_i} \theta(b_{ij}, h_i^{(l)}, h_j^{(l)}) + \sum_{j \neq i} \phi(r_{ij}, h_i^{(l)}, h_j^{(l)}) \tag{2}$$

Where, $r_{ij} = ||x_i - x_j||^2$ is the inter-atomic distance between atoms $i$ and $j$. $\phi(\cdot)$ is a learnable function of the interatomic distances and node features, shared between all node pairs.

*Equivariant GNN:*

$$h_i^{(l+1)} = h_i^{(l)} W^{(l)} + \sum_{j \in \mathcal{N}_i} \theta(b_{ij}, h_i^{(l)}, h_j^{(l)}) + \sum_{j \neq i} \phi(r_{ij}^{(l)}, h_i^{(l)}, h_j^{(l)}) \tag{3}$$

$$x^{(l+1)} = x^{(l)} + \sum_{j \neq i} \left( \frac{x_i^{(l)} - x_i^{(l)}}{r_{ij}^{(l)} + \xi} \right) \psi(r_{ij}^{(l)}, h_i^{(l)}, h_j^{(l)}) \tag{4}$$

Where, $r_{ij}^{(l)} = ||x_i^{(l)} - x_j^{(l)}||^2$ is the inter-atomic distance between atoms $i$ and $j$ at the $l$-th layer. $\xi$ is a small constant for numerical stability. $\psi(\cdot)$ is a learnable function of the inter-atomic distances and node features, shared between all node pairs.

As we can see Eq. 3 is equivalent to Eq. 2 but with a per-layer coordinate update. Furthermore, Eq. 2 is equivalent to Eq. 1 but with an additional term dependent on the pairwise distances $r_{ij}$.

### 8.3.1 Readout Function

The readout function, $\mathcal{R}$ of a GNN aggregates the node-level information on the graph and combines them to get a graph-level output. The readout function can be any permutation invariant function such that, $\mathcal{R} : \mathbb{R}^{|\mathcal{V} \times F|} \to \mathbb{R}^K$, where $F$ is the per-vertex feature dimension, and $K$ is the output dimension ($K = 1$ in the case of regression). The flexibility in the readout function can be used to provide target-specific inductive bias such as using a summing over the vertices for extensive properties while taking the mean output for the vertices for intensive properties.
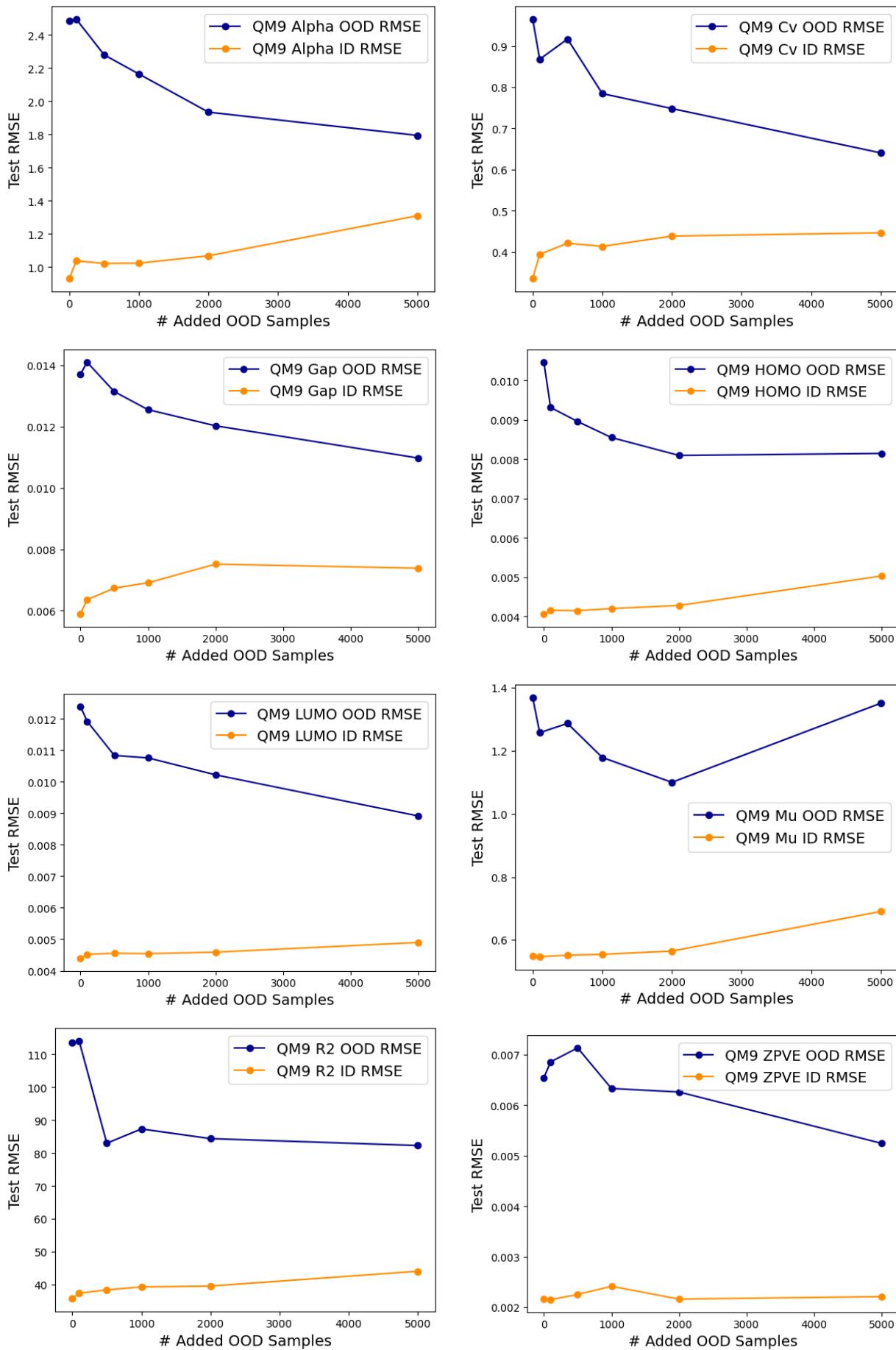
MACE and ET use modified readout functions for some properties, such as $\mu$, while we are using the unmodified readout function. We have not had success modifying the readout function as described in their publication, but we are working with the authors to replicate their results. We plan on investigating this further.
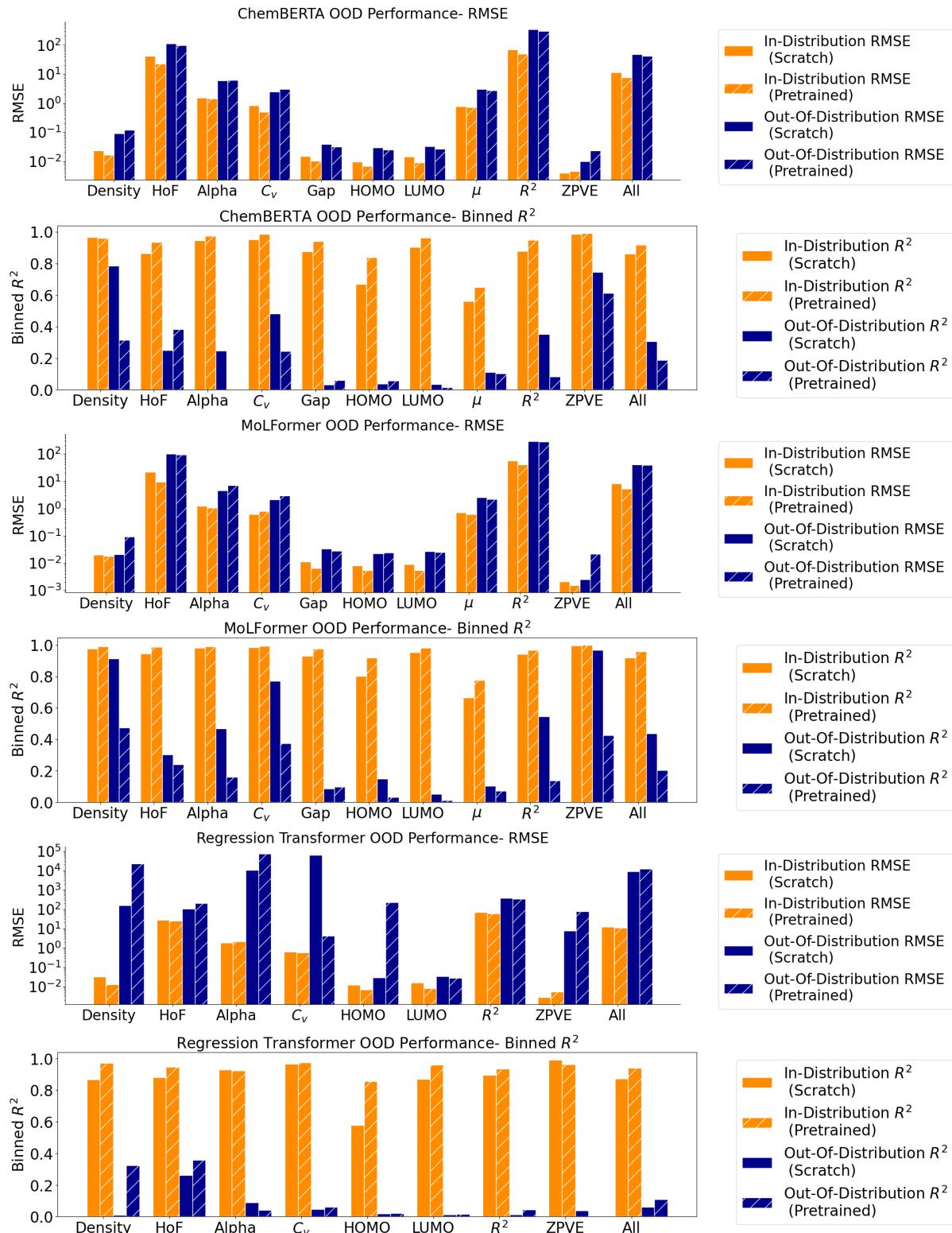
## 8.4 Transformer Fine-tuning Details

ChemBERTa and MolFormer models are pre-trained with a masked language modeling (MLM) task and the Regression Transformer is pretrained with a permutation language modeling (PLM) task. During MLM pretraining, a predetermined fraction of the SMILES string of the molecule is masked and then predicted by the model. The Regression Transformer foundation model uses a PLM pretraining task, which seeks to autoregressively predict masked tokens from a permuted sequence of both SMILES and property tokens.

For Regression Transformer and ChemBERTa, the models without pretraining are initialized with random weights, whereas the MoLFormer model without pretraining is loaded directly from the provided checkpoint saved at the beginning (0th iteration) of pretraining. For all three models, the pretrained models are initialized from the provided model checkpoints, before finetuning on each of the 10 downstream OOD tasks. Both the pretrained and scratch models are fine-tuned according to the same learning schedule hyperparameters.

## 8.5 Additional Plots

**Figure 11.** OOD Performance of chemical foundation models (ChemBERTA, MoLFormer and Regression Transformer) with and without pretraining. The performance of Regression Transformer on the QM9 dipole moment and HOMO-LUMO Gap properties are omitted due to the inability of the scratch Regression Transformer model to converge on these properties.
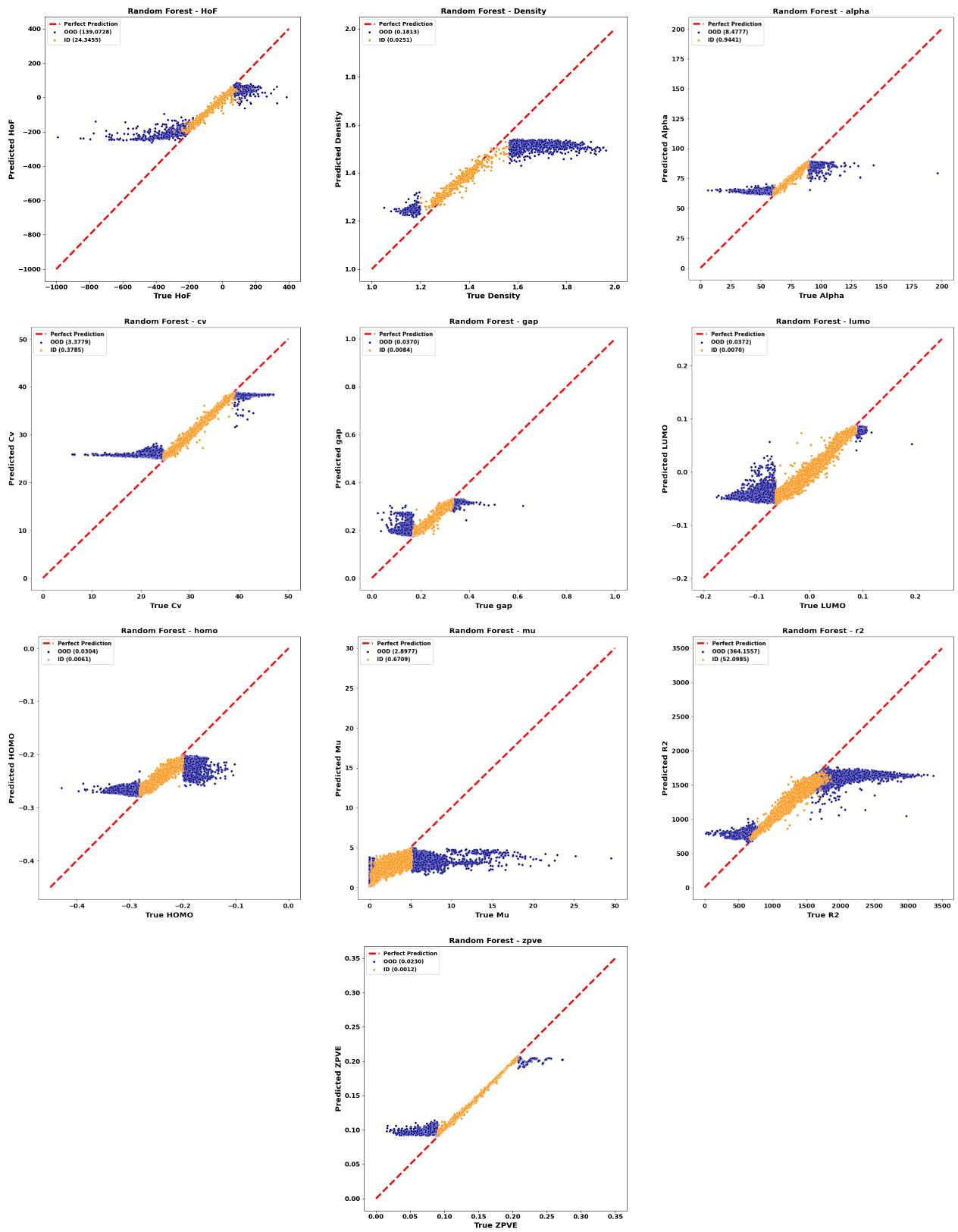
## 8.6 $R^2$ results

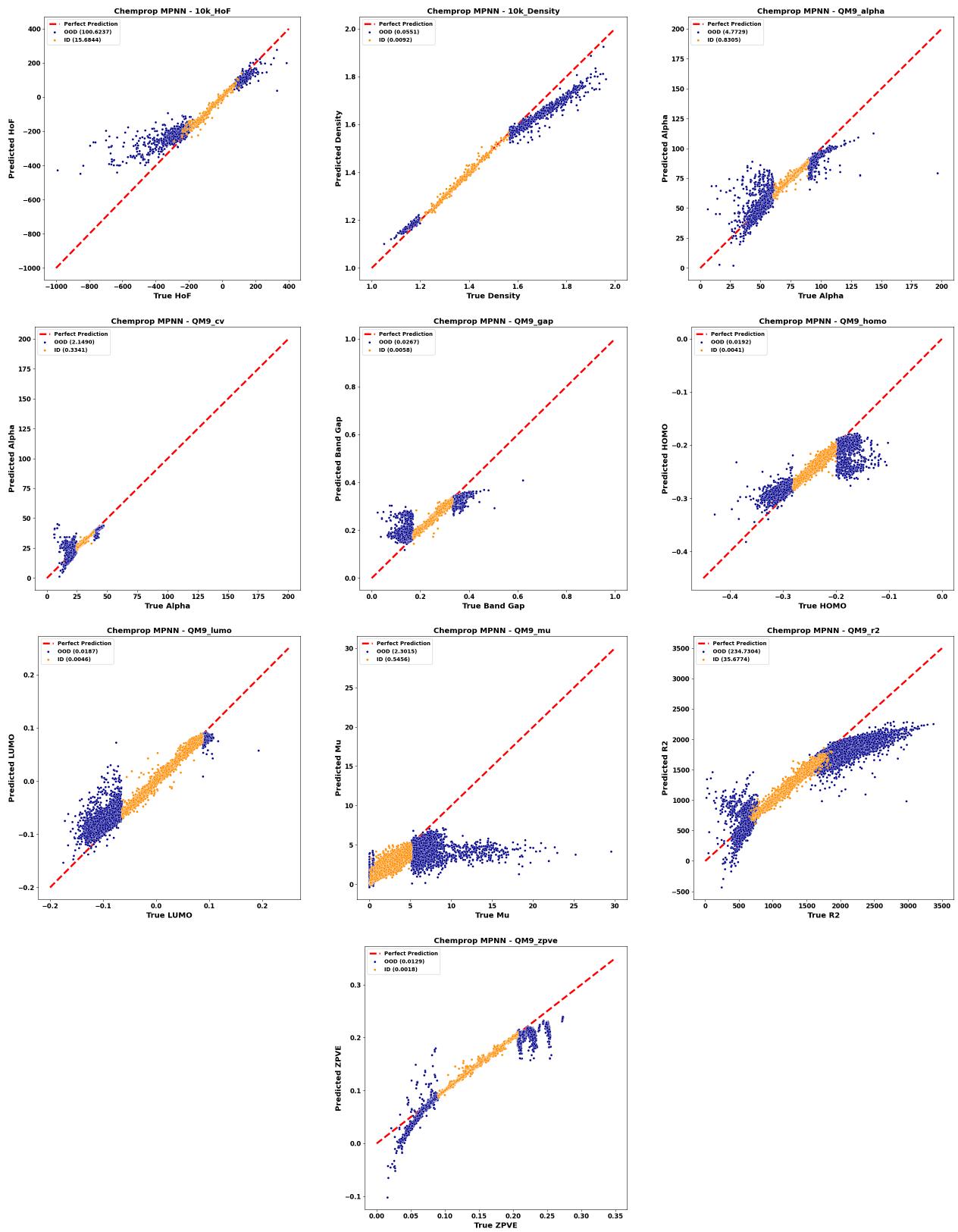| Model | Split | HoF | Density | HOMO | LUMO | GAP | ZPVE | $R^2$ | $\alpha$ | $\mu$ | $C_v$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Random Forest | ID | 0.920 | 0.910 | 0.875 | 0.970 | 0.958 | 0.998 | 0.933 | 0.978 | 0.682 | 0.987 |
| | OOD | 0.152 | 0.018 | 0.022 | 0.012 | 0.041 | 0.039 | 0.077 | 0.098 | 0.098 | 0.010 |
| RT | ID | 0.947 | 0.971 | 0.856 | 0.960 | 0.923 | 0.963 | 0.934 | 0.925 | 0.475 | 0.975 |
| | OOD | 0.358 | 0.324 | 0.020 | 0.015 | 0.026 | 0.003 | 0.045 | 0.040 | 0.086 | 0.061 |
| ChemBERTa | ID | 0.935 | 0.962 | 0.840 | 0.963 | 0.942 | 0.991 | 0.948 | 0.976 | 0.650 | 0.987 |
| | OOD | 0.385 | 0.315 | 0.058 | 0.017 | 0.062 | 0.612 | 0.085 | 0.004 | 0.103 | 0.244 |
| MolFormer | ID | 0.988 | 0.992 | 0.921 | 0.983 | 0.977 | 0.999 | 0.970 | 0.993 | 0.777 | 0.995 |
| | OOD | 0.242 | 0.475 | 0.034 | 0.014 | 0.098 | 0.428 | 0.139 | 0.160 | 0.074 | 0.375 |
| Chemprop | ID | 0.963 | 0.985 | 0.941 | 0.987 | 0.980 | 0.996 | 0.968 | 0.983 | 0.790 | 0.990 |
| | OOD | 0.537 | 0.897 | 0.197 | 0.128 | 0.077 | 0.609 | 0.333 | 0.341 | 0.096 | 0.460 |
| EGNN | ID | 0.989 | 0.989 | 0.932 | 0.988 | 0.977 | 0.998 | 0.990 | 0.992 | 0.842 | 0.994 |
| | OOD | 0.883 | 0.888 | 0.227 | 0.081 | 0.034 | 0.701 | 0.402 | 0.208 | 0.079 | 0.469 |
| IGNN | ID | 0.980 | 0.987 | 0.920 | 0.984 | 0.972 | 0.998 | 0.978 | 0.984 | 0.821 | 0.986 |
| | OOD | 0.855 | 0.879 | 0.220 | 0.135 | 0.035 | 0.715 | 0.347 | 0.212 | 0.093 | 0.363 |
| TGNN | ID | 0.969 | 0.965 | 0.895 | 0.946 | 0.919 | 0.999 | 0.930 | 0.986 | 0.713 | 0.990 |
| | OOD | 0.819 | 0.833 | 0.214 | 0.240 | 0.006 | 0.991 | 0.127 | 0.552 | 0.082 | 0.931 |
| MACE | ID | 0.995 | 0.509 | 0.553 | 0.916 | 0.827 | 0.996 | 0.997 | 0.998 | 0.902 | 0.999 |
| | OOD | 0.956 | 0.219 | 0.033 | 0.197 | 0.004 | 0.979 | 0.840 | 0.827 | 0.058 | 0.992 |
| Graphormer (3D) | ID | 0.981 | 0.988 | 0.942 | 0.988 | 0.981 | 1.000 | 0.972 | 0.995 | 0.720 | 0.997 |
| | OOD | 0.829 | 0.917 | 0.181 | 0.044 | 0.083 | 0.504 | 0.410 | 0.372 | 0.091 | 0.608 |
| ET | ID | 0.983 | 0.992 | 0.975 | 0.995 | 0.989 | 1.000 | 0.947 | 0.998 | 0.294 | 0.999 |
| | OOD | 0.903 | 0.942 | 0.272 | 0.302 | 0.090 | 0.940 | 0.484 | 0.628 | 0.046 | 0.844 |
| ModernBERT | ID | 0.942 | 0.976 | 0.849 | 0.967 | 0.944 | 0.999 | 0.947 | 0.977 | 0.650 | 0.989 |
| | OOD | 0.620 | 0.888 | 0.353 | 0.303 | 0.329 | 0.990 | 0.731 | 0.629 | 0.075 | 0.917 |

**Table 4.** Batched $R^2$ scores of all models on OOD and ID tasks. Best performing **ID** and **OOD** models are highlighted in **Black** and **Blue** respectively. The worst performing **ID** and **OOD** models are highlighted in **Orange** and **Red** respectively. The graph-based and hybrid models provide the best scores across nearly all tasks for OOD and ID splits.
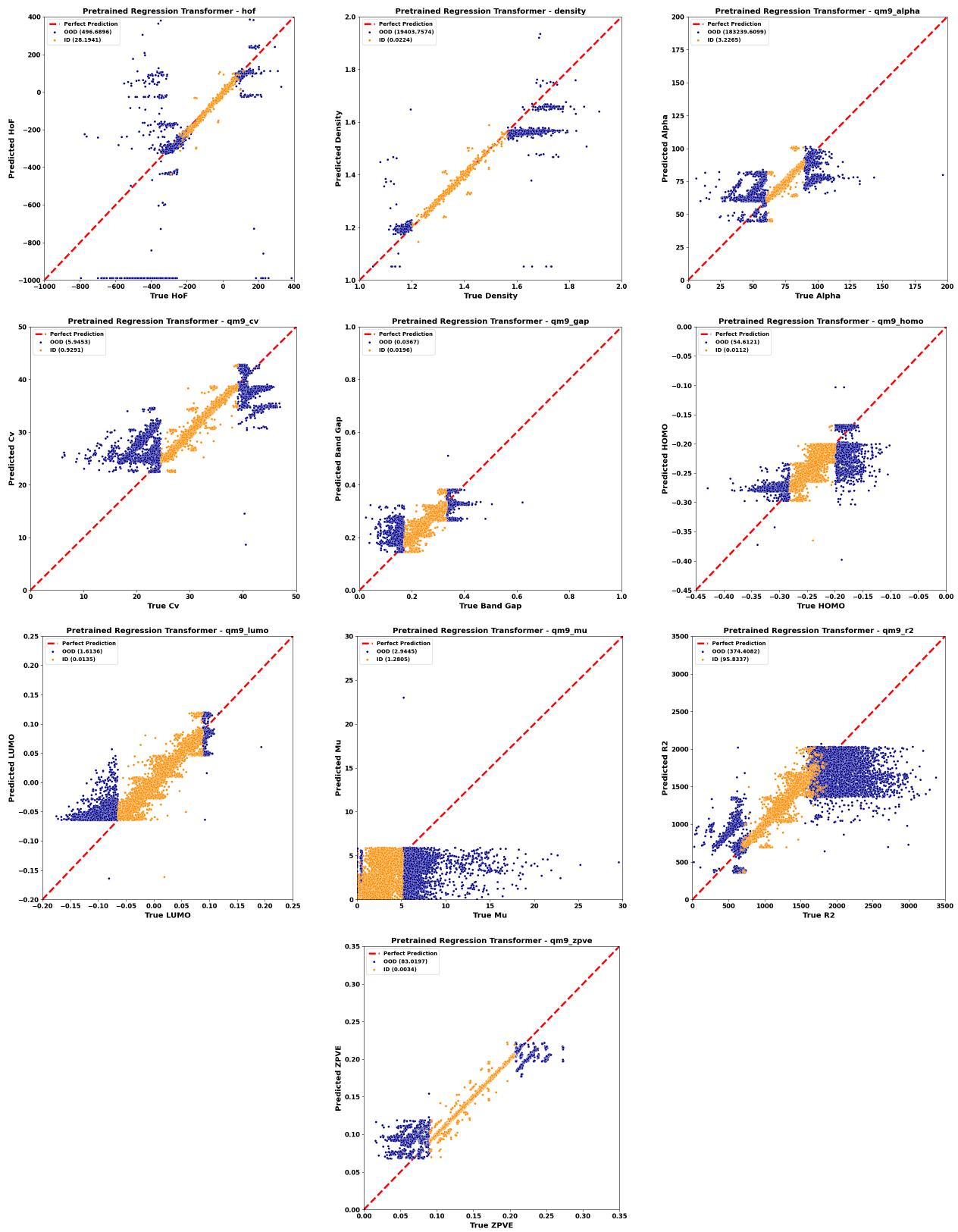
# 9 Parity Plots

As there are more than 150 plots, we provide the parity plots for all of our experiments in a compressed layout in the following section, intended for observing the prediction trends for the model. We also upload the higher resolution images, as well as the actual predictions, including the training/fine-tuning code for all models, to our repository.
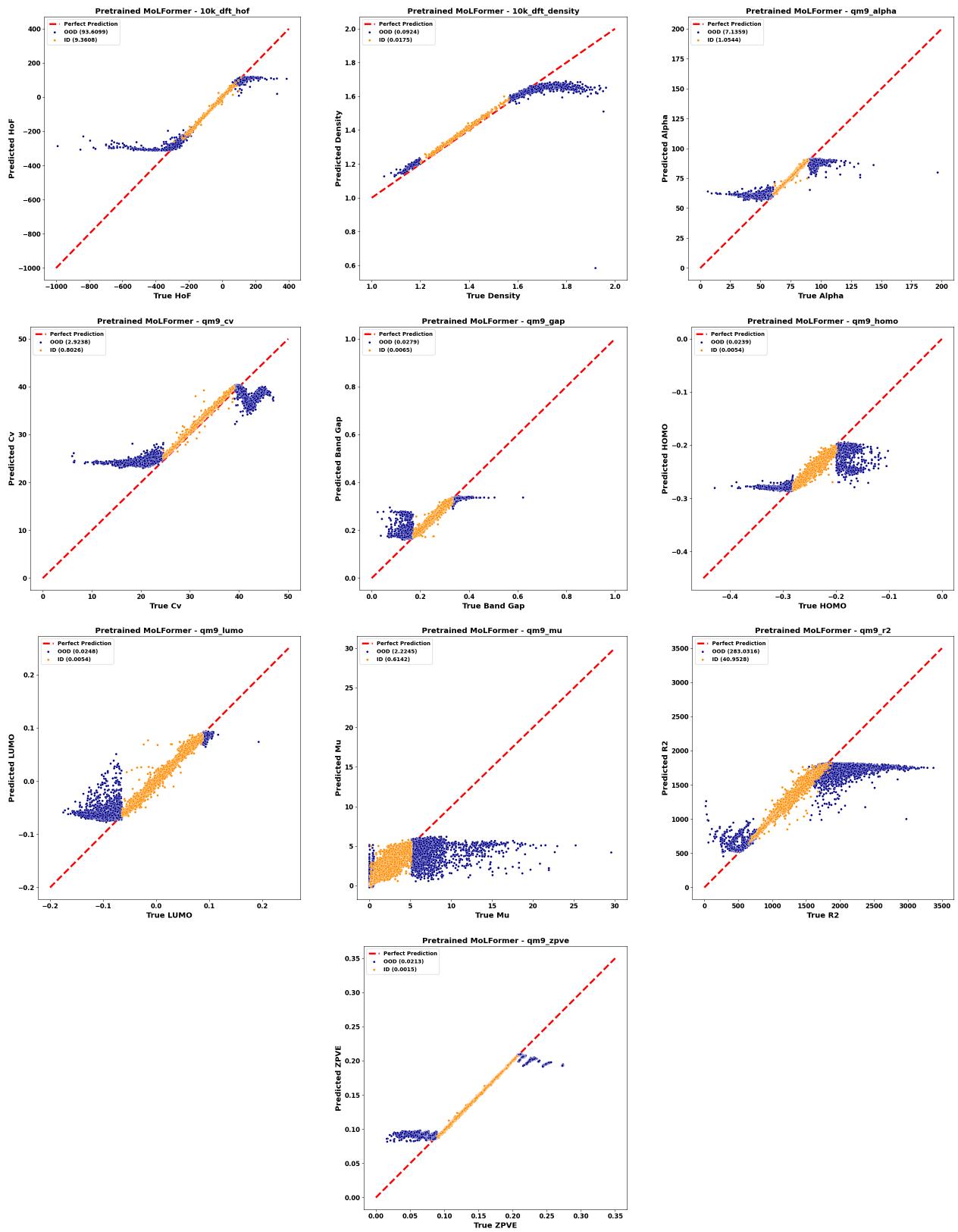
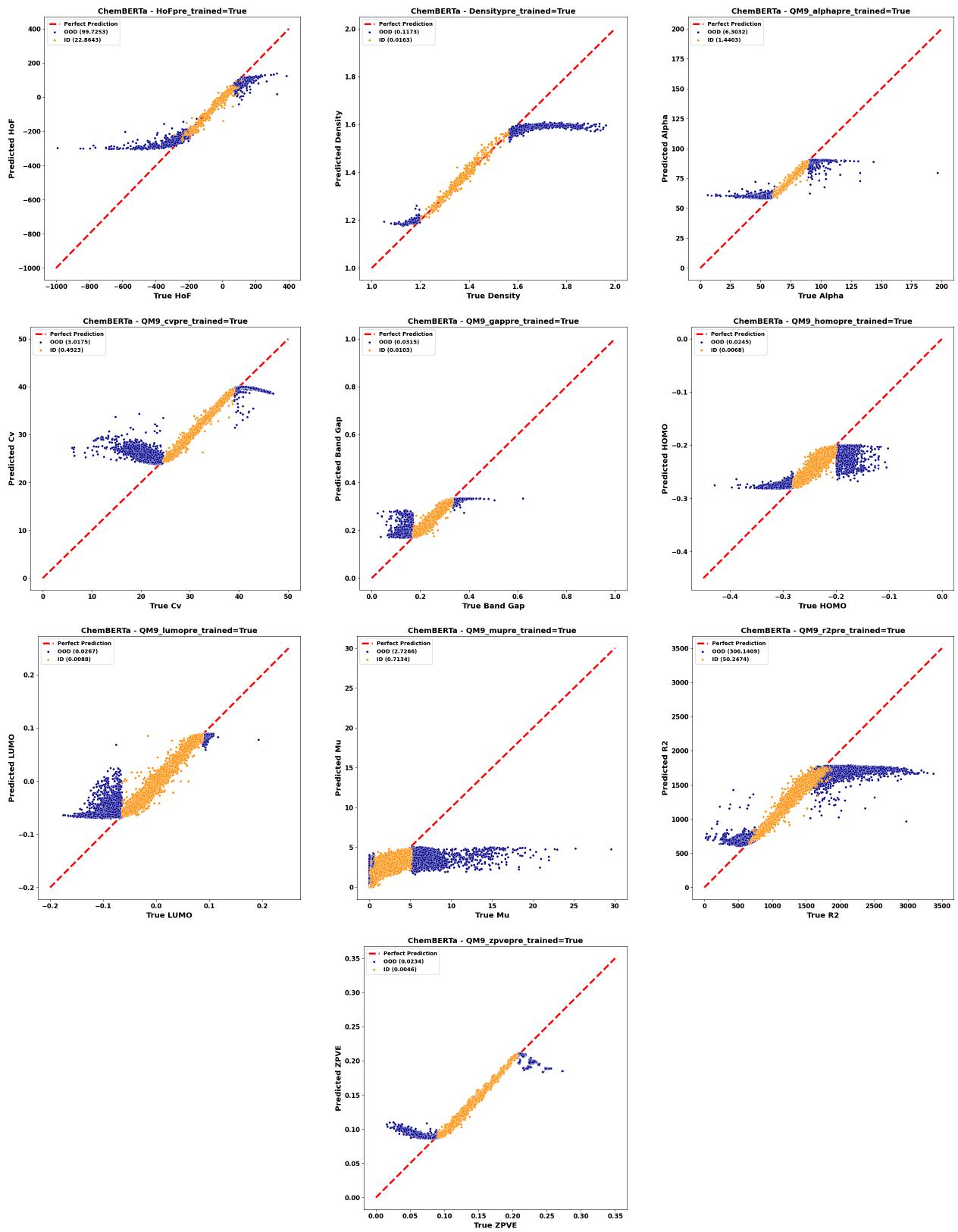**Figure 12.** Parity Plots for Random Forest on 10K and QM9 OOD tasks.

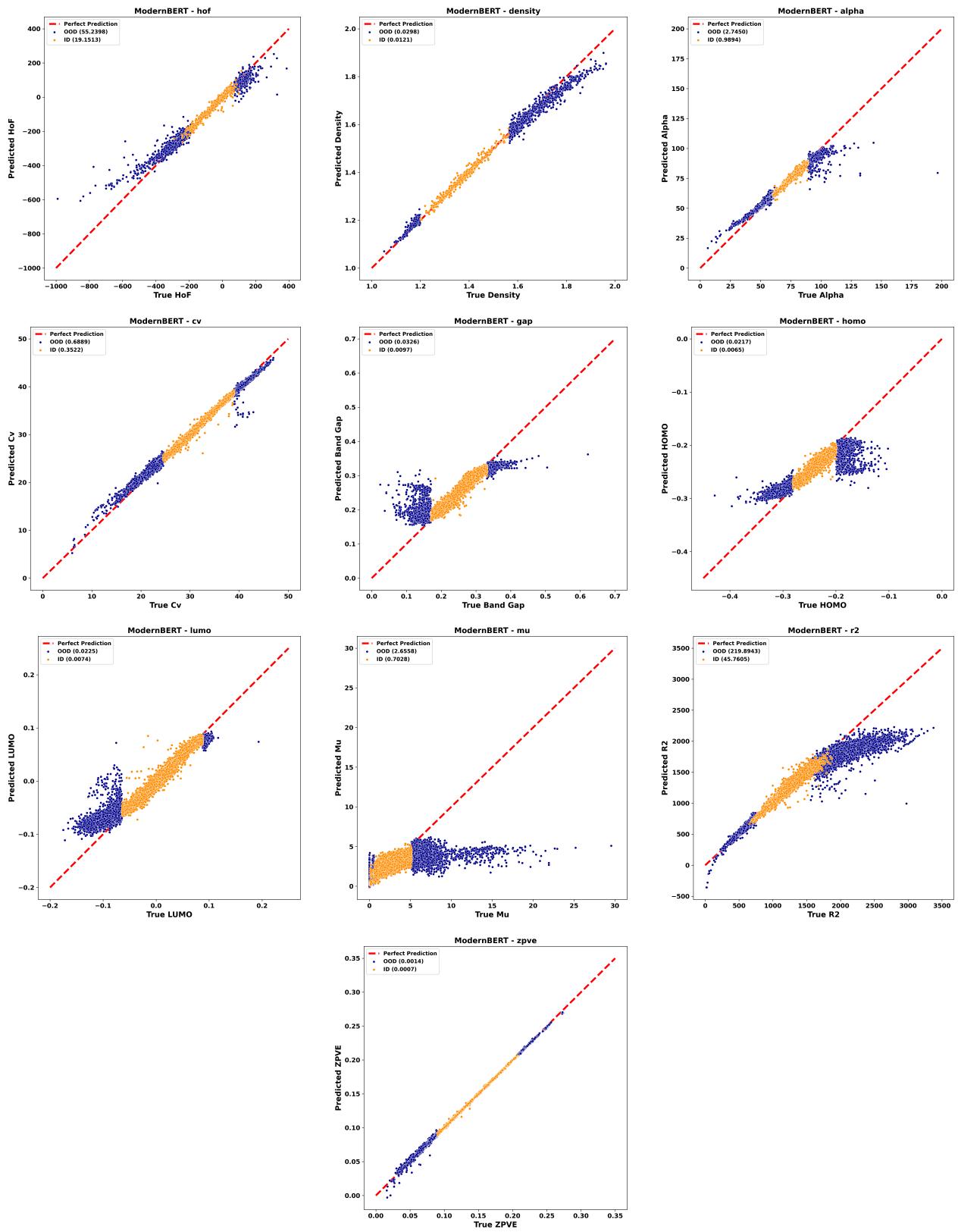**Figure 13.** Parity Plots for Chemprop on 10K and QM9 OOD tasks.

**Figure 14.** Parity Plots for Regression Transformer (with Pretraining) on 10K and QM9 OOD tasks.
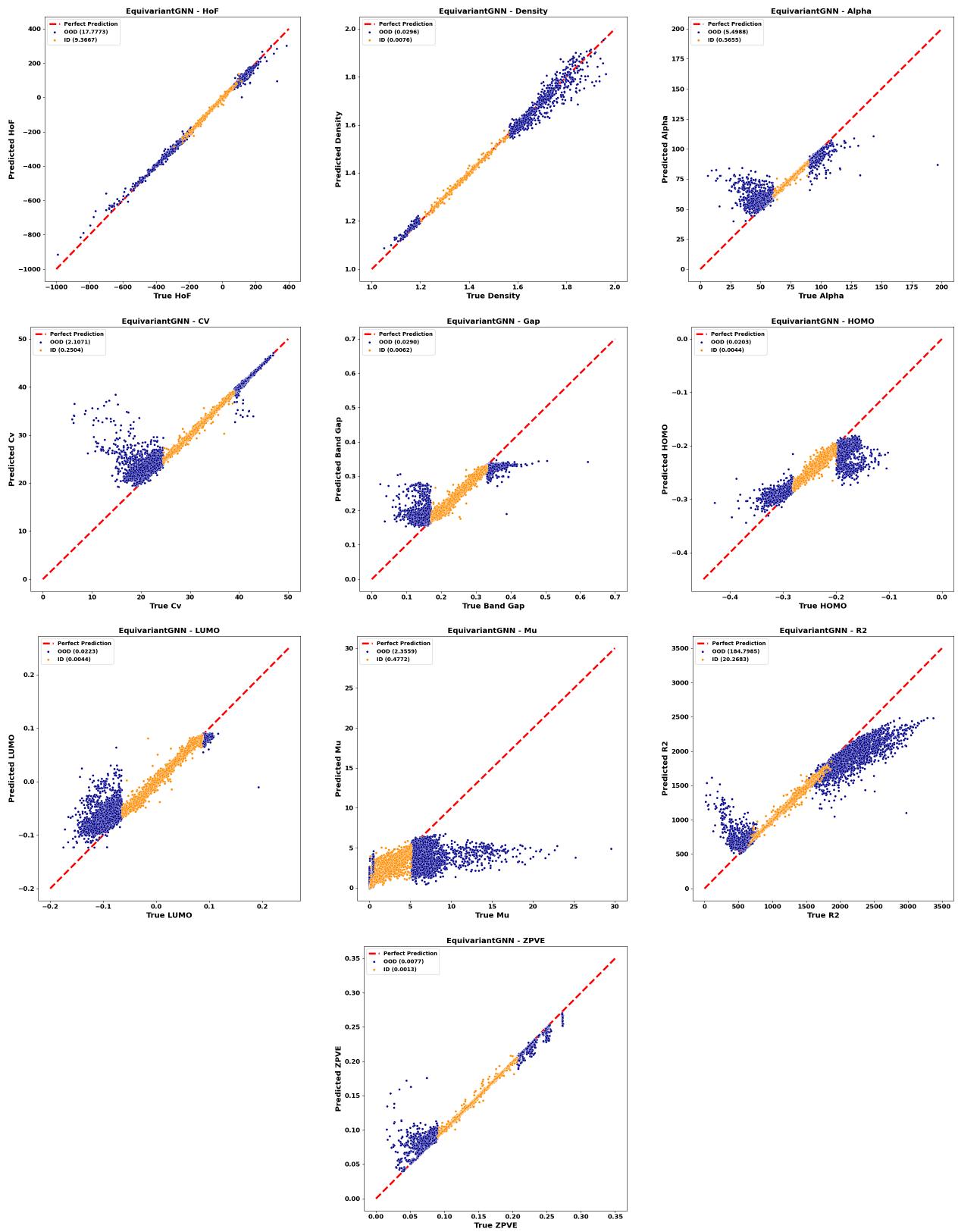
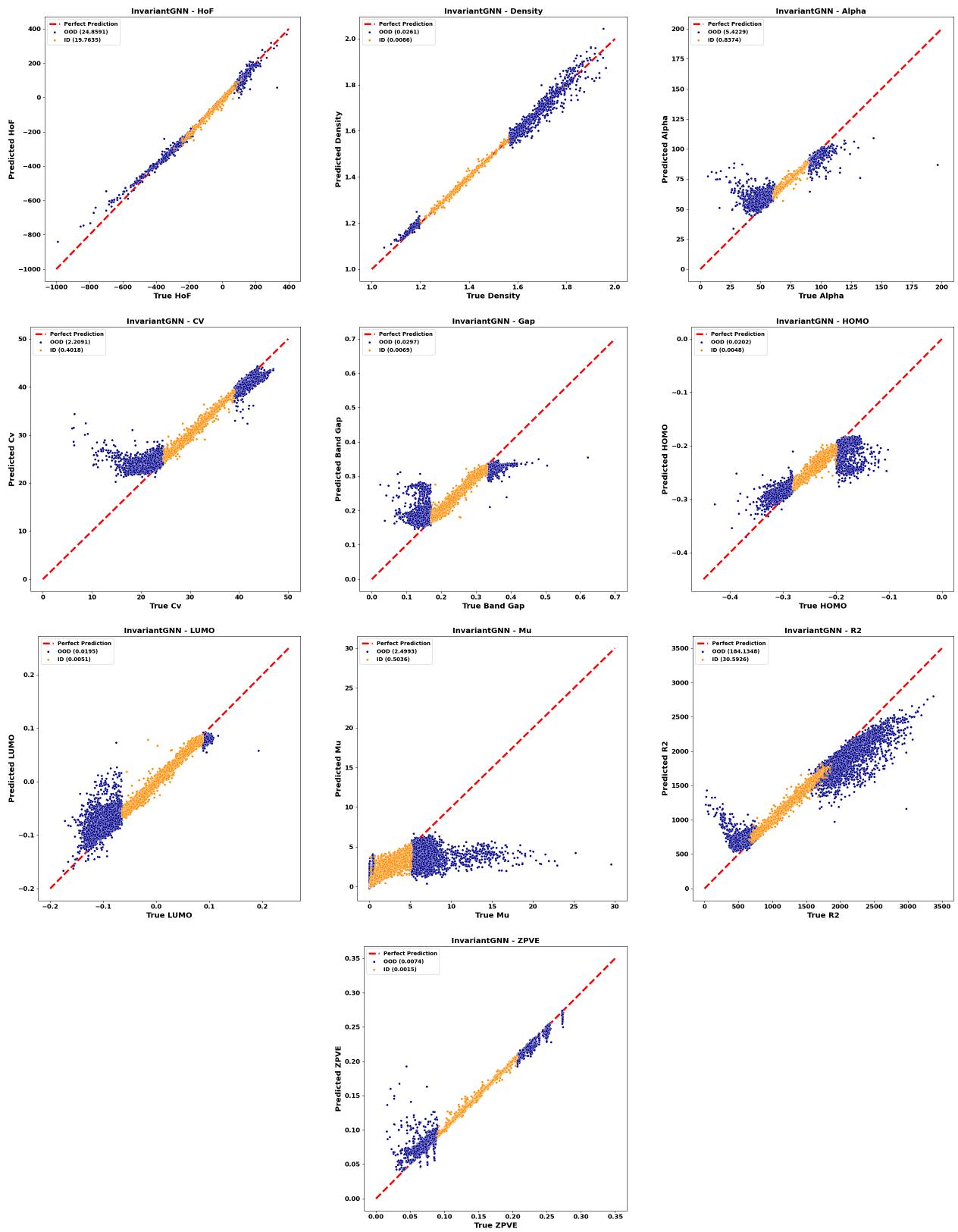**Figure 15.** Parity Plots for MoLFormer (with Pretraining) on 10K and QM9 OOD tasks.

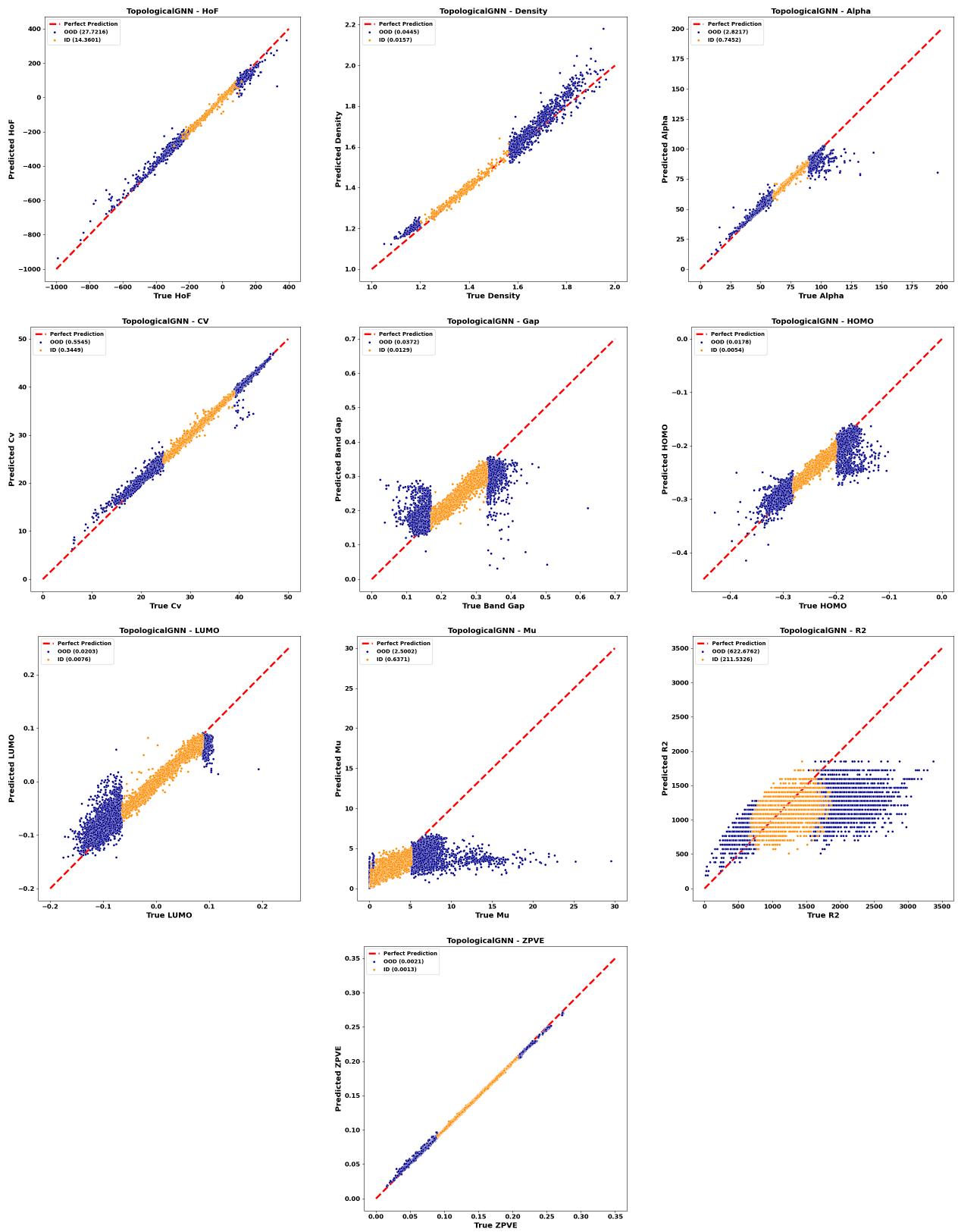**Figure 16.** Parity Plots for ChemBERTa (with Pretraining) on 10K and QM9 OOD tasks.

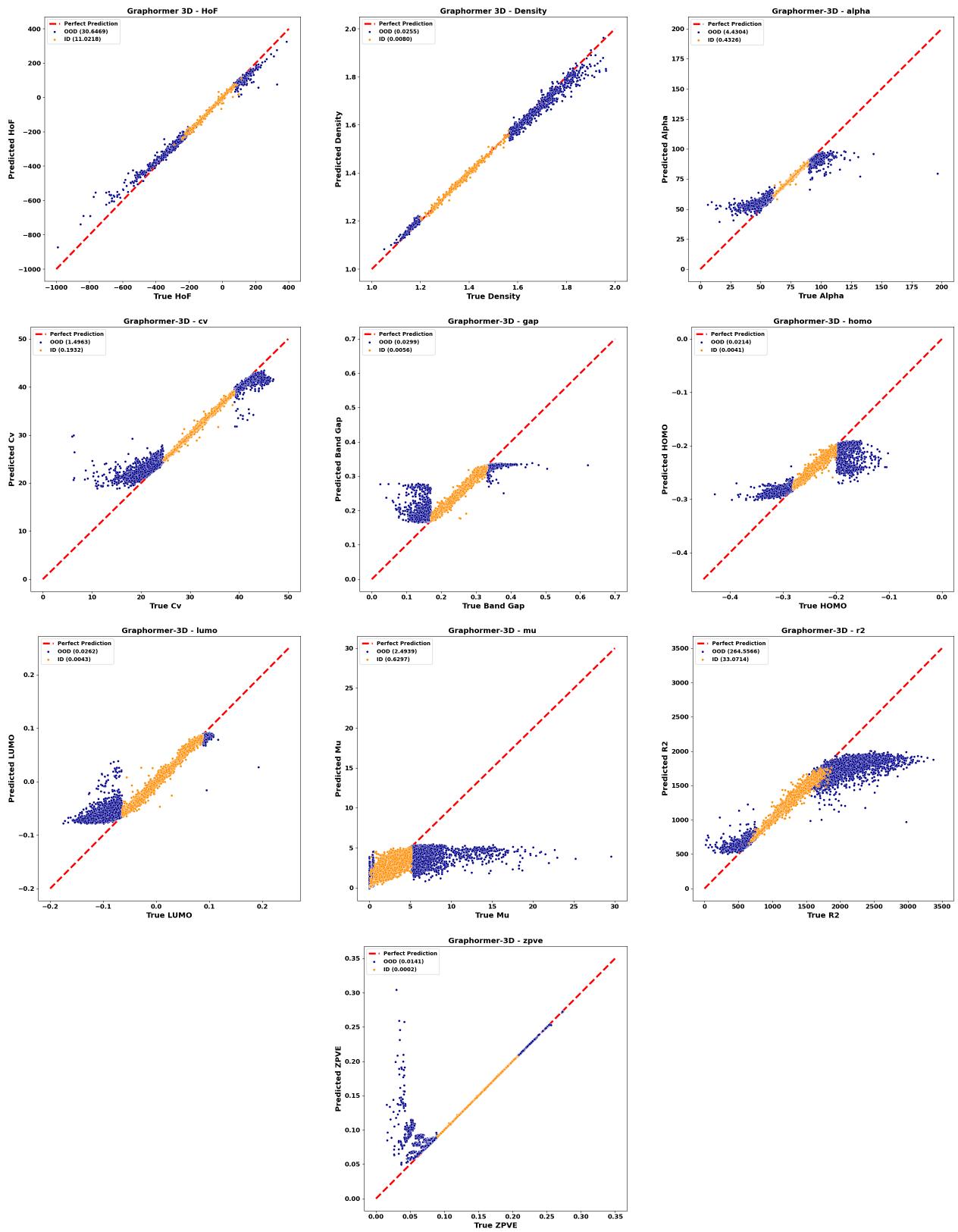**Figure 17.** Parity Plots for ModernBERT on 10K and QM9 OOD tasks.

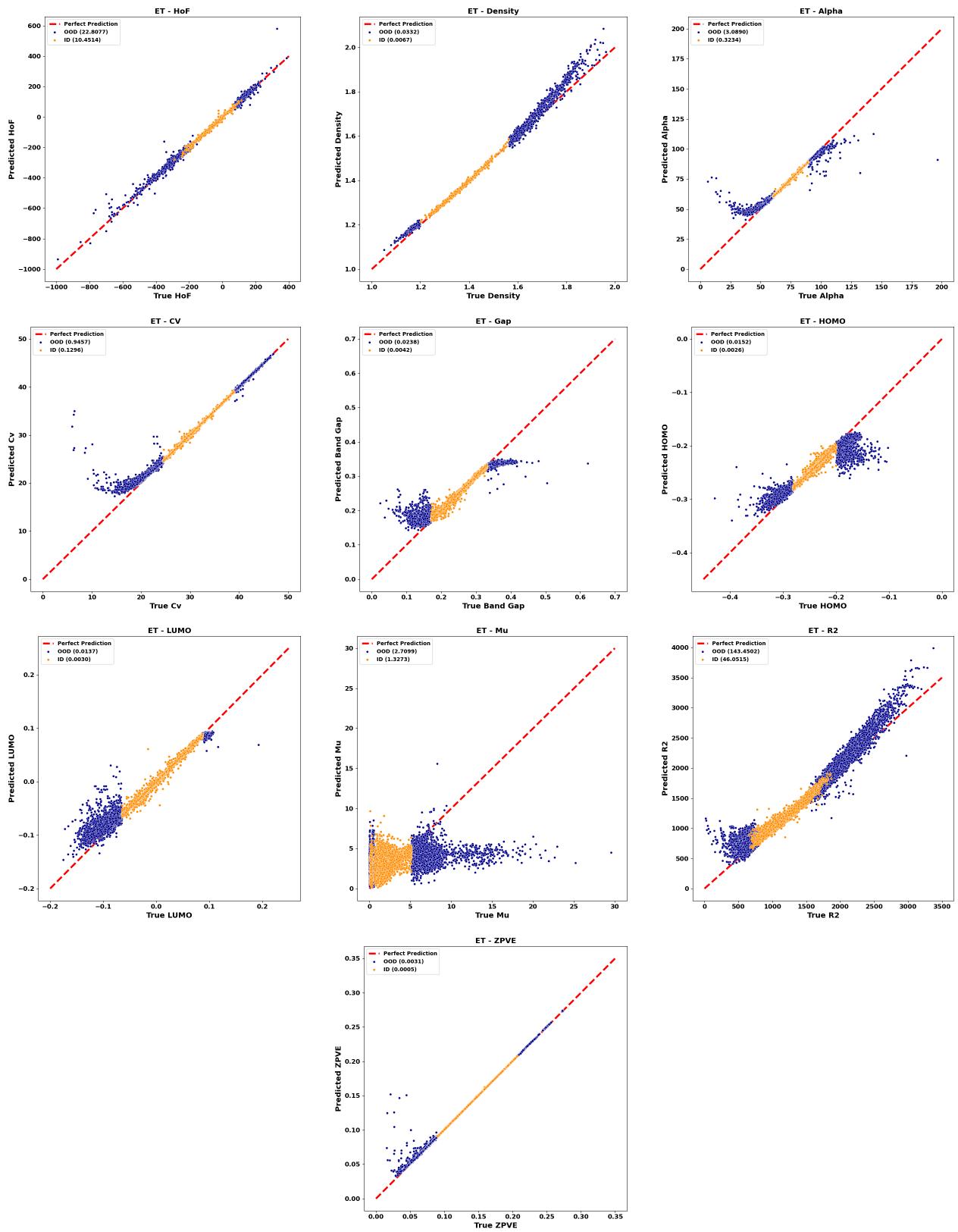**Figure 18.** Parity Plots for EGNN on 10K and QM9 OOD tasks.

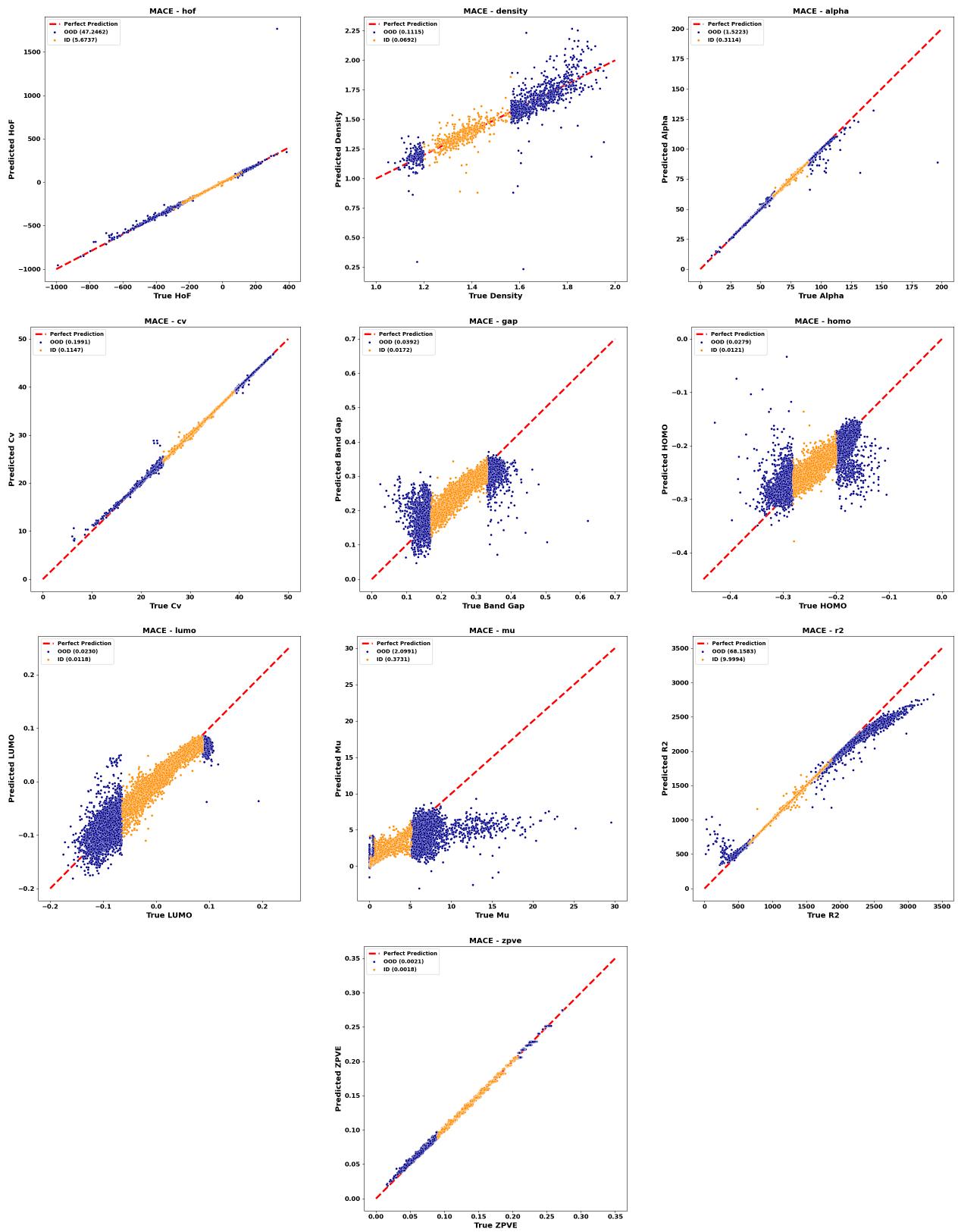**Figure 19.** Parity Plots for IGNN on 10K and QM9 OOD tasks.

**Figure 20.** Parity Plots for TGNN on 10K and QM9 OOD tasks.

**Figure 21.** Parity Plots for Graphormer(3D) on 10K and QM9 OOD tasks.

**Figure 22.** Parity Plots for TorchMD-ET on 10K and QM9 OOD tasks.

**Figure 23.** Parity Plots for MACE on 10K and QM9 OOD tasks.