
Machine Learning Force Fields with Data Cost Aware Training

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Machine learning force fields (MLFF) have been proposed to accelerate molecular
2 dynamics (MD) simulation, which finds widespread applications in chemistry and
3 biomedical research. Even for the most data-efficient MLFFs, reaching chemical
4 accuracy can require hundreds of frames of force and energy labels generated by
5 expensive quantum mechanical algorithms, which may scale as $O(n^3)$ to $O(n^7)$,
6 with n proportional to the number of basis functions. To address this issue, we
7 propose a multi-stage computational framework – ASTEROID, which lowers the
8 data cost of MLFFs by leveraging a combination of cheap inaccurate data and
9 expensive accurate data. The motivation behind ASTEROID is that inaccurate
10 data, though incurring large bias, can help capture the sophisticated structures
11 of the underlying force field. Therefore, we first train a MLFF model on a large
12 amount of inaccurate training data, employing a bias-aware loss function to prevent
13 the model from overfitting to the potential bias of this data. We then fine-tune the
14 obtained model using a small amount of accurate training data, which preserves the
15 knowledge learned from the inaccurate training data while significantly improving
16 the model’s accuracy. Moreover, we propose a variant of ASTEROID based on
17 score matching for the setting where the inaccurate training data are unlabeled.
18 Extensive experiments on MD datasets validate the efficacy of ASTEROID.

19 1 Introduction

20 Molecular dynamics (MD) simulation is a key technology driving scientific discovery in fields such
21 as chemistry, biophysics, and materials science [Alder and Wainwright, 1960, McCammon et al.,
22 1977]. By simulating the dynamics of molecules, important macro statistics such as the folding
23 probability of a protein [Tuckerman, 2010] or the density of new materials [Varshney et al., 2008]
24 can be estimated. These macro statistics are an essential part of many important applications such as
25 structure-driven drug design [Hospital et al., 2015] and battery development [Leung and Budzien,
26 2010]. Most MD simulation techniques share a common iterative structure: MD simulations calculate
27 the forces on each atom in the molecule, and use these forces to move the molecule forward to the
28 next state.

29 The fundamental challenge of MD simulation is how to efficiently calculate the forces at each
30 iteration. An exact calculation requires solving the Schrödinger equation, which is not feasible
31 for many-body systems [Berezin and Shubin, 2012]. Instead approximation methods such as the
32 Lennard-Jones potential [Johnson et al., 1993], Density Functional Theory (DFT, Kohn [2019]),
33 or Coupled Cluster Single-Double-Triple (CCSD(T), Scuseria et al. [1988]) are used. CCSD(T)
34 is seen as the gold-standard for force calculation, but is computationally expensive. In particular,

35 CCSD(T) has complexity $\mathcal{O}(n^7)$ with respect to the number of basis functions used along with a
 36 huge storage requirement [Chen et al., 2020]. To accelerate MD simulations while maintaining high
 37 accuracy, machine learning based force fields (MLFFs) have been proposed. MLFFs take a molecular
 38 configuration as input and then predict the forces on each atom in the molecule, consequently speeding
 39 up the force calculation step.

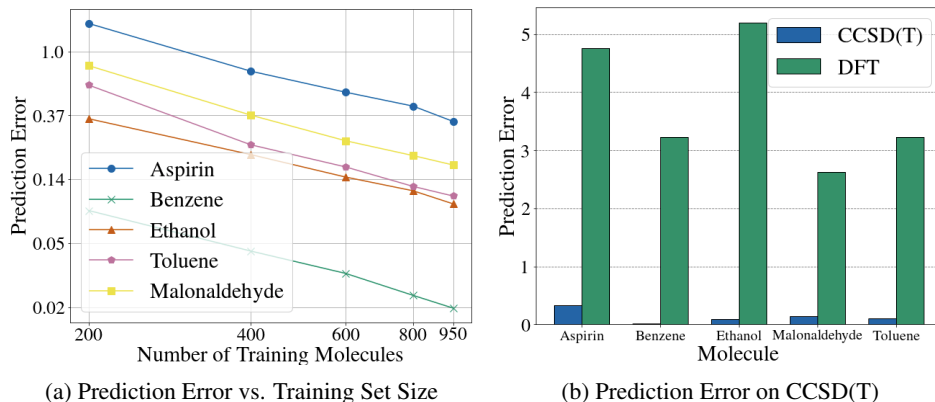


Figure 1: (a) Log-log plot of the number of training points versus the prediction error for deep force fields (b) Prediction error on CCSD labeled molecules for force fields trained on large amounts of DFT reference forces (100,000 configurations) and moderate amounts of CCSD reference forces (1000 configurations). In both cases the model architecture used is GemNet [Gasteiger et al., 2021].

40 Most recently, deep learning techniques for force fields have been developed, resulting in highly
 41 accurate force fields parameterized by large neural networks [Gasteiger et al., 2021, Batzner et al.,
 42 2022]. Despite their empirical success, these methods suffer from a critical drawback: *in order to*
 43 *train state-of-the-art machine learning force field models, a large amount of costly training data must*
 44 *be generated.* For example, to train a model at the CCSD(T) level of accuracy, at least a thousand
 45 CCSD(T) calculations must be done to construct the training set. This is computationally expensive
 46 due to the method’s $\mathcal{O}(n^7)$ cost.

47 A natural solution to this problem is to train on fewer data points. However, if the number of training
 48 points is decreased, the accuracy of the learned force fields quickly deteriorates. In our experiments,
 49 we empirically find that the prediction error and the number of training points roughly follow a
 50 power law relationship, with prediction error $\propto (\text{Number of Training Points})^{-1}$ [Müller et al., 1996,
 51 Cortes et al., 1993]. This can be seen in Figure 1a, where prediction error and train set size are
 52 observed to have a linear relationship with a slope of -1 when plotted on a log scale.

53 Another option is to train the force field model on less accurate but computationally cheap reference
 54 forces calculated using DFT [Kohn, 2019] or empirical force field methods [Johnson et al., 1993].
 55 However, these algorithms introduce undesirable bias into the force labels, meaning that the trained
 56 models will have poor performance. This phenomenon can be seen in Figure 1b, where models
 57 trained on large quantities of DFT reference forces are shown to perform poorly relative to force
 58 fields trained on moderate quantities of CCSD(T) reference forces. Therefore current methodologies
 59 are not sufficient for training force field models in low resource settings, as training on either small
 60 amounts of accurate data (i.e. from CCSD(T)) or large amounts of inaccurate data (i.e. from DFT or
 61 empirical force fields) will result in inaccurate force fields.

62 To address this issue, we propose to use both large amounts of inaccurate force field data and small
 63 amounts of accurate data to reduce the data generation cost needed to achieve highly accurate force
 64 fields. Our motivation is that computationally cheap data, though incurring large bias, can help
 65 capture the sophisticated structures of the underlying force field. Moreover, if treated properly, we
 66 can further reduce the bias of the obtained model by taking advantage of the accurate data.

67 Specifically, we propose a multi-stage computational framework – **data cosST aware tRaining of**
 68 **fOrce fIeIDs (ASTEROID)**. In the first stage, small amounts of accurate data are used to identify the

69 bias of force labels in a large but inaccurate dataset. In the second stage, the model is trained on the
70 large inaccurate dataset with a bias-aware loss function. This loss function generates smaller weights
71 for data points with larger bias, suppressing the effect of label noise on training. The inaccurately
72 trained model serves as a warm start for the third stage, where it is fine-tuned on the small and
73 accurate dataset. Together, these stages allow the model to learn from many molecular configurations
74 while incorporating highly accurate force data, significantly outperforming conventional methods
75 trained with similar data generation budgets.

76 Beyond using cheap labeled data to boost model performance, we also develop a method for the case
77 where a large amount of unlabeled molecular configurations are cheaply available [Smith et al., 2017,
78 Köhler et al., 2022]. Without labels, we cannot adopt the supervised learning approach. Instead,
79 we draw a connection to score matching, which learns the gradient of the log density function with
80 respect to each data point (called the score) [Hyvärinen, 2005]. In the context of molecular dynamics,
81 we notice that if the log density function is proportional to the energy of each molecule, then the
82 score function with respect to a molecule’s position is equal to the force on the molecule. Based on
83 this insight, we show that the supervised force matching problem can be tackled in an unsupervised
84 manner. This unsupervised approach can then be incorporated into the ASTEROID framework,
85 improving performance when limited data is available.

86 We demonstrate the effectiveness of our framework with extensive experiments on different force
87 field data sets and downstream simulation tasks. We use two popular model architectures, GemNet
88 [Gasteiger et al., 2021] and EGNN [Satorras et al., 2021], and verify the performance of our method
89 in a variety of settings. These experiments show that ASTEROID can lead to significant gains when
90 either DFT reference forces or empirical force field forces are viewed as inaccurate data and CCSD(T)
91 configurations are used as accurate data. In addition, we show that we can learn accurate forces via
92 the connection to score matching, and that using this objective in the second stage of training can
93 improve performance on both DFT and CCSD(T) datasets.

94 2 Background

95 \diamond **Machine Learning Force Fields.** Recent years have seen a surge of interest in MLFFs. Much
96 of this work has focused on developing machine learning architectures that have physically correct
97 equivariances, resulting in large graph neural networks that can generate highly accurate force and
98 energy predictions [Gasteiger et al., 2021, Satorras et al., 2021, Batzner et al., 2022]. Two popular
99 architectures are EGNN and GemNet. Both models are translation invariant, rotationally equivariant,
100 and permutation equivariant. EGNN is a smaller model and is often used when limited resources
101 are available. The GemNet architecture is significantly larger and more refined than the EGNN
102 architecture, modeling various types of inter-atom interactions. GemNet is therefore more powerful
103 and can achieve state-of-the-art performance, but requires more resources to train.

104 It has been observed that modern MLFFs often cannot achieve sufficient test accuracy to be reliable
105 for MD simulations [Stocker et al., 2022]. Critically, the accuracy of deep force fields such as GemNet
106 and EGNN is highly dependent on the size and quality of the training dataset. With limited training
107 data, MLFFs cannot achieve the required accuracy for usefulness, preventing their application in
108 settings where data is expensive to generate (e.g. large molecules). The amount of resources needed
109 to train is therefore a key bottleneck preventing the widespread use of MLFFs.

110 \diamond **Data Generation Cost.** The training data for MLFFs can be generated by a variety of force
111 calculation methods. These methods exhibit an accuracy cost tradeoff: accurate reference forces
112 from methods such as CCSD(T) require high computational costs to generate reference forces, while
113 inaccurate reference forces from methods such as DFT and empirical force fields can be generated
114 fairly quickly. Concretely, CCSD(T) is highly accurate but has $\mathcal{O}(n^7)$ complexity, DFT is less
115 accurate with complexity $\mathcal{O}(n^3)$, and empirical force fields are inaccurate but could have complexity
116 as low as $\mathcal{O}(n)$ [Lin et al., 2019, Ratcliff et al., 2017]. CCSD(T) is typically viewed as the gold
117 standard for calculating reference forces, but its computational costs often make it impractical for MD
118 simulation (it has been estimated that “a nanosecond-long MD trajectory for a single ethanol molecule
119 executed with the CCSD(T) method would take roughly a million CPU years on modern hardware”)

120 [Chmiela et al., 2018]. Due to this large expense, MLFF training data is typically generated first with
 121 MD simulations driven by DFT or empirical force fields. These simulations generate a large number
 122 of molecular configurations, and then CCSD(T) reference forces are computed for a small portion of
 123 these configurations. Therefore, a large amount of inaccurately labeled molecular configurations are
 124 often available along with the accurate CCSD(T) labeled data.

125 3 ASTEROID

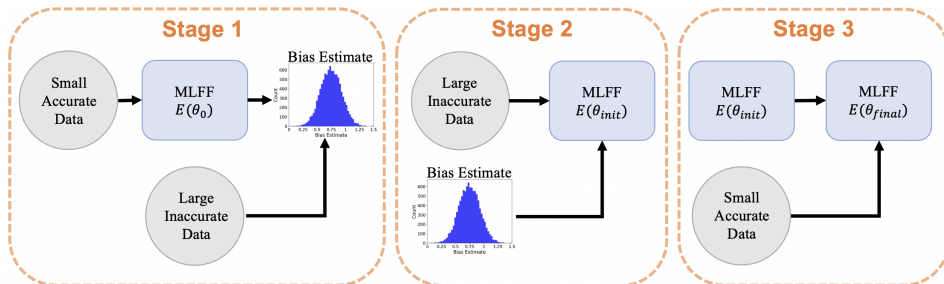


Figure 2: Asteroid workflow diagram.

126 To reduce the data generation cost needed to train MLFFs, we propose a multi-stage training
 127 framework, ASTEROID, to learn from a combination of both cheaply available inaccurate data and
 128 more expensive accurate data.

129 **Preliminaries.** For a molecule with k atoms, we denote a configuration (the positions of its atoms in
 130 3D) of this molecule as $x \in \mathbb{R}^{3k}$, its respective energy as $E(x) \in \mathbb{R}$, and its force as $F(x) \in \mathbb{R}^{3k}$.
 131 We denote the accurately labeled data as $\mathcal{D}_A = \{(x_1^a, e_1^a, f_1^a), \dots, (x_N^a, e_N^a, f_N^a)\}$ and the inaccurately
 132 labeled data as $\mathcal{D}_I = \{(x_1^n, e_1^n, f_1^n), \dots, (x_M^n, e_M^n, f_M^n)\}$, where (x_i^a, e_i^a, f_i^a) represents the position,
 133 potential energy, and force of the i th accurately labeled molecule, (similarly (x_j^n, e_j^n, f_j^n) for the j th
 134 inaccurately labeled data). Conventional methods train a force field model $E(\cdot; \theta)$ with parameters θ
 135 on the accurate data by minimizing the loss

$$\min_{\theta} \mathcal{L}(\mathcal{D}_A, \theta) = \frac{(1-\rho)}{3N} \sum_{i=1}^N \ell_f(f_i^a, \nabla_x E(x_i^a; \theta)) + \frac{\rho}{N} \sum_{i=1}^N \ell_e(e_i^a, E(x_i^a; \theta)), \quad (1)$$

136 where ℓ_f is the loss function for the force prediction, and ℓ_e is the loss function for the energy
 137 prediction. Here the force is denoted by $\nabla_x E(x; \theta)$, i.e., the gradient of the energy $E(x; \theta)$ w.r.t. to
 138 the input x . In practice, most of the emphasis is placed on the force prediction, e.g. $\rho = 0.001$.

139 3.1 Bias Identification

140 The goal of ASTEROID is to leverage cheap MD simulation data to boost MLFF accuracy. However,
 141 the approximation algorithms used to generate cheap data \mathcal{D}_I introduce a large amount of bias into
 142 some force labels f^n , which may significantly hurt accuracy. Motivated by this phenomenon, we aim
 143 to identify the most biased force labels so that we can avoid overfitting the bias during training. To do
 144 so, we use small amounts of accurately labeled data \mathcal{D}_A to identify the levels of bias in the inaccurate
 145 dataset \mathcal{D}_I . Specifically, we train a force field model by minimizing $\mathcal{L}(\mathcal{D}_A, \theta)$ (Eq. 1), the loss over
 146 the accurate data, to get parameters θ_0 . Although the resulting model $E(\cdot; \theta_0)$ will not necessarily
 147 have good prediction performance because of the limited amount of training data, it can still help
 148 estimate the bias of the inaccurate data. For every configuration x_j^n in the inaccurate dataset \mathcal{D}_I , we
 149 suspect it to have a large bias if there is a large discrepancy between its force label f_j^n and the force
 150 label predicted by the accurately trained model $\nabla_x E(x_j^n; \theta_0)$. We can therefore use this discrepancy
 151 as a surrogate for bias, i.e. $B(x_j^n) = \|\nabla_x E(x_j^n; \theta_0) - f_j^n\|_1$.

152 3.2 Bias-Aware Training with Inaccurate Data

153 In the second stage of our framework, we train a force field model $E(\cdot; \theta_{\text{init}})$ from scratch on large
 154 amounts of *inaccurately labeled data* \mathcal{D}_I . Although this data can effectively capture the intrinsic

155 problem structure, the high levels of bias on some data points may propagate to the final model and
 156 harm generalization performance. To avoid over-fitting to the biased force labels, we use a bias-aware
 157 loss function that weighs the inaccurate data according to their bias. In particular, we use the weights
 158 $w_j = \exp(-B(x_j^n)/\gamma)$ for configuration x_j^n , where γ is a hyperparameter to be tuned. In this way,
 159 low-bias points are given higher importance and high-bias points are treated more carefully. We then
 160 minimize the bias-aware loss function

$$\min_{\theta} \mathcal{L}_w(\mathcal{D}_I, \theta) = (1 - \rho) \sum_{i=1}^M w_i \cdot \ell_f(f_i^n, \nabla_x E(x_i^n; \theta)) + \rho \sum_{i=1}^M w_i \cdot \ell_e(e_i^n, E(x_i^n; \theta)) \quad (2)$$

161 to get parameters θ_{init} , resulting in the initial estimate of the MLFF $E(\cdot; \theta_{\text{init}})$.

162 3.3 Fine-Tuning over Accurate Data

163 The model $E(\cdot; \theta_{\text{init}})$ contains information useful to the force prediction problem, but may still
 164 contain bias because it is trained on inaccurately labeled data \mathcal{D}_I . Therefore, we further refine it
 165 using accurately labeled data \mathcal{D}_A . Specifically, we use $E(\cdot; \theta_{\text{init}})$ as initialization for our final stage,
 166 in which we fine-tune the model over the accurate data by minimizing $\mathcal{L}(\mathcal{D}_A, \theta_{\text{final}})$ (Eq. 1). The full
 167 ASTEROID framework is illustrated in Figure 2.

168 4 ASTEROID for Unlabeled Data

In several settings, molecular configurations are generated without force labels, either because they
 are not generated via MD simulation (e.g. normal mode sampling, Smith et al. [2017]) or because
 the forces are not stored during the simulation [Köhler et al., 2022]. Although these unlabeled
 configurations may be cheaply available, they are not generated for the purpose of learning force
 fields and have not been used in existing literature. Here, we show that the unlabeled configurations
 can be used to obtain an initial estimate of the force field, which can then be further fine-tuned on
 accurate data. More specifically, we consider a molecular system where the number of particles,
 volume, and temperature are constant (NVT ensemble). Let x refer to a molecule’s configuration and
 $E(x)$ refer to the corresponding potential energy. It is known that x follows a Boltzmann distribution,
 i.e.

$$p(x) = \frac{1}{Z} \exp\left(-\frac{1}{k_{\beta} T} E(x)\right),$$

169 where Z is a normalizing constant, T is the temperature, and k_{β} is the Boltzmann constant. In practice,
 170 configurations generated using normal mode sampling [Unke et al., 2021] or via a sufficiently long
 171 NVT MD simulation follow a Boltzmann distribution.

172 Recall that we model the energy $E(x)$ as $E(x; \theta)$, and the force can be calculated as $F(x; \theta) =$
 173 $\nabla_x E(x; \theta)$. It follows from Hyvärinen [2005] that we can learn the score function of the Boltzmann
 174 distribution using score matching, where the score function is defined as the gradient of the log density
 175 function $\nabla_x \log p(x)$. In our case, we observe that the force on a configuration x is proportional
 176 to the score function, i.e., $F(x) \propto \nabla_x \log p(x)$. Therefore, we can use score matching to learn the
 177 forces by minimizing the unsupervised loss

$$L(\theta) = \mathbb{E}_{p(x)} \left[\frac{1}{\beta} \text{Tr}[\nabla_x F(x; \theta)] + \frac{1}{2} \|F(x; \theta)\|^2 \right], \quad (3)$$

178 where $\beta = \frac{1}{k_{\beta} T}$. A derivation can be found in Appendix A.5. Although this objective allows us
 179 to solve the force matching problem in an unsupervised manner, the unsupervised loss is difficult
 180 to optimize in practice. To reduce the cost of solving Eq. 3, we adopt sliced score matching [Song
 181 et al., 2020]. Sliced score matching takes advantage of random projections to significantly reduce the
 182 cost of solving Eq. 3, allowing us to apply score matching to large neural models such as GemNet.

183 In our experiments, we find that score matching does not match the accuracy of CCSD(T) force labels.
 184 Instead, we can think of score-matching as a form of inaccurate training. We therefore use score
 185 matching as an alternative to stages one and two of the ASTEROID framework. That is, we minimize
 186 Eq. 3 to get θ_{init} , after which the model is fine-tuned on the accurate data.

187 5 Experiments

188 For our main experiments, we evaluate ASTEROID on MLFF datasets and downstream MD simula-
189 tion tasks. For ASTEROID, we consider three settings: using DFT data to enhance CCSD(T) training,
190 using empirical force field data to enhance CCSD(T) training, and using unlabeled configurations to
191 enhance CCSD(T) training. In each setting, we evaluate the performance of ASTEROID and standard
192 training over a variety of data generation budgets.

193 5.1 Datasets and Models

194 For the CCSD(T) data, we use MD17@CCSD, which contains 1,000 configurations labeled at the
195 CCSD(T) and CCSD level of accuracy for five molecules [Chmiela et al., 2017]. For DFT data, we
196 use the MD17 dataset, which contains molecular configurations labeled at the DFT level of accuracy
197 [Chmiela et al., 2017]. For the empirical force field data, we generate 100,000 configurations for each
198 molecule using the OpenMM empirical force field software [Eastman et al., 2017]. For the unlabeled
199 datasets, we use MD17 with the force labels removed.

200 The MD17 datasets do not release the computational cost of data generation, but when we replicate
201 their experiments, we find that CCSD(T) labels cost roughly 40 times more than DFT labels. However,
202 the difference in cost between CCSD(T) and DFT labels may change drastically depending on the
203 implementation of each method. Therefore we evaluate the performance of ASTEROID when
204 CCSD(T) force labels are 20, 40, and 80 times more expensive than DFT force labels. Note that the
205 cost of empirical force labels is essentially negligible (more than 10,000 times cheaper) compared to
206 CCSD(T) labels [Folmsbee and Hutchison, 2021].

207 In each setting, we compare standard training with 250, 450, 650, or 850 CCSD(T) training samples
208 with ASTEROID. For ASTEROID, we use either 1000, 2000, or 4000 DFT datapoints (corresponding
209 to cost ratios of 20:1, 40:1, and 80:1 for DFT and CCSD(T) labels), and 200, 400, 600, or 800
210 CCSD(T) data points. The computational budget of standard training and ASTEROID are therefore
211 equivalent. A validation set of size 50 and a test set of size 500 are used in all experiments.

212 We implement our method on GemNet and EGNN. For GemNet we use the same model parameters
213 as Gasteiger et al. [2021]. For EGNN, we use a 5-layer model and an embedding size of 128. When
214 training with inaccurate data, we train with a batch size of 16 and stop training when the loss stabilizes.
215 In the fine-tuning stage, we use a batch size of 10 and train for a maximum of 2000 epochs. To tune
216 the bias aware loss parameter γ , we search in the set $\{0.1, 0.5, 1.0, 2.0\}$ and select the model with
217 the lowest validation loss. Comprehensive experimental details are deferred to Appendix A.6.

218 5.2 Enhancing Force Fields with DFT

219 We display the results for using DFT data to enhance CCSD(T) training in Figure 3 for GemNet
220 and Figure 4 for EGNN. From these figures, we can see that ASTEROID can outperform standard
221 training for all amounts of data and cost ratios. Using larger amounts of inaccurate data can
222 significantly reduce prediction error, but the 20:1 cost ratio already has large performance gains
223 over standard training. When applied to GemNet in low resource settings, ASTEROID reduces
224 the average prediction error by 39.4% and improves sample efficiency by a factor of 2. For EGNN,
225 ASTEROID improves prediction error by 56% and increases sample efficiency by more than 3 times.
226 The large performance increase for EGNN may be due to the fact that the EGNN architecture has
227 less inductive bias than GemNet, and therefore may struggle to learn the structures of the underlying
228 force field with only a small amount of data.

229 5.3 Enhancing Force Fields with Empirical Force Calculation

230 We present the results for empirical force field in Table 1. Additional results for GemNet can be
231 found in Appendix A.7. Again we find that ASTEROID significantly outperforms the supervised
232 baseline, improving prediction accuracy by 36% for GemNet and by 17% for EGNN. The good
233 performance on empirical force fields indicates that ASTEROID is relatively robust to the label noise
234 on the inaccurate data.

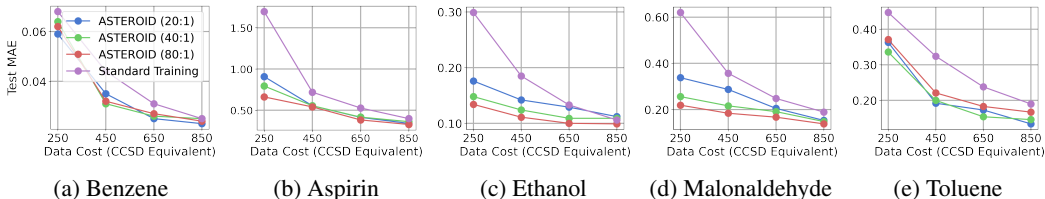


Figure 3: Main results for GemNet when DFT data is viewed as inaccurate. The ratio refers to the number of DFT calculations that are equivalent to one CCSD(T) calculation. The results are measure in kcal/mol/Å, averaged across dimensions and atoms.

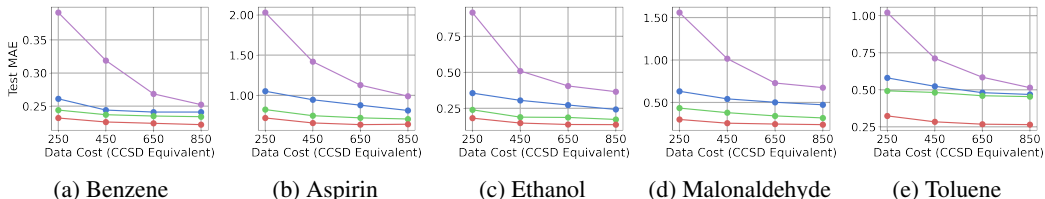


Figure 4: Main results for EGNN when DFT data is viewed as inaccurate.

Table 1: Test MAE of ASTEROID with empirical force field data. The results are measure in kcal/mol/Å, averaged across dimensions and atoms. The training set for the fine-tuning stage contains 200 molecules labeled at the CCSD(T) level. ‘‘Malo.’’ refers to malonaldehyde and ‘‘Standard Tr.’’ refers to standard training.

	Aspirin	Benzene	Malonaldehyde	Toluene	Ethanol
GemNet					
Standard Training	1.554	0.083	0.801	0.591	0.348
ASTEROID	0.843	0.048	0.516	0.337	0.301
EGNN					
Standard Training	1.897	0.297	1.466	0.777	0.840
ASTEROID	1.314	0.268	1.341	0.664	0.637

235 5.4 Enhancing Force Fields with unlabeled Molecules

236 We first verify that our proposed score matching approach can learn the forces on unlabeled molecules
 237 by comparing the prediction accuracy of models trained by score matching with models trained on
 238 supervised data (DFT and empirical force fields). We measure prediction accuracy on CCSD(T)
 239 datasets and show the results in Figure 7. Surprisingly, we find that the prediction error of score
 240 matching is between that of DFT and empirical force fields. This indicates that relatively accurate
 241 force predictions can be obtained by only solving the unsupervised loss in Eq. 3.

242 Next we apply ASTEROID to settings where unlabeled data is available by fine-tuning the model
 243 obtained from score matching. We present the results in Table 2, where we find that ASTEROID
 244 can improve prediction accuracy by 18% for GemNet and 4% for EGNN. If unsupervised data can
 245 be generated cheaply (i.e. through normal mode sampling), then our approach can be used to boost
 246 the performance of MLFFs with little additional cost.

247 6 Discussion

248 **Related Work.** There are several works which we compare ASTEROID with.

249 \diamond Δ -ML [Ramakrishnan et al., 2015, Bogojeski et al., 2020], learns the difference between inaccurate
 250 (DFT) and accurate (CCSD(T)) force predictions, therefore speeding up MD simulation while

Table 2: Accuracy of ASTEROID with unlabeled molecular configurations. The results are measure in kcal/mol/Å, averaged across dimensions and atoms. The training set for the fine-tuning stage contains 200 CCSD(T) labeled molecules.

	Aspirin	Benzene	Malonaldehyde	Toluene	Ethanol
GemNet					
Standard Training	1.554	0.083	0.801	0.591	0.348
ASTEROID	0.928	0.093	0.629	0.475	0.314
EGNN					
Standard Training	1.897	0.297	1.466	0.777	0.840
ASTEROID	1.756	0.305	1.382	0.740	0.823

251 maintaining high accuracy. However, this approach requires a DFT calculation to be done during
 252 inference, greatly increasing inference time compared to ASTEROID or standard MLFFs [Folmsbee
 253 and Hutchison, 2021].

254 \diamond **ANI-1ccx** [Smith et al., 2019, Deringer et al., 2020] train an MLFF on a huge DFT dataset
 255 comprised of many molecules, and then finetune on many CCSD(T) labeled molecules with a goal of
 256 learning a general MLFF. Notably, the method from Smith et al. [2017] only trains on equilibrium
 257 states and may not work well for MD trajectory data. To compare ANI-1ccx with ASTEROID, we
 258 evaluate the provided model checkpoint in the zero-shot setting (as in [Smith et al., 2019]) and when
 259 finetuned on each MD17 molecule. Note that the data generation cost of ANI-1ccx is much more
 260 expensive than ASTEROID, using 2,500 times more CCSD(T) data and 500 times more DFT data.

261 \diamond **sGDML** [Chmiela et al., 2019] is a kernel-based MLFF method that can perform well when limited
 262 training data is available by incorporating relevant physical constraints into the MLFF.

263 As can be seen in Table 3, ASTEROID trained MLFFs can achieve lower test errors than all of the
 264 baselines except Δ -ML. However, since Δ -ML requires a DFT calculation during inference, MD
 265 simulation with Δ -ML will take 100 to 1000 times longer than with ASTEROID [Folmsbee and
 266 Hutchison, 2021, Gasteiger et al., 2021]. Therefore ASTEROID results in the most useful force
 267 fields out of all the baselines, while having a smaller or equivalent data generation cost.

Table 3: Accuracy of ASTEROID compared with competitive baselines with a data budget of 250 CCSD(T) points. FT refers to fine-tuning ANI-1ccx on MD17@CCSD. The model is GemNet.

	Aspirin	Benzene	Malonaldehyde	Toluene	Ethanol
ANI-1	1.897	0.297	1.466	0.777	0.840
ANI-1 (FT)	1.314	0.268	1.341	0.664	0.637
Δ -ML	0.801	–	0.182	0.350	–
sGDML	1.727	0.097	0.923	0.478	0.902
ASTEROID	0.908	0.059	0.338	0.306	0.176

268 **Ablation.** We conduct a detailed ablation study in Appendix A.3, which shows ASTEROID is fairly
 269 robust to hyperparameter selection.

270 **MD Simulation** We show the results of MD simulation results in Appendix A.2, where we observe
 271 ASTEROID can result in stable simulation.

272 **Asteroid with Multiple Molecules.** We try ASTEROID on multiple molecules simultaneously in
 273 Appendix A.1. We find mixed results, indicating this could be an exciting direction to explore further.

274 **References**

- 275 Berni Julian Alder and Thomas Everett Wainwright. Studies in molecular dynamics. ii. behavior of a
276 small number of elastic spheres. *The Journal of Chemical Physics*, 33(5):1439–1451, 1960.
- 277 J Andrew McCammon, Bruce R Gelin, and Martin Karplus. Dynamics of folded proteins. *nature*,
278 267(5612):585–590, 1977.
- 279 Mark Tuckerman. *Statistical mechanics: theory and molecular simulation*. Oxford university press,
280 2010.
- 281 Vikas Varshney, Soumya S Patnaik, Ajit K Roy, and Barry L Farmer. A molecular dynamics study of
282 epoxy-based networks: cross-linking procedure and prediction of molecular and material properties.
283 *Macromolecules*, 41(18):6837–6842, 2008.
- 284 Adam Hospital, Josep Ramon Goñi, Modesto Orozco, and Josep L Gelpí. Molecular dynamics simu-
285 lations: advances and applications. *Advances and applications in bioinformatics and chemistry:*
286 *AABC*, 8:37, 2015.
- 287 Kevin Leung and Joanne L Budzien. Ab initio molecular dynamics simulations of the initial stages of
288 solid–electrolyte interphase formation on lithium ion battery graphitic anodes. *Physical Chemistry*
289 *Chemical Physics*, 12(25):6583–6586, 2010.
- 290 Feliks Aleksandrovich Berezin and Mikhail Shubin. *The Schrödinger Equation*, volume 66. Springer
291 Science & Business Media, 2012.
- 292 J Karl Johnson, John A Zollweg, and Keith E Gubbins. The lennard-jones equation of state revisited.
293 *Molecular Physics*, 78(3):591–618, 1993.
- 294 Walter Kohn. Density functional theory. *Introductory Quantum Mechanics with MATLAB: For Atoms,*
295 *Molecules, Clusters, and Nanocrystals*, 2019.
- 296 Gustavo E Scuseria, Curtis L Janssen, and Henry F Schaefer Iii. An efficient reformulation of the
297 closed-shell coupled cluster single and double excitation (ccsd) equations. *The Journal of Chemical*
298 *Physics*, 89(12):7382–7387, 1988.
- 299 Jiu-Li Chen, Tao Sun, Yi-Bo Wang, and Weizhou Wang. Toward a less costly but accurate calculation
300 of the ccsd (t)/cbs noncovalent interaction energy. *Journal of Computational Chemistry*, 41(13):
301 1252–1260, 2020.
- 302 Johannes Gasteiger, Florian Becker, and Stephan Günnemann. Gemnet: Universal directional
303 graph neural networks for molecules. *Advances in Neural Information Processing Systems*, 34:
304 6790–6802, 2021.
- 305 Simon Batzner, Albert Musaelian, Lixin Sun, Mario Geiger, Jonathan P Mailoa, Mordechai Kornbluth,
306 Nicola Molinari, Tess E Smidt, and Boris Kozinsky. E (3)-equivariant graph neural networks for
307 data-efficient and accurate interatomic potentials. *Nature communications*, 13(1):1–11, 2022.
- 308 K-R Müller, Michael Finke, Noboru Murata, Klaus Schulten, and Shun-ichi Amari. A numerical
309 study on learning curves in stochastic multilayer feedforward networks. *Neural Computation*, 8(5):
310 1085–1106, 1996.
- 311 Corinna Cortes, Lawrence D Jackel, Sara Solla, Vladimir Vapnik, and John Denker. Learning curves:
312 Asymptotic values and rate of convergence. *Advances in neural information processing systems*, 6,
313 1993.
- 314 Justin S Smith, Olexandr Isayev, and Adrian E Roitberg. Ani-1, a data set of 20 million calculated
315 off-equilibrium conformations for organic molecules. *Scientific data*, 4(1):1–8, 2017.
- 316 Jonas Köhler, Yaoyi Chen, Andreas Krämer, Cecilia Clementi, and Frank Noé. Force-matching
317 coarse-graining without forces. *arXiv preprint arXiv:2203.11167*, 2022.

- 318 Aapo Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of*
319 *Machine Learning Research*, 6(4), 2005.
- 320 Victor Garcia Satorras, Emiel Hoogetboom, and Max Welling. E (n) equivariant graph neural networks.
321 In *International conference on machine learning*, pages 9323–9332. PMLR, 2021.
- 322 Sina Stocker, Johannes Gasteiger, Florian Becker, Stephan Günnemann, and Johannes T Margraf.
323 How robust are modern graph neural network potentials in long and hot molecular dynamics
324 simulations? *Machine Learning: Science and Technology*, 3(4):045010, 2022.
- 325 Lin Lin, Jianfeng Lu, and Lexing Ying. Numerical methods for kohn–sham density functional theory.
326 *Acta Numerica*, 28:405–539, 2019.
- 327 Laura E Ratcliff, Stephan Mohr, Georg Huhs, Thierry Deutsch, Michel Masella, and Luigi Genovese.
328 Challenges in large scale quantum mechanical calculations. *Wiley Interdisciplinary Reviews:*
329 *Computational Molecular Science*, 7(1):e1290, 2017.
- 330 Stefan Chmiela, Huziel E Saucedo, Klaus-Robert Müller, and Alexandre Tkatchenko. Towards exact
331 molecular dynamics simulations with machine-learned force fields. *Nature communications*, 9(1):
332 1–10, 2018.
- 333 Oliver T Unke, Stefan Chmiela, Huziel E Saucedo, Michael Gastegger, Igor Poltavsky, Kristof T
334 Schütt, Alexandre Tkatchenko, and Klaus-Robert Müller. Machine learning force fields. *Chemical*
335 *Reviews*, 121(16):10142–10186, 2021.
- 336 Yang Song, Sahaj Garg, Jiaxin Shi, and Stefano Ermon. Sliced score matching: A scalable approach
337 to density and score estimation. In *Uncertainty in Artificial Intelligence*, pages 574–584. PMLR,
338 2020.
- 339 Stefan Chmiela, Alexandre Tkatchenko, Huziel E Saucedo, Igor Poltavsky, Kristof T Schütt, and
340 Klaus-Robert Müller. Machine learning of accurate energy-conserving molecular force fields.
341 *Science advances*, 3(5):e1603015, 2017.
- 342 Peter Eastman, Jason Swails, John D Chodera, Robert T McGibbon, Yutong Zhao, Kyle A Beauchamp,
343 Lee-Ping Wang, Andrew C Simmonett, Matthew P Harrigan, Chaya D Stern, et al. Openmm 7:
344 Rapid development of high performance algorithms for molecular dynamics. *PLoS computational*
345 *biology*, 13(7):e1005659, 2017.
- 346 Dakota Folmsbee and Geoffrey Hutchison. Assessing conformer energies using electronic structure
347 and machine learning methods. *International Journal of Quantum Chemistry*, 121(1):e26381,
348 2021.
- 349 Raghunathan Ramakrishnan, Pavlo O Dral, Matthias Rupp, and O Anatole Von Lilienfeld. Big data
350 meets quantum chemistry approximations: the δ -machine learning approach. *Journal of chemical*
351 *theory and computation*, 11(5):2087–2096, 2015.
- 352 Mihail Bogojeski, Leslie Vogt-Maranto, Mark E Tuckerman, Klaus-Robert Müller, and Kieron Burke.
353 Quantum chemical accuracy from density functional approximations via machine learning. *Nature*
354 *communications*, 11(1):1–11, 2020.
- 355 Justin S Smith, Benjamin T Nebgen, Roman Zubatyuk, Nicholas Lubbers, Christian Devereux,
356 Kipton Barros, Sergei Tretiak, Olexandr Isayev, and Adrian E Roitberg. Approaching coupled
357 cluster accuracy with a general-purpose neural network potential through transfer learning. *Nature*
358 *communications*, 10(1):1–8, 2019.
- 359 Volker L Deringer, Miguel A Caro, and Gábor Csányi. A general-purpose machine-learning force
360 field for bulk and nanostructured phosphorus. *Nature communications*, 11(1):1–11, 2020.

- 361 Stefan Chmiela, Huziel E Sauceda, Igor Poltavsky, Klaus-Robert Müller, and Alexandre Tkatchenko.
362 sgdml: Constructing accurate and data efficient molecular force fields using machine learning.
363 *Computer Physics Communications*, 240:38–45, 2019.
- 364 Ask Hjorth Larsen, Jens Jørgen Mortensen, Jakob Blomqvist, Ivano E Castelli, Rune Christensen,
365 Marcin Dułak, Jesper Friis, Michael N Groves, Bjørk Hammer, Cory Hargus, et al. The atomic
366 simulation environment—a python library for working with atoms. *Journal of Physics: Condensed*
367 *Matter*, 29(27):273002, 2017.

368 A Appendix

369 A.1 ASTEROID with Multiple Molecules

Table 4: Accuracy of ASTEROID when the inaccurate data is comprised of multiple molecules.

	Aspirin	Benzene	Malonaldehyde	Toluene	Ethanol
Standard Training	1.554	0.074	0.776	0.566	0.351
ASTEROID (Multi)	0.716	0.05	0.480	0.237	0.269
ASTEROID	0.908	0.059	0.338	0.306	0.176

370 To explore the generality of ASTEROID, we further investigate the setting where the inaccurate
371 data for ASTEROID is comprised of multiple molecules. After training such a general purpose (but
372 inaccurate) MLFF, we separately fine-tune the MLFF on each of the MD17 molecules labeled at the
373 CCSD(T) level of accuracy. This setting is very intriguing, since it means that only one network
374 must be pre-trained per molecule. This approach could potentially allow for a large reduction in the
375 memory requirement and pre-training time of ASTEROID.

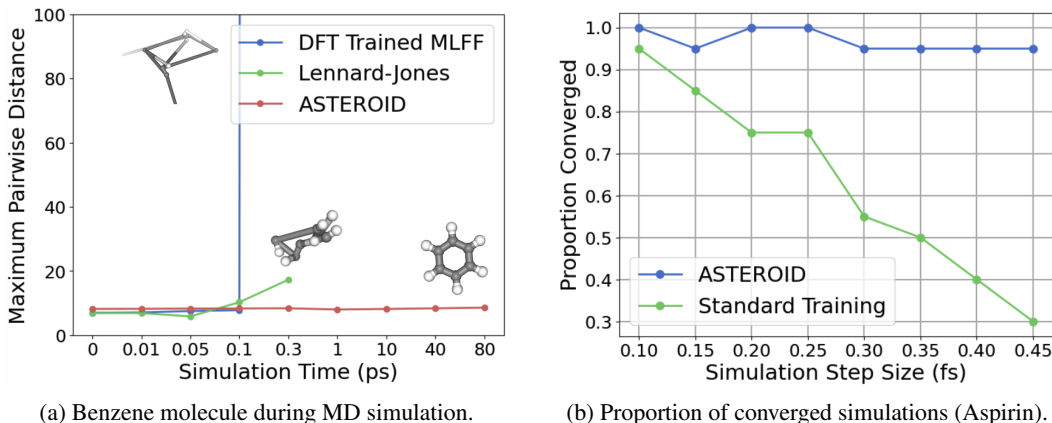
376 The results for a total budget of 250 CCSD(T) data points can be seen in Table 4. From Table 4
377 we can see that training ASTEROID over multiple molecules can significantly reduce the test error
378 compared to standard training. On Aspirin, Benzene, and Malonaldehyde, ASTEROID trained over
379 multiple molecules can perform better than ASTEROID for just a single molecule, likely due to the
380 fact that these molecules all share common structures. However for Malonaldehyde and Ethanol,
381 training over multiple molecules harms performance. Given the mixed performance and the simplicity
382 of single molecule pre-training, it is expected that single molecule pre-training would be favored in
383 most scenarios.

384 A.2 MD simulation

385 It has been observed that low test errors are not sufficient for obtaining stable MD simulation
386 dynamics [Stocker et al., 2022]. To ensure that ASTEROID can be used for MD simulations, we
387 evaluate the performance of MLFFs trained by ASTEROID in downstream MD simulation tasks.
388 First, we demonstrate that ASTEROID-trained MLFFs can produce stable dynamics, while MLFFs
389 trained on DFT data and empirical force fields diverge. Using the Atomic Simulation Environment
390 (ASE) [Larsen et al., 2017], we simulate the behavior of a benzene molecule using forces calculated
391 by a MLFF trained with ASTEROID, an MLFF trained on DFT data only, and the Lennard-Jones
392 empirical force field. We simulate the molecule with Langevin dynamics, where the steps size is 0.5
393 femtoseconds, the temperature is 500K, the friction coefficient is 0.002, and the maximum number of
394 time steps is 10000. The results of these simulations can be seen in Figure 5a, where ASTEROID is
395 able to produce stable dynamics. On the other hand, the error compounding of the DFT trained MLFF
396 and the Lennard-Jones potential results in diverged simulations and unlikely molecular configurations.

397 We also compare the MD simulations generated using ASTEROID with those generated using
398 standard training, where both MLFFs are trained with a data budget equivalent to 250 CCSD(T)
399 points. Inspired by Stocker et al. [2022], we run MD simulations with varying step sizes on the aspirin
400 molecule to evaluate robustness. In Figure 5b we plot the proportion of simulations that converge
401 with varying simulation step sizes. We define a simulation as converged if the maximum pairwise
402 distance between atoms remains within a specified threshold. For each step size, we report the result
403 over 20 Langevin dynamics simulations, each with a length of one picosecond. The ASTEROID
404 framework is able to maintain steady performance across step sizes, and almost all the simulations
405 converge. In contrast, the simulations powered by standard MLFFs fail with larger step sizes.

406 Figures 5a and 5b show the advantages of ASTEROID go beyond reducing test error and allow for
407 stable simulations to be run over 3 times as fast as standard MLFFs. Interestingly, Stocker et al.
408 [2022] find that to train robust MLFFs, much more training data than the amount needed for low test



(a) Benzene molecule during MD simulation.

(b) Proportion of converged simulations (Aspirin).

Figure 5: MD simulation analysis for ASTEROID.

409 error should be used. ASTEROID provides a cost-efficient way to increase the size of the dataset,
 410 therefore enhancing robustness at a low data cost.

411 A.3 Analysis

412 \diamond **Ablation Study** We study the effectiveness of each component of ASTEROID. Specifically, we
 413 investigate the importance of bias-aware training (BAT) and fine-tuning (FT) when compared with
 414 standard training. The results for Gemnet can be seen in Table 5. As shown in Table 5, each of
 415 ASTEROID’s components is effective and complementary to one another. We find that bias-aware
 416 training is most helpful with GemNet, where it reduces test error by 6.5% on average, possibly due
 417 to the fact that GemNet has more capacity to overfit harmful data points than EGNN.

Table 5: Ablation study for ASTEROID on Gemnet. The inaccurate data is DFT labeled configurations and the accurate dataset contains 200 CCSD(T) labeled configurations. “AST.” refers to ASTEROID.

	Aspirin	Benzene	Malonaldehyde	Toluene	Ethanol
Standard Training	1.554	0.074	0.776	0.566	0.351
AST. w/o FT	4.670	3.252	2.726	3.342	5.107
AST. w/o BAT	1.095	0.064	0.347	0.309	0.183
ASTEROID	0.908	0.059	0.338	0.306	0.176

418 \diamond **Sensitivity** We also investigate the sensitivity of ASTEROID to the hyperparameters γ . We use a
 419 data budget of 250 CCSD(T) points. From Figure 6a we can see that ASTEROID is robust to the
 420 choice of hyperparameters, outperforming standard training in every setting.

421 \diamond **Size of inaccurate data.** To demonstrate that ASTEROID can exploit varying amounts of inaccurate
 422 data, we plot the performance of ASTEROID with different cost ratios. This can be seen in Figure 6b,
 423 where a budget of 250 CCSD(T) points is used. ASTEROID performs best when large amounts of
 424 inaccurate data are available but still increases the accuracy by 20% when the cost ratio is small.

425 A.4 Accuracy of Score Matching

426 A.5 Derivation of Score Matching for Forces

For a given molecule with conformations x_1, \dots, x_n , let us denote energy as $E(x)$. Then the the Boltzmann/Equilibrium distribution for the molecule is given by

$$p(x) = \frac{1}{Z} \exp(-\beta E(x)),$$

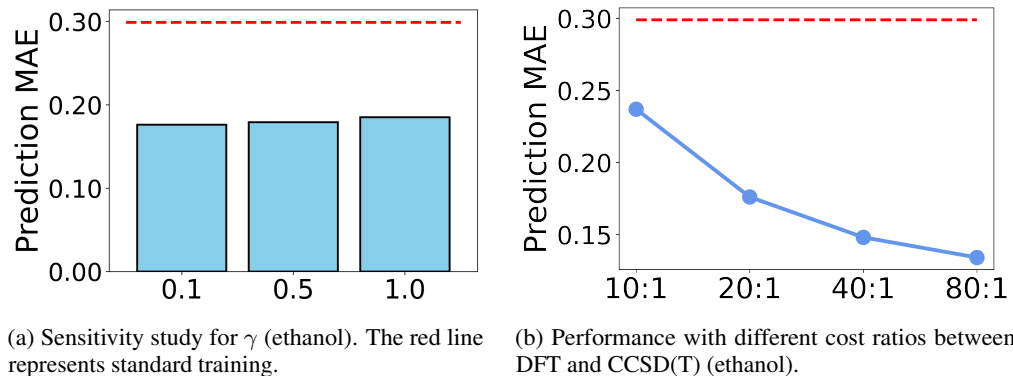


Figure 6: Ablation and sensitivity studies for ASTEROID.

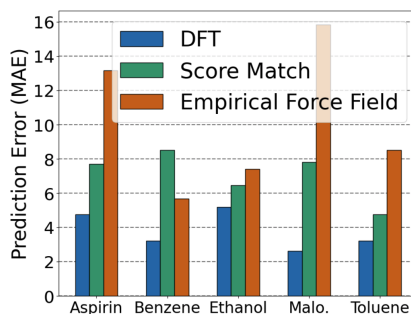


Figure 7: Prediction errors of models tested on CCSD(T) data. Models are not fine-tuned on the CCSD(T) data.

where Z is a normalizing constant, $\beta = \frac{1}{k_{\beta}T}$, k_{β} is the Boltzmann constant, and T is the temperature under which the simulation is run. Then we can see that the force on a conformation x is equivalent to the score, i.e. $F(x) = -\nabla_x E(x) = \frac{1}{\beta} \nabla_x \log p(x)$. Therefore learning the force $F(x)$ is equivalent to learning the score $\frac{1}{\beta} \nabla_x \log p(x)$. Suppose we parameterize the MLFF to directly predict the force as $F_{\theta}(x)$. Then the force matching loss can be written as

$$\mathcal{L}(\theta) = \frac{1}{2} E_{x \sim p(x)} \|F_{\theta}(x) - F(x)\|_2^2 = \frac{1}{2} E_{x \sim p(x)} \|F(x)\|_2^2 - E_{x \sim p(x)} [\langle F_{\theta}(x), F(x) \rangle] + \frac{1}{2} E_{x \sim p(x)} \|F_{\theta}(x)\|_2^2.$$

427 The middle term can then be expanded as

$$\begin{aligned}
E_{x \sim p(x)} [\langle F_\theta(x), F(x) \rangle] &= \int_x p(x) \langle F_\theta(x), F(x) \rangle dx && \text{Integration over } x. \\
&= \int_x p(x) \sum_{i=1}^d \left(\frac{1}{\beta} \frac{d \log p(x)}{dx_i} F_\theta(x)_i \right) dx && \text{Expansion of inner product.} \\
&= \frac{1}{\beta} \sum_{i=1}^d \int_x \frac{dp(x)}{dx_i} F_\theta(x)_i dx && \text{Simplify and move summation.} \\
&= \frac{1}{\beta} \sum_{i=1}^d \int_{x_{i-}} \int_{x_i} F_\theta(x)_i dp(x) dx_{i-} && \text{Integrate over } x_i. \\
&= \frac{1}{\beta} \sum_{i=1}^d \int_{x_{i-}} \left(F_\theta(x)_i dp(x) \Big|_{-\infty}^{+\infty} - \int_{x_i} p(x) dF_\theta(x)_i \right) dx_{i-} && \text{Partial integration.} \\
&= -\frac{1}{\beta} \sum_{i=1}^d \int_{x_{i-}} \int_{x_i} p(x) \frac{dF_\theta(x)_i}{dx_i} dx_i dx_{i-} && \text{Normality assumption.} \\
&= -\frac{1}{\beta} \sum_{i=1}^d E_{x \sim p(x)} \left[\frac{dF_\theta(x)_i}{dx_i} \right] = -\frac{1}{\beta} E_{x \sim p(x)} [\text{Tr} [\nabla_x F_\theta(x)]] .
\end{aligned}$$

Therefore we have the loss

$$\mathcal{L}(\theta) = E_{x \sim p(x)} \left[\frac{1}{\beta} \text{Tr} [\nabla_x F_\theta(x)] + \frac{1}{2} \|F_\theta(x)\|_2^2 \right].$$

428 The first term in the loss disappears as it is not dependent on θ .

429 A.6 Experimental Details

430 In this section, we go over the experimental details.

431 **GemNet Training Details.** To train the bias identification method, we train a freshly initialized
432 model with a batch size of 10 on the accurate dataset for 2000 epochs. To train the inaccurate model,
433 we train a freshly initialized model with the bias aware loss function and batch size 16 over the
434 inaccurate dataset. Finally, to finetune the inaccurately trained model, we train a model with a batch
435 size of 10 on the accurate dataset for 2000 epochs. In each stage of training, we use the following
436 hyperparameters:

- 437 • Evaluation Interval: 1 epoch
- 438 • Decay steps: 1200000
- 439 • Warmup steps: 10000
- 440 • Decay patience: 50000
- 441 • Decay cooldown: 50000

442 The rest of the parameters are the same as used in Gasteiger et al. [2021].

443 **EGNN Training Details.** The EGNN training setup is similar to GemNet. To train the bias
444 identification method, we train a freshly initialized model with a batch size of 10 on the accurate
445 dataset for 2000 epochs. To train the inaccurate model, we train a freshly initialized model with
446 the bias aware loss function and batch size 32 over the inaccurate dataset. Finally to finetune the
447 inaccurately trained model, we train a with a batch size of 10 on the accurate dataset for 2000 epochs.
448 In each stage of training we use the following hyperparameters:

- 449 • Evaluation Interval: 1 epoch

- 450 • Learning rate: 10^{-4} for inaccurate training, 10^{-5} for finetuning
- 451 • num_layers: 5
- 452 • embedding_size: 128

453 A.7 Additional Results

454 Here we include additional results for ASTEROID when empirical force field data is viewed as
 455 inaccurate. For the baseline model we use GemNet. The ASTEROID framework again leads to
 456 consistent gains across all amounts of data.

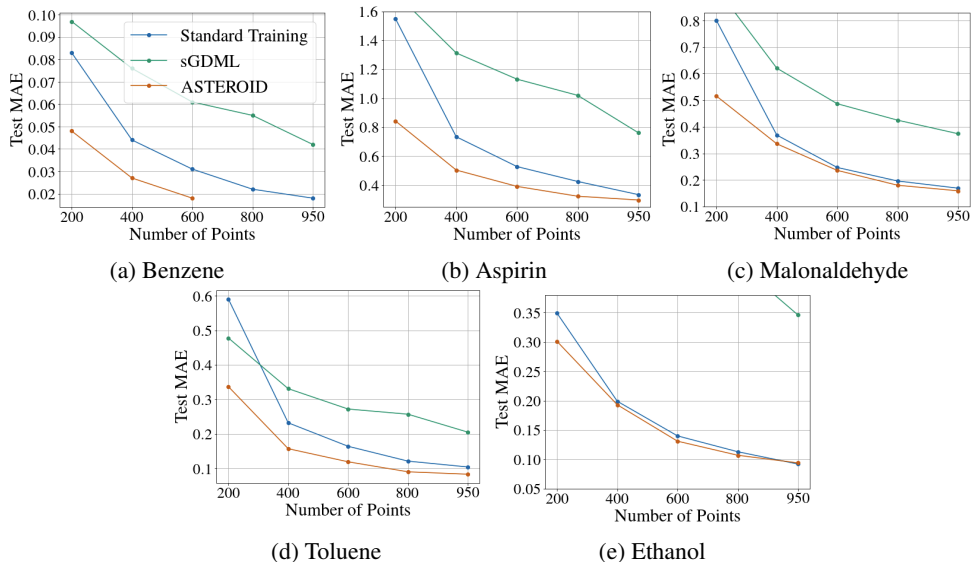


Figure 8: Main results for GemNet when empirical force field data is viewed as inaccurate.

457 A.8 Baseline Implementations

458 \diamond **ANI-1ccx.** In order to have a fair comparison with Smith et al. [2019], we consider two ANI-1ccx
 459 based baselines. In the first baseline, we take the provided ANI-1ccx checkpoint and analyze its zero-
 460 shot performance on the MD17 dataset. For the second ANI-1ccx baseline, we finetune ANI-1ccx
 461 separately on each molecule in MD17 until the validation loss has converged.

462 \diamond **Δ -ML for GemNet.** For a fair comparison with ASTEROID, we implement Δ -ML on GemNet
 463 and the MD17 molecules. Given a molecular configuration x , it’s corresponding DFT force labels f^i ,
 464 and the CCSD(T) force labels f^a , optimize the supervised loss

$$\min_{\theta} \mathcal{L}_{\Delta}(x, \theta) = \ell_f(f^a, f^i + \nabla_x E(x; \theta)). \quad (4)$$

465 We then optimize this loss over all train configurations for a given molecule, using an energy loss
 466 similar to (1). During inference, we predict the CCSD(T) force labels as $f^i + \nabla_x E(x; \theta)$, which
 467 requires the DFT force label to be computed.

468 Since the mapping between DFT labeled configurations and MD-17 labeled configurations is not
 469 explicitly given, we must find it ourselves. For every point in the CCSD(T) dataset, we find the
 470 closest point to it in the DFT dataset. For each of the molecules listed in Section 6, the difference
 471 between the CCSD(T) configuration and closest DFT configuration is 1×10^{-5} . For Benzene and
 472 Ethanol, we find that such a mapping is not available.

473 **A.9 ASTEROID Toy Example**

474 We have added a new result using a two-layer MLP with 128 hidden units each and synthetic data.
475 This experiment shows that ASTEROID can significantly improve generalization error in more general
476 settings. In this experiment, we generate a biased dataset of 2000 points according to $Y = AX + b$,
477 where where $X \sim N(0, 1)$ has dimension 16, b is the bias, and A is a randomly generated Gaussian
478 matrix of dimension 16×16 . The bias b is chosen uniformly from the set $[0, 2, 4, 8, 16]$. We also
479 generate varying levels of accurate data according to $Y = AX$, where $X \sim N(0, 1)$. We then
480 evaluate the test MAE of ASTEROID and standard training over a variety of accurate data sizes. We
481 find that ASTEROID significantly outperforms standard training.

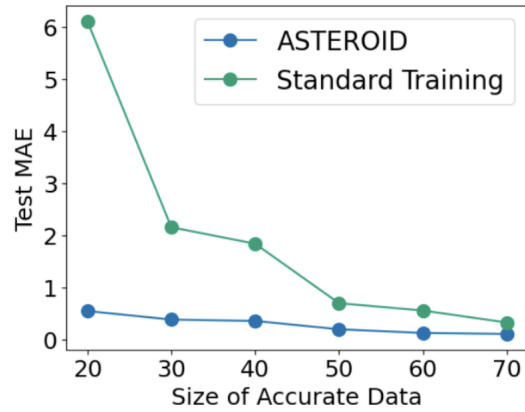


Figure 9: Asteroid toy example.