

The Open Catalyst 2020 (OC20) Dataset and Community Challenges

Lowik Chanussot,^{†,‡} Abhishek Das,^{†,‡} Siddharth Goyal,^{†,‡} Thibaut Lavril,^{†,‡}
 Muhammed Shuaibi,^{†,¶} Morgane Riviere,[‡] Kevin Tran,[¶] Javier Heras-Domingo,[¶]
 Caleb Ho,[‡] Weihua Hu,[§] Aini Palizhati,[¶] Anuroop Sriram,[‡] Brandon Wood,^{||}
 Junwoong Yoon,[¶] Devi Parikh,^{‡,⊥} C. Lawrence Zitnick,^{*,‡} and Zachary Ulissi^{*,#,¶}

[†]*Co-first Author*

[‡]*Facebook AI Research*

[¶]*Department of Chemical Engineering, Carnegie Mellon University*

[§]*Computer Science Department, Stanford University*

^{||}*National Energy Research Scientific Computing Center (NERSC)*

[⊥]*School of Interactive Computing, Georgia Tech*

[#]*Scott Institute for Energy Innovation, Carnegie Mellon University*

E-mail: zitnick@fb.com; zulissi@andrew.cmu.edu

Abstract

Catalyst discovery and optimization is key to solving many societal and energy challenges including solar fuels synthesis, long-term energy storage, and renewable fertilizer production. Despite considerable effort by the catalysis community to apply machine learning models to the computational catalyst discovery process, it remains an open challenge to build models that can generalize across both elemental compositions of surfaces and adsorbate identity/configurations, perhaps because datasets have been smaller in catalysis than related fields. To address this we developed the OC20 dataset, consisting of 1,281,040 Density Functional Theory (DFT) relaxations ($\sim 264,890,000$ single point evaluations) across a wide swath of materials, surfaces, and adsorbates (nitrogen, carbon, and oxygen chemistries). We supplemented this dataset with randomly perturbed structures, short timescale molecular dynamics, and electronic structure analyses. The dataset comprises three central tasks indicative of day-to-day catalyst modeling and comes with pre-defined train/validation/test splits to facilitate direct comparisons with future model development efforts. We applied three state-of-the-art graph neural network models (CGCNN, SchNet, DimeNet++) to each of these tasks as baseline demonstrations for the community to build on. In almost every task, no upper limit on model size was identified, suggesting that even larger models are likely to improve on initial results. The dataset and baseline models are both provided as open resources, as well as a public leader board to encourage community contributions to solve these important tasks.

Keywords

Catalysis, renewable energy, datasets, machine learning, graph convolutions, force field

Introduction

Advancements to renewable energy processes are needed urgently to address climate change and energy scarcity around the world.^{1,2} These include the generation of electricity through fuel cells, fuel generation from renewable resources, and the production of ammonia for fertilization. Catalysis plays a key role in each of these by enabling new reactions and improving process efficiencies.³⁻⁵ Unfortunately, discovering or optimizing catalysts remains a time-intensive process. The space of possible catalyst materials that can be synthesized or engineered is vast and modeling their full complexity under reaction conditions remains elusive. Simulation tools such as Density Functional Theory (DFT)⁶ have greatly expanded our field’s ability to develop reaction mechanisms for specific materials, rationalize experimental measurements, and suggest more active or selective structures for experimental testing. Despite steady growth in computing resources from Moore’s law, the computational complexity of DFT remains a limiting factor in the large-scale exploration of new catalysts.^{7,8} Given its societal importance, finding computationally efficient methods for molecular simulations is of utmost necessity. One potentially promising approach is the use of efficient Machine Learning (ML) models trained with data produced from computationally expensive models, such as DFT.

Indeed, the application of Artificial Intelligence and Machine Learning (AI/ML) to molecular simulations has increased in popularity recently, due to its ability to efficiently model complex functions in data-rich domains. There have been a number of demonstrations from domain scientists for specific challenges such as reaction network elucidation,⁹⁻¹¹ thermochemistry prediction,¹²⁻²⁰ structure optimization,²¹⁻²⁵ accelerating individual calculations,²⁶⁻²⁹ and integration with characterization³⁰ (see recent reviews for a more thorough discussion³¹⁻⁴⁴). Most of these tasks are variations on the same fundamental problem: modeling heterogeneous catalysis. The dataset developed seeks to target a specific subclass of

this problem, periodic slab models. Such modeling involves predicting the energy and forces of various configurations of adsorbate molecules at inorganic interfaces.

Unfortunately, modeling of heterogeneous catalysts entails all the known difficulties of modeling both organic and inorganic chemistry. In organic chemistry modeling involves an overwhelming space of molecules and reactions and many similar, low-energy conformers. Inorganic chemistry involves a large diversity in elements, coordination environments, lattice structures, and long-range interactions. The result is a complex space of compositions and chemistries for which computationally efficient modeling methods are needed for thorough exploration.

A critical factor in building ML models is the data used for training. Despite the importance of heterogeneous catalysis, datasets for it remain smaller than those in other related fields^{45,46} due to additional complexity and higher computational cost. Much of the progress in applying AI/ML in heterogeneous catalysis has been driven by increasingly large and diverse datasets of electronic structure calculations. In the past few years there has been a push towards larger datasets in catalysis, going from $O(100)$ ⁴⁷⁻⁵¹ to $O(1,000)$ ⁵²⁻⁵⁴ then $O(100,000)$ ^{15,55,56} relaxations. Most focus on relaxed adsorption energies of simple adsorbates with smaller datasets of transition state calculations. State-of-the-art ML methods are still improving as data is added to these datasets, so there is no indication that we have saturated the performance of these models. Further, models trained on these datasets have shown limited ability to generalize, which suggests that the models are not yet learning fundamental physical representations. As has been shown in other ML tasks,⁵⁷⁻⁵⁹ we expect that significantly larger datasets will lead to improved accuracy and better generalization.

In this paper, we present the Open Catalyst 2020 (OC20) dataset, (Figure 1) which comprises over 1.2 million DFT relaxations of molecular adsorptions onto surfaces (*ca.* 250 million single-point calculations) across a substantially larger structure and chemistry space

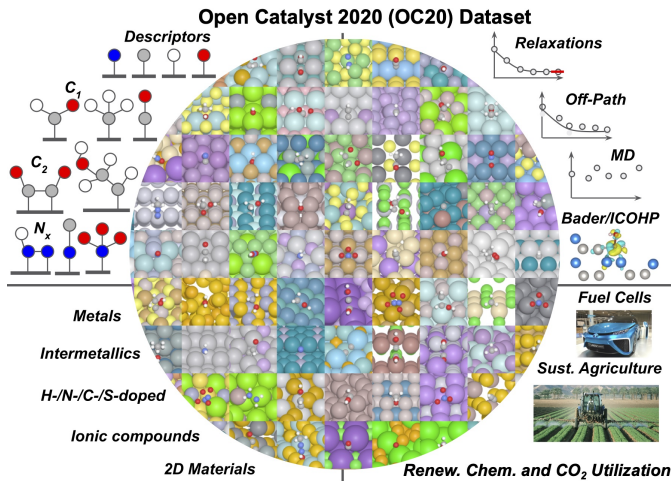


Figure 1: Adsorbates, materials, calculations, and impact areas of the OC20 dataset. Images are a random sample of the dataset.

than previously realized. We envision OC20 to serve as a crucial stepping stone in the development of ML models for practical catalysis applications.

While a dataset of this magnitude will lead to significant improvements in ML models, this is still an extremely sparse sampling of all possibilities. We consider 82 different adsorbates (small adsorbates, C_1/C_2 compounds, and N/O-containing intermediates) that are relevant for renewable energy and environmental applications. Relaxations are performed on randomly sampled low-Miller-index facets of stable materials from the Materials Project,⁶⁰ resulting in surfaces from 55 different elements and mixtures thereof. For each of the calculations, we include relaxation trajectories, Bader charges, and LOBSTER^{61,62}-calculated orbital information. To aid in training more robust models, we additionally compute short, high-temperature *ab initio* Molecular Dynamics (MD) trajectories on a randomly sampled subset of the relaxed states. We also randomly perturb the atomic positions in a subset of the structures along the relaxation pathways and perform single point DFT calculations for these perturbed/rattled structures. We recognize that OC20 addresses a simplified version of heterogeneous catalysis - single adsorbates on idealized structures. Although useful as a first step to informing reaction pathways, the reality in-

volves a number of additional complexities that impact catalyst performance, including reaction conditions, solvation effects, kinetics, etc. While we believe OC20’s approximations to be a reliable step forward, it is important to understand the limits of models developed from this dataset. Future work that incorporates more of the complexities mentioned will undoubtedly benefit from the developments related to OC20. The dataset is publicly available at <http://opencatalystproject.org>. We also plan to upload the dataset to other open systems (e.g. NOMAD or Zenodo) for long-term availability.

In addition to generating and sharing the dataset, we propose three related domain challenges as an open competition: (1) predict the energy and force for a given state, (2) predict a nearby relaxed state given an initial starting state, and (3) predict the relaxed adsorption energy given an initial state. The dataset is split into train/validation/test splits indicative of common situations in catalysis: predicting these properties for a previously unseen adsorbate, for a previously unseen crystal structure or composition, or both. To bootstrap research and the competition, we also provide an open software repository (<https://github.com/Open-Catalyst-Project/ocp>) containing a set of baseline models, data loaders, and training scripts for each of these tasks. While we focus on a subset of tasks, we believe that models capable of solving these tasks on the OC20 dataset will also be able to address a large number of related catalysis problems.

Tasks

Our goal is to improve the efficiency with which inorganic and organic interfaces can be simulated for use in catalysis. Since the primary computational bottlenecks are the DFT calculations used to compute a structure’s forces and energy, we focus on the general challenge of efficient DFT approximation. We focus on structure relaxation – a fundamental calculation in catalysis used in determining a structure’s activity and selectivity. We define three related

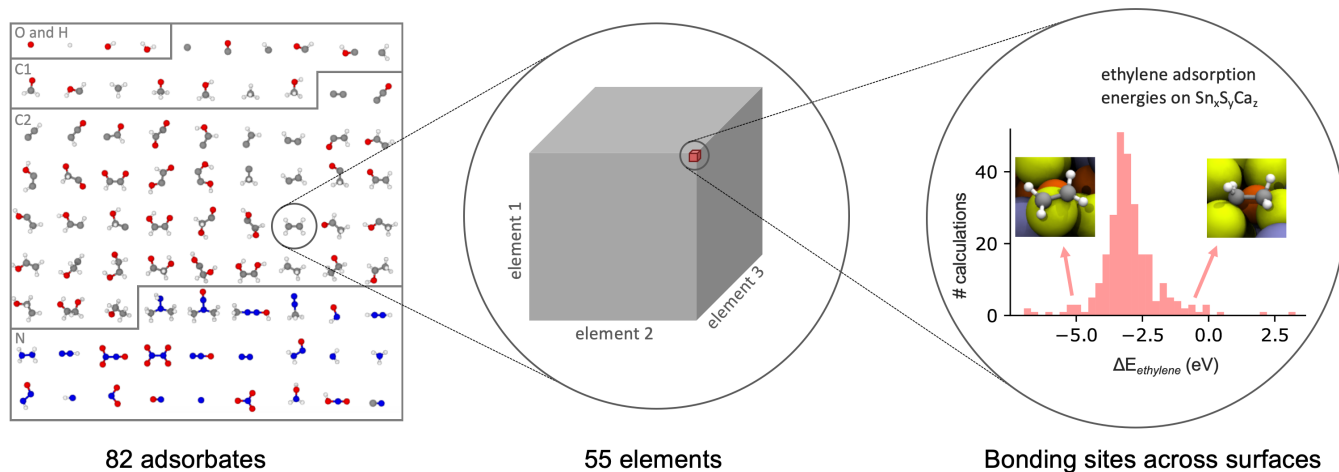


Figure 2: The adsorbates used to generate the Open Catalyst Dataset contain oxygen, hydrogen, C_1 , C_2 , and nitrogen molecules useful for renewable energy applications. Adsorbates that contain both carbon and nitrogen were counted both as C_X adsorbates and as nitrogen-containing adsorbates. For each adsorbate, up to 55^3 different catalyst compositions were considered, with up to dozens of adsorption energy calculations per adsorbate-composition pairing.

tasks, in that success in one task may aid other tasks. These are not the only possibilities for this dataset, and future tasks may be added with additional data generation and input from the community.

In all our tasks, the structure contains a surface and adsorbate. The surface is defined by a unit cell that is periodic in all directions with a vacuum layer of at least 20\AA applied in the z direction. Initial structures are heuristically determined. Ground truth data is computed for all tasks using DFT. Dataset details and evaluation metrics are provided in following sections.

Structure to Energy and Forces ($S2EF$) is to take the positions of the atoms as input and predict the energy and per-atom forces as calculated by DFT. For the purposes of this manuscript, energy refers to adsorption energy unless otherwise noted. The adsorption energy is defined as the energy of the combined surface and adsorbate system (relaxed or not) minus the energy of the relaxed slab and the relaxed gas phase adsorbate molecule. The force is defined as the negative gradient of the energy with respect to the atomic positions.

This is our most general task and has the broadest applicability across catalysis and related fields. It is essentially identical to existing challenges in developing machine learning

potentials.⁶³ However, the inclusion of both inorganic and organic materials and the dataset size make this challenge unique.

Initial Structure to Relaxed Structure ($IS2RS$) takes as input an initial structure and predicts the atomic positions in their final, relaxed state. Traditional relaxations are performed through an iterative process that estimates the atomic forces using DFT, which are in turn used to update atom positions until convergence. This very computationally expensive process typically requires hundreds of DFT calculations to converge.

If the $IS2RS$ task is approached using ML approximations of DFT to estimate atomic forces ($S2EF$ task), evaluation on the $IS2RS$ task may help determine whether models built for $S2EF$ are sufficiently accurate for practical applications. Alternatively, it may be possible to predict the relaxed structure directly, without estimating a structure’s energy or forces (Figure 3(B)), as many of the changes during relaxation (say due to particular initial guess strategies) are systematic. These direct $IS2RS$ approaches may lead to even further improvements in computational efficiency.

Initial Structure to Relaxed Energy ($IS2RE$) task is to take the initial structure as input and predict the structure’s energy in the

relaxed state. This is the most common task in catalysis, as the relaxed energies are often correlated with catalyst activity and selectivity, and the energies are important parameters for detailed microkinetic models. Similar to *IS2RS*, this task may be approached by estimating the relaxed structure and energy by iteratively applying *S2EF*, or by directly regressing the energy from the initial structure without estimating the intermediate or relaxed structures.

The OC20 Dataset

The OC20 dataset is constructed to provide both training and evaluation data for our three previously defined tasks involving DFT approximation and structure relaxation. Modern machine learning models, especially those employing deep learning, require sufficiently large datasets to learn accurate models. For training, we provide 640,081 relaxations across a wide variety of surfaces and adsorbates. The intermediate structures and their corresponding energy and forces are provided for each relaxation resulting in over 133 million training structures. To potentially aid in training and to provide additional information for the catalysis community, we performed DFT calculations on rattled and *ab initio* Molecular Dynamics (MD) data. We also computed Bader charges and LOBSTER analyses (over 1.8 million examples each) as these computed properties may be useful for models by explaining why the energies are what they are.

Dataset Generation

The dataset is constructed in four stages: 1) adsorbate selection, 2) surface selection, 3) initial structure generation, and 4) structure relaxation. We describe each of these four stages in turn, followed by a description of the additional data provided with the main dataset. All source code to generate the configurations are provided in the Open Catalyst Dataset repository (<https://github.com/Open-Catalyst-Project/Open-Catalyst-Dataset>).

Adsorbate Selection

Adsorbates are sampled randomly from a set of 82 molecules that are chosen for their utility to renewable energy applications. As shown in Figure 2, this includes adsorbates that contain only oxygen or hydrogen, C₁ molecules, C₂ molecules, and nitrogen-containing molecules. We enumerated the oxygen and hydrogen molecules for their ubiquitous presence in water-solvated electrochemical reactions. C₁ and C₂ molecules are important for solar fuel synthesis, while nitrogen-containing molecules have applicability in solar fuel and solar chemical synthesis. Note that some of the C₂ molecules have two binding sites; we refer to these as bidentate adsorbates. The list of all 82 adsorbates is provided in the Supplementary Information.

Surface Selection

Surfaces are sampled in three stages. First, the number of elements is selected with a 5% chance of choosing a unary material, 65% chance for a binary material, and a 30% chance for a ternary material. Greater emphasis is given to binary and ternary materials because these sets contain a wider variety of understudied materials. Next, a stable bulk material is randomly selected from the 11,451 materials in the Materials Project⁶⁰ with the number of elements chosen in the first step. Finally, all symmetrically distinct surfaces from the material with Miller indices less than or equal to 2 are enumerated, including possibilities for different absolute positions of surface plane. From this list of surfaces one is randomly selected. The surface atoms were replicated to a depth of at least 7 Å and a width of at least 8 Å.

Pymatgen⁶⁴ was used to search over all bulk materials in the Materials Project with non-positive formation energies and energies-above-lower-hulls of at most 0.1 eV/atom. The enumeration of symmetrically distinct surfaces was also performed using pymatgen.⁶⁴ Elements for the bulk materials were chosen from a set of 55 elements comprising reactive nonmetals, alkali metals, alkaline earth metals, metalloids, transition metals, and post-transition metals.

Note that DFT was used to re-optimize the bulk structures prior to surface enumeration to ensure differences between the DFT settings used in the Materials Project and OC20 did not induce unintended stress or strain effects. Any bulks that we could not successfully relax were omitted from this dataset.

Initial Structure Generation

The initial structures are generated by placing the selected adsorbates on the selected surfaces using CatKit⁶⁵ and the atomic simulation environment (ASE).⁶⁶ Surface atoms are identified by their positions above the center-of-mass, their z -distance within 2 Å of the upper-most atom, and by their under-coordination relative to the bulk atoms. Atomic coordination environments were calculated using pymatgen’s Voronoi tessellation algorithm.⁶⁴ Next, we manually tagged the adsorbates’ binding atoms for both mono- and bi-dentate adsorbates. Finally, we gave the surface structure, adsorbate, the identified surface atoms, and identified adsorbate binding sites to CatKit.⁶⁵ CatKit used this information to enumerate a list of symmetrically distinct adsorption sites along with suggested per-site orientations for the adsorbates. From this list, an adsorption configuration is randomly selected. The sites selected are not necessarily the most stable adsorption site on each surface. Since one of our goals is to calculate adsorption energies, we generate two sets of inputs for each system: (1) the adsorbate placed over the catalyst atoms, and (2) just the catalyst atoms without the adsorbate. This resulted in a total of 1,919,165 and 616,124 unique inputs for (1) and (2), respectively, which were later filtered and segregated into suitable train, validation, and test validation splits as described later in this section.

Structure Relaxation

All structure relaxations were performed using the Vienna Ab Initio simulation Package (VASP)^{67–71} until all per-atom forces are less than 0.03 eV/Å. Calculations were allowed up to 144 hours (12 cores) for the relaxation. Sys-

tems that timed out before reaching the specified force threshold were set aside for the S2EF task. All intermediate structures, energies, and forces are stored for future training and evaluation. During the relaxations only adsorbate and surface atoms (as defined during the generation above) were allowed to move; subsurface atoms were maintained at fixed positions. This was done to avoid unrealistic structure deformations and to simulate a semi-infinite condition with bulk material far below the catalyst surface. Given the intended scale of OC20, the careful consideration of DFT settings was a non-trivial challenge. Relaxations generally followed previous high-throughput catalysis efforts with reasonable trade-offs between accuracy for surface chemistry and computational cost¹⁶ (VASP,^{67–71} RPBE,⁷² no spin polarization, etc). The choices made for DFT were a result of several important considerations: ensuring calculations were representative, concerns associated with inconsistent cutoffs/settings, and representative of typical numerical/convergence issues the computational chemistry field faces. The assumptions made were necessary to achieve the dataset’s scale. Detecting small numerical or convergence errors is a non-trivial problem that could be improved with this dataset. Most importantly, we anticipate models and methods that solve the S2EF, IS2RE, or IS2RS tasks for this dataset are very likely to solve future challenges for future surface science datasets with different DFT modeling choices.

System DFT energies were referenced to represent adsorption energies. Adsorption energies were calculated according to the Equation below, where E_{sys} is the DFT energy of the combined surface (i.e. slab) and adsorbate — this energy can be from both relaxed and intermediate structures. The reference energies for each system, E_{slab} and E_{gas} are the DFT energy of the relaxed surface and adsorbate molecule respectively. The value of E_{gas} for each adsorbate was computed as a linear combination of N₂, H₂O, CO, and H₂ resulting in the atomic energies found in the supplementary.

$$E_{ad} = E_{sys} - E_{slab} - E_{gas}$$

Table 1: Size of train/validation splits (number of structures for *S2EF* and initial structures for *IS2RS* and *IS2RE*). The structures for *S2EF* are sampled from 640,081 relaxations for train, and from 30k-70k relaxations for each validation and test split. Subsplits of validation and test are roughly the same size, but are exclusive of each other. Subsplits include sampling from the same distribution as training (In Domain), unseen adsorbates (Out of Domain (OOD) Adsorbate), unseen element compositions for catalysts (OOD Catalyst), and unseen adsorbates and catalysts (OOD Both). Test sizes are similar.

Task	Train	In Domain	OOD Adsorbate	OOD Catalyst	OOD Both
S2EF	133,934,018	987,036	999,838	987,343	997,922
IS2RS	460,328	24,733	24,961	24,738	24,971
IS2RE	460,328	24,733	24,961	24,738	24,971

Resulting trajectories were further analyzed for per-atom force criterion, numerical issues, or catastrophic reconstructions as described below in the Train, Validation, and Test Splits section.

MD and Rattled Calculations

The intermediate structures from the relaxations may result in a dataset biased towards structures with lower energies. To learn robust models, training samples with higher forces and greater configurational diversity may be needed. We adopted two strategies for generating additional training data: (1) partial MD in VASP⁶⁷⁻⁷¹ and (2) normally-distributed random position perturbation methods colloquially known in molecular simulations as “rattling.”

MD calculations simulate the atomic interactions when heat is added to the system. Partial MD calculations were carried out on previously relaxed structures with random initial velocities generated from a Maxwell-Boltzmann distribution at a temperature of 900 K. We integrated the MD trajectories over 80 fs or 320 fs with integration steps of 2 fs in the NVE ensemble. Time-scales were selected to allow systems to explore local configurations while minding computational costs.

To diversify the distribution of single-point structures in the dataset, we “rattled” some of the structures by adding random displacements to the atomic positions with ASE.⁶⁶ For each relaxation, 20% of the images in the trajectories were selected for rattling. The atomic displacements were sampled from a heuristically-

generated normal distribution with a $\mu = 0$ and $\sigma = 0.05$. Single point DFT calculations were then performed on the rattled structures.

Similar to the relaxations, only the top surface atom layers were allowed to move in both the MD and rattled calculations with the rest of the atom positions held fixed. All calculations were performed at the same theoretical level and energy/forces convergence criteria as in the relaxation calculations. Approximately 950 thousand MD (*ca.* 64 million single-point energies/forces) and 30 million rattled calculations were carried out.

Bader Charges and LOBSTER Analyses

We performed electronic structure calculations for general use by the catalysis research field. These calculations (i.e., Bader charges^{61,73,74} and LOBSTER^{75,76} analyses) were carried out on relaxed structures and also on randomly selected snapshots from both MD and rattled trajectories. Bader charge analyses provides charge density maxima at each atomic center and the Bader volume for each atom through the zero-flux partitioning method.⁶² LOBSTER enables chemical-bonding analysis based on periodic DFT outputs.⁷⁵ LOBSTER calculates atom-projected densities of states (pDOS) or projected crystal orbital Hamilton population (pCOHP) curves, among others. Literature has demonstrated that such electronic structure information can provide valuable insights to the theoretical and the ML communities.⁷⁷⁻⁷⁹

Dataset profile

Approximately 872,000 adsorption energies were calculated successfully. Of these, 3.7% were calculations on unary catalysts; 61.4% were on binaries; and 34.9% were on ternaries. Among these calculations, 28.9% of them had reactive nonmetal elements in the catalyst; 8.1% of them had alkali metals; 10.2% had alkaline earth metals; 26.4% had metalloids; 81.3% had transition metals; and 37.2% had post-transition metals. Considering adsorbates: 6.6% of the calculations had adsorbates containing only oxygen or hydrogen; 25.2% of the calculations had C_1 adsorbates; 44.4% had C_2 adsorbates; and 29.0% had nitrogen-containing adsorbates.

Despite this dataset’s large size compared to previous catalytic datasets, it still very small in comparison to the number of potential calculations. Of the $\binom{55}{3} + \binom{55}{2} + \binom{55}{1} = 27,775$ possible compositions, only 5,243 (18.9%) of them were successfully sampled here. Of the compositions sampled, there were an average of 249 successful adsorption calculations for each. Additionally: if we compare the number of sites we sampled here to rough estimates of the number of sites we could have sampled given our constraints on adsorbates, surfaces, and bulks, then we find that we performed *ca.* 0.07% of the possible calculations. This severe sparsity in the data compared to its large scale emphasizes the need for surrogate models.

Train, Validation and Test Splits

We split our dataset into training, validation, and testing sets. The training set is used to learn model parameters; the validation set is used to tune model hyperparameters and to perform ablation studies; and the test set is used to report model performance.

A careful choice of validation and test splits can help evaluate a model’s performance on both interpolative and extrapolative tasks. Interpolative evaluation tests the ability to model variations of the training data, and is performed by sampling examples from the same distribution as the training dataset. Extrapolative eval-

uation tests a model’s performance on unseen tasks, e.g., new materials or adsorbates. In the context of catalytic development, we strive to extrapolate beyond data we have already seen so that we can discover new materials and search spaces.^{80,81}

We explore extrapolation along two dimensions; new adsorbates and new catalyst compositions. Adsorbate extrapolation is performed by holding out 14 adsorbates from the training dataset sampled from all types (O, H, C1, C2, and N) of adsorbates. Similarly for catalyst compositions, a subset of element combinations for catalysts is held out from the training dataset. These were sampled from the 1,485 binary and 26,235 ternary material combinations of the 55 elements used in the dataset. No surfaces with unary materials are in the extrapolative subsplits for validation and testing. A full list of the adsorbates materials in train and validation splits are in the SI.

We used four subsplits for each of the validation and test sets by considering all combinations of potential extrapolations (Table 1). These include In-Domain (sampled from the training distribution), Out-of-Domain Adsorbate (OOD Adsorbate), OOD Catalyst, and OOD Both (both unseen adsorbate and unseen catalyst compositions). As shown in Table 1, each subsplits in validation and testing contains *ca.* 25,000 relaxations. For the *S2EF* task we randomly select a one million structure subset from the relaxations in each subsplit. Note that the extrapolative subsplits of our validation set are completely exclusive to the extrapolative subsplits in the test set, e.g., the adsorbates in the validation adsorbate subsplit are unique from the adsorbates in the test adsorbate subsplit. This helps ensure overfitting to the test set does not occur during hyperparameter tuning on the validation set.

Baseline GNN Models

We evaluate our tasks using a set of baseline models that are representative of the current state-of-the-art. The set of models we evaluate is by no means comprehensive, but

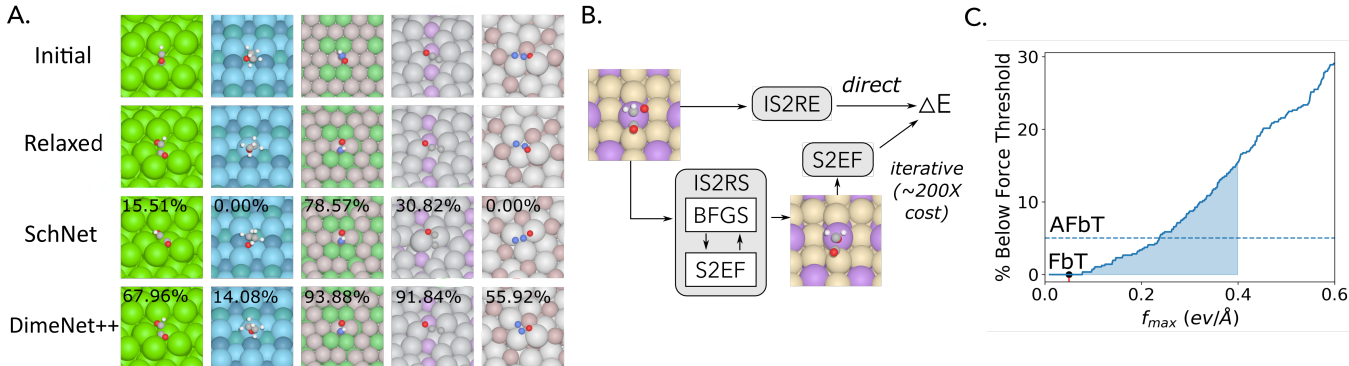


Figure 3: Demonstration of baselines SchNet and DimeNet++ models for solving the *IS2RE*, *S2EF*, and *IS2RS* tasks and the inter-relationships. (A) Snapshots of five representative initial adsorbate configurations before DFT relaxations, the same adsorbates after DFT relaxation, and the relaxed structures as relaxed by SchNet and DimeNet++ after fitting the *S2EF* task. ADwT metrics are overlaid on the model snapshots. (B) Three ways to predict the relaxed energy: directly through *IS2RE*, indirectly through *IS2RS*, and confirmation of the relaxed structure with a single DFT single-point. (C) SchNet force-only performance as characterized by the percentage of structures within the desired max force threshold of 0.05 eV/Å (FbT) and average percentage of force below threshold (AFbT) of 0.4 eV/Å (shaded area).

they demonstrate what is feasible with current models. Code and pretrained models for our baseline ML approaches implemented in PyTorch Geometric^{82,83} are publicly available at the Open Catalyst Project (<http://opencatalystproject.org>).

Our baseline ML approaches are all based on Graph Neural Networks (GNNs)⁸⁴ that operate over a graph structure containing nodes and edges. In our domain, the nodes represent atoms and edges represent the relationship between neighboring atoms. At each node, an atom embedding is iteratively updated based on messages passed along the edges. During this message-passing phase, GNNs employ neural networks to learn the atomic representations,^{85,86} and unlike traditional descriptor-based models do not require hand-crafting. Node embeddings are initialized based on the atom’s properties, such as their atomic number, group number, electronegativity, atomic volume, etc.⁸⁷ Outputs for the GNN may be computed from individual node (atom) embeddings for node-specific information (per-atom forces), or over the pooled node embeddings for system outputs (structure energy).

We benchmark three recent GNN methods: Crystal Graph Convolutional Neural Network

(CGCNN),⁸⁷ SchNet⁸⁸ and DimeNet++.^{89,90} CGCNN is one of the first approaches to use GNNs on periodic crystal systems and uses a diverse set of features as input to the node embeddings. The original model encoded edge information using the discretized distances between atoms. SchNet proposed using continuous edge filters, which allows for the computation of per-atom forces through partial derivatives of the structure’s energy with respect to the atom positions. To allow CGCNN to compute per-atom forces in the same manner, we updated the distance encoding to use gaussian basis functions but without the envelope distance function used in SchNet in our experiments. Finally, to not only encode distance information but also angular information between triplets of atoms, DimeNet introduced the use of directional message passing. DimeNet++, an extension to DimeNet, replaces the Bilinear layer with a Hadamard product and additional multi-layer perceptrons; providing reported speed improvements of 8x and a 10% accuracy boost on QM9.⁹¹

For all approaches, graph edges were determined by a nearest neighbor search limited by a cutoff radius of 6Å, retaining up to the 50 nearest neighbors. When computing distances, pe-

riodic boundary conditions were taken into consideration. Atoms were tagged as three types, slab (fixed), surface (free), and adsorbate (free), to allow loss functions to emphasize free atoms over fixed atoms. The number of hidden channels is 128, 1024, 192 for CGCNN, SchNet and DimeNet++ respectively unless stated otherwise; resulting in 3.6 million (CGCNN), 7.4 million (SchNet) and 1.8 million (DimeNet++) parameters. Model sizes were chosen so that runtimes were roughly equivalent. Note the size of the models was increased from their original implementations to account for OC20’s larger size. Model hyperparameters and additional modifications can be found in the supplementary.

Since both the computed energies and forces are evaluated, the baseline loss function^{27,90} uses the following form:

$$\mathcal{L} = \lambda_E \sum_i |E_i - E_i^{DFT}| + \lambda_F \sum_{i,j} \frac{1}{N_i} |F_{i,j} - F_{i,j}^{DFT}|,$$

where λ_E and λ_F are empirical parameters, E_i is the energy of image i , and $F_{i,j}$ is the force of the j th free atom in image i , and N_i is the number of free atoms in image i . For the *IS2RE* task, in which only the energy is evaluated, only the first term of the loss function is used ($\lambda_F = 0$).

All of the models are ML-based as there are currently no physical models that operate over such a large composition space with reasonable accuracy and elemental parameterizations. In particular, the recently developed GFNO-xTB method⁹² is parameterized for all of the elements in this dataset and is fast enough (approx 1,000X faster than DFT) to compete on these benchmarks and preliminary results are reported in the SI. However, since the method was not fit for inorganic surfaces and the xTB code⁹³ is still under active development for periodic boundary conditions, the results were excluded from the summaries here. We hope that the release of our dataset will inspire future efforts on parameterizing tight-binding DFT codes or reactive force field methods for these materials.

Experiments

We begin by describing the metrics used to evaluate our three tasks, followed by the results of our baseline models.

Evaluation Metrics

For each task, we define evaluation metrics to track the progress in the field, as well as to measure the practical utility of the approaches. All ground truth values are computed using DFT. Our evaluation metrics are as follows:

S2EF: The *S2EF* task has three metrics: the Mean Absolute Error (MAE) for energy, MAE for forces on free atoms and a combined metric. Our combined metric, Energy and Forces within Threshold (EFwT), is designed to measure the practical usefulness of a model for replacing DFT by evaluating whether both the computed energy and forces are close to the ground truth.

Energy MAE: Mean Absolute Error between the computed energy and the ground truth energy.

Force MAE: Mean Absolute Error between the computed per-atom forces and the ground truth forces. Errors are only computed for free catalyst and adsorbate atoms.

Force cosine: Mean cosine of the angle between the computed per-atom forces and the ground-truth forces. Similar to MAE, these are only computed for free atoms.

EFwT: The percentage of structures in which the computed energy is within $\epsilon = 0.02$ eV of the ground truth energy, and the maximum error in per-atom forces is below $\alpha = 0.03$ eV/Å. Both these criteria must be met for the structure to be labeled as “correct”.

IS2RS: Several methods exist for determining the accuracy of relaxed structures predicted by ML models. The simplest is to measure the distance between the predicted 3D positions of the atoms and those of the ground truth. However, small changes in position can lead to significant changes in the per-atom forces and a

structure’s energy. For this reason, a better measure of a proposed relaxed structure is the magnitude of its per-atom forces as measured by a single point DFT calculation. If the proposed relaxed structure represents a true local energy minimum, the forces should be close to zero.

ADwT: The Average DwT (Distance within Threshold) across thresholds ranging from $\beta = 0.01\text{\AA}$ to $\beta = 0.5\text{\AA}$ in increments of 0.001\AA . DwT is computed as the percentage of structures with an atom position MAE below the threshold. MAE is only computed for free catalyst and adsorbate atom positions while taking into account periodic boundary conditions. We use ADwT as opposed to the MAE on 3D atom positions, since ADwT is robust to outliers and better indicates the percentage of relaxations that are likely to be successful.

FbT: The percentage of relaxed structures with maximum DFT calculated per-atom force magnitudes below a threshold of $\alpha = 0.05\text{ eV/\AA}$. Force magnitudes of only free catalyst and adsorbate atoms are used. A value of $\alpha = 0.05\text{ eV/\AA}$ represents a practical threshold by which DFT relaxations are commonly assumed to have converged. To ensure that the ML relaxations find a relaxed structure that isn’t significantly different from the ground truth relaxed structures, e.g., the adsorbate moves to a different binding site, an additional filtering step is applied. We filter on the atom position MAE (free catalyst and adsorbate atoms) with a threshold of $\beta = 0.5\text{\AA}$. Thus, to be considered correct, a relaxed structure must meet both the FbT and the DwT criterion.

AFbT: The Average FbT (Forces below Threshold) over a range of thresholds ranging from $\alpha = 0.01\text{ eV/\AA}$ to $\alpha = 0.4\text{ eV/\AA}$ in increments of 0.001 eV/\AA , Figure 3(C). This metric measures progress over a wider range of thresholds, which may be important for early algorithm development that may need thresholds more lenient than $\alpha = 0.05\text{ eV/\AA}$ to see improvement. Similar to FbT, the

relaxed structures must also meet the same DwT criterion with $\beta = 0.5\text{\AA}$.

Note that FbT and AFbT require the computation of single point DFT calculations, which are computationally expensive. For this reason, a random subset of 500 relaxed structures are chosen from the validation and test set splits (2000 total for each) for evaluating these metrics. If a DFT calculation fails to converge within 60 electronic steps or a wall time of 2 hrs, the system is assumed to be incorrect with forces beyond the thresholds for both FbT and AFbT.

Table 2: Predicting energy and forces from a structure (*S2EF*) as evaluated by Mean Absolute Error (MAE) of the energies, forces MAE, and the percentage of Energies and Forces within Threshold (EFwT). Results reported for models training on the entire training dataset.

<i>S2EF</i> Test				
Model	ID	OOD Ads	OOD Cat	OOD Both
		Energy MAE [eV] ↓		
Median baseline	2.0596	2.4188	2.0110	2.5460
CGCNN ⁸⁷	0.5105	0.6321	0.5202	0.7681
SchNet ⁸⁸	0.4421	0.4858	0.5279	0.7057
SchNet ⁸⁸ - force-only	34.0689	33.7670	35.2701	38.4607
SchNet ⁸⁸ - energy-only	0.3975	0.4533	0.5626	0.7241
DimeNet++ ^{89,90}	0.4579	0.4701	0.5056	0.6489
DimeNet++ ^{89,90} - force-only	28.2214	28.9404	28.8636	34.9118
DimeNet++ ^{89,90} - energy-only	0.3585	0.4022	0.5041	0.6549
DimeNet++ ^{89,90} -Large - force-only	29.3504	30.0338	30.0074	36.7665
		Force MAE [eV/\AA] ↓		
Median baseline	0.0808	0.0801	0.0787	0.0978
CGCNN ⁸⁷	0.0683	0.0728	0.0670	0.0851
SchNet ⁸⁸	0.0493	0.0529	0.0509	0.0655
SchNet ⁸⁸ - force-only	0.0442	0.0469	0.0459	0.0591
SchNet ⁸⁸ - energy-only	0.5794	0.5974	0.5852	0.6463
DimeNet++ ^{89,90}	0.0442	0.0458	0.0444	0.0559
DimeNet++ ^{89,90} - force-only	0.0331	0.0341	0.0340	0.0417
DimeNet++ ^{89,90} - energy-only	0.3399	0.3395	0.3395	0.3643
DimeNet++ ^{89,90} -Large - force-only	0.0280	0.0289	0.0312	0.0371
		Force cosine ↑		
Median baseline	0.0000	0.0000	0.0000	0.0000
CGCNN ⁸⁷	0.1541	0.1369	0.1492	0.1444
SchNet ⁸⁸	0.3184	0.2954	0.2956	0.2987
SchNet ⁸⁸ - force-only	0.3595	0.3391	0.3279	0.3403
SchNet ⁸⁸ - energy-only	0.0845	0.0798	0.0804	0.0830
DimeNet++ ^{89,90}	0.3628	0.3476	0.3465	0.3684
DimeNet++ ^{89,90} - force-only	0.4870	0.4717	0.4607	0.4954
DimeNet++ ^{89,90} - energy-only	0.1066	0.0959	0.1048	0.1015
DimeNet++ ^{89,90} -Large - force-only	0.5638	0.5502	0.5115	0.5516
		EFwT ↑		
Median baseline	0.00%	0.00%	0.00%	0.00%
CGCNN ⁸⁷	0.01%	0.00%	0.01%	0.00%
SchNet ⁸⁸	0.11%	0.04%	0.06%	0.01%
SchNet ⁸⁸ - force-only	0.00%	0.00%	0.00%	0.00%
SchNet ⁸⁸ - energy-only	0.00%	0.00%	0.00%	0.00%
DimeNet++ ^{89,90}	0.10%	0.03%	0.05%	0.01%
DimeNet++ ^{89,90} - force-only	0.00%	0.00%	0.00%	0.00%
DimeNet++ ^{89,90} - energy-only	0.00%	0.00%	0.00%	0.00%
DimeNet++ ^{89,90} -Large - force-only	0.00%	0.00%	0.00%	0.00%

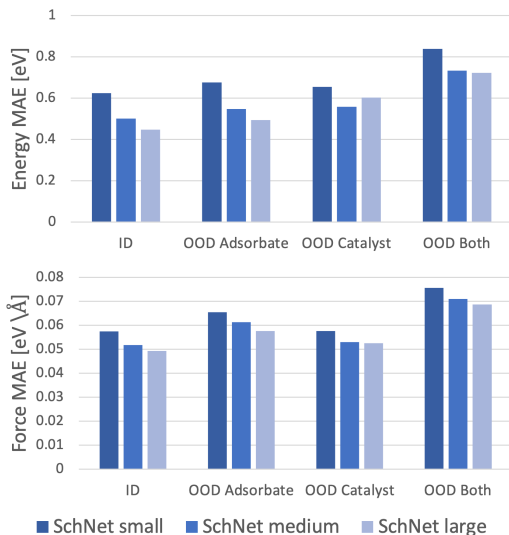


Figure 4: Predicting Structure to Energy and Forces ($S2EF$) as evaluated by Mean Absolute Error (MAE) of the energies and forces. The small, medium and large SchNet models have the following sizes: Small: 256 hidden, 4 message-passing layers, 1,316,097 params, Medium: 1024 hidden, 3 message-passing layers, 5,704,193 params, Large: 1024 hidden, 4 message-passing layers, 7,396,353 params. Results reported for models trained on the entire training dataset.

$IS2RE$: Similar to the $S2EF$ task we propose two metrics for $IS2RE$. The first measures the MAE on the computed and ground truth energy. The second measures the energies within a threshold (EwT) of the ground truth, which once again measures the percentage of estimated energies that are likely to be practically useful.

Energy MAE: Mean Absolute Error between the computed relaxed energy and the ground truth relaxed energy.

EwT: The percentage of computed relaxed energies within $\epsilon = 0.02$ eV of the ground truth relaxed energy.

While our evaluation metrics focus on accuracy, it is important to note that methods should also be significantly faster than conventional DFT. As a rough benchmark, we desire energy and force estimates at approximately 10

Table 3: Predicting relaxed structure from initial structure ($IS2RS$) as evaluated by Average Distance within Threshold (ADwT), Forces below Threshold (FbT), and Average Forces below Threshold (AFbT). All values in percentages, higher is better. Results reported for structure to force models trained on the All training dataset. The initial structure was used as a naive baseline (IS baseline). FbT and AFbT metrics are only computed when ADwT metrics are greater than 20.26%.

$IS2RS$ Test				
Model	ID	OOD Ads	OOD Cat	OOD Both
ADwT \uparrow				
IS baseline	21.37%	19.09%	21.42%	26.28%
SchNet ⁸⁸	15.92%	12.83%	14.63%	14.78%
SchNet ⁸⁸ – force-only	32.47%	28.59%	30.94%	35.09%
DimeNet++ ^{89,90}	30.62%	26.66%	30.01%	32.29%
DimeNet++ ^{89,90} – force-only	48.73%	45.19%	48.54%	53.17%
DimeNet++ ^{89,90} -Large – force-only	52.43%	48.47%	50.91%	54.85%
FbT \uparrow				
IS baseline	0.00%	0.00%	0.00%	0.00%
SchNet ⁸⁸	-	-	-	-
SchNet ⁸⁸ – force-only	0.00%	0.00%	0.00%	0.00%
DimeNet++ ^{89,90}	0.00%	0.20%	0.00%	0.00%
DimeNet++ ^{89,90} – force-only	0.61%	0.20%	0.00%	0.20%
DimeNet++ ^{89,90} -Large – force-only	1.02%	0.40%	0.00%	0.20%
AFbT \uparrow				
IS baseline	0.06%	0.34%	0.21%	0.00%
SchNet ⁸⁸	-	-	-	-
SchNet ⁸⁸ – force-only	5.31%	2.82%	2.66%	2.73%
DimeNet++ ^{89,90}	3.60%	3.01%	2.61%	2.33%
DimeNet++ ^{89,90} – force-only	17.42%	14.67%	14.12%	14.46%
DimeNet++ ^{89,90} -Large – force-only	25.58%	20.73%	20.05%	20.62%

ms which would significantly improve the applicability of DFT. Significantly faster than this (closer in speed to classical force fields) would open up even more interesting applications. We ask that users self-report timing results, but we are not going to make that a core part of the challenge as computation time can likely be further optimized for the best models and with hardware acceleration.

Leaderboard

To ensure consistent and fair evaluation, a public leaderboard is available on the Open Catalyst Project webpage (<http://opencatalystproject.org>). Results on any of the tasks’ test datasets may be uploaded for evaluation. Ground truth test data is not publicly released to reduce potential overfitting. Evaluation on the test set may only be done through the leaderboard. Ablation stud-

Table 4: Predicting relaxed state energy from initial structure (*IS2RE*) as evaluated by Mean Absolute Error (MAE) of the energies and the percentage of Energies within a Threshold (EwT) of the ground truth energy. Results reported for models trained on the All training dataset.

Model	Approach	<i>IS2RE</i> Test							
		Energy MAE [eV] ↓				EwT ↑			
		ID	OOD Ads	OOD Cat	OOD Both	ID	OOD Ads	OOD Cat	OOD Both
Median baseline	-	1.7489	1.8911	1.7107	1.6807	0.75%	0.69%	0.83%	0.78%
CGCNN ⁸⁷	Direct	0.6135	0.9155	0.6211	0.8506	3.41%	1.93%	3.11%	1.99%
SchNet ⁸⁸	Direct	0.6372	0.7342	0.6611	0.7035	2.96%	2.33%	2.95%	2.22%
DimeNet++ ^{89,90}	Direct	0.5605	0.7252	0.5750	0.6613	4.26%	2.06%	4.10%	2.42%
SchNet ⁸⁸	Relaxation	0.7088	0.7741	0.7665	0.8055	4.23%	2.63%	3.52%	2.52%
SchNet ⁸⁸ – force-only + energy-only	Relaxation	0.7066	0.7420	0.7966	0.7493	4.18%	2.98%	3.39%	2.70%
DimeNet++ ^{89,90}	Relaxation	0.6687	0.6864	0.6858	0.6835	4.29%	3.36%	3.79%	3.51%
DimeNet++ ^{89,90} – force-only + energy-only	Relaxation	0.5112	0.5744	0.5922	0.6130	6.14%	4.29%	5.10%	3.84%
DimeNet++ ^{89,90} – large force-only + energy-only	Relaxation	0.5022	0.5430	0.5780	0.6117	6.58%	4.34%	5.09%	3.93%

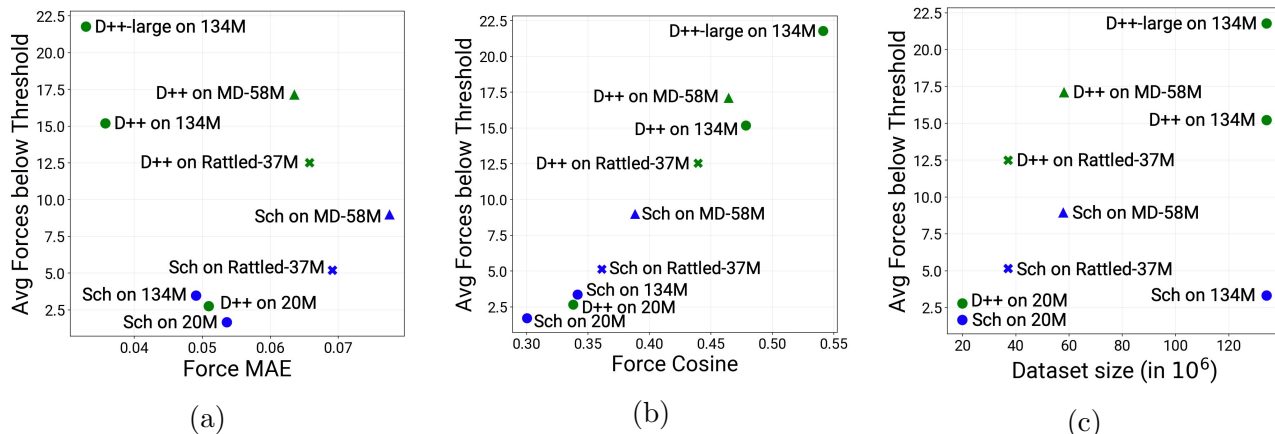


Figure 5: Results of force-only SchNet (denoted by ‘Sch’) and DimeNet++ (‘D++’) *S2EF* models trained on *S2EF-20M*, *S2EF-100M*, *S2EF-20M* + Rattled (‘Rattled-37M’) and *S2EF-20M* + MD (‘MD-58M’) dataset splits used to drive relaxations from given initial structures (*IS2RS*). We plot *IS2RS* AFbT performance against *S2EF* force cosine, *S2EF* force MAE and number of training samples for the different variants. 5a,5b: *IS2RS* AFbT seems to correlate better with *S2EF* force cosine than *S2EF* force MAE, especially when analyzing models trained on Rattled-37M or MD-58M data. 5c: Further, both DimeNet++ and SchNet achieve higher AFbT when trained on MD-58M than *S2EF-134M*. Additional MD data seems to offer a stronger learning signal than additional *S2EF* data.

ies and hyper-parameter tuning may be done and reported on using the validation datasets.

Results

To provide baselines for the OC20 dataset, we report results using three state-of-the-art approaches: CGCNN,⁸⁷ SchNet,⁸⁸ and DimeNet++.^{89,90} Details of the models’ implementations can be found in the Baselines Section.

S2EF: Results on CGCNN,⁸⁷ SchNet,⁸⁸ and DimeNet++^{89,90} are evaluated. All approaches predict structure energies in their forward pass and per-atom forces by the negative gradient of the predicted energy with respect to atomic positions.⁹⁴ Across most metrics DimeNet++ performs the best, with SchNet marginally outperforming DimeNet++ and CGCNN on EFwT. SchNet outperforms CGCNN across all metrics. Since tradeoffs exist in the prediction of energy and forces, we trained

three variants of SchNet and DimeNet++ with $\{\lambda_E, \lambda_F\} = \{1, 30\}, \{0, 100\}, \{100, 1\}$ for SchNet/DimeNet++, SchNet/DimeNet++ force-only and SchNet/DimeNet++ energy-only respectively. As expected, the energy-only model performs best on energy MAE, while the force-only performs best on force MAE. DimeNet++ and SchNet both provide a balance between the two and the best results on EFwT. All approaches perform badly on the EFwT metric; indicating that the results are still far from being practically useful. Table 2 and Figure 4 show results across subsplits. As expected, the In Domain (ID) achieves the best results and the OOD Both performs the worst. However, results are not dramatically different between In Domain, OOD Adsorbate and OOD Catalyst, which shows some generalization to new adsorbates and catalysts. Increases in training data sizes results in significant improvements, Figure 6(A). The rate and amount of improvement varies based on the model. Finally, wider and deeper models are shown to improve accuracies in Figure 4. Both increased depth (Medium to Large) and width (Small to Medium) show improvements.

IS2RS: For *IS2RS*, we use our *S2EF* baselines to drive ML relaxations from the given initial structures to estimate the relaxed structures using L-BGFS,⁹⁵ examples are shown in Figure 3(A). Table 3 shows that DimeNet++ outperforms SchNet in the ADwT and AFbT metrics. However, the FbT metrics indicate both methods do not produce relaxed structures with forces below thresholds used in practice. Since only the computed forces are used for the IS2RS task and not the energies, it is not surprising that the DimeNet++ force-only model performs the best. It was trained using only force losses and performs significantly better on AFbT and ADwT, but still is near zero when measured by FbT. A plot of FbT across thresholds from 0.01 to 0.6 for SchNet is shown in Figure 3(C). Both methods show better generalization to new adsorbates vs new catalyst material compositions. Similar to *S2EF* improved results are found with more training data, especially for DimeNet++ and SchNet, Figure 6(B). Experiments using the additional

rattled and MD data are shown in Figure 5. Interestingly, the force cosine metric appears to better correlate with AFbT scores than force MAE. A discussion on these results may be found in the supplementary.

IS2RE: For *IS2RE* we explore two pathways for computing the relaxed energy from the initial state, Figure 3(B). The first directly computes the relaxed energy given the initial state. The same model architectures are used as the *S2EF* task, but with new weights learned. The second approach uses models trained on the *S2EF* task to perform ML relaxations from which the resulting energy is returned. Note that the ML relaxation approach is about 200 times more expensive to compute, since energies needs to be computed at each relaxation step.

As shown in Table 4, the hybrid relaxation approaches outperformed the direct across all metrics. The percentage of predicted energies within a tight threshold (EwT) ranged from 2% to 6%; indicating that accuracies are still below practical usefulness. Generalization to new catalyst compositions performed better than new adsorbates. As shown in Figure 6(C), larger dataset sizes could significantly improve performance. The best direct-based approach, DimeNet++, was evaluated via the relaxation-based approach. The use of DimeNet++ force-only to perform the relaxation, followed by DimeNet++ energy-only to compute the relaxed energy significantly outperformed the use of a single model (optimized for EFwT) to compute both. Best metrics were achieved using the large DimeNet++ force-only model, followed by DimeNet++ energy-only.

Outlook and Future Directions

The baseline models in this work give significant insights into the complexity of day-to-day challenges in catalysis and what it will take to achieve generalizable models. Motivated by previous efforts,⁹⁸ we analyzed model performance for increasing dataset sizes to illustrate the differences between catalysis and

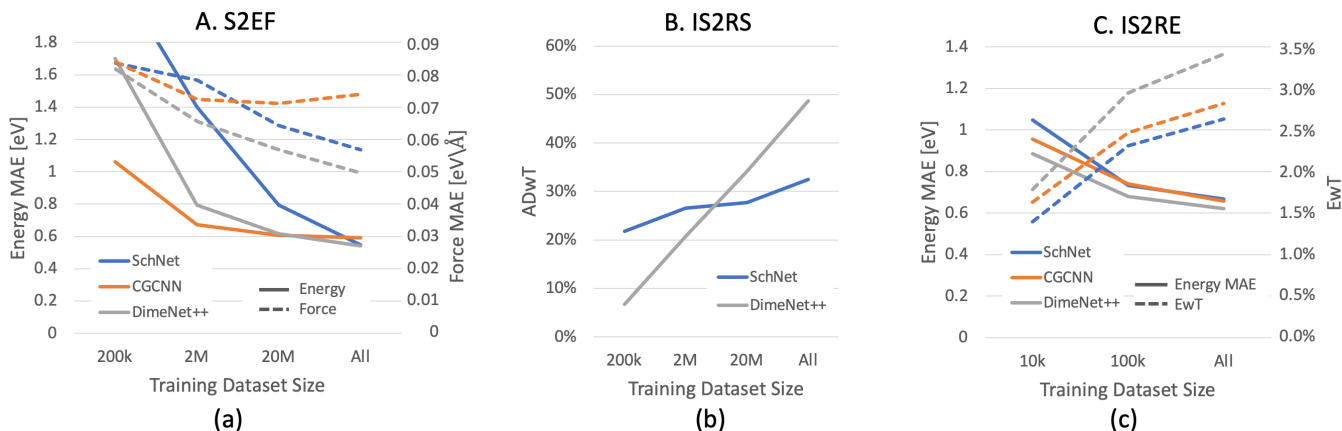


Figure 6: (A) Predicting energy and forces from a structure ($S2EF$) as evaluated by Mean Absolute Error (MAE) of the energies and forces. (B) Predicting relaxed structure from initial structure ($IS2RS$) as evaluated by Average Distance within Threshold (ADwT) using force-only models. (C) Predicting relaxed state energy from initial structure ($IS2RE$) as evaluated by Mean Absolute Error (MAE) of the energies and the percentage of Energies within a Threshold (EwT, $\epsilon = 0.02$ eV) of the ground truth energy. Results reported for $S2EF$ and $IS2RS$ trained on 200k, 2M, 20M and All dataset sizes. Results reported for $IS2RE$ trained on 10k, 100k, and All dataset sizes. $S2EF$ and $IS2RE$ values averaged across validation subsplits. $IS2RS$ values evaluated on the test in-domain (ID) subsplit.

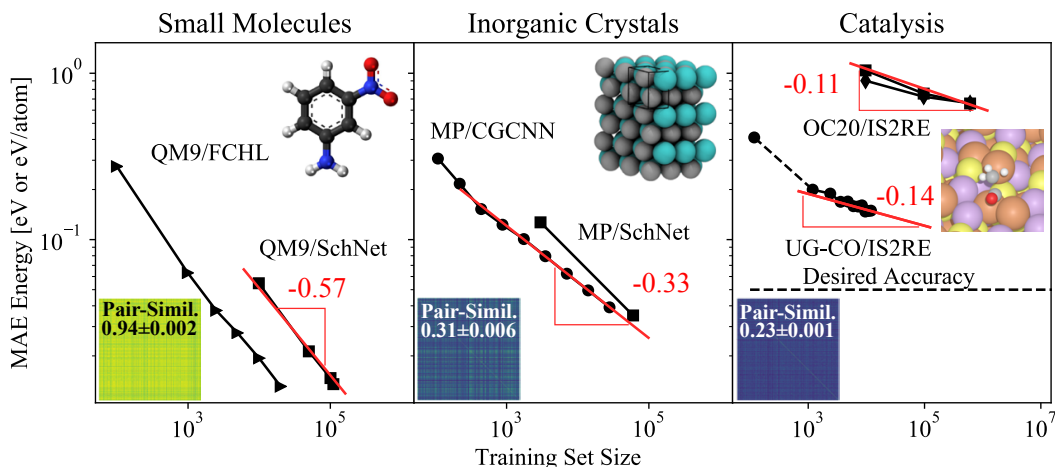


Figure 7: Model performance versus dataset size across three related atomistic domains. Insets are pairwise similarity for selected structures from the respective dataset using GraphDot (see the SI for details) (0/dark-blue/not-similar to 1/yellow/identical).^{96,97} (left) Results⁶³ for FCHL/SchNet models trained on the QM9 small molecule dataset (slope -0.57). (middle) Models^{87,88} trained on Materials Project formation energies (slope -0.33, more difficult). (right) Results for catalysis including a literature dataset for CO adsorbates¹⁶ and this work (slope -0.11 to -0.14, most difficult). Note that reaching the desired accuracy will require several orders of magnitude more data with current models.

related efforts—e.g., materials sciences or small molecule property prediction. Figure 7(left) and Figure 7(middle) show the performance of GNN models similar to the baseline models in this work on datasets for small molecules

(QM9) and materials (formation energies from the materials project). The scaling of model accuracy with respect to dataset size is related to the effective dimensionality of the task and the effective representation in the model. Compar-

ing DimeNet++ performance across all three tasks shows that the aggressive scaling for small molecules is reduced for inorganic materials, and further reduced for surfaces. Focusing on results from this study in Figure 7(right) shows that the scaling is similar for the same baseline models trained on the OC20 dataset and a related literature dataset of CO adsorption energies (see the SI). Importantly, this suggests that achieving the desired accuracy using the current baseline models would require a dataset nearly 10 orders of magnitude larger than the current dataset. This implies that this problem will not be solved through brute-force methods alone, and that significantly improved ML representations are also necessary. This is an exciting opportunity for the broader community.

For the computer science and ML communities, we expect that this dataset will provide unique challenges and spur innovation in atomistic simulations. Many state-of-the-art methods for organic and inorganic materials are based on graph convolutional networks,⁸⁷ which have seen rapid progress. With the above perspective, we expect that additional creative solutions will be necessary to fully solve these tasks. While they have not been demonstrated for inorganic materials, physics-informed tensor representations for small molecules may be helpful.^{99–102} Element embeddings and representations will be important to scale across materials. Incorporation of lower-level physics-based potentials is welcomed and encouraged. This includes the use of related datasets (organic molecules or inorganic materials) for pre-training or learning priors. Incorporating other electronic features in the training set, such as charge distribution to correctly localize effects is also an opportunity to effectively reduce the dimensionality of the problem.

Note that the size of this dataset is larger by 2 orders of magnitude than previous catalyst DFT dataset efforts.^{16,103} Along with the potential for more accurate ML models, it provides practical challenges to training atomistic machine learning models at scale, similar to software engineering challenges in image recognition and NLP.^{104,105} The largest baseline models with *ca.* 10 million parameters were trained

on upwards of 32 GPUs at a time, so we encourage the catalysis community to take advantage of these GPU-enabled resources. This is well-timed with the wave of large GPU-enabled supercomputers that are well-suited to these challenges, such as Perlmutter (DOE NERSC) or Summit (DOE OLCF), among many others.

The baseline models in this work represent the state-of-the-art for deep learning methods to predict thermochemistry for small molecules on inorganic surfaces. Solving this challenge with future model development efforts would enable a new generation of computational chemistry methods. In particular, on-the-fly thermochemistry for reaction intermediates would enable reaction mechanism prediction across materials or composition space. Accelerated methods would also enable the more routine use of more accurate computational methods (e.g. hybrid, exact-exchange, or RPA calculations) by focusing these efforts on the most promising and pre-relaxed structures. A solution to the S2EF task would enable transition state calculations, kinetic approximations, vibrational frequency calculations, and the more routine use of long timescale molecular dynamics for studying these systems. Sensitivity analyses will be necessary to understand the level of accuracy needed for models to be practically relevant for varying applications. Given the sparsity and breadth of OC20, the availability of relevant experimental data will also be a crucial challenge in the next stage of validating model results with experiments. The potential applicability of the OC20 dataset is not just catalysis, but also has implications for areas where organic and inorganic materials interact, such as water quality remediation, geochemistry, advanced manufacturing, and durable energy materials.

Supporting Information Available

The supporting information contains details on the precise DFT calculation methods, the adsorption energy reference energies, the adsorbates and their assumed binding config-

urations, details on graph construction, description of the graph similarity metrics, a few sample GFN0-XTB relaxations, the precise train/test/validation splits, details on the modified CGCNN/SchNet/DimeNet++ implementations, results on the Rattled/MD experiments, hyperparameters for baseline models, a list of adsorbates in OC20, and full results on the validation splits. The full open dataset is provided at <http://opencatalystproject.org> in accessible extxyz format, and the baseline models are provided as an open source repository at <https://github.com/Open-Catalyst-Project/ocp>.

Acknowledgement This research used resources of the National Energy Research Scientific Computing Center (NERSC), a U.S. Department of Energy Office of Science User Facility operated under Contract No. DE-AC02-05CH11231. B.W. acknowledges support from the NERSC Early Science Application Program. The authors acknowledge very helpful discussions with John Kitchin (CMU), Dionisios Vlachos (UD), and Philippe Sautet (UCLA) on the dataset construction and usage.

References

- (1) Newell, R. G.; Raimi, D.; Villanueva, S.; Prest, B. *Global Energy Outlook 2020: Energy Transition or Energy Addition? With Commentary on Implications of the COVID-19 Pandemic*; 2020.
- (2) *Annual Energy Outlook 2020*; U.S. Energy Information Administration, 2020.
- (3) Nørskov, J. K.; Studt, F.; Abild-Pedersen, F.; Bligaard, T. *Fundamental concepts in heterogeneous catalysis*; John Wiley & Sons, 2014; pp 1–4.
- (4) She, Z. W.; Kibsgaard, J.; Dickens, C. F.; Chorkendorff, I.; Nørskov, J. K.; Jaramillo, T. F. Combining theory and experiment in electrocatalysis: Insights into materials design. *Science* **2017**, *355*, p. N/A.
- (5) Nørskov, J. K.; Bligaard, T. *The catalyst genome*. 2013.
- (6) Sholl, D. S.; Steckel, J. A. *Density Functional Theory*; John Wiley & Sons, Inc.: Hoboken, NJ, USA, 2009; pp 1–31.
- (7) Matera, S.; Schneider, W. F.; Heyden, A.; Savara, A. Progress in Accurate Chemical Kinetic Modeling, Simulations, and Parameter Estimation for Heterogeneous Catalysis. *ACS Catalysis* **2019**, *9*, 6624–6647.
- (8) Medford, A. J.; Kunz, M. R.; Ewing, S. M.; Borders, T.; Fushimi, R. Extracting Knowledge from Data through Catalysis Informatics. *ACS Catalysis* **2018**, *8*, 7403–7429.
- (9) Ulissi, Z. W.; Medford, A. J.; Bligaard, T.; Nørskov, J. K. To address surface reaction network complexity using scaling relations machine learning and DFT calculations. *Nature Communications* **2017**, *8*, 1–7.
- (10) Li, B.; Rangarajan, S. Designing compact training sets for data-driven molecular property prediction through optimal exploitation and exploration. *Molecular Systems Design & Engineering* **2019**, *4*, 1048–1057.
- (11) Gu, G. H.; Plechac, P.; Vlachos, D. G. Thermochemistry of gas-phase and surface species via LASSO-assisted subgraph selection. *Reaction Chemistry & Engineering* **2018**, *3*, 454–466.
- (12) Liu, Y.; Hong, W.; Cao, B. Machine learning for predicting thermodynamic properties of pure fluids and their mixtures. *Energy* **2019**, *188*, 116091.
- (13) Jirasek, F.; Alves, R. A.; Damay, J.; Vandermeulen, R. A.; Bamler, R.; Bortz, M.; Mandt, S.; Kloft, M.; Hasse, H. Machine Learning in Thermodynamics: Prediction of Activity Coefficients by Matrix Completion. *Journal of Physical Chemistry Letters* **2020**, *11*, 981–985.

- (14) Kauwe, S. K.; Graser, J.; Vazquez, A.; Sparks, T. D. Machine learning prediction of heat capacity for solid inorganics. *Integrating Materials and Manufacturing Innovation* **2018**, *7*, 43–51.
- (15) Tran, K.; Neiswanger, W.; Yoon, J.; Zhang, Q.; Xing, E.; Ulissi, Z. W. Methods for comparing uncertainty quantifications for material property predictions. *Machine Learning: Science and Technology* **2020**, *1*, 025006.
- (16) Tran, K.; Ulissi, Z. W. Active learning across intermetallics to guide discovery of electrocatalysts for CO₂ reduction and H₂ evolution. *Nature Catalysis* **2018**, *1*, 696–703.
- (17) Gu, G. H.; Noh, J.; Kim, S.; Back, S.; Ulissi, Z.; Jung, Y. Practical Deep-Learning Representation for Fast Heterogeneous Catalyst Screening. *The Journal of Physical Chemistry Letters* **2020**, *11*, 3185–3191, PMID: 32191473.
- (18) Back, S.; Yoon, J.; Tian, N.; Zhong, W.; Tran, K.; Ulissi, Z. W. Convolutional neural network of atomic surface structures to predict binding energies for high-throughput screening of catalysts. *The Journal of Physical Chemistry Letters* **2019**, *10*, 4401–4408.
- (19) Kim, M.; Yeo, B. C.; Park, Y.; Lee, H. M.; Han, S. S.; Kim, D. Artificial Intelligence to Accelerate the Discovery of N₂ Electroreduction Catalysts. *Chemistry of Materials* **2019**, *32*, 709–720.
- (20) Esterhuizen, J. A.; Goldsmith, B. R.; Linic, S. Theory-Guided Machine Learning Finds Geometric Structure-Property Relationships for Chemisorption on Sub-surface Alloys. *Chem* **2020**, *6*, 3100 – 3117.
- (21) Sun, G.; Sautet, P. Metastable Structures in Cluster Catalysis from First-Principles: Structural Ensemble in Reaction Conditions and Metastability Triggered Reactivity. *Journal of the American Chemical Society* **2018**, *140*, 2812–2820.
- (22) Li, B.; Huang, C.; Li, X.; Zheng, S.; Hong, J. Non-iterative structural topology optimization using deep learning. *CAD Computer Aided Design* **2019**, *115*, 172–180.
- (23) Hayashi, K.; Ohsaki, M. Reinforcement Learning and Graph Embedding for Binary Truss Topology Optimization Under Stress and Displacement Constraints. *Frontiers in Built Environment* **2020**, *6*, 59.
- (24) Schmidt, J.; Marques, M. R.; Botti, S.; Marques, M. A. Recent advances and applications of machine learning in solid-state materials science. *npj Computational Materials* **2019**, *5*, 1–36.
- (25) Aksoz, Z.; Preisinger, C. *Impact: Design With All Senses*; Springer International Publishing, 2020; pp 18–31.
- (26) Boes, J. R.; Kitchin, J. R. Neural network predictions of oxygen interactions on a dynamic Pd surface. *Molecular Simulation* **2017**, *43*, 346–354.
- (27) Khorshidi, A.; Peterson, A. A. Amp: A modular approach to machine learning in atomistic simulations. *Computer Physics Communications* **2016**, *207*, 310–324.
- (28) Peterson, A. A. Acceleration of saddle-point searches with machine learning. *The Journal of Chemical Physics* **2016**, *145*, 074106.
- (29) Sun, G.; Sautet, P. Toward fast and reliable potential energy surfaces for metallic Pt clusters by hierarchical delta neural networks. *Journal of Chemical Theory and Computation* **2019**, *15*, 5614–5627.
- (30) Timoshenko, J.; Frenkel, A. I. “Inverting” X-ray absorption spectra of catalysts by machine learning in search for activity descriptors. *ACS Catalysis* **2019**, *9*, 10192–10211.

- (31) Kitchin, J. R. Machine learning in catalysis. *Nature Catalysis* **2018**, *1*, 230–232.
- (32) Li, Z.; Wang, S.; Xin, H. Toward artificial intelligence in catalysis. *Nature Catalysis* **2018**, *1*, 641–642.
- (33) Medford, A. J.; Kunz, M. R.; Ewing, S. M.; Borders, T.; Fushimi, R. Extracting knowledge from data through catalysis informatics. *ACS Catalysis* **2018**, *8*, 7403–7429.
- (34) Goldsmith, B. R.; Esterhuizen, J.; Liu, J.-X.; Bartel, C. J.; Sutton, C. Machine learning for heterogeneous catalyst design and discovery. *AIChE Journal* **2018**, *64*, 2311–2323.
- (35) Schlexer Lamoureux, P.; Winther, K. T.; Garrido Torres, J. A.; Streibel, V.; Zhao, M.; Bajdich, M.; Abild-Pedersen, F.; Bligaard, T. Machine learning for computational heterogeneous catalysis. *ChemCatChem* **2019**, *11*, 3581–3601.
- (36) Li, H.; Zhang, Z.; Liu, Z. Application of artificial neural networks for catalysis: a review. *Catalysts* **2017**, *7*, 306.
- (37) Tabor, D. P.; Roch, L. M.; Saikin, S. K.; Kreisbeck, C.; Sheberla, D.; Montoya, J. H.; Dwaraknath, S.; Aykol, M.; Ortiz, C.; Tribukait, H., et al. Accelerating the discovery of materials for clean energy in the era of smart automation. *Nature Reviews Materials* **2018**, *3*, 5–20.
- (38) Chen, C.; Zuo, Y.; Ye, W.; Li, X.; Deng, Z.; Ong, S. P. A Critical Review of Machine Learning of Energy Materials. *Advanced Energy Materials* **2020**, *10*, 1903242.
- (39) Artrith, N. Machine learning for the modeling of interfaces in energy storage and conversion materials. *Journal of Physics: Energy* **2019**, *1*, 032002.
- (40) Gu, G. H.; Noh, J.; Kim, I.; Jung, Y. Machine learning for renewable energy materials. *Journal of Materials Chemistry A* **2019**, *7*, 17096–17117.
- (41) Toyao, T.; Maeno, Z.; Takakusagi, S.; Kamachi, T.; Takigawa, I.; Shimizu, K.-i. Machine Learning for Catalysis Informatics: Recent Applications and Prospects. *ACS Catalysis* **2019**, *10*, 2260–2297.
- (42) Gu, G. H.; Choi, C.; Lee, Y.; Situmorang, A. B.; Noh, J.; Kim, Y.-H.; Jung, Y. Progress in Computational and Machine-Learning Methods for Heterogeneous Small-Molecule Activation. *Advanced Materials* **2020**, *32*, 1907865.
- (43) Schlexer Lamoureux, P.; Winther, K. T.; Garrido Torres, J. A.; Streibel, V.; Zhao, M.; Bajdich, M.; Abild-Pedersen, F.; Bligaard, T. Machine Learning for Computational Heterogeneous Catalysis. 2019.
- (44) Takahashi, K.; Takahashi, L.; Miyazato, I.; Fujima, J.; Tanaka, Y.; Uno, T.; Satoh, H.; Ohno, K.; Nishida, M.; Hirai, K.; Ohyama, J.; Nguyen, T. N.; Nishimura, S.; Taniike, T. The Rise of Catalyst Informatics: Towards Catalyst Genomics. 2019.
- (45) Curtarolo, S.; Setyawan, W.; Wang, S.; Xue, J.; Yang, K.; Taylor, R. H.; Nelson, L. J.; Hart, G. L. W.; Sanvito, S.; Buongiorno-Nardelli, M.; Mingo, N.; Levy, O. AFLOWLIB.ORG: A distributed materials properties repository from high-throughput ab initio calculations. *Computational Materials Science* **2012**, *58*, 227–235.
- (46) Kirklin, S.; Saal, J. E.; Meredig, B.; Thompson, A.; Doak, J. W.; Aykol, M.; Rühl, S.; Wolverton, C. The Open Quantum Materials Database (OQMD): Assessing the accuracy of DFT formation energies. *npj Computational Materials* **2015**, N/A.

- (47) Calle-Vallejo, F.; Martínez, J. I.; García-Lastra, J. M.; Sautet, P.; Loffreda, D. Fast Prediction of Adsorption Properties for Platinum Nanocatalysts with Generalized Coordination Numbers. *Angew. Chem. Int. Ed.* **2014**, *53*, 8316–8319.
- (48) Abild-Pedersen, F.; Greeley, J.; Studt, F.; Rossmeisl, J.; Munter, T.; Moses, P. G.; Skulason, E.; Bligaard, T.; Nørskov, J. K. Scaling Properties of Adsorption Energies for Hydrogen-containing Molecules on Transition-metal Surfaces. *Phys. Rev. Lett.* **2007**, *99*, 016105.
- (49) Ma, X.; Xin, H. Orbitalwise Coordination Number for Predicting Adsorption Properties of Metal Nanocatalysts. *Phys. Rev. Lett.* **2017**, *118*, 036101.
- (50) Noh, J.; Back, S.; Kim, J.; Jung, Y. Active learning with Non-ab Initio Input Features toward Efficient CO₂ Reduction Catalysts. *Chem. Sci.* **2018**, *9*, 5152–5159.
- (51) Andersen, M.; Levchenko, S. V.; Scheffler, M.; Reuter, K. Beyond Scaling Relations for the Description of Catalytic Materials. *ACS Catalysis* **2019**, *9*, 2752–2759.
- (52) Dickens, C. F.; Montoya, J. H.; Kulkarni, A. R.; Bajdich, M.; Nørskov, J. K. An Electronic Structure Descriptor for Oxygen Reactivity at Metal and Metal-oxide Surfaces. *Surf. Sci.* **2019**, *681*, 122–129.
- (53) Li, Z.; Wang, S.; Chin, W. S.; Achenie, L. E.; Xin, H. High-throughput Screening of Bimetallic Catalysts Enabled by Machine Learning. *J. Mater. Chem. A* **2017**, *5*, 24131–24138.
- (54) Batchelor, T. A.; Pedersen, J. K.; Winther, S. H.; Castelli, I. E.; Jacobsen, K. W.; Rossmeisl, J. High-entropy Alloys as a Discovery Platform for Electrocatalysis. *Joule* **2019**, *3*, 834–845.
- (55) Mamun, O.; Winther, K. T.; Boes, J. R.; Bligaard, T. High-throughput calculations of catalytic properties of bimetallic alloy surfaces. *Scientific data* **2019**, *6*, 76.
- (56) Winther, K. T.; Hoffmann, M. J.; Boes, J. R.; Mamun, O.; Bajdich, M.; Bligaard, T. Catalysis-Hub.org, an open electronic structure database for surface reactions. *Scientific Data* **2019**, N/A.
- (57) Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; Fei-Fei, L. Imagenet: A large-scale hierarchical image database. 2009 IEEE conference on computer vision and pattern recognition. 2009; pp 248–255.
- (58) Panayotov, V.; Chen, G.; Povey, D.; Khudanpur, S. Librispeech: an asr corpus based on public domain audio books. 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). 2015; pp 5206–5210.
- (59) Antol, S.; Agrawal, A.; Lu, J.; Mitchell, M.; Batra, D.; Lawrence Zitnick, C.; Parikh, D. Vqa: Visual question answering. Proceedings of the IEEE international conference on computer vision. 2015; pp 2425–2433.
- (60) Jain, A.; Ong, S. P.; Hautier, G.; Chen, W.; Richards, W. D.; Dacek, S.; Cholia, S.; Gunter, D.; Skinner, D.; Ceder, G.; Persson, K. The Materials Project: A materials genome approach to accelerating materials innovation. *APL Materials* **2013**, *1*, 011002.
- (61) Henkelman, G.; Arnaldsson, A.; Jónsson, H. A fast and robust algorithm for Bader decomposition of charge density. *Computational Materials Science* **2006**, *36*, 354–360.
- (62) Bader, R.; Bader, R. *Atoms in Molecules: A Quantum Theory*; International series of monographs on chemistry; Clarendon Press, 1994; pp 13–52.

- (63) von Lilienfeld, O. A.; Müller, K.-R.; Tkatchenko, A. Exploring chemical compound space with quantum-based machine learning. *Nature Reviews Chemistry* **2020**, 1–12.
- (64) Ong, S. P.; Richards, W. D.; Jain, A.; Hautier, G.; Kocher, M.; Cholia, S.; Gunter, D.; Chevrier, V. L.; Persson, K. A.; Ceder, G. Python Materials Genomics (pymatgen): A robust, open-source python library for materials analysis. *Computational Materials Science* **2013**, 68, 314–319.
- (65) Boes, J. R.; Mamun, O.; Winther, K.; Bligaard, T. Graph theory approach to high-throughput surface adsorption structure generation. *The Journal of Physical Chemistry A* **2019**, 123, 2281–2285.
- (66) Larsen, A. H. et al. The atomic simulation environment—a Python library for working with atoms. *Journal of Physics: Condensed Matter* **2017**, 29, 273002.
- (67) Kresse, G.; Hafner, J. Ab initio molecular-dynamics simulation of the liquid-metal–amorphous-semiconductor transition in germanium. *Physical Review B* **1994**, 49, 14251–14269.
- (68) Kresse, G.; Furthmüller, J. Efficiency of ab-initio total energy calculations for metals and semiconductors using a plane-wave basis set. *Computational Materials Science* **1996**, 6, 15–50.
- (69) Kresse, G.; Furthmüller, J. Efficient iterative schemes for ab initio total-energy calculations using a plane-wave basis set. *Physical Review B* **1996**, 54, 11169–11186.
- (70) “The calculations in this work have been performed using the ab-initio total-energy and molecular-dynamics package VASP (Vienna ab-initio simulation package) developed at the Institut für Materialphysik of the Universität Wien^{2,3}”.
- (71) Kresse, G.; Joubert, D. From ultrasoft pseudopotentials to the projector augmented-wave method. *Physical Review B* **1999**, 59, 1758.
- (72) Perdew, J. P.; Burke, K.; Ernzerhof, M. Generalized gradient approximation made simple. *Physical Review Letters* **1996**, 77, 3865.
- (73) Tang, W.; Sanville, E.; Henkelman, G. A grid-based Bader analysis algorithm without lattice bias. *Journal of Physics: Condensed Matter* **2009**, 21, 084204.
- (74) Sanville, E.; Kenny, S. D.; Smith, R.; Henkelman, G. Improved grid-based algorithm for Bader charge allocation. *Journal of Computational Chemistry* **2007**, 28, 899–908.
- (75) Nelson, R.; Ertural, C.; George, J.; Deringer, V. L.; Hautier, G.; Dronskowski, R. LOBSTER: Local orbital projections, atomic charges, and chemical-bonding analysis from projector-augmented-wave-based density-functional theory. *Journal of Computational Chemistry* **2020**, 41, 1931–1940.
- (76) Deringer, V. L.; Tchougréeff, A. L.; Dronskowski, R. Crystal orbital Hamilton population (COHP) analysis as projected from plane-wave basis sets. *The Journal of Physical Chemistry A* **2011**, 115, 5461–5466.
- (77) Gong, S.; Xie, T.; Zhu, T.; Wang, S.; Fadel, E. R.; Li, Y.; Grossman, J. C. Predicting charge density distribution of materials using a local-environment-based graph convolutional network. *Physical Review B* **2019**, 100, 184103.
- (78) Nagai, R.; Akashi, R.; Sugino, O. Completing density functional theory by machine learning hidden messages from molecules. *npj Computational Materials* **2020**, 6, 1–8.

- (79) Chandrasekaran, A.; Kamal, D.; Batra, R.; Kim, C.; Chen, L.; Ramprasad, R. Solving the electronic structure problem with machine learning. *npj Computational Materials* **2019**, *5*, 1–7.
- (80) Kim, Y.; Kim, E.; Antono, E.; Meredig, B.; Ling, J. Machine-learned metrics for predicting the likelihood of success in materials discovery. **2019**, 1–13.
- (81) Meredig, B.; Antono, E.; Church, C.; Hutchinson, M.; Ling, J.; Paradiso, S.; Blaiszik, B.; Foster, I.; Gibbons, B.; Hattrick-Simpers, J.; Mehta, A.; Ward, L. Can machine learning identify the next high-temperature superconductor? Examining extrapolation performance for materials discovery. *Molecular Systems Design and Engineering* **2018**, *3*, 819–825.
- (82) Fey, M.; Lenssen, J. E. Fast graph representation learning with PyTorch Geometric. *arXiv preprint arXiv:1903.02428* **2019**, p. N/A.
- (83) Paszke, A.; Gross, S.; Massa, F.; Lerer, A.; Bradbury, J.; Chanan, G.; Killeen, T.; Lin, Z.; Gimelshein, N.; Antiga, L., et al. Pytorch: An imperative style, high-performance deep learning library. Advances in neural information processing systems. 2019; pp 8026–8037.
- (84) Hamilton, W. L.; Ying, R.; Leskovec, J. Representation learning on graphs: Methods and applications. *IEEE Data Engineering Bulletin* **2017**, p. N/A.
- (85) Behler, J. Perspective: Machine learning potentials for atomistic simulations. *The Journal of Chemical Physics* **2016**, *145*, 170901.
- (86) Bartók, A. P.; Payne, M. C.; Kondor, R.; Csányi, G. Gaussian approximation potentials: The accuracy of quantum mechanics, without the electrons. *Physical Review Letters* **2010**, *104*, 136403.
- (87) Xie, T.; Grossman, J. C. Crystal graph convolutional neural networks for an accurate and interpretable prediction of material properties. *Physical Review Letters* **2018**, *120*, 145301.
- (88) Schütt, K.; Kindermans, P.-J.; Felix, H. E. S.; Chmiela, S.; Tkatchenko, A.; Müller, K.-R. SchNet: A continuous-filter convolutional neural network for modeling quantum interactions. Advances in neural information processing systems. 2017; pp 991–1001.
- (89) Klicpera, J.; Giri, S.; Margraf, J. T.; Günnemann, S. Fast and Uncertainty-Aware Directional Message Passing for Non-Equilibrium Molecules. *arXiv preprint arXiv:2011.14115* **2020**,
- (90) Klicpera, J.; Groß, J.; Günnemann, S. Directional Message Passing for Molecular Graphs. International Conference on Learning Representations (ICLR). 2020.
- (91) Ramakrishnan, R.; Dral, P. O.; Rupp, M.; von Lilienfeld, O. A. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific Data* **2014**, *1*.
- (92) Pracht, P.; Caldeweyher, E.; Ehlert, S.; Grimme, S. A Robust Non-Self-Consistent Tight-Binding Quantum Chemistry Method for large Molecules. *ChemRxiv* **2019**, p. N/A.
- (93) Bannwarth, C.; Caldeweyher, E.; Ehlert, S.; Hansen, A.; Pracht, P.; Seibert, J.; Spicher, S.; Grimme, S. Extended tight-binding quantum chemistry methods. *Wiley Interdisciplinary Reviews: Computational Molecular Science* **2020**, e01493.
- (94) Pukrittayakamee, A.; Malshe, M.; Hagan, M.; Raff, L.; Narulkar, R.; Bukkapatnum, S.; Komanduri, R. Simultaneous fitting of a potential-energy surface and its corresponding force fields using feed-forward neural networks. *The Journal of chemical physics* **2009**, *130*, 134101.

- (95) Liu, D. C.; Nocedal, J. On the limited memory BFGS method for large scale optimization. *Mathematical programming* **1989**, *45*, 503–528.
- (96) Tang, Y.; Selvitopi, O.; Popovici, D. T.; Buluç, A. A High-Throughput Solver for Marginalized Graph Kernels on GPU. 2020 IEEE International Parallel and Distributed Processing Symposium (IPDPS). 2020; pp 728–738.
- (97) Tang, Y.-H.; de Jong, W. A. Prediction of atomization energy using graph kernel and active learning. *The Journal of Chemical Physics* **2019**, *150*, 044107.
- (98) Huang, B.; Symonds, N. O.; von Lilienfeld, O. A. The fundamentals of quantum machine learning. *arXiv preprint arXiv:1807.04259* **2018**, p. N/A.
- (99) Miller, B. K.; Geiger, M.; Smidt, T. E.; Noé, F. Relevance of Rotationally Equivariant Convolutions for Predicting Molecular Properties. *arXiv preprint arXiv:2008.08461* **2020**, p. N/A.
- (100) Bratholm, L. A.; Gerrard, W.; Anderson, B.; Bai, S.; Choi, S.; Dang, L.; Hanchar, P.; Howard, A.; Huard, G.; Kim, S., et al. A community-powered search of machine learning strategy space to find NMR property prediction models. *arXiv preprint arXiv:2008.05994* **2020**, p. N/A.
- (101) Anderson, B.; Hy, T. S.; Kondor, R. Cormorant: Covariant molecular neural networks. *Advances in Neural Information Processing Systems*. 2019; pp 14537–14546.
- (102) Nigam, J.; Pozdnyakov, S.; Ceriotti, M. Recursive evaluation and iterative contraction of N -body equivariant features. *The Journal of Chemical Physics* **2020**, *153*, 121101.
- (103) Hummelshøj, J. S.; Abild-Pedersen, F.; Studt, F.; Bligaard, T.; Nørskov, J. K. CatApp: A web application for surface chemistry and heterogeneous catalysis. *Angewandte Chemie - International Edition* **2012**, *51*, 272–274.
- (104) He, K.; Zhang, X.; Ren, S.; Sun, J. Deep residual learning for image recognition. *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016; pp 770–778.
- (105) Radford, A.; Wu, J.; Child, R.; Luan, D.; Amodei, D.; Sutskever, I. Language models are unsupervised multitask learners. *OpenAI Blog* **2019**, *1*, 9.
- (106) Blöchl, P. E. Projector augmented-wave method. *Physical Review B* **1994**, *50*, 17953.
- (107) Hammer, B.; Hansen, L. B.; Nørskov, J. K. Improved adsorption energetics within density-functional theory using revised Perdew-Burke-Ernzerhof functionals. *Physical Review B* **1999**, *59*, 7413.
- (108) Monkhorst, H. J.; Pack, J. D. Special points for Brillouin-zone integrations. *Physical Review B* **1976**, *13*, 5188.
- (109) Teter, M. P.; Payne, M. C.; Allan, D. C. Solution of Schrödinger’s equation for large systems. *Physical Review B* **1989**, *40*, 12255.
- (110) Press, W. H.; Flannery, B. P.; Teukolsky, S. A.; Vetterling, W. T. *Numerical Recipes*, Cambridge: Cambridge Univ. 1986.
- (111) Feynman, R. P. Forces in molecules. *Physical Review* **1939**, *56*, 340.

DFT Relaxations

DFT calculations were performed with the *Vienna Ab Initio Simulation Package* (VASP)^{67–71} with periodic boundary conditions and the projector augmented wave (PAW) pseudopotentials.^{71,106} The external electrons were expanded in plane waves with kinetic energy cut-offs of 350 eV. Exchange and correlation effects were taken into account via the generalized gradient approximation⁷² and the revised Perdew-Burke-Ernzerhof (RPBE) functional, because of its improved description of the energetics of atomic and molecular bonding to surfaces.¹⁰⁷ Bulk and surface calculations were performed considering a K-point mesh for the Brillouin zone derived from the unit cell parameters as an on-the-spot method, employing the Monkhorst-Pack grid.¹⁰⁸ The ionic degrees of freedom were relaxed using a Conjugate Gradient minimization.^{109,110} The relaxation was terminated when either the Hellmann-Feynman forces¹¹¹ were less than 0.03 eV/Å or the relaxation required more than 200 steps in a single uninterrupted VASP call. This limit was reset each time the calculation was checkpointed allowing some relaxations to exceed this 200 steps. The final distribution of residual forces is shown in Figure 8 in the SI. Relaxations still converging after approximately 5,000 core-hours were terminated and not included in the dataset. For the electronic degrees of freedom, the energy convergence criteria was fixed to 10^{-4} eV, where no spin magnetism or dispersion corrections were included.

Adsorption Energy

$$E_{ad} = E_{sys} - E_{slab} - E_{gas}$$

Gas phase references, E_{gas} , for each adsorbate was computed as a linear combination of N₂, H₂O, CO, and H₂ resulting in the atomic energies from Table 5.

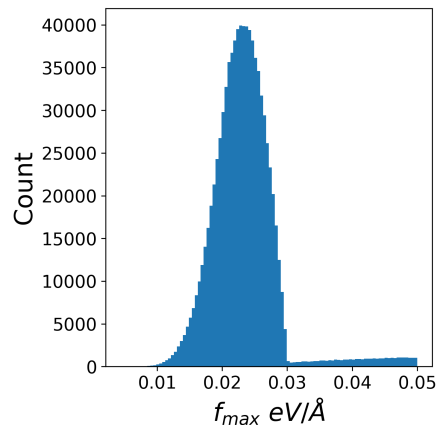


Figure 8: The distribution of max-absolute forces, f_{max} , for systems that converged and completed successfully. Systems in which $f_{max} > 0.05$ eV/Å were excluded from all tasks except S2EF.

Table 5: The per atom energy of individual adsorbate atoms used to calculate the gas phase reference energy for an adsorbate molecule

Adsorbate atom	Energy (eV)
H	-3.477
O	-7.204
C	-7.282
N	-8.083

Computational Workflow

An illustration of the workflow used to sample from the dataset and perform calculations is shown in Figure 9.

Graph Construction

Given a set of atoms in the 3D unit cell that is periodically repeated, we construct a radius graph where nodes represent the atoms and edges represent nearby interaction between pairs of atoms. Specifically, we draw a directed edge from atom j to atom i if atom j is within the cutoff distance from atom i , and vice versa. This means that the edges are always bidirectional. Furthermore, since the nodes are periodically repeated, two atoms may have multiple directed edges if they lie within the cutoff

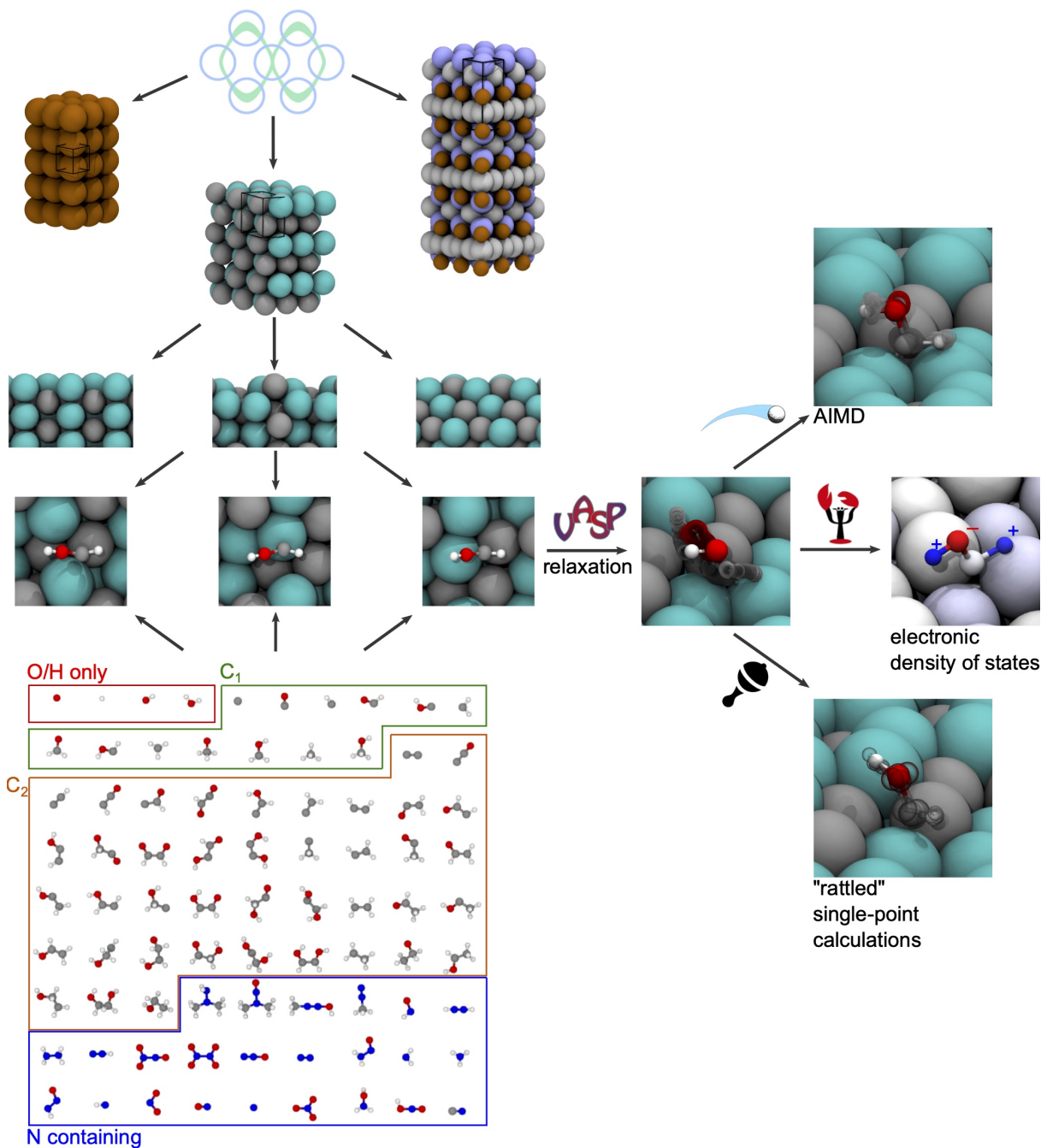


Figure 9: The workflow used to generate the Open Catalyst Dataset. Stable materials were downloaded from The Materials Project⁶⁰ and paired with heuristically chosen adsorbates to create adsorption structures. These structures were randomly sampled for DFT relaxation and then subsequent AIMD, electronic structure analysis, and single-point rattling calculations.

distance in multiple repeated cells. If an atom i has more than one edge to an atom j , each edge represents atom j in a different cell, resulting in unique relative distances and edge features, Figure 10. From the atom-centric view, the above directed multi-graph representation of the atomic system precisely captures the local 3D structure surrounding each atom, taking periodic boundary condition into account.

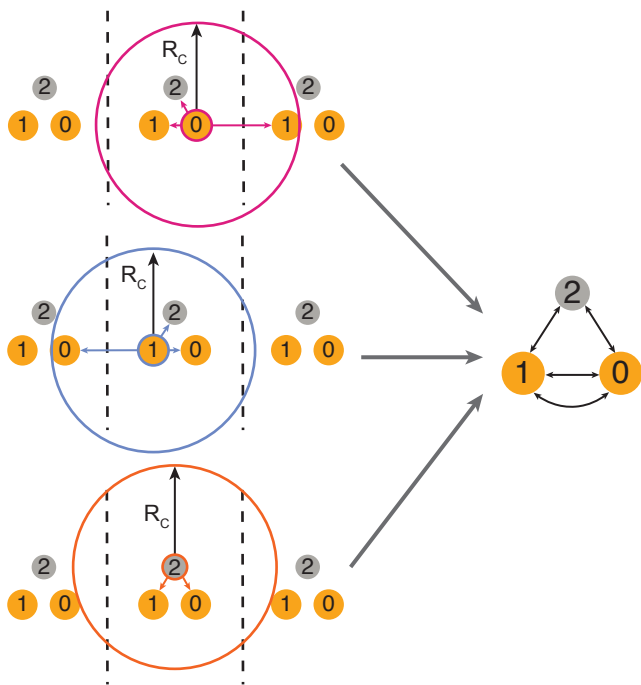


Figure 10: A simple example of constructing a radius graph with periodic boundary conditions. The graph on the right represents all edges assuming each atom as the center node individually (shown on the left).

Graph Pairwise Similarity

The mean pairwise similarity (mps) between a collection of graphs gives an indication of the diversity present in a given dataset and is comparable between different datasets. Pairwise similarity was computed as the mean of the elements in the upper triangle of the similarity matrix (\mathbf{K}) without the diagonal elements included (Equation below). The similarity matrix was calculated using graphs and the molecular kernel from the GraphDot package (<https://graphdot.readthedocs.io/en/latest/>),

details of these methods are provided by Tang et al.⁹⁷ Mean pairwise similarity values range from 1, where all graphs are the same and decay to 0. The mean pairwise similarity can be compared between datasets if the graph and the kernel parameters are consistent. For the results in Figure 6 of the main text, we randomly sampled 1000 systems (N) from a 10,000 subsample of each respective dataset and computed the mean pairwise similarity, this was repeated six times to collect statistics. Random subsampling was done to keep the similarity matrix the same size across datasets and to decrease the computational cost. For the similarity matrix calculation the adjacency length scale used to convert atomic structures to graphs was set to 6 Å and the molecular kernel edge length scale was set to 18 Å, nearly identical results were achieved with 2 Å and 5 Å respectively. All other parameters were set to default values.

$$\text{mps} = \frac{1}{N(N-1)/2} \sum_{i,j} \mathbf{K}_{ij}$$

where $i < j$

Baseline Models Implementation

All proposed baseline models were implemented using PyTorch Geometric. Several implementation changes, however, were necessary to make such models relevant to our dataset and tasks. We outline the modifications below:

SchNet

- Periodic boundary conditions (PBCs) were incorporated into the PyTorch Geometric implementation of SchNet.

DimeNet++

- PBCs were incorporated into the PyTorch Geometric implementation of DimeNet++.

CGCNN

- Similar to SchNet, a Gaussian basis function was incorporated to the edge features. Although not contained within the original CGCNN implementation, a significant performance increase was observed.
- In order to make force predictions, a gradient call was included in the forward pass with respect to positions. The original CGCNN implementation was only concerned with energy predictions.

Hyperparameters for Baseline Models

Model hyperparameters for the ‘All’ splits of the IS2RE and S2EF tasks are provided in Tables 6, 7, and 8. Hyperparameters of the remaining splits can be found in the corresponding repo: <https://github.com/Open-Catalyst-Project/ocp/tree/master/configs>.

Table 6: CGCNN⁸⁷ hyperparameters on the All split of the IS2RE and S2EF tasks.

Hyperparameters	IS2RE	S2EF
Size of atom embeddings	384	512
Size of fully connected layers	512	128
Number of fully connected layers	4	3
Number of graph convolutional layers	6	3
Number of Gaussians used for smearing	100	100
Cutoff distance for interatomic interactions	6	6
Batch size (per gpu)	16	24
Initial learning rate	0.001	0.0005
Learning rate gamma	0.1	0.1
Learning rate milestones	[5, 9, 13]	[3, 5, 7]
Warmup epochs	3	2
Warmup factor	0.2	0.2
Max epochs	20	20
Force coefficient	N/A	10

Table 7: SchNet⁸⁸ hyperparameters on the All split of the IS2RE and S2EF tasks.

Hyperparameters	IS2RE	S2EF
Number of hidden channels	384	1024
Number of filters	128	256
Number of interaction blocks	4	5
Number of Gaussians used for smearing	100	200
Cutoff distance for interatomic interactions	6	6
Global aggregation	add	add
Batch size (per gpu)	64	20
Initial learning rate	0.001	0.0001
Learning rate gamma	0.1	0.1
Learning rate milestones	[10, 15, 20]	[3, 5, 7]
Warmup epochs	3	2
Warmup factor	0.2	0.2
Max epochs	30	15
Force coefficient	N/A	30

Table 8: DimeNet++^{89,90} hyperparameters on the All split of the IS2RE and S2EF tasks.

Hyperparameters	IS2RE	S2EF
Number of hidden channels	256	192
Output block embedding size	192	192
Number of interaction blocks	3	3
Number of radial basis functions	6	6
Number of spherical harmonics	7	7
Number of residual layers before skip connection	1	1
Number of residual layers after skip connection	2	2
Number of linear layers in output blocks	3	3
Cutoff distance for interatomic interactions	6	6
Batch size (per GPU)	4	8
Initial learning rate	0.0001	0.0001
Learning rate gamma	0.1	0.1
Learning rate milestones	[4, 8, 12]	[2, 3, 4]
Warmup epochs	2	2
Warmup factor	0.2	0.2
Max epochs	20	7
Force coefficient	N/A	50

IS2RE Performance of Baseline Models on Previous Datasets

The MAE metrics of the baseline models for the *IS2RE* task are significantly higher than have been reported in recent studies applying ML models to predict adsorption energies.^{15,17,18} There are three key differences in this work. First, the dataset here is larger, more diverse, sparser, and more uniformly sampled than previous datasets making this task more challenging. Second, we are using a more difficult definition of the *IS2RE* task - predict the final energy directly from the initial structure, rather than a clean representation of the final structure.¹⁸ Finally, the baseline models themselves are somewhat different (both implementation, and details of the training and precise form).

To test that the baselines models were consistent with previous efforts, we applied all three models to the *IS2RE* task for a literature dataset of CO adsorption energies,^{16,18} show in Table 9. Our results are consistent, and often better, than previously reported validation accuracy for a CGCNN-based model at approximately 0.190 eV MAE on the literature dataset. This is far lower than the 0.57 eV MAE for our baseline models trained only on the CO subset of the OC20 dataset. This suggests that the dataset diversity is the dominant factor in this variation, and further emphasizes that a uniformly sampled dataset can be more difficult to fit than one obtained through an active learning process that emphasizes high-performing catalysts.

Adsorbates Included

The full list of adsorbates is indicated in Table 10. This list was constructed by considering the four monatomic species and adding common intermediates for renewable energy challenges. The number of possible organic molecules is combinatorially large, so this is not a comprehensive list. Larger molecules (e.g. C3) are also relevant but have an even larger num-

ber of possible configurations. Most adsorbates were mono-dentate (binding through a single adsorbate atom), but larger molecules known to bind in bi-dentate configurations were initialized that way. The atoms considered for either mono-dentate or bi-dentate adsorption location is indicated by *.

Train/Test/Validation Splits

The following adsorbates were reserved for validation subsplits: *CH, *CHO, *COCH₂OH, *COH, *NH₂, *NH₂N(CH₃)₂, and *ONOH. Asterisks represent the binding atoms. The following adsorbates were reserved for the test subsplits: *CH₂*CH₂, *CO, *COHCH₂, *NHN₂, *NNCH₃, *OCHCH₂, and *ONNO₂.

Tight Binding Baseline

Obtaining reasonable energies, forces, and relaxed structures from tight binding codes is an enticing possibility because of the low computational cost compared to DFT; however, tight binding calculations on systems for catalysis remain a challenge, as demonstrated by SI Figure 11. We performed tight binding calculations on 100 random systems from the validation set with extended tight binding (xTB) and the atomic simulation environment (ASE)⁶⁶ interface using the GFN0 parameters.⁹² All xTB calculations were carried out in accordance to our DFT procedures with a few notable differences. For the combined systems, i.e. an adsorbate on a surface, all surface atoms were fixed during the relaxation. Relaxations with xTB featured a BFGS optimizer instead of conjugated gradient, but the convergence criteria remained the same as other DFT calculations, f_{max} of 0.03 eV/Å or a maximum of 200 steps except for adsorbate references where f_{max} was 0.05 eV/Å. Additionally, the surface energies used for the computation of adsorption energies were approximated with single point energies. We did not allow surfaces to relax because of unphysical behavior during optimization, which we likely attribute to periodic boundary conditions (PBCs). We are aware that the xTB

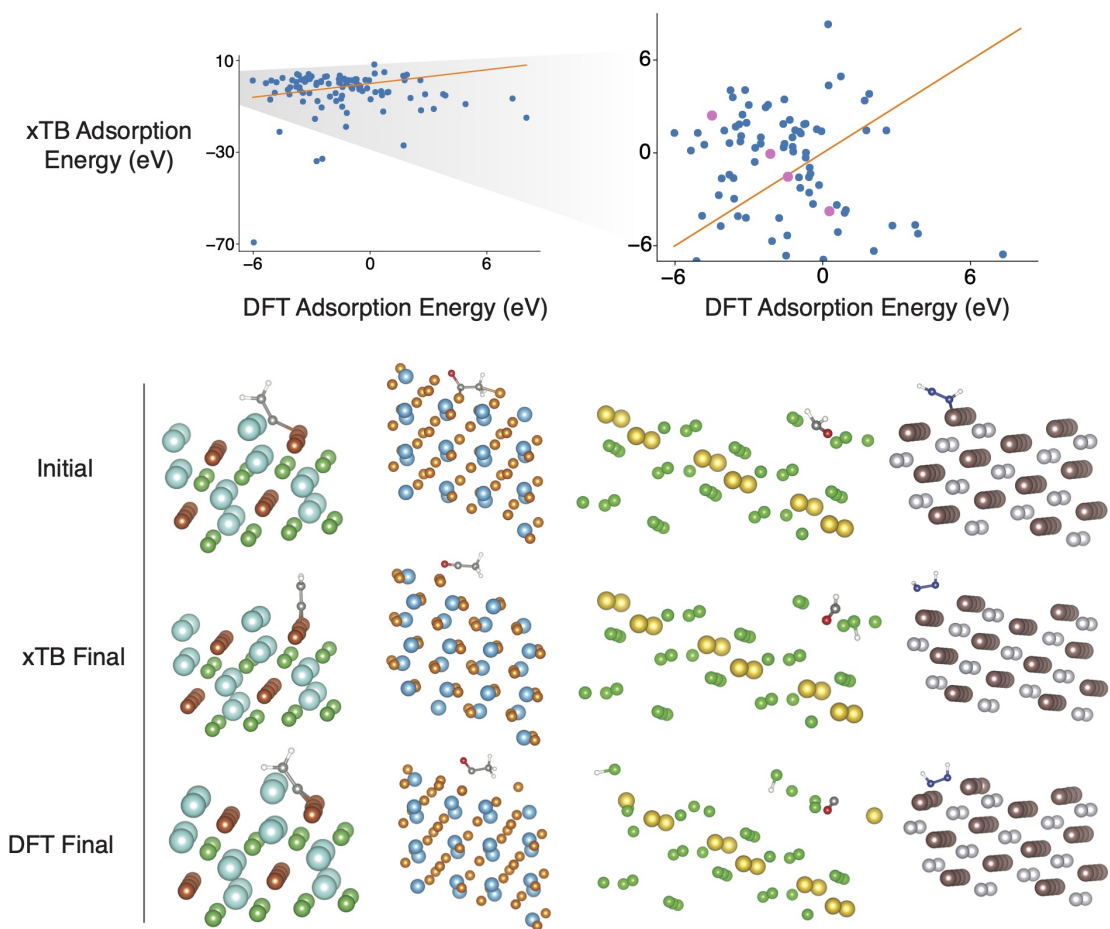


Figure 11: Top: A parity plot comparing xTB adsorption energies with DFT adsorption energies and an inset that limits xTB values to a range similar to that of DFT. Bottom: Initial and final structures corresponding to the pink markers in the plot above organized from left to right.

Table 9: Benchmark of our baseline models’ implementations on a literature CO dataset^{16,18} as evaluated by Energy MAE.

Model	Validation
	Energy MAE [eV] ↓
Previous Work ^{16,18}	0.190
CGCNN ⁸⁷	0.174
SchNet ⁸⁸	0.170
DimeNet++ ^{89,90}	0.149

code was designed for non-periodic systems and that incorporation of PBCs is an ongoing effort. Overall, the speed of the xTB was impressive and we look forward to future developments related to systems with PBCs.

Additional Data: Rattled & Molecular Dynamics

Off-equilibrium data was additionally generated to diversify the structures in the dataset. Two approaches were used to generate this additional data: structural perturbations (“rattled”) and molecular dynamics.

Rattled. Structures along the relaxation path way were sampled, perturbed via random atomic position displacements, and evaluated with DFT. For each relaxation, 20% of the intermediate structures were sampled for rattling. Atomic displacements were sampled from a normal distributions with $\mu = 0$ and $\sigma = 0.05$. Approximately 30 million single-point calculations were carried out. Upon filtering, 17M *S2EF* data points were used for training.

Molecular Dynamics. Short time-scale molecular dynamics simulations were performed on previously relaxed structures. Simulations took place at 900K for 80 or 320 fs with an integration step size of 2 fs in the NVE ensemble. Approximately 64 million single-point calculations were carried out. Upon filtering, 38M *S2EF* data points were used for training.

Performance of baseline models. We report *S2EF* and *IS2RS* results for SchNet⁸⁸ and DimeNet++⁸⁹ models optimized for force-prediction in Table 11. Consistent with results in the main paper, we find that DimeNet++

outperforms SchNet (lower Force MAE, higher Force cosine, higher AFbT). Compared to training only on *S2EF* data, training on MD data seems to provide a complementary learning signal and leads to better sample efficiency – both DimeNet++ and SchNet trained on *S2EF-20M* + MD (58M training samples) outperform corresponding models trained on *S2EF-All* (134M training samples) as per AFbT. Finally, *IS2RS* AFbT seems to correlate better with *S2EF* Force cosine than *S2EF* Force MAE, especially when comparing models trained on Rattled or MD data.

Results on Validation splits

Full results on the validation splits are shown in Tables 13, 14, and 12 for the *S2EF*, *IS2RS*, and *IS2RE* tasks respectively.

Table 10: Adsorbates considered in OC20

Adsorbate class	# of adsorbates	Adsorbates
O/H Only	4	*H, *O, *OH, *OH ₂
C ₁	13	*C, *CO, *CH, *CHO, *COH, *CH ₂ , *CH ₂ *O, *CHOH, *CH ₃ , *OCH ₃ , *CH ₂ OH, *CH ₄ , *OHCH ₃
C ₂	41	*C*C, *CCO, *CCH, *CHCO, *CCHO, *COCHO, *CCHOH, *CCH ₂ , *CH*CH, CH ₂ *CO, *CHCHO, *CH*COH, *COCH ₂ O, *CHO*CHO, *COHCHO, *COHCOH, *CCH ₃ , *CHCH ₂ , *COCH ₃ , *OCHCH ₂ , *COHCH ₂ , *CHCHOH, *CCH ₂ OH, *CHOCHOH, *COCH ₂ OH, *COHCHOH, *CH ₂ *CH ₂ , *OCHCH ₃ , *COHCH ₃ , *CHOHCH ₂ , *CHCH ₂ OH, *OCH ₂ CHOH, *CHOCH ₂ OH, *COHCH ₂ OH, *CHOHCHOH, *CH ₂ CH ₃ , *OCH ₂ CH ₃ , *CHOHCH ₃ , *CH ₂ CH ₂ OH, *CHOHCH ₂ OH, *OHCH ₂ CH ₃
Nitrogen-based	24	*NH ₂ N(CH ₃) ₂ , *ONN(CH ₃) ₂ , *OHNNCH ₃ , *NNCH ₃ , *ONH, *NHNH, *NHN ₂ , *N*NH, *ONNO ₂ , *NO ₂ NO ₂ , *N*NO, *N ₂ , *ONNH ₂ , *NH ₂ , *NH ₃ , *NONH, *NH, *NO ₂ , *NO, *N, *NO ₃ , *OHNH ₂ , *ONOH, *CN

Table 11: *S2EF* and IS2RS results of force-only SchNet and DimeNet++ models on *S2EF*, MD, and Rattled data.

Model	Training Data	# Samples	S2EF Test		IS2RS Test		
			Force MAE	Force cosine	ADwT	FbT	AFbT
SchNet ⁸⁸	S2EF-20M	20M	0.0535	0.3006	27.68%	0.00%	1.68%
SchNet ⁸⁸	S2EF-All	134M	0.0490	0.3417	31.78%	0.00%	3.38%
SchNet ⁸⁸	S2EF-20M + Rattled	37M	0.0691	0.3619	36.70%	0.10%	5.14%
SchNet ⁸⁸	S2EF-20M + MD	58M	0.0775	0.3885	41.10%	0.15%	8.97%
DimeNet++ ^{89,90}	S2EF-20M	20M	0.0509	0.3382	34.37%	0.00%	2.67%
DimeNet++ ^{89,90}	S2EF-All	134M	0.0357	0.4787	48.91%	0.25%	15.17%
DimeNet++ ^{89,90}	S2EF-20M + Rattled	37M	0.0658	0.4395	43.94%	0.05%	12.51%
DimeNet++ ^{89,90}	S2EF-20M + MD	58M	0.0635	0.4644	47.69%	0.15%	17.09%
DimeNet++ ^{89,90} -large	S2EF-All	134M	0.0313	0.5443	51.67%	0.40%	21.74%

Table 12: Predicting relaxed state energy from initial structure (*IS2RE*) as evaluated by Mean Absolute Error (MAE) of the energies and the percentage of Energies within a Threshold (EwT) of the ground truth energy. Results reported for trained on the All training dataset.

Model	Approach	Energy MAE [eV] ↓				EwT ↑			
		ID	OOD Ads	OOD Cat	OOD Both	ID	OOD Ads	OOD Cat	OOD Both
Median baseline	-	1.7466	1.7647	1.7283	1.5640	0.78%	0.80%	0.83%	0.91%
CGCNN ⁸⁷	Direct	0.6203	0.7426	0.6001	0.6708	3.36%	2.11%	3.53%	2.29%
SchNet ⁸⁸	Direct	0.6465	0.7074	0.6475	0.6626	2.96%	2.22%	3.03%	2.38%
DimeNet++ ^{89,90}	Direct	0.5636	0.7127	0.5612	0.6492	4.25%	2.48%	4.40%	2.56%
SchNet ⁸⁸	Relaxation	0.7150	0.7395	0.8010	0.8197	4.03%	3.09%	3.87%	2.72%
SchNet ⁸⁸ – force-only + energy-only	Relaxation	0.7110	0.7574	0.8316	0.8075	4.33%	2.88%	3.63%	2.57%

Table 13: Predicting energy and forces from a structure (*S2EF*) as evaluated by Mean Absolute Error (MAE) of the energies, force MAE, force cosine, and the percentage of Energies and Forces within Threshold (EFwT). Results reported for models trained on the entire training dataset (*S2EF-All*).

<i>S2EF</i> Validation				
Model	ID	OOD Ads	OOD Cat	OOD Both
Energy MAE [eV] ↓				
Median baseline	2.0715	2.2275	2.0558	2.3313
CGCNN ⁸⁷	0.5041	0.5986	0.5252	0.7308
SchNet ⁸⁸	0.4468	0.4973	0.5453	0.7047
SchNet ⁸⁸ – force-only	34.0183	33.4238	34.2519	38.1693
SchNet ⁸⁸ – energy-only	0.4011	0.4727	0.5607	0.7165
DimeNet++ ^{89,90}	0.4545	0.5093	0.5184	0.6753
DimeNet++ ^{89,90} – force-only	28.2095	28.4266	28.8740	35.0468
DimeNet++ ^{89,90} – energy-only	0.3599	0.4500	0.5412	0.7108
DimeNet++ ^{89,90} -Large – force-only	29.3524	29.4825	29.9799	36.6944
Force MAE [eV/Å] ↓				
Median baseline	0.0810	0.0799	0.0798	0.0942
CGCNN ⁸⁷	0.0684	0.0746	0.0679	0.0852
SchNet ⁸⁸	0.0493	0.0574	0.0520	0.0685
SchNet ⁸⁸ – force-only	0.0442	0.0514	0.0465	0.0618
SchNet ⁸⁸ – energy-only	0.5810	0.6254	0.5875	0.6562
DimeNet++ ^{89,90}	0.0443	0.0508	0.0445	0.0589
DimeNet++ ^{89,90} – force-only	0.0331	0.0366	0.0343	0.0436
DimeNet++ ^{89,90} – energy-only	0.3410	0.3322	0.3425	0.3502
DimeNet++ ^{89,90} -Large – force-only	0.0281	0.0318	0.0315	0.0396
Force Cosine ↑				
Median baseline	0.0000	0.000	0.000	0.000
CGCNN ⁸⁷	0.1550	0.1320	0.1456	0.1338
SchNet ⁸⁸	0.3185	0.2862	0.2973	0.2854
SchNet ⁸⁸ – force-only	0.3604	0.3296	0.3294	0.3266
SchNet ⁸⁸ – energy-only	0.0841	0.0695	0.0807	0.0699
DimeNet++ ^{89,90}	0.3632	0.3401	0.3512	0.3556
DimeNet++ ^{89,90} – force-only	0.4877	0.4747	0.4599	0.4849
DimeNet++ ^{89,90} – energy-only	0.1064	0.0855	0.1043	0.0880
DimeNet++ ^{89,90} -Large – force-only	0.5640	0.5500	0.5106	0.5390
EFwT ↑				
Median baseline	0.00%	0.01%	0.01%	0.01%
CGCNN ⁸⁷	0.01%	0.00%	0.00%	0.01%
SchNet ⁸⁸	0.13%	0.00%	0.10%	0.00%
SchNet ⁸⁸ – force-only	0.00%	0.00%	0.00%	0.00%
SchNet ⁸⁸ – energy-only	0.00%	0.00%	0.00%	0.00%
DimeNet++ ^{89,90}	0.09%	0.00%	0.09%	0.00%
DimeNet++ ^{89,90} – force-only	0.00%	0.00%	0.00%	0.00%
DimeNet++ ^{89,90} – energy-only	0.00%	0.00%	0.00%	0.00%
DimeNet++ ^{89,90} -Large – force-only	0.00%	0.00%	0.00%	0.00%

Table 14: Predicting relaxed structure from initial structure (*IS2RS*) as evaluated by Average Distance within Threshold (ADwT). All values in percentages, higher is better. Results reported for structure to energy-force (S2EF) models trained on the All training dataset. The initial structure was used as a naive baseline (IS baseline). Note that metrics requiring expensive DFT calculations – FbT and AFbT – are only computed for test splits, not val.

<i>IS2RS</i> Validation				
Model	ID	OOD Ads	OOD Cat	OOD Both
			ADwT \uparrow	
IS baseline	21.18%	23.49%	20.25%	28.29%
SchNet ⁸⁸	15.53%	16.57%	14.50%	17.29%
SchNet ⁸⁸ – force-only	32.41%	33.33%	30.02%	37.48%
DimeNet++ ^{89,90}	30.40%	30.77%	29.94%	34.89%
DimeNet++ ^{89,90} – force-only	49.05%	46.91%	46.54%	55.23%

Changelog

This section tracks the changes to this document since the original release.

v1. Initial version.

v2.

- DimeNet⁹⁰ results replaced with DimeNet++.^{89,90} DimeNet++ is more memory-efficient and performs slightly better.
- Force cosine similarity added as an additional *S2EF* metric. It correlates better with downstream *IS2RS* AFbT.
- 81 systems removed from the original 1.28M systems due to convergence issues later discovered. Models were not re-trained due to the negligible amount of data ($\sim 0.00675\%$).
- Rattled/MD data experiments added.

v3. Included additional VASP⁶⁷⁻⁷¹ citations.

v4.

- Some systems removed or modified due to subtle convergence and trajectory stitching issues later discovered, affecting a very small proportion of the data. More details can be found at <https://github.com/Open-Catalyst-Project/ocp/blob/master/DATASET.md#version-2-feb-2021>. S2EF models were not retrained due to the negligible amount of data affected.
- IS2RE models retrained and metrics re-evaluated.
- IS2RS - ADwT metrics re-evaluated.

v5.

- A bug was resolved in the IS2RE via relaxation approach. Metrics have been updated, results now outperforming the direct-based approaches.
- Some systems removed from the validation and test splits due to errors found in their initial placements and/or improper split classification ($\sim 0.167\%$ IS2RE, $\sim 0.043\%$ S2EF).
- S2EF, IS2RE, and IS2RS metrics were re-evaluated with the updated splits.