## Supplementary Information

## A Implementation Details

#### A.1 Graph Neural Networks

Our implementation of CGCNN and TorchMD-Net utilizes the MatDeepLearn package 36 as a backbone. Graph representations are calculated through the algorithms released by the Open Catalyst Project (OCP) 38.

We report the main hyperparameters for CGCNN and TorchMD-Net below.

Tabl	e A.1: CGCNN hyperparameters for	pre-training on MP Force
	Parameter	Value or description
	Dimension 1	100
	Dimension 2	150
	Pre-graph convolution FC layers	1
	Graph convolution	4
	Post-graph convolution FC layers	3
	Activation function	silu
	Dropout rate	0
	Pooling	Global add pool
	Pool order	Early

# Table A.2: TorchMD-Net hyperparameters for pre-training on MP Forces Parameter Value or description

	1
Hidden channels	128
Number of filters	128
Number of layers	8
Number of RBF	50
RBF type	expnorm
Trainable RBF	True
Activation	silu
Attention activation	silu
Number of heads	8
Distance influence	both
Number of post layers	2
Post hidden channels	128
Pooling	Global add pool
Pool order	Early

#### A.2 Training Settings

Pre-training datasets, namely MP Forces and MP Forces Relaxed, are split on a train:test:val ratio of 0.8:0.15:0.05. All finetuning datasets are split on a train:test:val ratio of 0.6:0.2:0.2 to ensure consistency and fair comparison. Models are pre-trained for 100 epochs and finetuned for 200 epochs for all datasets.

For graph representation generation, we used a cutoff radius of 8.0 and a maximum neighbor count of 250.

In the case of pre-training via denoising, the amount of noise added to atomic positions is generated according to a normal distribution of zero mean and a standard deviation of 0.1.

Several different GPU models are used in this work, including NVIDIA A100 80GB, NVIDIA A100 40GB, and

NVIDIA A40 48GB. Pre-training of CGCNN and TorchMD-Net on MP Forces takes approximately 10 and 60 hours respectively on a single A100 80GB GPU.

## **B** Errors on Main Results

Here we present tables detailing the standard errors associated with the main results. Each entry in the main results tables represents the average of 5 runs. To quantify the uncertainty, we calculate the error as the standard deviation divided by the square root of 5:

Standard Error =  $\frac{\text{Standard Deviation}}{\sqrt{\# \text{ of Measurements}}}$ 

		JDFT	Phonons	Dielectric	GVRH	KVRH	Perovskites	2D	MOF	Surface	MP gap	MP form
	Baseline	$\pm 4.37$	$\pm 4.18$	$\pm 0.0169$	$\pm 0.00124$	$\pm 0.00118$	$\pm 0.000294$	$\pm 0.00366$	$\pm 0.0104$	$\pm 0.00121$	$\pm 0.00351$	$\pm 0.000516$
Forces	1:0:0	$\pm 5.46$	$\pm 3.86$	$\pm 0.0535$	$\pm 0.00173$	$\pm 0.00058$	$\pm 0.000741$	$\pm 0.00326$	$\pm 0.0029$	$\pm 0.000717$	$\pm 0.00198$	$\pm 0.00216$
	0:1:0 0:1:1	$\pm 3.32 \\ \pm 5.82$	$\pm 4.40 \\ \pm 5.40$	$\pm 0.00634 \\ \pm 0.0277$	$\pm 0.00524 \\ \pm 0.00146$	$\pm 0.0178 \\ \pm 0.00180$	$\pm 0.000529$ $\pm 0.000709$	$\pm 0.00546 \\ \pm 0.00569$	$\pm 0.00229 \\ \pm 0.00339$	$\pm 0.000964$ $\pm 0.00137$	$\pm 0.00545 \\ \pm 0.00258$	$\pm 0.00296$ $\pm 0.00181$
	1:1:1 1:500:500	$\pm 5.14 \pm 3.49$	$\pm 2.21 \\ \pm 3.20$	$\pm 0.0129$ $\pm 0.0265$	$\pm 0.00246$ $\pm 0.000812$	$\pm 0.000986$ $\pm 0.00126$	$\pm 0.000511$ $\pm 0.000298$	$\pm 0.00481$ $\pm 0.00546$	$\pm 0.00101$ $\pm 0.00235$	$\pm 0.00110$ $\pm 0.000662$	$\pm 0.00262$ $\pm 0.00180$	$\pm 0.00143$ $\pm 0.00336$
Derivative-based denoising		$\pm 2.17$	±7.74	$\pm 0.0343$	$\pm 0.00207$	$\pm 0.00109$	$\pm 0.000422$	$\pm 0.00616$	$\pm 0.0513$	$\pm 0.00209$	$\pm 0.00365$	$\pm 0.000975$
Prediction head denoising		±3.87	$\pm 2.39$	$\pm 0.0360$	$\pm 0.00212$	$\pm 0.0132$	$\pm 0.000628$	$\pm 0.00585$	$\pm 0.00483$	$\pm 0.000874$	$\pm 0.00268$	$\pm 0.00350$

#### Table B.1: Standard errors associated with Table 2

Table B.2: Standard errors associated with Table 3

		JDFT	Phonons	Dielectric	GVRH	KVRH	Perovskites	2D	MOF	Surface	MP gap	MP form
	Baseline	$\pm 2.32$	$\pm 4.34$	$\pm 0.0356$	$\pm 0.00132$	$\pm 0.0021$	$\pm 0.0000842$	$\pm 0.00421$	$\pm 0.00138$	$\pm 0.00117$	$\pm 0.00181$	$\pm 0.000339$
Forces	1:0:0 0:1:0	$\pm 5.46$ $\pm 3.32$	$\pm 3.86$ $\pm 4.40$	$\pm 0.0535$ $\pm 0.00634$	$\pm 0.00173$ $\pm 0.00524$	$\pm 0.00058$ $\pm 0.0178$	$\pm 0.000741$ $\pm 0.000529$	$\pm 0.00326$ $\pm 0.00546$	$\pm 0.0029$ $\pm 0.00229$	$\pm 0.000717$ $\pm 0.000964$	$\pm 0.00198$ $\pm 0.00545$	$\pm 0.00216$ $\pm 0.00296$
	0:1:1 1:1:1 1:500:500	$\pm 5.82 \\ \pm 5.14 \\ \pm 3.49$	$\pm 5.40 \\ \pm 2.21 \\ \pm 3.20$	$\pm 0.0277$ $\pm 0.0129$ $\pm 0.0265$	$\pm 0.00146$ $\pm 0.00246$ $\pm 0.000812$	$\pm 0.00180$ $\pm 0.000986$ $\pm 0.00126$	$\pm 0.000709$ $\pm 0.000511$ $\pm 0.000298$	$\pm 0.00569$ $\pm 0.00481$ $\pm 0.00546$	$\pm 0.00339$ $\pm 0.00101$ $\pm 0.00235$	$\pm 0.00137$ $\pm 0.00110$ $\pm 0.000662$	$\pm 0.00258$ $\pm 0.00262$ $\pm 0.00180$	$\pm 0.00181$ $\pm 0.00143$ $\pm 0.00336$
Derivative-based denoising		$\pm 3.25$	$\pm 7.66$	$\pm 0.0333$	$\pm 0.00122$	$\pm 0.000833$	$\pm 0.000259$	$\pm 0.00831$	$\pm 0.00223$	$\pm 0.00086$	$\pm 0.00207$	$\pm 0.000333$
Prediction head denoising		$\pm 3.66$	±12.7	$\pm 0.0206$	$\pm 0.00198$	$\pm 0.00052$	$\pm 0.000558$	$\pm 0.00761$	$\pm 0.0039$	$\pm 0.000587$	$\pm 0.00109$	$\pm 0.000261$

## C Parity Plots of Forces For Pre-trained Models



Figure C.1: Parity plots of forces prediction versus ground-truth labels. Axes are trimmed to [-100, +100]. a) CGCNN pre-trained models. b) TorchMD-Net pre-trained models. c) CGCNN pre-trained for different number of epochs on ratio  $\lambda_{\text{energy}} : \lambda_{\text{forces}} : \lambda_{\text{stress}} = 0:1:0$ . d) CGCNN pre-trained for different number of epochs on ratio  $\lambda_{\text{energy}} : \lambda_{\text{forces}} : \lambda_{\text{stress}} = 1:500:500$ .

### **D** Embedding Visualization

To understand the nature of the benefit to finetuning performance from pre-training with energies and forces, we visualized the node-level embeddings for the CGCNN and TorchMD-Net models trained under the different ratios. Specifically, principal component analysis (PCA) is applied to 50,000 uniformly selected atom embeddings from the training split in the MP Forces dataset for a visualization in the two-dimensional space. It should be noted that both the derivative-based denoising and prediction head denoising models are trained on the smaller MP Forces Relaxed dataset. Consequently, it is possible that these two denoising models have not been exposed to some of the 50,000 atoms in the dataset.

We observe that the embeddings for CGCNN prediction head denoising is dissimilar from the rest. Our hypothesis is that this divergence stems from the fact that this particular model hasn't effectively learnt during the pre-training phase, as supported by its training loss barely decreasing and plateauing very early on.



Figure D.1: Visualization of atom embedding space wherein atoms are categorized into S, P, D and F blocks on the periodic table. Principal component analysis (PCA) is applied to all atomic representations in the training split of the MP Forces dataset, and 50,000 uniformly selected atoms are chosen to be plotted. a) CGCNN pre-trained models. b) TorchMD-Net pre-trained models.