
A Large Encoder-Decoder Family of Foundation Models For Chemical Language

Eduardo Soares
IBM Research Brazil
Rio de Janeiro, RJ, Brazil
eduardo.soares@ibm.com

Victor Shirasuna
IBM Research Brazil
São Paulo, SP, Brazil
victor.shirasuna@ibm.com

Emilio Vital Brazil
IBM Research Brazil
Rio de Janeiro, RJ, Brazil
evital@br.ibm.com

Renato Cerqueira
IBM Research Brazil
Rio de Janeiro, RJ, Brazil
rcerq@br.ibm.com

Dmitry Zubarev
IBM Research Almaden
San Jose, CA, USA
dmitry.zubarev@ibm.com

Kristin Schmidt
IBM Research Almaden
San Jose, CA, USA
schmidkr@us.ibm.com

Abstract

Large-scale pre-training methodologies for chemical language models represent a breakthrough in cheminformatics. These methods excel in tasks such as property prediction and molecule generation by learning contextualized representations of input tokens through self-supervised learning on large unlabeled corpora. Typically, this involves pre-training on unlabeled data followed by fine-tuning on specific tasks, reducing dependence on annotated datasets and broadening chemical language representation understanding. This paper introduces a large encoder-decoder chemical foundation models pre-trained on a curated dataset of 91 million SMILES samples sourced from PubChem, which is equivalent to 4 billion of molecular tokens. The proposed foundation model supports different complex tasks, including quantum property prediction, and offer flexibility with two main variants (289M and $8 \times 289M$). Our experiments across multiple benchmark datasets validate the capacity of the proposed model in providing state-of-the-art results for different tasks. We also provide a preliminary assessment of the compositionality of the embedding space as a prerequisite for the reasoning tasks. We demonstrate that the produced latent space is separable compared to the state-of-the-art with few-shot learning capabilities.

1 Introduction

Understanding molecular properties is crucial for accelerating discoveries in different fields, including drug development and materials science [1]. Traditional methods rely on labor-intensive trial-and-error experiments, which are both costly and time-consuming [2]. However, recent advances in deep learning have enabled the use of foundation models to predict molecular properties and generate molecule candidates [3, 4, 5], marking significant progress in scientific exploration.

The introduction of large-scale pre-training methodologies for chemical language models (LMs) represents a significant advancement in cheminformatics [6]. These methodologies have demonstrated impressive results in challenging molecular tasks such as predicting properties and generating molecules [7]. The success of these models can be attributed to their ability to learn contextualized representations of input tokens through self-supervised learning on large unlabeled corpora [8]. This methodological approach typically involves two phases: pre-training on unlabeled data followed by fine-tuning on specific downstream task [9]. By reducing the reliance on annotated datasets, this approach has broadened our understanding of chemical language representations [10].

Simplified Molecular-Input Line Entry System, SMILES, provide natural graphs that encode the connectivity information from the line annotations of molecular structures [11]. SMILES defines a character string representation of a molecule by performing a depth-first pre-order spanning tree traversal of the molecular graph, generating symbols for each atom, bond, tree-traversal decision, and broken cycles [12]. Therefore, the resulting character string corresponds to a flattening of a spanning tree of the molecular graph. SMILES is widely adopted for molecular property prediction as SMILES is generally more compact than other methods of representing structure, including graphs [13]. There are billions of SMILES available on different open-sources repositories [14]. However, most SMILES sequences do not belong to well-defined molecules [15]. Alternative string-based representations exist, such as SELFIES. However, focusing on molecular optimization tasks on the learned representation space, suggested no obvious shortcoming of SMILES with respect to SELFIES in terms of optimization ability and sample efficiency [16]. The quality of the pre-training data plays a more important role on the outcome of the foundation model [4, 17].

Towards this direction, we present a novel family of molecular encoder-decoder foundation models, denoted as SMI-TED289M. Our SMI-TED289M encoder-decoder foundation model was obtained using a transformer-based molecular tokens encoder model aligned with an encoder-decoder mechanism trained on a large corpus of 91 million carefully curated molecules from PubChem [18], resulting in 4 billion molecular tokens. Our main contributions are:

- We pre-train a large-scale family of encoder-decoder molecular open-source foundation models, denoted as SMI-TED289M, on over 91 million molecules carefully curated from PubChem [18], which is equivalent to 4 billion of molecular tokens.
- A molecular dataset for pre-training of chemical foundation models, 91 million molecules carefully curated from PubChem [18].
- Our SMI-TED289M family of foundation models encompasses two distinct configurations: base, which has 289 million parameters; and the Mixture-of-SMI-TED-Experts, SMI-TED8x289M, characterized by a composition of $8 \times 289M$ parameters. The source code is available at: <https://github.com/IBM/materials>.
- We perform extensive experimentation on several classification and regression tasks from 11 benchmark datasets, covering quantum mechanical, physical, biophysical, and physiological property prediction of small molecules. We also evaluate the reconstruction capacity of our SMI-TED289M considering the MOSES benchmarking dataset [19]. Furthermore, a study investigating the embedding created by SMI-TED289M and few-shot learning is also provided, indicating compositionality of the learned molecular representations.

Our results section demonstrates state-of-the-art performance of SMI-TED289M on different tasks, molecular properties prediction, molecule reconstruction, and an efficient metric for molecular latent space. Compositionality of the latent space suggests strong potential for chemical reasoning tasks. The SMI-TED289M family consists of two main variants (289M, and $8 \times 289M$), offering flexibility and scalability for different scientific applications.

2 Overview of the proposed approach

This section presents an overview of the proposed SMI-TED289M foundation model for small molecules. Here, we outline the process of collecting, curating, and pre-processing the pre-train data. Additionally, we describe the token encoder process and the SMILES encoder-decoder process. Finally, we explain the Mixture-of-SMI-TED-Experts approach used to scale the base model. Fig. 1 illustrates the general architecture of the base model.

2.1 Pre-training Data

The pretraining data originated from the PubChem data repository, a public database containing information on chemical substances and their biological activities [18]. Initially, 113 million SMILES strings were collected from PubChem. These molecular strings underwent deduplication and canonicalization processes to ensure uniqueness [20]. Subsequently, a molecular transformation was conducted to verify the validity of the molecules derived from the unique SMILES strings, resulting in a set of 91 million unique and valid molecules.

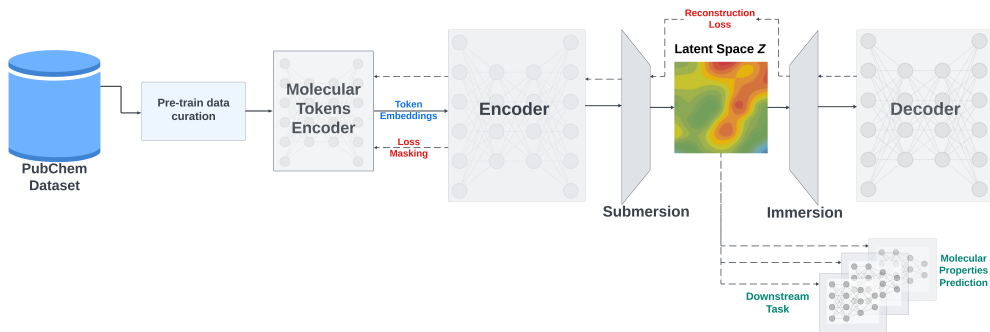


Figure 1: This figure illustrates the general architecture of the base SMI-TED289M model.

To construct the vocabulary, we employed the molecular tokenizer proposed by [21]. All 91 million molecules curated from PubChem were utilized in the tokenization process, resulting in a set of 4 billion molecular tokens. The unique tokens extracted from the resulting output provided a vocabulary of 2988 tokens plus 5 special tokens. In comparison, MoLFormer, trained on 1 billion samples with minimal curation, presented a vocabulary of 2362 tokens using the same tokenization process [7]. This suggests an improvement in the vocabulary model due to our curation process.

2.2 Model Architecture

We conduct training for SMI-TED289M model employing a deep-bidirectional-transformers-based encoder [22] for tokens and an encoder-decoder architecture to compose SMILES. The hyperparameters of SMI-TED289M base model are detailed in Table 1

Table 1: SMI-TED289M base architecture specificity.

Hidden size	Attention heads	Layers	Dropout	Normalization
768	12	12	0.2	LayerNorm

Vocab size	# SMILES	# Mol tokens	# Encoder	# Decoder	Total params
2993	91M	4T	47M	242M	289M

To optimize the relative encoding through position-dependent rotations R_m of the query and keys at position m , the SMI-TED289M uses a modified version of the RoFormer [23] attention mechanism. These rotations can be implemented as pointwise multiplications and do not significantly increase computational complexity as shown in Eq. (1).

$$Attention_m(Q, K, V) = \frac{\sum_{n=1}^N \langle \varphi(R_m q_m), \varphi(R_n k_n) \rangle v_n}{\sum_{n=1}^N \langle \varphi(R_m q_m), \varphi(R_n k_n) \rangle} \quad (1)$$

where Q, K, V are the query, key, and value respectively, and φ is a random feature map.

We start with a sequence of tokens extracted from SMILES, each embedded in a 768-dimensional space. The encoder-decoder layer is designed to process molecular token embeddings, represented as $\mathbf{x} \in \mathbb{R}^{D \times L}$, where D denotes the maximum number of tokens and L represents the embedding space dimension. We limited D at 202 tokens, as 99.4% of molecules in the PubChem dataset contain fewer tokens than this threshold.

In encoder-only models, a mean pooling layer is typically employed to represent tokens as SMILES in the latent space. However, this approach is limited by the lack of a natural inversion process for the mean pooling operation. To overcome this limitation, we aim to construct a latent space representation for SMILES by submersing the \mathbf{x} in a latent space, denoted as \mathbf{z} , as described in Eq. 2.

$$\mathbf{z} = (\text{LayerNorm}(\text{GELU}(\mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)))\mathbf{W}_2, \quad (2)$$

where $\mathbf{z} \in \mathbb{R}^L$, $\mathbf{W}_1 \in \mathbb{R}^{D \times L}$, $\mathbf{b}_1 \in \mathbb{R}^L$, $\mathbf{W}_2 \in \mathbb{R}^{L \times L}$, with L denoting the latent space size (specifically, $L = 768$) and D representing the original feature space size (namely, $D = 202$). Subsequently, we can immerse \mathbf{z} back by calculating Eq. 3.

$$\hat{\mathbf{x}} = (\text{LayerNorm}(\text{GELU}(\mathbf{z}\mathbf{W}_3 + \mathbf{b}_3)))\mathbf{W}_4 \quad (3)$$

where $\hat{\mathbf{x}} \in \mathbb{R}^{D \times L}$, $\mathbf{W}_3 \in \mathbb{R}^{L \times L}$, $\mathbf{b}_3 \in \mathbb{R}^L$, $\mathbf{W}_4 \in \mathbb{R}^{L \times D}$.

A language layer (decoder) is used to process $\hat{\mathbf{x}}$, where it applies non-linearity and normalization, and projects the resulting vector into a set of logits over the vocabulary, which can then be used to predict the next token in the molecular [24].

2.3 Pre-training strategies

Pre-training of SMI-TED289M was performed for 40 epochs through the entire curated PubChem dataset with a fixed learning rate of 1.6e-4 and a batch size of 288 molecules on a total of 24 NVIDIA V100 (16G) GPUs parallelized into 4 nodes using DDP and *torch run*. It involves two distinct phases: i) Learning of token embeddings through a masking process; ii) Subsequently, the token embeddings are mapped into a common latent space that encapsulates the entire SMILES string. This latent space not only facilitates the representation of the SMILES but also enables the reconstruction of both individual tokens and complete SMILES strings. Consequently, the pre-training process involves two separate loss functions: one for the token embeddings, which is based on the masking process, and another for the encoder-decoder layer, which focuses on the reconstruction of tokens. Two pre-training strategies are employed:

- In phase 1, the token encoder is initially pre-trained using 95% of the available samples, while the remaining 5% is reserved for training the encoder-decoder layer. This partitioning is necessary as the token embeddings may encounter convergence difficulties in the initial epochs, which could adversely affect the training of the encoder-decoder layer.
- In phase 2, once the token embeddings layer has achieved convergence, the pre-training process is expanded to utilize 100% of the available samples for both phases. This approach leads to an enhancement in the performance of the encoder-decoder layer, particularly in terms of token reconstruction.

For encoder pre-training we use the masked language model method defined in [22]. Initially 15% of the tokens are selected for possible learning. From that selection, 80% of the tokens are randomly selected and replaced with the [MASK] token, 10% of the tokens are randomly selected to be replaced with a random token, while the remaining 10% of the tokens will be unchanged.

The adoption of different pre-training strategies has proven instrumental in enhancing the efficiency of our model, as evidenced by improvements observed in the loss functions. For detailed insights into the loss functions and pre-training methodologies, refer to the Supplementary Materials.

2.4 Mixture-of-SMI-TED-Experts

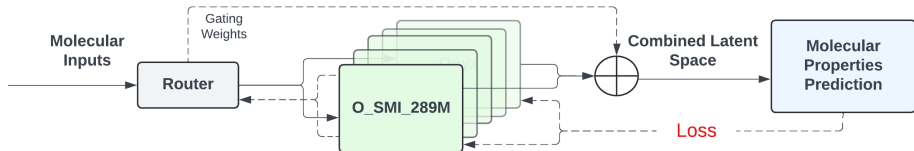


Figure 2: Mixture-of-SMI-TED-Experts for downstream tasks.

The Mixture-of-SMI-TED-Experts, SMI-TED8x289M comprises a set of n “expert networks” labeled as E_1, E_2, \dots, E_n , augmented through a gating network denoted as G , tasked with generating a sparse n -dimensional embedding space optimized for a downstream task as illustrated by Fig. 2.

Here, we map each SMILES into tokens and then convert the input tokens to the latent space. A mean pooling method is applied to all token embeddings in order to produce a meaningful embedding of the molecule. The architecture is equipped with a router module responsible for determining the n

experts that will be activated, refining the adaptability and specialization of the system. Let $G(x)$ and $E_i(\hat{x})$ denote the output of the gating network and the output of the i -th expert network, respectively, for a given input \hat{x} of SMILES and x , which is the embeddings space, following a similar notation as proposed in [25]. The resulting output y is defined as follows:

$$y = \sum_{i=1}^n G(x)_i E_i(\hat{x})$$

The resulting embedding space y is used to train a task-specific feed-forward network, where the loss function is chosen according to the studied downstream task. The optimization process refines the parameters of $G(x)$. If the gating vector is sparse, we can use softmax over the Top-K logits of a linear layer [25].

$$G(x) := \text{Softmax}(\text{TopK}(x \cdot Wg))$$

where $(\text{TopK}(\ell))_i := \ell_i$ if ℓ_i is among the TopK coordinates of logits $\ell \in \mathbb{R}^n$ and $(\text{TopK}(\ell))_i := \infty$ otherwise. The router layer retains only the top k values, setting the remaining values to $-\infty$ (which effectively assigns corresponding gate values as 0). This sparsity-inducing step serves to optimize computational efficiency [26]. Here, we define SMI-TED8x289M as $n = 8$ and $k = 2$, which means that SMI-TED8x289M is composed by $8 \times$ SMI-TED289M models, which 2 models are activated through the router each round.

3 Experiments

To evaluate the effectiveness of our proposed methodology, we conducted experiments using a set of 11 datasets sourced from MoleculeNet [27] as demonstrated in Table 2. Specifically, we evaluated 6 datasets for classification task and 5 datasets for regression tasks. To ensure an unbiased assessment, we maintained consistency with the original benchmark by adopting identical train/validation/test splits for all tasks [27]. We also conducted the experiments considered 10 different seeds for all the tests in other to guarantee the robustness of the approach. Details are provided in the Supplementary Materials.

Table 2: Evaluated datasets description

Dataset	Description	# compounds	# tasks	Metric
BBBP	Blood brain barrier penetration dataset	2039	1	ROC-AUC
HIV	Ability of small molecules to inhibit HIV replication	41127	1	ROC-AUC
BACE	Binding results for a set of inhibitors for β - secretase 1	1513	1	ROC-AUC
Clintox	Clinical trial toxicity of drugs	1478	2	ROC-AUC
SIDER	Drug side effect on different organ classes	1427	27	ROC-AUC
Tox21	Toxicity measurements on 12 different targets	7831	12	ROC-AUC
QM9	12 quantum mechanical calculations	133885	12	Average MAE
QM8	12 excited state properties of small molecules	21786	12	Average MAE
ESOL	Water solubility dataset	1128	1	RMSE
FreeSolv	Hydration free energy of small molecules in water	642	1	RMSE
Lipophilicity	Octanol/water distribution coefficient of molecules	4200	1	RMSE

To assess the reconstruction/decoder capacity of SMI-TED289M we considered the MOSES benchmarking dataset [19]. The MOSES dataset contains 1,936,962 molecular structures. For experiments, we consider the split proposed by [19], where the dataset was divided into a training, test and scaffold test sets containing around 1.6M, 176k, and 176k molecules respectively. The scaffold test set contains unique Bemis-Murcko scaffolds that were not present in the training and test sets. We use this set to assess how well the model can generate previously unobserved scaffolds. An evaluation of the embedding space of SMI-TED289M is also provided, it uses the compositional molecules to evaluate the capability of the model to generate metric latent spaces.

4 Results and Discussion

In this section, we present the analysis of results obtained using SMI-TED289M for different experiments conducted with various versions of the base model. We include: i) A study comparing

frozen and fine-tuned versions of SMI-TED289M; and a comparison with the State-of-the-Art (SOTA) on different benchmarking datasets for classification and regression molecular prediction tasks; ii) An evaluation of SMI-TED8x289M for molecular properties prediction; iii) An evaluation of the Decoder module considering the MOSES benchmarking dataset; iv) A study comparing the latent space of SMI-TED289M based on compositional molecules metrics.

4.1 Comparison with SOTA on benchmarking tasks

Results for classification tasks: The analysis investigates the comparative efficacy of SMI-TED289M in its fine-tuned and frozen states versus state-of-the-art algorithms for molecular properties classification, as demonstrated in Table 3.

Table 3: Methods and Performance for the classification tasks of MoleculeNet benchmark datasets

Method	Dataset					
	BBBP	ClinTox	HIV	BACE	SIDER	Tox21
GraphMVP [28]	72.4 ± 1.6	79.1 ± 2.8	77.0 ± 1.2	81.2 ± 0.9	63.9 ± 1.2	75.9 ± 0.5
GEM [29]	72.4 ± 0.4	90.1 ± 1.3	80.6 ± 0.9	85.6 ± 1.1	67.2 ± 0.4	78.1 ± 0.1
GROVER _{Large} [30]	69.5 ± 0.1	76.2 ± 3.7	68.2 ± 1.1	81.0 ± 1.4	65.4 ± 0.1	73.5 ± 0.1
ChemBerta [31]	64.3	90.6	62.2	-	-	-
ChemBerta2 [32]	71.94	90.7	-	85.1	-	-
Galatica 30B [33]	59.6	82.2	75.9	72.7	61.3	68.5
Galatica 120B [33]	66.1	82.6	74.5	61.7	63.2	68.9
Uni-Mol [34]	72.9 ± 0.6	91.9 ± 1.8	80.8 ± 0.3	85.7 ± 0.2	65.9 ± 1.3	79.6 ± 0.5
MolFM [34]	72.9 ± 0.1	79.7 ± 1.6	78.8 ± 1.1	83.9 ± 1.1	64.2 ± 0.9	77.2 ± 0.7
MolFormer [35]	73.6 ± 0.8	91.2 ± 1.4	80.5 ± 1.65	86.3 ± 0.6	65.5 ± 0.2	80.46 ± 0.2
SMI-TED289M (Frozen Weights)	91.46 ± 0.47	93.49 ± 0.85	80.51 ± 1.34	85.58 ± 0.92	66.01 ± 0.88	81.53 ± 0.45
SMI-TED289M (Fine-tuned)	92.26 ± 0.57	94.27 ± 1.83	76.85 ± 0.89	88.24 ± 0.50	65.68 ± 0.45	81.85 ± 1.42

Table 3 displays the performance of different advanced methods on different benchmarking datasets used for molecule classification tasks. SMI-TED289M consistently shows superior performance in four out of six datasets. Interestingly, using SMI-TED289M with its initial settings provided comparable results to SOTA methods available. However, fine-tuning SMI-TED289M further enhances its performance across all datasets. This indicates SMI-TED289M potential for accurate molecule classification, with potential for further optimization through fine-tuning. Detailed results for all the experiments are presented in the Supplementary Materials due to limit of pages.

Results for regression tasks: Next, we applied SMI-TED289M for prediction of chemical properties. The performance results across five challenging regression benchmarks, namely QM9, QM8, ESOL, FreeSolv, and Lipophilicity, are summarized in Table 4.

Table 4: Methods and Performance for the regression tasks of MoleculeNet benchmark datasets.

Method	Dataset				
	QM9	QM8	ESOL	FreeSolv	Lipophilicity
D-MPNN [36]	3.241 ± 0.119	0.0143 ± 0.0022	0.98 ± 0.26	2.18 ± 0.91	0.65 ± 0.05
N-Gram [37]	2.51 ± 0.19	0.0320 ± 0.003	1.074 ± 0.107	2.688 ± 0.085	0.812 ± 0.028
PretrainGNN [38]	-	-	1.100 ± 0.006	2.764 ± 0.002	0.739 ± 0.003
GROVER _{Large} [30]	-	-	0.895 ± 0.017	2.272 ± 0.051	0.823 ± 0.010
ChemBERTa-2 [32]	-	-	0.89	-	0.80
SPMM [35]	-	-	0.818 ± 0.008	1.907 ± 0.058	0.692 ± 0.008
MolCLR _{GIN} [39]	2.357 ± 0.118	0.0174 ± 0.0013	1.11 ± 0.01	2.20 ± 0.20	0.65 ± 0.08
Hu et al. [40]	4.349 ± 0.061	0.0191 ± 0.0003	1.22 ± 0.02	2.83 ± 0.12	0.74 ± 0.00
MolFormer [35]	1.5894 ± 0.0567	0.0102	0.880 ± 0.028	2.342 ± 0.052	0.700 ± 0.012
SMI-TED289M (Frozen Weights)	7.4883 ± 0.0659	0.0179 ± 0.0004	0.7045 ± 0.0344	1.668 ± 0.0616	0.6499 ± 0.012
SMI-TED289M (Fine-tuned)	1.3246 ± 0.0157	0.0095 ± 0.0001	0.6112 ± 0.0096	1.2233 ± 0.0029	0.5522 ± 0.0194

Results presented in Table 4 indicates that SMI-TED289M presents superior results when compared to the state-of-the-art, outperforming its competitors in all the 5 datasets considered. To fine-tune SMI-TED289M is important to achieve state-of-the-art results in regression datasets, due to the complexity of such tasks. Table 4 elucidates the superiority of SMI-TED289M over the QM9 dataset. The QM9 dataset is composed by 12 tasks regarding to the quantum properties of molecules. A detailed overview over the results for QM9 are depicted in the next subsection. Detailed results for all experiments are in the Supplementary Materials of this paper.

A deeper analysis over the QM9 benchmark: In this subsection, we provide a deeper analysis over the results for the QM9 dataset. Table 5 details the results of the SOTA approaches each property

that composes QM9. Our comparative analysis extends to benchmarking the proposed encoder-decoder foundation model against state-of-the-art models derived from three distinct categories: (i) Graph-based, (ii) Geometry-based, and (iii) SMILES-based methodologies for prediction of molecular properties. The included baselines models are: 123-gnn [41], a multitask neural net encoding the Coulomb Matrix (CM) [42], and its GNN variant as in the deep tensor neural net (DTNN) [43].

Table 5: Comparing state-of-the-art models performance over the QM9 dataset. **Blue** and **Orange** indicates best and second-best performing model, respectively.

Measure	Graph-based			Geometry-based			SMILES-based	
	A-FP	123-gnn	GC	CM	DTNN	MPNN	MoLFormer-XL	This paper
α	0.49	0.27	1.37	0.85	0.95	0.89	0.33	0.27
C_v	0.25	0.09	0.65	0.39	0.27	0.42	0.14	0.12
G	0.89	0.05	3.41	2.27	2.43	2.02	0.34	0.11
gap	0.0052	0.0048	0.01126	0.0086	0.0112	0.0066	0.0038	0.0036
H	0.89	0.04	3.41	2.27	2.43	2.02	0.25	0.09
ϵ_{homo}	0.0036	0.0034	0.0072	0.0051	0.0038	0.0054	0.0029	0.0027
ϵ_{lumo}	0.0041	0.0035	0.0092	0.0064	0.0051	0.0062	0.0027	0.0026
μ	0.451	0.476	0.583	0.519	0.244	0.358	0.361	0.384
$\langle R^2 \rangle$	26.84	22.90	35.97	46.00	17.00	28.5	17.06	14.72
U_0	0.898	0.0427	3.41	2.27	2.43	2.05	0.3211	0.0850
U	0.89	0.111	3.41	2.27	2.43	2.00	0.25	0.0905
ZPVE	0.00207	0.0002	0.00299	0.00207	0.0017	0.00216	0.0003	0.0002
Avg MAE	2.6355	1.9995	4.3536	4.7384	2.3504	3.1898	1.5894	1.3246
Avg std MAE	0.0854	0.0658	0.1683	0.1281	0.1008	0.1108	0.0567	0.0157

Table 5 compares existing SOTA models in predicting quantum properties of molecules. The evaluation demonstrates that the proposed encoder-decoder foundation model outperforms current models in predicting 7 out of 12 quantum properties, and achieves either the best or second-best results in 11 out of 12 tasks.

However, when comparing with MoLFormer-XL, a model showing the second-best average error rate, it is noted that MoLFormer-XL’s performance is influenced by its results on a specific property $\langle R^2 \rangle$. Although MoLFormer-XL performs well in average error rate, 123-gnn performs better in a larger number of tasks. In comparison, the proposed SMI-TED289M maintains consistent performance across all tasks, suggesting its robustness in predicting complex molecular properties.

4.2 Mixture-of-SMI-TED-Experts perform studies

This study compare the results of MoE-SMI-TED against a single SMI-TED289M models (frozen and fine-tuned). SMI-TED8x289M is composed by $8 \times 289M$ fine-tuned models for each specific task, we set $k = 2$, which means that 2 models are activated every step. The results for this study are shown in Table 6, which considers classification and regression tasks for molecular properties. Results refers to the best run of each version.

Table 6: SMI-TED8x289M and single SMI-TED289M models for molecular properties prediction.

Method	Dataset								
	BBBP \uparrow	ClinTox \uparrow	HIV \uparrow	BACE \uparrow	SIDER \uparrow	Tox21 \uparrow	ESOL \downarrow	FreeSolv \downarrow	Lipo \downarrow
SMI-TED289M - Frozen	92.27	95.02	81.81	87.18	67.11	82.22	0.6784	1.5832	0.6311
SMI-TED289M - Fine-Tuned	93.07	97.97	79.09	89.33	65.97	83.72	0.6024	1.2167	0.5413
SMI-TED8x289M	93.72	95.62	80.42	89.84	68.08	84.07	0.5566	1.1181	0.5376

Table 6 summarizes the performance metrics for each model across the different datasets. The results from the study indicate that SMI-TED8x289M consistently achieves higher performance metrics compared to single SMI-TED289M models (Frozen and Fine-Tuned) models across different tasks, especially in regression tasks where it improved results in all scenarios. These findings suggest that the MoE approach effectively leverages specialized sub-models to capture diverse patterns in the data, leading to improved accuracy in molecular property predictions. The mixture-of-experts approach serves as an efficient solution to scale single models and enhance performance for various tasks due to its ability to allocate specific tasks to different experts, optimizing single model’s overall predictive capabilities.

4.3 Decoder evaluation over MOSES benchmarking dataset

Next, we compared SMI-TED289M with different baseline models, such as the character-level recurrent neural network (CharRNN) [19], SMILES variational autoencoder (VAE) [19], junction tree VAE (JT-VAE) [44], latent inceptionism on molecules (LIMO) [45], MolGen-7b [46], and GP-MoLFormer [47]. All baseline performances are reported on their corresponding test set consisting of 176k molecules. Standard metrics for evaluating model-generated molecules are reported in Table 7. All metrics are computed using MOSES.

Table 7: MOSES benchmarking dataset evaluation.

Metric	Frag \uparrow	Scaf \uparrow	SNN \uparrow	IntDiv \uparrow	FCD \downarrow
CharRNN	0.9998	0.9242	0.6015	0.8562	0.0732
VAE	0.9984	0.9386	0.6257	0.8558	0.0990
JT-VAE	0.9965	0.8964	0.5477	0.8551	0.3954
LIMO	0.6989	0.0079	0.2464	0.9039	26.78
MolGen-7b	0.9999	0.6538	0.5138	0.8617	0.0435
GP-MoLFormer	0.9998	0.7383	0.5045	0.8655	0.0591
SMI-TED289M	0.9999	0.9999	0.9998	0.8565	1.1532

When compared to baselines, SMI-TED289M is equally performant in generating unique, valid, and novel molecules that share high cosine similarity with the corresponding reference molecules at the fragment (Frag) level, consistent with low Fréchet ChemNet Distance (FCD). At the same time, SMI-TED289M generates molecules with high internal diversity (IntDiv), i.e., average pairwise dissimilarity. The scaffold cosine similarity (Scaf) and similarity to the nearest neighbor in the test set (SNN) of SMI-TED289M is superior to the baselines demonstrating that SMI-TED289M is effective in generating molecules of varying structures and quality compared to baseline methods.

4.4 Latent space study

We conducted an experiment to investigate the structure of the latent space created by Large Language Models in the context of Chemistry. Molecular structures are composable from fragments, motifs, and functional groups. The composability of structure often translates into compositionality of structure-property relations, which is exemplified by powerful group contribution methods in chemical sciences. Compositionality of the learnt representation, however, does not follow automatically from the structure of the data and requires some combination of the learning architecture and learning constraints to emerge. Our approach was to utilize simple chemical structures that can be easily understood by humans, allowing us to anticipate relationships between elements, and examine the latent space for similar patterns. We constructed a dataset consisting of six families of carbon chains: $\mathcal{F} = \{CC, CO, CN, CS, CF, CP\}$. For each family, we generated a sequence of molecules by incrementally adding carbon atoms to the end of the SMILES string, up to a maximum of ten carbon atoms. For example, the family CO consists of $\{CO, CCO, \dots, CCCCCCCCCCO\}$. According to the domain expert’s intuition consistent with the theory of chemical structure, in a metric space, such sequences should exhibit a hierarchical distance structure, where the distance between consecutive elements is smaller than the distance between elements with a larger difference in carbon count, i.e., $|\overline{C_n \mathcal{F}_i} - \overline{C_{n+1} \mathcal{F}_i}| < |\overline{C_n \mathcal{F}_i} - \overline{C_{n+2} \mathcal{F}_i}|$. Here, n represents the number of carbon atoms, and $\overline{\text{SMILE}}$ denotes the projection of the SMILE string onto the embedding space.

First, we generated the embeddings for two different encoders, the MoLFormer and SMI-TED289M, and used the t-SNE [48] projection technique to generate pictures (Fig. 3) for visually inspecting the spaces. It is worth noting that the SMI-TED289M generated an embedding space that creates a nice separation of each family and respects the hierarchical distance structure, almost creating a linear relationship between each family. To quantify this relationship, we created a dataset of triples of SMILES, $\mathcal{T} = \{(C_n \mathcal{F}_{CC}, C_k \mathcal{F}_i, C_{n+k} \mathcal{F}_i) \mid 0 < n \leq 4, 0 < k \leq 5\}$, for the six families \mathcal{F}_i , resulting in six sub-datasets with 20 elements each, e.g., $(CC, CCO, CCCCCO)$ is one element of the subset of type CO where $n = 1, k = 2$. Then, we randomly selected one triple from each subset to feed a linear regression calculating α, β , and B_0 such that $\alpha \cdot \overline{C_n \mathcal{F}_{CC}} + \beta \cdot \overline{C_k \mathcal{F}_i} + B_0 = \overline{C_{n+k} \mathcal{F}_i}$. We validated the linearity using the remaining 114 elements. The linear regression on the MoLFormer embeddings resulted in $R^2 = 0.55$ and $MSE = 0.237$, while on our model embeddings, it resulted in $R^2 = 0.99$ and $MSE = 0.002$.

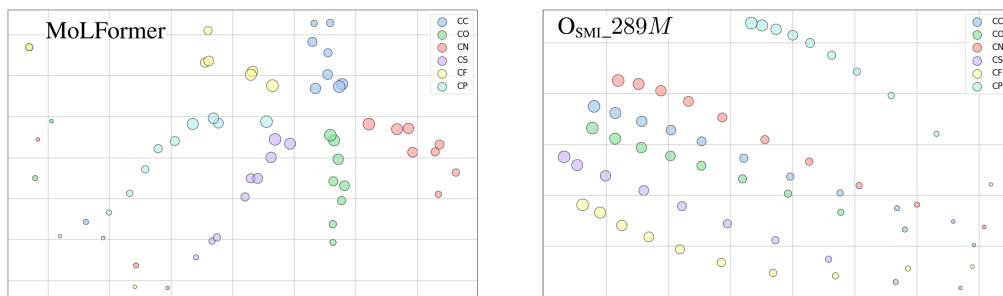


Figure 3: The figure shows the t-SNE projection of 60 small molecule embeddings. Color distinguishes between families, and point size represents the number of carbon atoms in the chain. Left: MoLFormer embeddings; Right: SMI-TED289M embeddings.

We evaluated our encoder-decoder model using a few-shot learning process, where we input a few examples of triples, such as those mentioned earlier, to calculate α , β , and B_0 . We then use these parameters to generate embeddings for subsequent SMILES pairs and recreate the SMILES strings. To validate our approach, we tested the process on the same dataset of triples. We calculated the molecule similarity between the expected and generated results using the Tanimoto score (TS) [49]. We repeated this test with different combinations of input triples, yielding similar results. For example, when using the input triples $[CC+CCCS = CCCCCS, CCCCC+CCCS = CCCCCCS]$ and querying all pairs in our subsets, we obtained a mean TS of 0.52. The top two similar results were $CC + CCCCCS = CCCCCS$ with TS = 0.92 and $CC + CCCCCO = CCCCCO$ with TS = 0.92, while the bottom two results were $CCCCC + CF = F[PH3+]F$ with TS = 0.06 and $CCCC + CF = F[PH3+]F$ with TS = 0.07.

Historically, group contribution was introduced in supervised learning context of structure-property relations. Our simple tests indicate that SMI-TED289M derived an equivalent of group contribution method purely from self-supervised learning of molecular structure. Signs of the emergence of compositionality of the learned molecular representations suggest strong potential of SMI-TED289M for reasoning applications. Further studies consistent with methodologies of compositionality analysis in natural languages are required to make stronger statements.

5 Conclusion

This paper introduces the SMI-TED289M family of chemical foundation models, which are pre-trained on a curated dataset of 91 million SMILES samples from PubChem, amounting to 4 billion molecular tokens. The SMI-TED289M family includes two configurations: the base model with 289 million parameters and the MoE SMI-TED8x289M model, which consists of $8 \times 289M$ parameters.

The performance of these models was evaluated through an extensive experimentation on different tasks, including molecular properties classification and prediction. Our approach achieved state-of-the-art results in most tasks, particularly in predicting molecular quantum mechanics, where it achieved the best or second-best results in 11 out of 12 tasks of the QM9 dataset.

We also investigated the structure of the latent space created by these language-based foundation models, using simple chemical structures for clarity. SMI-TED289M generated an embedding space that creates a nice separation of each family and respects the hierarchical distance structure, almost creating a linear relationship between each family. The encoder-decoder model’s capabilities in few-shot learning were assessed by generating embeddings from a few example triples and using them to recreate SMILES strings, achieving a Tanimoto score of 0.92 in the best case.

The family of chemical foundation models presented in this paper offers flexibility and scalability for different scientific applications. The source code is available at: <https://github.com/IBM/materials>.

References

- [1] J. Pan, “Large language model for molecular chemistry,” *Nature Computational Science*, vol. 3, no. 1, pp. 5–5, 2023.

- [2] K. M. Jablonka, P. Schwaller, A. Ortega-Guerrero, and B. Smit, "Leveraging large language models for predictive chemistry," *Nature Machine Intelligence*, pp. 1–9, 2024.
- [3] D. Flam-Shepherd, K. Zhu, and A. Aspuru-Guzik, "Language models can learn complex molecular distributions," *Nature Communications*, vol. 13, no. 1, p. 3293, 2022.
- [4] H. Wang, T. Fu, Y. Du, W. Gao, K. Huang, Z. Liu, P. Chandak, S. Liu, P. Van Katwyk, A. Deac *et al.*, "Scientific discovery in the age of artificial intelligence," *Nature*, vol. 620, no. 7972, pp. 47–60, 2023.
- [5] M. Wen, E. W. C. Spotte-Smith, S. M. Blau, M. J. McDermott, A. S. Krishnapriyan, and K. A. Persson, "Chemical reaction networks and opportunities for machine learning," *Nature Computational Science*, vol. 3, no. 1, pp. 12–24, 2023.
- [6] A. V. Sadybekov and V. Katritch, "Computational approaches streamlining drug discovery," *Nature*, vol. 616, no. 7958, pp. 673–685, 2023.
- [7] J. Ross, B. Belgodere, V. Chenthamarakshan, I. Padhi, Y. Mroueh, and P. Das, "Large-scale chemical language representations capture molecular structure and properties," *Nature Machine Intelligence*, vol. 4, no. 12, pp. 1256–1264, 2022.
- [8] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill *et al.*, "On the opportunities and risks of foundation models," *arXiv preprint arXiv:2108.07258*, 2021.
- [9] S. Yang, O. Nachum, Y. Du, J. Wei, P. Abbeel, and D. Schuurmans, "Foundation models for decision making: Problems, methods, and opportunities," *arXiv preprint arXiv:2303.04129*, 2023.
- [10] T. Guo, B. Nan, Z. Liang, Z. Guo, N. Chawla, O. Wiest, X. Zhang *et al.*, "What can large language models do in chemistry? a comprehensive benchmark on eight tasks," *Advances in Neural Information Processing Systems*, vol. 36, pp. 59 662–59 688, 2023.
- [11] Z. Li, M. Jiang, S. Wang, and S. Zhang, "Deep learning methods for molecular representation and property prediction," *Drug Discovery Today*, vol. 27, no. 12, p. 103373, 2022.
- [12] L. Wei, N. Fu, Y. Song, Q. Wang, and J. Hu, "Probabilistic generative transformer language models for generative design of molecules," *Journal of Cheminformatics*, vol. 15, no. 1, p. 88, 2023.
- [13] H. Öztürk, A. Özgür, P. Schwaller, T. Laino, and E. Ozkirimli, "Exploring chemical space using natural language processing methodologies for drug discovery," *Drug Discovery Today*, vol. 25, no. 4, pp. 689–705, 2020.
- [14] B. I. Tingle, K. G. Tang, M. Castanon, J. J. Gutierrez, M. Khurelbaatar, C. Dandarchuluun, Y. S. Moroz, and J. J. Irwin, "Zinc 22 a free multi-billion-scale database of tangible compounds for ligand discovery," *Journal of chemical information and modeling*, vol. 63, no. 4, pp. 1166–1176, 2023.
- [15] D. S. Wigh, J. M. Goodman, and A. A. Lapkin, "A review of molecular representation in the age of machine learning," *Wiley Interdisciplinary Reviews: Computational Molecular Science*, vol. 12, no. 5, p. e1603, 2022.
- [16] W. Gao, T. Fu, J. Sun, and C. Coley, "Sample efficiency matters: a benchmark for practical molecular optimization," *Advances in neural information processing systems*, vol. 35, pp. 21 342–21 357, 2022.
- [17] S. Takeda, A. Kishimoto, L. Hamada, D. Nakano, and J. R. Smith, "Foundation model for material science," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 13, 2023, pp. 15 376–15 383.
- [18] S. Kim, J. Chen, T. Cheng, A. Gindulyte, J. He, S. He, Q. Li, B. A. Shoemaker, P. A. Thiessen, B. Yu *et al.*, "Pubchem 2023 update," *Nucleic acids research*, vol. 51, no. D1, pp. D1373–D1380, 2023.
- [19] D. Polykovskiy, A. Zhebrak, B. Sanchez-Lengeling, S. Golovanov, O. Tatanov, S. Belyaev, R. Kurbanov, A. Artamonov, V. Aladinskiy, M. Veselov *et al.*, "Molecular sets (moses): a benchmarking platform for molecular generation models," *Frontiers in pharmacology*, vol. 11, p. 565644, 2020.
- [20] E. Heid, J. Liu, A. Aude, and W. H. Green, "Influence of template size, canonicalization, and exclusivity for retrosynthesis and reaction prediction applications," *Journal of Chemical Information and Modeling*, vol. 62, no. 1, pp. 16–26, 2021.
- [21] P. Schwaller, T. Laino, T. Gaudin, P. Bolgar, C. A. Hunter, C. Bekas, and A. A. Lee, "Molecular transformer: a model for uncertainty-calibrated chemical reaction prediction," *ACS central science*, vol. 5, no. 9, pp. 1572–1583, 2019.

- [22] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *North American Chapter of the Association for Computational Linguistics*, 2019. [Online]. Available: <https://api.semanticscholar.org/CorpusID:52967399>
- [23] J. Su, Y. Lu, S. Pan, A. Murtadha, B. Wen, and Y. Liu, "Roformer: Enhanced transformer with rotary position embedding," *arXiv preprint arXiv:2104.09864*, 2021.
- [24] J. Ferrando, G. I. Gállego, I. Tsiamas, and M. R. Costa-jussà, "Explaining how transformers use context to build predictions," *arXiv preprint arXiv:2305.12535*, 2023.
- [25] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean, "Outrageously large neural networks: The sparsely-gated mixture-of-experts layer," *arXiv preprint arXiv:1701.06538*, 2017.
- [26] A. Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary, C. Bamford, D. S. Chaplot, D. d. l. Casas, E. B. Hanna, F. Bressand *et al.*, "Mixtral of experts," *arXiv preprint arXiv:2401.04088*, 2024.
- [27] Z. Wu, B. Ramsundar, E. N. Feinberg, J. Gomes, C. Geniesse, A. S. Pappu, K. Leswing, and V. Pande, "Moleculenet: a benchmark for molecular machine learning," *Chemical science*, vol. 9, no. 2, pp. 513–530, 2018.
- [28] S. Liu, H. Wang, W. Liu, J. Lasenby, H. Guo, and J. Tang, "Pre-training molecular graph representation with 3d geometry," *arXiv preprint arXiv:2110.07728*, 2021.
- [29] X. Fang, L. Liu, J. Lei, D. He, S. Zhang, J. Zhou, F. Wang, H. Wu, and H. Wang, "Geometry-enhanced molecular representation learning for property prediction," *Nature Machine Intelligence*, vol. 4, no. 2, pp. 127–134, 2022.
- [30] Y. Rong, Y. Bian, T. Xu, W. Xie, Y. Wei, W. Huang, and J. Huang, "Self-supervised graph transformer on large-scale molecular data," *Advances in Neural Information Processing Systems*, vol. 33, pp. 12 559–12 571, 2020.
- [31] S. Chithrananda, G. Grand, and B. Ramsundar, "Chemberta: large-scale self-supervised pretraining for molecular property prediction," *arXiv preprint arXiv:2010.09885*, 2020.
- [32] W. Ahmad, E. Simon, S. Chithrananda, G. Grand, and B. Ramsundar, "Chemberta-2: Towards chemical foundation models," *arXiv preprint arXiv:2209.01712*, 2022.
- [33] R. Taylor, M. Kardas, G. Cucurull, T. Scialom, A. Hartshorn, E. Saravia, A. Poulton, V. Kerkez, and R. Stojnic, "Galactica: A large language model for science," *arXiv preprint arXiv:2211.09085*, 2022.
- [34] G. Zhou, Z. Gao, Q. Ding, H. Zheng, H. Xu, Z. Wei, L. Zhang, and G. Ke, "Uni-mol: a universal 3d molecular representation learning framework," *ChemRxiv preprint*, 2023.
- [35] J. Chang and J. C. Ye, "Bidirectional generation of structure and properties through a single molecular foundation model," *Nature Communications*, vol. 15, no. 1, p. 2323, 2024.
- [36] K. Yang, K. Swanson, W. Jin, C. Coley, P. Eiden, H. Gao, A. Guzman-Perez, T. Hopper, B. Kelley, M. Mathea *et al.*, "Analyzing learned molecular representations for property prediction," *Journal of chemical information and modeling*, vol. 59, no. 8, pp. 3370–3388, 2019.
- [37] S. Liu, M. F. Demirel, and Y. Liang, "N-gram graph: Simple unsupervised representation for graphs, with applications to molecules," *Advances in neural information processing systems*, vol. 32, 2019.
- [38] W. Hu, B. Liu, J. Gomes, M. Zitnik, P. Liang, V. Pande, and J. Leskovec, "Strategies for pre-training graph neural networks," *arXiv preprint arXiv:1905.12265*, 2019.
- [39] Y. Wang, J. Wang, Z. Cao, and A. Barati Farimani, "Molecular contrastive learning of representations via graph neural networks," *Nature Machine Intelligence*, vol. 4, no. 3, pp. 279–287, 2022.
- [40] Z. Hu, Y. Dong, K. Wang, K.-W. Chang, and Y. Sun, "Gpt-gnn: Generative pre-training of graph neural networks," in *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, 2020, pp. 1857–1867.
- [41] C. Morris, M. Ritzert, M. Fey, W. L. Hamilton, J. E. Lenssen, G. Rattan, and M. Grohe, "Weisfeiler and leman go neural: Higher-order graph neural networks," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, no. 01, 2019, pp. 4602–4609.
- [42] M. Rupp, A. Tkatchenko, K.-R. Müller, and O. A. Von Lilienfeld, "Fast and accurate modeling of molecular atomization energies with machine learning," *Physical review letters*, vol. 108, no. 5, p. 058301, 2012.

- [43] K. T. Schütt, F. Arbabzadah, S. Chmiela, K. R. Müller, and A. Tkatchenko, "Quantum-chemical insights from deep tensor neural networks," *Nature communications*, vol. 8, no. 1, p. 13890, 2017.
- [44] W. Jin, R. Barzilay, and T. Jaakkola, "Junction tree variational autoencoder for molecular graph generation," in *International conference on machine learning*. PMLR, 2018, pp. 2323–2332.
- [45] P. Eckmann, K. Sun, B. Zhao, M. Feng, M. K. Gilson, and R. Yu, "Limo: Latent inceptionism for targeted molecule generation," *Proceedings of machine learning research*, vol. 162, p. 5777, 2022.
- [46] Y. Fang, N. Zhang, Z. Chen, L. Guo, X. Fan, and H. Chen, "Domain-agnostic molecular generation with self-feedback," *arXiv preprint arXiv:2301.11259*, 2023.
- [47] J. Ross, B. Belgodere, S. C. Hoffman, V. Chenthamarakshan, Y. Mroueh, and P. Das, "Gp-molformer: A foundation model for molecular generation," *arXiv preprint arXiv:2405.04912*, 2024.
- [48] L. van der Maaten and G. Hinton, "Visualizing high-dimensional data using t-sne," *Journal of Machine Learning Research*, vol. 9, no. nov, pp. 2579–2605, 2008, pagination: 27.
- [49] A. H. Lipkus, "A proof of the triangle inequality for the tanimoto distance," *Journal of Mathematical Chemistry*, vol. 26, no. 1, pp. 263–265, Oct 1999.
- [50] T. Chen, T. He, M. Benesty, V. Khotilovich, Y. Tang, H. Cho, K. Chen, R. Mitchell, I. Cano, T. Zhou *et al.*, "Xgboost: extreme gradient boosting," *R package version 0.4-2*, vol. 1, no. 4, pp. 1–4, 2015.
- [51] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, "Optuna: A next-generation hyperparameter optimization framework," in *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, 2019, pp. 2623–2631.

A Supplementary Materials

A.1 Detailed results - frozen weights

Here, we provide the detailed results for every experiment conducted in this paper. First, we present the detailed results for the experiments considering frozen weights of SMI-TED289M for both, classification and regression tasks, considering the MoleculeNet benchmarking dataset. For SMI-TED289D frozen weights, we considered XGBoost [50] as learner, and Optuna [51] for hyper-parameters optimization. Table 8 illustrates the results for the classification tasks using for 10 different seeds, and considering frozen weights.

Table 8: Classification results for 10 different seeds considering SMI-TED289 frozen weights.

SEED	ROC-AUC \uparrow					
	BBBP	HIV	BACE	SIDER	Clintox	Tox21
0	91.66	81.68	85.05	67.46	93.62	80.90
10	91.17	79.66	84.59	66.43	93.92	81.15
20	91.30	81.69	84.56	66.21	94.40	82.00
30	91.33	81.81	86.02	64.79	93.73	81.55
40	91.22	81.00	85.51	65.88	92.85	82.00
50	91.89	81.80	86.68	64.99	95.02	82.22
60	90.67	80.21	84.72	66.18	92.03	81.68
70	91.94	79.69	86.26	65.86	92.99	81.18
80	91.19	77.69	85.25	65.05	92.95	81.60
90	92.27	79.91	87.18	67.11	93.41	81.04
Average	91.46	80.51	85.58	66.00	93.49	81.53
Std	0.47	1.34	0.92	0.88	0.85	0.45

Table 9 elucidates the results for the regression tasks using for 10 different seeds, and considering frozen weights. Similar to the classification tasks, here we also use XGBoost as learner and Optuna for hyper-parameters optimization.

Table 9: Regression results for 10 different seeds considering SMI-TED289M frozen weights.

SEED	RMSE↓			MAE↓	
	ESOL	FreeSolv	Lipophilicity	QM8	QM9
0	0.6846	1.6248	0.6681	0.0184	7.4126
10	0.6784	1.7022	0.6400	0.0180	7.4956
20	0.6886	1.5832	0.6528	0.0174	7.6201
30	0.6880	1.7418	0.6311	0.0177	7.4845
40	0.7100	1.6443	0.6603	0.0185	7.5486
50	0.6933	1.6495	0.6515	0.0181	7.5118
60	0.6793	1.6285	0.6477	0.0182	7.5056
70	0.6884	1.7482	0.6411	0.0177	7.4128
80	0.7746	1.7468	0.6410	0.0179	7.4774
90	0.7599	1.6104	0.6654	0.0174	7.4135
Average	0.7045	1.6680	0.6499	0.0179	7.4883
Std	0.0344	0.0616	0.0120	0.0004	0.0659

A.2 Detailed results - Fine-tuning

To fine-tune SMI-TED289M, we used a fully connected network with 2 layers. Table 10 provides a detailed overview of the hyper-parameters considered for the fine-tuning of SMI-TED289M. We used a single V100 NVIDIA (16G) GPU for the task. Detailed results considering SMI-TED289M for both, classification and regression tasks using the MoleculeNet benchmarking dataset are illustrated in Table 11 and Table 12. We run each task for 10 different seeds to guarantee the robustness of the results.

Table 10: SMI-TED289M fine-tuning architecture specificity.

Hidden size	Attention heads	Layers	Dropout	Normalization
768	12	12	0.2	LayerNorm

Learning rate	# batch	# epochs	# tokens	# GPUs	Total params
3e-5	32	500	202	1 NVIDIA V100 (32G)	289M

Table 11 presents the results BBBP, HIV, BACE, SIDER, Clintox, Tox21 datasets. For these classifications tasks, ROC-AUC has been defined as evaluation metric as in the MoleculeNet. We run each seed for 500 epochs.

Table 11: Classification results for 10 different seeds considering SMI-TED289M fine-tuning.

SEED	ROC-AUC↑					
	BBBP	HIV	BACE	SIDER	Clintox	Tox21
0	92.42	76.76	88.02	65.88	96.55	81.87
10	92.20	76.89	87.82	66.12	91.86	82.20
20	92.48	75.72	88.63	65.05	94.95	80.58
30	92.17	76.52	87.82	65.97	97.97	83.72
40	91.94	77.01	88.32	65.30	92.90	83.08
50	91.29	79.09	88.63	66.51	93.95	83.27
60	93.07	76.49	89.33	65.49	94.32	80.26
70	92.84	76.52	87.91	65.22	93.41	79.41
80	92.74	76.33	87.80	65.71	92.85	81.44
90	91.49	77.20	88.08	65.59	93.96	82.65
Average	92.26	76.85	88.24	65.68	94.27	81.85
Std	0.57	0.89	0.50	0.45	1.83	1.42

Results for ESOL, FreeSolv, Lipophilicity, QM8, and QM9 are presented in Table 12. As for classification tasks, we also run each regression task for 10 different seeds, each one considering 500 epochs.

Table 12: Prediction results for 10 different seeds considering SMI-TED289M fine-tuning.

SEED	RMSE↓			MAE↓	
	ESOL	FreeSolv	Lipophilicity	QM8	QM9
0	0.6110	1.2258	0.5426	0.0092	1.2814
10	0.6110	1.2230	0.5375	0.0095	1.3371
20	0.6024	1.2230	0.5561	0.0094	1.3245
30	0.6124	1.2258	0.5472	0.0095	1.3291
40	0.6024	1.2258	0.5435	0.0095	1.3338
50	0.6024	1.2230	0.5413	0.0096	1.3302
60	0.6355	1.2167	0.5611	0.0099	1.3265
70	0.6116	1.2230	0.5513	0.0094	1.3293
80	0.6124	1.2258	0.5381	0.0095	1.3290
90	0.6110	1.2212	0.6029	0.0094	1.3249
Average	0.6112	1.2233	0.5522	0.0095	1.3246
Std	0.0096	0.0029	0.0194	0.0002	0.0157

QM9 and QM8 datasets contains 12 different metrics referring to the quantum properties of the molecules. Table 13 presents the results for the QM9 metrics: α , C_v , G , gap , H , ϵ_{homo} , ϵ_{lumo} , μ , $\langle R^2 \rangle$, U_0 , U , $ZPVE$. Table 13 also show the avg MAE and avg std MAE. For each seed we considered 500 epochs.

Table 13: Prediction results over SMI-TED289M fine-tuning for QM9 dataset considering 10 different seeds.

SEED	QM9												Average
	α	C_v	G	gap	H	ϵ_{homo}	ϵ_{lumo}	μ	$\langle R^2 \rangle$	U_0	U	$ZPVE$	
0	0.2266	0.0893	0.1503	0.0035	0.0873	0.0025	0.0024	0.3859	14.2478	0.0919	0.0890	0.0002	1.2814
10	0.2898	0.1283	0.1276	0.0037	0.1126	0.0027	0.0025	0.3850	14.7824	0.1005	0.1093	0.0007	1.3371
20	0.2826	0.1226	0.0937	0.0036	0.0871	0.0026	0.0025	0.3846	14.7603	0.0737	0.0804	0.0005	1.3245
30	0.2827	0.1249	0.1270	0.0036	0.1088	0.0026	0.0026	0.3842	14.7041	0.1010	0.1069	0.0010	1.3291
40	0.2880	0.1351	0.1219	0.0043	0.1099	0.0035	0.0032	0.3853	14.7624	0.0935	0.0971	0.0019	1.3338
50	0.2832	0.1241	0.1042	0.0036	0.0816	0.0027	0.0025	0.3845	14.8141	0.0794	0.0814	0.0007	1.3302
60	0.2835	0.1263	0.0964	0.0036	0.0870	0.0027	0.0025	0.3850	14.7702	0.0785	0.0819	0.0007	1.3265
70	0.2873	0.1284	0.1014	0.0036	0.0864	0.0026	0.0027	0.3845	14.7972	0.0758	0.0810	0.0006	1.3293
80	0.2866	0.1270	0.0844	0.0036	0.0843	0.0027	0.0025	0.3842	14.8097	0.0752	0.0875	0.0007	1.3290
90	0.2829	0.1257	0.0957	0.0036	0.0874	0.0027	0.0025	0.3848	14.7414	0.0809	0.0907	0.0006	1.3249
Average	0.2793	0.1232	0.1103	0.0037	0.0932	0.0027	0.0026	0.3848	14.7190	0.0850	0.0905	0.0008	1.3246
Std	0.0187	0.0124	0.0205	0.0002	0.0120	0.0003	0.0002	0.0005	0.1688	0.0106	0.0107	0.0004	0.0157

Table 14 illustrates the results for the QM8 metrics: E1-CAM, E1-CC2, E1-PBE0, E2-CAM, E2-CC2, E2-PBE0, f1-CAM, f1-CC2, f1-PBE0, f2-CAM, f2-CC2, f2-PBE0. We also show the results for the average MAE and average std MAE. For both tasks, QM8 and QM9, our proposed SMI-TED289M demonstrated better results when compared to the state-of-the-art methods. To demonstrate the robustness and reliability of our approach we extensively evaluated it over 10 different seeds, considering 500 epochs for each seed.

Table 14: Prediction results over SMI-TED289M fine-tuning for QM8 dataset considering 10 different seeds.

SEED	QM8												Average
	E1-CAM	E1-CC2	E1-PBE0	E2-CAM	E2-CC2	E2-PBE0	f1-CAM	f1-CC2	f1-PBE0	f2-CAM	f2-CC2	f2-PBE0	
0	0.0040	0.0037	0.0037	0.0041	0.0050	0.0046	0.0081	0.0097	0.0078	0.0188	0.0226	0.0182	0.0092
10	0.0040	0.0039	0.0038	0.0043	0.0051	0.0053	0.0085	0.0100	0.0083	0.0195	0.0231	0.0186	0.0095
20	0.0040	0.0038	0.0037	0.0042	0.0050	0.0051	0.0084	0.0100	0.0082	0.0194	0.0231	0.0183	0.0094
30	0.0040	0.0038	0.0038	0.0043	0.0051	0.0053	0.0085	0.0100	0.0083	0.0195	0.0229	0.0185	0.0095
40	0.0041	0.0039	0.0039	0.0042	0.0051	0.0052	0.0084	0.0100	0.0081	0.0194	0.0230	0.0185	0.0095
50	0.0040	0.0039	0.0039	0.0043	0.0051	0.0053	0.0086	0.0100	0.0084	0.0195	0.0231	0.0185	0.0096
60	0.0043	0.0042	0.0042	0.0046	0.0054	0.0056	0.0091	0.0103	0.0085	0.0200	0.0235	0.0189	0.0099
70	0.0040	0.0038	0.0037	0.0042	0.0050	0.0050	0.0083	0.0101	0.0081	0.0193	0.0230	0.0186	0.0094
80	0.0040	0.0038	0.0038	0.0043	0.0051	0.0053	0.0084	0.0100	0.0083	0.0197	0.0230	0.0187	0.0095
90	0.0040	0.0038	0.0038	0.0042	0.0051	0.0051	0.0085	0.0101	0.0082	0.0194	0.0228	0.0183	0.0094
Average	0.0040	0.0039	0.0038	0.0043	0.0051	0.0052	0.0085	0.0100	0.0082	0.0194	0.0230	0.0185	0.0095
Std	0.0001	0.0001	0.0002	0.0001	0.0001	0.0003	0.0003	0.0001	0.0002	0.0003	0.0002	0.0002	0.0001