
Molecule Attention Transformer

Łukasz Maziarka^{1,2} Tomasz Danel^{1,2} Sławomir Mucha² Krzysztof Rataj¹ Jacek Tabor²
Stanisław Jastrzębski^{3,4}

Abstract

Designing a single neural network architecture that performs competitively across a range of molecule property prediction tasks remains largely an open challenge, and its solution may unlock a widespread use of deep learning in the drug discovery industry. To move towards this goal, we propose Molecule Attention Transformer (MAT). Our key innovation is to augment the attention mechanism in Transformer using inter-atomic distances and the molecular graph structure. Experiments show that MAT performs competitively on a diverse set of molecular prediction tasks. Most importantly, with a simple self-supervised pretraining, MAT requires tuning of only a few hyperparameter values to achieve state-of-the-art performance on downstream tasks. Finally, we show that attention weights learned by MAT are interpretable from the chemical point of view.

1. Introduction

The task of predicting properties of a molecule lies at the center of applications such as drug discovery or material design. In particular, estimated 85% drug candidates fail the clinical trials in the United States after a long and costly development process (Wong et al., 2018). Potentially, many of these failures could have been avoided by having correctly predicted a clinically relevant property of a molecule such as its toxicity or bioactivity.

Following the breakthroughs in image (Krizhevsky et al., 2012) and text classification (Vaswani et al., 2017), deep neural networks (DNNs) are expected to revolutionize other fields such as drug discovery or material design (Jr et al.,

2019). However, on many molecular property prediction tasks DNNs are outperformed by *shallow* models such as support vector machine or random forest (Korotcov et al., 2017; Wu et al., 2018). On the other hand, while DNNs can outperform shallow models on some tasks, they tend to be difficult to train (Ishiguro et al., 2019; Hu et al., 2019), and can require tuning of a large number of hyperparameters. We also observe both issues on our benchmark (see Section 4.2).

Making deep networks easier to train has been the central force behind their widespread use. In particular, one of the most important breakthroughs in deep learning was the development of initialization methods that allowed to train easily deep networks end-to-end (Goodfellow et al., 2016). In a similar spirit, our aim is to develop a deep model that is simple to use out-of-the-box, and achieves strong performance on a wide range of tasks in the field of molecule property prediction.

In this paper we propose the Molecule Attention Transformer (MAT). We adapt Transformer (Devlin et al., 2018) to chemical molecules by augmenting the self-attention with inter-atomic distances and molecular graph structure. Figure 1 shows the architecture. We demonstrate that MAT, in contrast to other tested models, achieves strong performance across a wide range of tasks (see Figure 2). Next, we show that self-supervised pre-training further improves performance, while drastically reducing the time needed for hyperparameter tuning (see Table 3). In these experiments we tuned only the learning rate, testing 7 different values. Finally, we find that MAT has interpretable attention weights. We share pretrained weights at <https://github.com/gmum/MAT>.

2. Related work

Molecule property prediction. Predicting properties of a candidate molecule lies at the heart of many fields such as drug discovery and material design. Broadly speaking, there are two main approaches to predicting molecular properties. First, we can use our knowledge of the underlying physics (Lipinski et al., 1997). However, despite recent advances (Schütt et al., 2017), current approaches remain prohibitively costly to accurately predict many properties of

¹Ardigen, Cracow, Poland.

²Jagiellonian University, Cracow, Poland.

³Molecule.one, Warsaw, Poland.

⁴New York University, New York, USA.

Correspondence to:

Łukasz Maziarka <lukasz.maziarka@ardigen.com>,
Stanisław Jastrzębski <staszek.jastrzebski@gmail.com>.

Preprint. Work in progress.

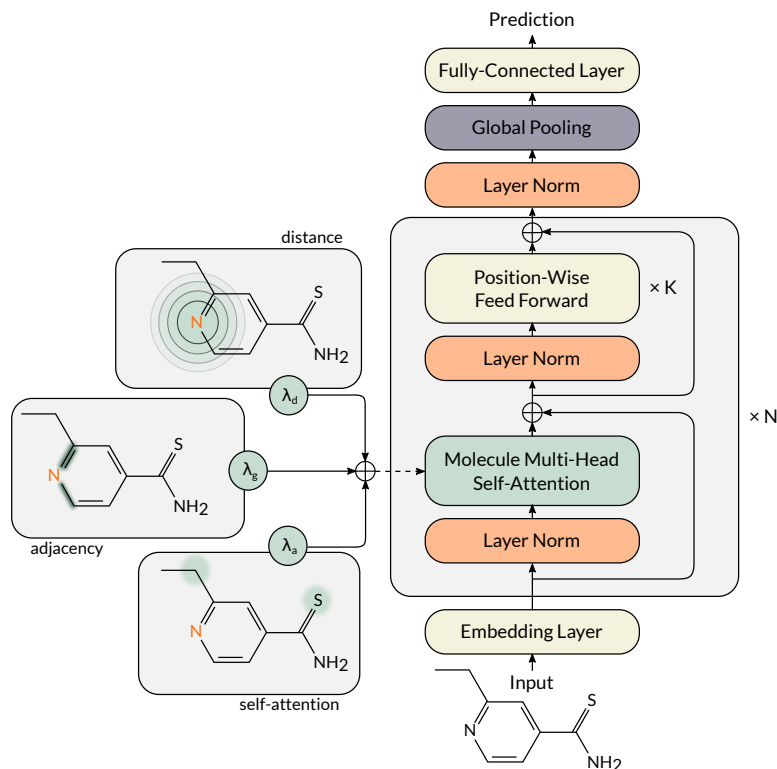


Figure 1. Molecule Attention Transformer architecture. We largely base our model on the Transformer encoder. In the first layer we embed each atom using one-hot encoding and atomic features. The main innovation is the Molecule Multi-Head Self-Attention layer that augments attention with distance and graph structure of the molecule. We implement this using a weighted (by λ_d , λ_g , and λ_a) element-wise sum of the corresponding matrices.

interest such as bioactivity. The second approach is to use existing data to train a predictive model (Haghighatlari & Hachmann, 2019). Here the key issue is the lack of large datasets. Even for the most popular drug targets, such as 5-HT1A (a popular target for depression), only thousands of active compounds are known. Promising direction is using hybrid approaches such as Wallach et al. (2015) or approaches leveraging domain knowledge and underlying physics to impose a strong prior such as Feinberg et al. (2018).

Deep learning for molecule property prediction. Deep learning has become a valuable tool for modeling molecules. During the years, the community has progressed from using handcrafted representations to representing molecules as strings of symbols, and finally to the currently popular approaches based on molecular graphs.

Graph convolutional networks in each subsequent layer gather information from adjacent nodes in the graph. In this way after N convolution layers each node has information from its N -edges distant neighbors. Using the graph structure improves performance in a range of molecule modeling tasks (Wu et al., 2018). Some of the most recent works

implement more sophisticated generalization methods for gathering the neighbor data. Veličković et al. (2017); Shang et al. (2018) propose to augment GCNs with an attention mechanism. Li et al. (2018) introduces a model that dynamically learns neighbourhood function in the graph.

In parallel to these advances, using the three-dimensional structure of the molecule is becoming increasingly popular. Perhaps the most closely related models are 3D Graph Convolutional Neural Network (3DGCN), Message Passing Neural Network (MPNN), and Adaptive Graph Convolutional Network (AGCN) (Cho & Choi, 2018; Gilmer et al., 2017; Li et al., 2018). 3DGCN and MPNN integrate graph and distance information in a single model, which enables them to achieve strong performance on tasks such as solubility prediction. In contrast to them, we additionally allow for a flexible neighbourhood based on self-attention.

Transformer, originally developed for natural language processing (Vaswani et al., 2017), has been recently applied to retrosynthesis in Karpov et al. (2019). They represent compounds as sentences using the SMILES notation (Weininger, 1988). In contrast to them, we represent compounds as a list of atoms, and ensure that models understand the structure of

the molecule by augmenting the self-attention mechanism (see Figure 1). Our ablation studies show it is a critical component of the model.

To summarize, methods related to our model have been proposed in the literature. Our contribution is unifying these ideas in a single model based on the state-of-the-art Transformer architecture that preserves strong performance across many chemical tasks.

How easy is it to use deep learning for molecule property prediction? DNNs performance is not always competitive to methods such as support vector machine or random forest. MoleculeNet is a popular benchmark for methods for molecule property prediction (Wu et al., 2018) that demonstrates this phenomenon. Similar results can be found in Withnall et al. (2019). We reproduce a similar issue on our benchmark. This makes using deep learning less applicable to molecule property prediction because in some cases practitioners might actually benefit from using other methods. Another issue is that graph neural networks, which are the most popular class of models for molecule property prediction, can be difficult to train. Ishiguro et al. (2019) show and try to address the problem that graph neural networks tend to underfit the training set. We also reproduce this issue on our benchmark (see also App. C).

There has been a considerable interest in developing easier to use deep models for molecule property prediction. Goh et al. (2017) pretrains a deep network that takes as an input an image of a molecule. Another studies highlight the need to augment feedforward (Mayr et al., 2018) and graph neural networks (Yang et al., 2019) with handcrafted representations of molecules. Hu et al. (2019) proposes pretraining methods for graph neural networks and shows this largely alleviates the problem of underfitting, present in these architectures (Ishiguro et al., 2019). We take inspiration from Hu et al. (2019) and use one of the three pretraining tasks proposed therein.

Concurrently, Wang et al. (2019); Honda et al. (2019) pre-train a vanilla Transformer (Devlin et al., 2018) that takes as input a text representation (SMILES) of a molecule. Honda et al. (2019) shows that decoding based approach improves data efficiency of the model. A similar approach, specialized to the task of drug-target interaction prediction, was concurrently proposed in Shin et al. (2019). In contrast to them, we adapt Transformer to chemical structures, which in our opinion is crucial for achieving strong empirical performance. We also use a domain-specific pretraining based on Wu et al. (2018). We further confirm importance of both approaches by comparing directly with Honda et al. (2019).

Self-attention based models. Arguably, the attention mechanism (Bahdanau et al., 2014) has been one of the

most important breakthroughs in deep learning. This is perhaps best illustrated by the wide-spread use of Transformer architecture in natural language processing (Vaswani et al., 2017; Devlin et al., 2018).

Multiple prior works have augmented self-attention in Transformer using domain-specific knowledge (Chen et al., 2018; Shaw et al., 2018; Bello et al., 2019; Guo et al., 2019). Guo et al. (2019) encourages Transformer to attend to adjacent words in a sentence, and Chen et al. (2018) encourages another attention-based model to focus on pairs of words in a sentence that are connected in an external knowledge base. Our novelty is applying this successive modeling idea to molecule property prediction.

3. Molecule Attention Transformer

As the rich literature on deep learning for molecule property prediction suggests, it is necessary for a model to be flexible enough to represent a range of possible relationships between atoms of a compound. Inspired by its flexibility and strong empirical performance, we base our model on the Transformer encoder (Vaswani et al., 2017; Devlin et al., 2018). It is worth noting that natural language processing has inspired important advances in cheminformatics (Segler et al., 2017; Gómez-Bombarelli et al., 2018), which might be due to similarities between the two domains (Jastrzębski et al., 2016).

Transformer. We begin by briefly introducing the Transformer architecture. On a high level, Transformer for classifications has N attention blocks followed by a pooling and a classification layer. Each attention block is composed of a multi-head self-attention layer, followed by a feed-forward block that includes a residual connection and layer normalization.

The multi-head self-attention is composed of H heads. Head i ($i = 1, \dots, H$) takes as input hidden state \mathbf{H} and computes first $\mathbf{Q}_i = \mathbf{H}\mathbf{W}_i^Q$, $\mathbf{K}_i = \mathbf{H}\mathbf{W}_i^K$, and $\mathbf{V}_i = \mathbf{H}\mathbf{W}_i^V$. These are used in the attention operation as follows:

$$\mathcal{A}^{(i)} = \rho \left(\frac{\mathbf{Q}_i \mathbf{K}_i^T}{\sqrt{d_k}} \right) \mathbf{V}_i, \quad (1)$$

Molecule Self-Attention. Using a naive Transformer architecture would require encoding of chemical molecules as sentences. Instead, inspired by Battaglia et al. (2018), we interpret the self-attention as a soft adjacency matrix between the elements of the input sequence. Following this line of thought, it is natural to augment the self-attention using information about the actual structure of the model. This allows us to avoid using linearized (textual) representation of molecule as input (Jastrzębski et al., 2016), which we expect to be a better inductive bias for the model.

More concretely, we propose the Molecule Self-Attention layer, which we describe in Equation 2. We augment the self-attention matrix as follows: let $\mathbf{A} \in \{0, 1\}^{N_{\text{atoms}} \times N_{\text{atoms}}}$ denote the graph adjacency matrix, and $\mathbf{D} \in \mathbb{R}^{N_{\text{atoms}} \times N_{\text{atoms}}}$ denote the inter-atomic distances. Let λ_a , λ_d , and λ_g denote scalars weighting the self-attention, distance, and adjacency matrices. We modify Equation 1 as follows:

$$\mathcal{A}^{(i)} = \left(\lambda_a \rho \left(\frac{\mathbf{Q}_i \mathbf{K}_i^T}{\sqrt{d_k}} \right) + \lambda_d g(\mathbf{D}) + \lambda_g \mathbf{A} \right) \mathbf{V}_i, \quad (2)$$

see also Figure 1. We denote λ_a , λ_d , and λ_g jointly as λ . We use as g either softmax (normalized over the rows), or an element-wise $g(d) = \exp(-d)$. Finally, the distance matrix \mathbf{D} is computed using RDKit package (Landrum, 2016).

Note that while we use only the adjacency and the distance matrices, MAT can be easily extended to include other types of information, e.g. forces between the atoms.

Molecule Attention Transformer. To define the model, we replace all self-attention layers in the original Transformer encoder by our Molecular Self Attention layers. We embed each atom as a 26 dimensional vector following (Coley et al., 2017), shown in Table 1. In the experiments, we treat λ_a , λ_d , and λ_g as hyperparameters and keep them frozen during training. Figure 1 illustrates the model.

Pretraining. We experiment with one of the two *node-level* pretraining tasks proposed in Hu et al. (2019), which involves predicting the masked input nodes. Consistently with Hu et al. (2019), we found it stabilizes learning (see Figure 6) and reduces the need for an extensive hyperparameter search (see Table 3). Given that MAT already achieves good performance using this simple pretraining task, we leave for future work exploring the other tasks proposed in Hu et al. (2019).

Other details. Inspired by Li et al. (2017); Clark et al. (2019), we add an artificial *dummy node* to the molecule. The dummy node is not connected by an edge to any other atom and the distance to any of them is set to 10^6 . Our motivation is to allow the model to skip searching for a molecular pattern if none is to find by putting higher attention on that distant node, which is similar to how BERT uses the separation token (Devlin et al., 2018; Clark et al., 2019). We confirm this intuition in Section 4.4 and Section 4.5.

Finally, the distance matrices are calculated from 3D conformers calculated using UFFOPTIMIZE MOLECULE function from the RDKit package (Landrum, 2016), and the default parameters (MAXITERS=200, VDWTHRESH=10.0, CONFID=-1, IGNOREINTERFRAGMENTINTERACTIONS=True). For each compound we use one pre-computed conformation.

We experimented with sampling more conformations for each compound, but did not observe a consistent boost in performance, however it is possible that using more sophisticated algorithms for compound 3D structure minimization could improve the results. We leave this for future work.

Table 1. Featurization used to embed atoms in MAT.

INDICES	DESCRIPTION
0 – 11	ATOMIC IDENTITY AS A ONE-HOT VECTOR OF B, N, C, O, F, P, S, CL, BR, I, DUMMY, OTHER
12 – 17	NUMBER OF HEAVY NEIGHBORS AS ONE-HOT VECTOR OF 0, 1, 2, 3, 4, 5
18 – 22	NUMBER OF HYDROGEN ATOMS AS ONE-HOT VECTOR OF 0, 1, 2, 3, 4
23	FORMAL CHARGE
24	IS IN A RING
25	IS AROMATIC

4. Experiments

We begin by comparing MAT to other popular models in the literature on a wide range of tasks. We find that with simple pretraining MAT outperforms other methods, while using a small budget for hyperparameter tuning.

In the rest of this section we try to develop understanding of what makes MAT work well. In particular, we find that individual heads in the multi-headed self-attention layers learn chemically interpretable functions.

4.1. Experimental settings

Comparing different models for molecule property prediction is challenging. Despite considerable efforts, the community still lacks a standardized way to compare different models. In our work, we use a similar setting to MoleculeNet (Wu et al., 2018).

Evaluation. Following recommendations of Wu et al. (2018) and the experimental setup of Podlowska & Kafel (2018), we use random split for FreeSolv, ESOL, and MetStab. For all the other datasets we use scaffold split, which assigns compounds that share the same molecular scaffolding to different subsets of the data (Bemis & Murcko, 1996). In regression tasks, the property value was standardized. Test performance is based on the model which gave best results in the validation setting. Each training was repeated 6 times, on different train/validation/test splits. All the other experimental details are reported in the Supplement.

Datasets. We run experiments on a wide range of datasets that represent typical tasks encountered in molecule mod-

eling. Below, we include a short description of these tasks, and a more detailed description is moved to App. A.

- **FreeSolv, ESOL.** Regression tasks used in Wu et al. (2018) for predicting water solubility in terms of the hydration free energy (FreeSolv) and log solubility in mols per litre (ESOL). The datasets have 642 and 1128 molecules, respectively.
- **Blood-brain barrier permeability (BBBP).** Binary classification task used in Wu et al. (2018) for predicting the ability of a molecule to penetrate the blood-brain barrier. The dataset has 2039 molecules.
- **Estrogen Alpha, Estrogen Beta.** The tasks are to predict whether a compound is active towards a given target (Estrogen- α , Estrogen- β) based on experimental data from the ChEMBL database (Gaulton et al., 2011). The datasets have 2398, and 1961 molecules, respectively.
- **MetStab_{high}, MetStab_{low}.** Binary classification tasks based on data from Podlowska & Kafel (2018) to predict whether a compound has high (over 2.32 h half-time) or low (lower than 0.6 h half-time) metabolic stability. Both datasets contain the same 2127 molecules.

4.2. Molecule Attention Transformer

Models. Similarly to Wu et al. (2018), we test a comprehensive set of baselines that span both shallow and deep models. We compare MAT to the following baselines: GCN (Duvenaud et al., 2015), Random Forest (RF) and Support Vector Machine with RBF kernel (SVM). We also test the following recently proposed models: Edge Attention-based Multi-Relational Graph Convolutional Network (EAGCN) (Shang et al., 2018), and Weave (Kearnes et al., 2016).

Hyperparameter tuning. For each method we extensively tune their hyperparameters using random search (Bergstra & Bengio, 2012). To ensure fair comparison, each model is given the same budget for hyperparameter search. We run two sets of experiments with budget of 150 and 500 evaluations. We include hyperparameter ranges in App. B.

Results. We evaluate models by their average rank according to the test set performance on the 7 datasets. Figure 2 reports ranks of all methods for the two considered hyperparameter budgets (150 and 500). Additionally, we report in Table 2 detailed scores on all datasets. We make three main observations.

First, graph neural networks (GCN, Weave, EAGCN) on average do not outperform the other models. The best graph

model achieves average rank 3.28 compared to 3.14 by RF. On the whole, performance of the deep models improves with larger hyperparameter search budget. This further corroborates the original motivation of our study. Indeed, using common deep learning methods for molecule property prediction is challenging in practice. It requires a large computational budget, and might still result in poor performance.

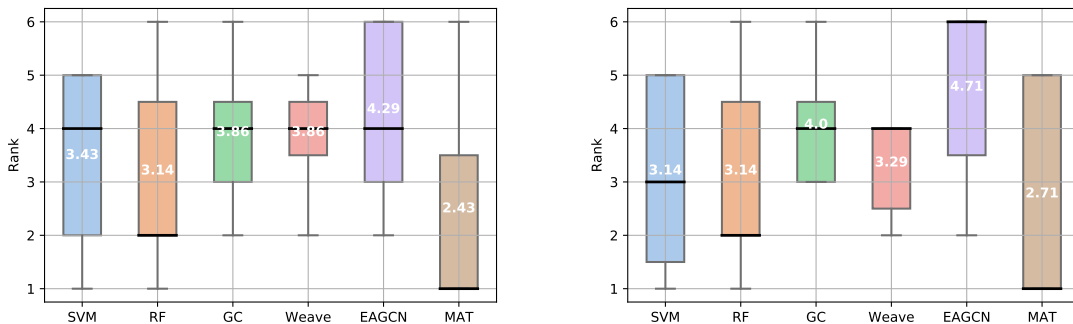
Second, MAT outperforms the other tested methods in terms of the average rank. MAT achieves average rank of 2.71 and 2.42 for 150 and 500 budgets, compared to 3.14 of RF, which is the second best performing model. This shows that architecture MAT is flexible enough and has the correct inductive bias to perform well on a wide range of tasks.

Examining performance of MAT across individual datasets, we observe that RF and SVM perform better on Estrogen- β , MetStab_{low}, and MetStab_{high}. Both RF and SVM use extended-connectivity fingerprint (Rogers & Hahn, 2010) as input representation, which encodes substructures in the molecule as features. Metabolic stability of a compound depends on existence of particular moieties, which are recognized by enzymes. Therefore a simple structure-based fingerprints perform well in such a setting. Wang et al. (2019); Mayr et al. (2018) show that using fingerprint as input representation improves performance of deep networks on related datasets. These two observations suggest that MAT could benefit from using fingerprints. Instead, we avoid using handcrafted representations, and investigate pretraining as an alternative in the next section. Though fingerprint-based models show excellent performance in all presented tasks, there are datasets on which they fail to match the performance of graph approaches. We observed this also on an energy prediction task (see the extension of our benchmark in App. C).

4.3. Pretrained Molecule Attention Transformer

Self-supervised pretraining has revolutionized natural language processing (Devlin et al., 2018) and has improved performance in molecule property prediction (Hu et al., 2019). We apply here node-level self-supervised pretraining from Hu et al. (2019) to MAT. The task is to predict features of masked out nodes. We refer the reader to App. D for more details.

Models. We compare MAT to the two following baselines. First, we apply the same pretraining to EAGCN, which we will refer to as “Pretrained EAGCN”. Second, we compare to a concurrent work by Honda et al. (2019). They pretrain a vanilla Transformer by decoding textual representation (SMILES) of molecules. We will refer to their method as “SMILES Transformer”.



(a) Hyperparameter search budget of 500 combinations.

(b) Hyperparameter search budget of 150 combinations.

Figure 2. The average rank across the 7 datasets in the benchmark. For each model we test 500 (left) or 150 (right) hyperparameter combinations. We split the data using random or scaffold split (according to the dataset description) 6 times into train/validation/test folds and use the mean metrics across the test folds to obtain the ranklists of models. Interestingly, *shallow* models (RF and SVM) outperform graph models (GCN, EAGCN and Weave).

Table 2. Test performances in the benchmark. For each model we test 500 (top) and 150 (bottom) hyperparameter combinations. On ESOL and FreeSolv we report RMSE (lower is better). The other tasks are evaluated using ROC AUC (higher is better). Experiments are repeated 6 times.

(a) Hyperparameter search budget of 500 combinations.

	BBBP	ESOL	Freesolv	ESTROGEN- α	ESTROGEN- β	METSTAB _{LOW}	METSTAB _{HIGH}
SVM	.707 \pm .000	.478 \pm .054	.461 \pm .077	.973 \pm .000	.778 \pm .000	.893 \pm .030	.890 \pm .029
RF	.725 \pm .006	.534 \pm .073	.523 \pm .097	.977 \pm .001	.797 \pm .007	.885 \pm .029	.888 \pm .030
GCN	.712 \pm .010	.357 \pm .032	.271 \pm .048	.975 \pm .003	.730 \pm .006	.881 \pm .031	.875 \pm .036
WEAVE	.701 \pm .016	.311 \pm .023	.311 \pm .072	.974 \pm .003	.769 \pm .023	.863 \pm .028	.882 \pm .043
EAGCN	.680 \pm .014	.316 \pm .024	.345 \pm .051	.961 \pm .011	.781 \pm .012	.883 \pm .024	.868 \pm .034
MAT (OURS)	.728 \pm .008	.285 \pm .022	.263 \pm .046	.979 \pm .003	.765 \pm .007	.862 \pm .038	.888 \pm .027

(b) Hyperparameter search budget of 150 combinations.

	BBBP	ESOL	Freesolv	ESTROGEN- α	ESTROGEN- β	METSTAB _{LOW}	METSTAB _{HIGH}
SVM	.723 \pm .000	.479 \pm .055	.461 \pm .077	.973 \pm .000	.772 \pm .000	.893 \pm .030	.890 \pm .029
RF	.721 \pm .003	.534 \pm .073	.524 \pm .098	.977 \pm .001	.791 \pm .012	.892 \pm .026	.888 \pm .030
GCN	.695 \pm .013	.369 \pm .032	.299 \pm .068	.975 \pm .003	.730 \pm .006	.884 \pm .033	.875 \pm .036
WEAVE	.702 \pm .009	.298 \pm .025	.298 \pm .049	.974 \pm .003	.769 \pm .023	.863 \pm .028	.885 \pm .042
EAGCN	.680 \pm .014	.322 \pm .052	.337 \pm .042	.961 \pm .011	.781 \pm .012	.859 \pm .024	.844 \pm .037
MAT (OURS)	.727 \pm .006	.290 \pm .019	.289 \pm .047	.979 \pm .003	.765 \pm .007	.861 \pm .029	.844 \pm .052

Hyperparameters. For all methods that use pre-training we reduce the hyperparameter grid to a minimum. We tune only the learning rate in $\{1e-3, 5e-4, 1e-4, 5e-5, 1e-5, 5e-6, 1e-6\}$. We set the other hyperparameters to reasonable defaults based on results from Section 4.2. For MAT and EAGCN, we follow (Devlin et al., 2018) and use the largest model that still fits the GPU memory. For SMILES Transformer we use pretrained weights provided by Honda et al. (2019).

Results. As in previous section, we compare the models based on their average rank on our benchmark. Figure 3 and Table 3 summarize the results.

We observe that Pretrained MAT achieves average rank of 1.57 and outperforms MAT (average rank of 2.14). Importantly, for Pretrained MAT we only tuned the learning rate by evaluating 7 different values. This is in stark contrast to the 500 hyperparameter combinations tested for MAT and EAGCN. To visualize this, in Figure 4 we plot the average test performance of all models as a function of the number of tested hyperparameter combinations. We also note that

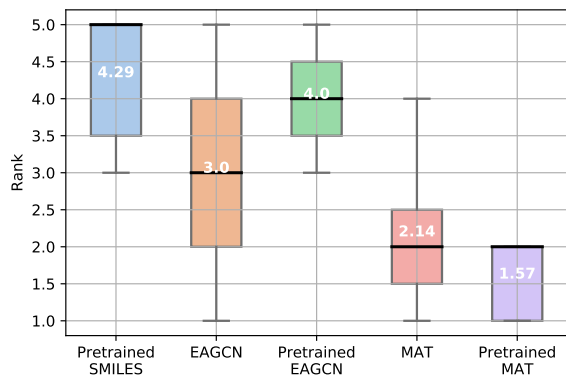


Figure 3. The average ranks across the 7 datasets in the benchmark. Pretrained MAT outperforms the other methods, despite a drastically smaller number of tested hyperparameters (7) compared to MAT and EAGCN (500).

Pretrained MAT is more competitive on the three datasets mentioned in the previous section.

We also find that Pretrained MAT outperforms the other two pretrained methods. Pretraining degrades the performance of EAGCN (average rank of 4.0), and SMILES Transformer achieves the worst average rank (average rank of 4.29). This suggests that both the architecture, and the choice of the pretraining task are important for the overall performance of the model.

4.4. Ablation studies

To better understand what contributes to the performance of MAT, we run a series of ablation studies on three representative datasets from our benchmark. We leave understanding how these choices interact with pretraining for future work.

For experiments in this section we generated additional splits for ESOL, FreeSolv and BBBP datasets (different than in Section 4.2). For each configuration we select the best hyperparameters settings using random search under a budget of 100 evaluations. Experiments are repeated 3 times.

Dummy node is not so dummy. MAT uses a *dummy node* that is disconnected from other atoms in the graph (Li et al., 2017). Our intuition is that such functionality can be useful to automatically adapt capacity on small datasets. By attending to the dummy node, the model can effectively choose to avoid changing the internal representation in a given layer. To examine this architectural choice, in Table 4 we compare MAT to a variant that does not include the dummy node. Results show that dummy node improves performance of the model.

Knowing molecular graph and distances between atoms improves performance.

Our key architectural innovation is integrating the molecule graph and inter-atomic distances with the self-attention layer in Transformer, as shown in Figure 1. To probe the importance of each of these sources of information, we removed them individually during training. Results in Table 5 suggest that keeping all sources of information results in the most stable performance across the three tasks, which is our primary goal. We also show that MAT can effectively use distance information in a toy task involving 3-dimensional distances between functional groups (see App.F).

Using a more complex featurization does not improve performance.

Many models for predicting molecule properties use additional edge features (Coley et al., 2017; Shang et al., 2018; Gilmer et al., 2017). In Table 6 we show that adding additional edge features does not improve MAT performance. This is certainly possible that a more comprehensive set of edge features or a better method to integrate them would improve performance, which we leave for future work. Procedure of using edge features is described in detail in App. E.

4.5. Analysis.

To understand MAT better, we investigate attention weights of the model, and the effect of pretraining on the learning dynamics.

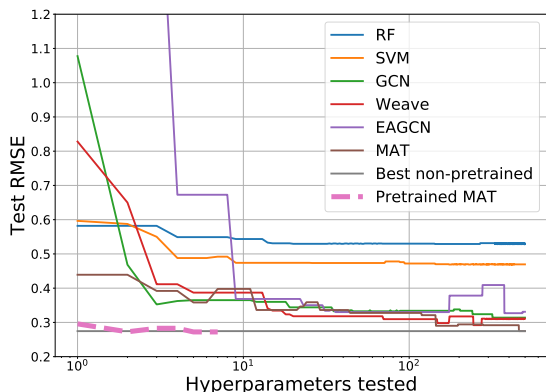
What is MAT looking at? In natural language processing, it has been shown that heads in Transformer seem to implement interpretable functions (Htut et al., 2019; Clark et al., 2019). Similarly, we investigate here the chemical function implemented by self-attention heads in MAT. We show patterns found in the model that was pretrained with the atom masking strategy (Hu et al., 2019), and then we verify our findings on a set of molecules extracted from the BBBP testing dataset.

Based on a manual inspection of attention matrices of MAT, we find two broad patterns: (1) many attention heads are almost fully focused on the dummy node, (2) many attention heads focus only on a few atoms. This seems consistent with observations about Transformer in Clark et al. (2019). We also notice that initial self-attention layers learn simple and easily interpretable chemical patterns, while subsequent layers capture more complex arrangements of atoms. In Figure 5 we exemplify attention patterns on a random molecule from the BBBP dataset.

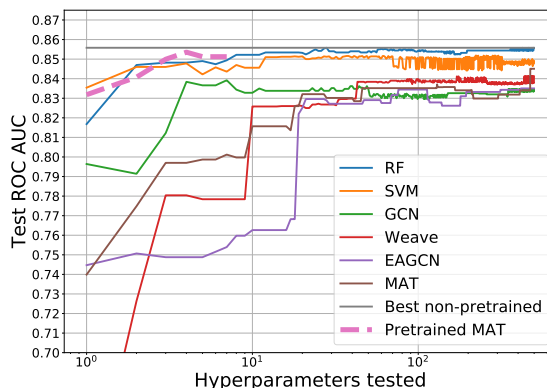
To quantify the above findings, we select six heads from the first layer that fit the second category and seem to implement six patterns: (i) focuses on 2-neighbouring aromatic carbons (not substituted); (ii) focuses on sulfurs; (iii) focuses on non-

Table 3. Test set performances of methods that use pretraining. Experiments are repeated 6 times. SMILES refers to SMILES Transformer from Honda et al. (2019).

	BBBP	ESOL	FREE SOLV	ESTROGEN- α	ESTROGEN- β	METSTAB _{LOW}	METSTAB _{HIGH}
MAT	.737 \pm .009	.278 \pm .020	.265 \pm .042	.998 \pm .000	.773 \pm .012	.862 \pm .025	.884 \pm .030
EAGCN	.687 \pm .023	.323 \pm .031	1.244 \pm .341	.994 \pm .002	.770 \pm .010	.861 \pm .029	.839 \pm .038
SMILES	.717 \pm .008	.356 \pm .017	.393 \pm .032	.953 \pm .002	.757 \pm .002	.860 \pm .038	.881 \pm .036



(a) Regression tasks.



(b) Classification tasks.

Figure 4. Test performance of all models as a function of the number of tested hyperparameter combinations (on a logarithmic scale). Figures show the aggregated mean RMSE for regression tasks (left) and the aggregated mean ROC AUC for classification tasks (right). Pretrained MAT requires tuning an order of magnitude less hyperparameters, and performs competitively on both sets of tasks.

Table 4. Test performance of MAT model variant without the dummy node (- DUMMY) compared to performance of the original MAT.

	BBBP	ESOL	FREE SOLV
MAT	.723 \pm .008	.286 \pm .006	.250 \pm .007
- DUMMY	.714 \pm .010	.317 \pm .014	.249 \pm .014

Table 5. Test performance of MAT with different sources of information removed (equivalent to setting the corresponding λ to zero).

	BBBP	ESOL	FREE SOLV
MAT	.723 \pm .008	.286 \pm .006	.250 \pm .007
- GRAPH	.716 \pm .009	.316 \pm .036	.276 \pm .034
- DISTANCE	.729 \pm .013	.281 \pm .001	.281 \pm .013
- ATTENTION	.692 \pm .001	.306 \pm .026	.329 \pm .014

Table 6. Test performance of MAT using additional edge features (+ EDGES F.), compared to vanilla MAT.

	BBBP	ESOL	FREE SOLV
MAT	.723 \pm .008	.286 \pm .006	.250 \pm .007
+ EDGES F.	.683 \pm .008	.314 \pm .014	.358 \pm .023

ring nitrogens; (iv) focuses on oxygen in carbonyl groups; (v) focuses on 3-neighbouring aromatic atoms (positions of aromatic ring substitutions) and on sulfur for different atoms; (vi) focuses on nitrogens in aromatic rings. We found that on the BBBP testing dataset the atoms corresponding to these definitions (queried with SMARTS expressions) have indeed higher attention weights assigned to them than other atoms. For each head, we calculated attention weights for all atoms in all molecules and compared those matching our hypothesis against the other atoms. Their distributions differ significantly ($p < 0.001$ in Kruskal-Wallis test) for all the patterns. The statistics and experimental details are summarized in App. G.

Effect of pretraining. Wu et al. (2018) observed that using pretraining stabilizes and speeds up training of graph convolutional models. We observe a similar effect in our case. Figure 6 reports training error of MAT and Pretrained MAT on the ESOL (left), and the FreeSolv (right) datasets. We use the learning rate that achieved the best generalization on each dataset in Sec. 4.3. The experiments are repeated 6 times. On both datasets, Pretrained MAT converges faster and has a lower variance of training error across repetitions. Mean standard deviation of training error for Pretrained MAT (MAT) is 0.027 (0.057) and 0.040 (0.076) for ESOL

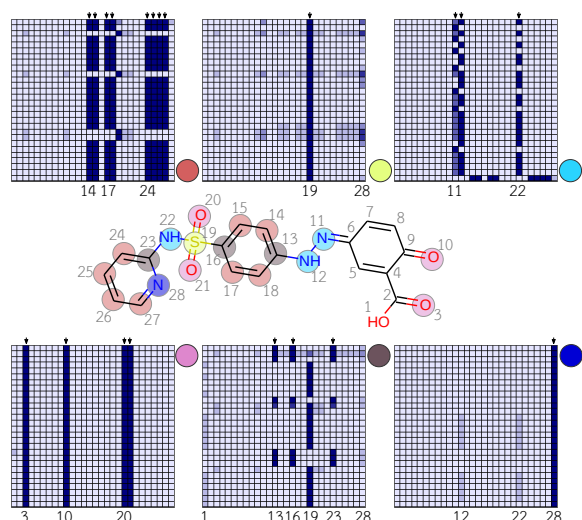


Figure 5. The heatmaps show selected self-attention weights from the first layer of MAT, on a random molecule from the BBBP dataset (center). The atoms, which these heads focus on, are marked with the same color as the corresponding matrix. The interpretation of the presented patterns is described in the text.

and FreeSolv, respectively.

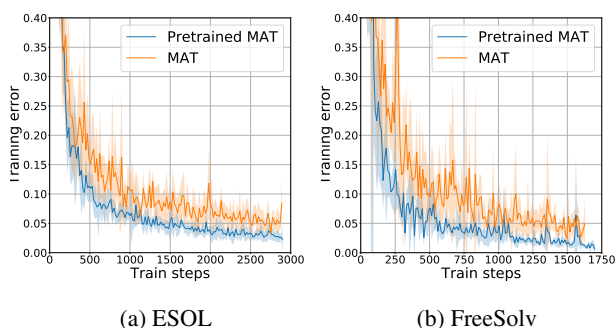


Figure 6. Training of MAT with (blue) and without (orange) pre-training, on ESOL (left) and FreeSolv (right). Pretraining stabilizes training (smaller variance of the training error) and improves convergence speed.

5. Conclusions.

In this work we propose Molecule Attention Transformer as a versatile architecture for molecular property prediction. In contrast to other tested models, MAT performs well across a wide range of molecule property prediction tasks. Moreover, inclusion of self-supervised pretraining further improves its performance, and drastically reduces the need for tuning of hyperparameters.

We hope that our work will widen adoption of deep learning in applications involving molecular property prediction,

as well as inspire new modeling approaches. One particularly promising avenue for future work is exploring better pretraining tasks for MAT.

References

- Bahdanau, D., Cho, K., and Bengio, Y. Neural machine translation by jointly learning to align and translate, 2014.
- Battaglia, P. W., Hamrick, J. B., Bapst, V., Sanchez-Gonzalez, A., Zambaldi, V. F., Malinowski, M., Tacchetti, A., Raposo, D., Santoro, A., Faulkner, R., Gülçehre, Ç., Song, F., Ballard, A. J., Gilmer, J., Dahl, G. E., Vaswani, A., Allen, K., Nash, C., Langston, V., Dyer, C., Heess, N., Wierstra, D., Kohli, P., Botvinick, M., Vinyals, O., Li, Y., and Pascanu, R. Relational inductive biases, deep learning, and graph networks. *CoRR*, abs/1806.01261, 2018.
- Bello, I., Zoph, B., Vaswani, A., Shlens, J., and Le, Q. V. Attention augmented convolutional networks. *CoRR*, abs/1904.09925, 2019.
- Bemis, G. W. and Murcko, M. A. The properties of known drugs. 1. molecular frameworks. *Journal of medicinal chemistry*, 39(15):2887–2893, 1996.
- Bergstra, J. and Bengio, Y. Random search for hyperparameter optimization. *J. Mach. Learn. Res.*, 13(1): 281–305, February 2012. ISSN 1532-4435.
- Chen, G., Chen, P., Hsieh, C.-Y., Lee, C.-K., Liao, B., Liao, R., Liu, W., Qiu, J., Sun, Q., Tang, J., et al. Alchemy: A quantum chemistry dataset for benchmarking ai models. *arXiv preprint arXiv:1906.09427*, 2019.
- Chen, Q., Zhu, X., Ling, Z.-H., Inkpen, D., and Wei, S. Neural natural language inference models enhanced with external knowledge. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2406–2417, Melbourne, Australia, July 2018. Association for Computational Linguistics.
- Cho, H. and Choi, I. S. Three-dimensionally embedded graph convolutional network (3DGCN) for molecule interpretation. *CoRR*, abs/1811.09794, 2018.
- Clark, K., Khandelwal, U., Levy, O., and Manning, C. D. What does bert look at? an analysis of bert’s attention. *arXiv preprint arXiv:1906.04341*, 2019.
- Coley, C. W., Barzilay, R., Green, W. H., Jaakkola, T. S., and Jensen, K. F. Convolutional embedding of attributed molecular graphs for physical property prediction. *Journal of Chemical Information and Modeling*, 57(8):1757–1772, 2017. doi: 10.1021/acs.jcim.6b00601. PMID: 28696688.

- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding, 2018. cite arxiv:1810.04805Comment: 13 pages.
- Duvenaud, D. K., Maclaurin, D., Iparraguirre, J., Bombarell, R., Hirzel, T., Aspuru-Guzik, A., and Adams, R. P. Convolutional networks on graphs for learning molecular fingerprints. In Cortes, C., Lawrence, N. D., Lee, D. D., Sugiyama, M., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 28*, pp. 2224–2232. Curran Associates, Inc., 2015.
- Feinberg, E. N., Sur, D., Wu, Z., Husic, B. E., Mai, H., Li, Y., Sun, S., Yang, J., Ramsundar, B., and Pande, V. S. Potentialnet for molecular property prediction. *ACS Central Science*, 4(11):1520–1530, 2018. doi: 10.1021/acscentsci.8b00507.
- Gaulton, A., Bellis, L. J., Bento, A. P., Chambers, J., Davies, M., Hersey, A., Light, Y., McGlinchey, S., Michalovich, D., Al-Lazikani, B., and Overington, J. P. ChEMBL: a large-scale bioactivity database for drug discovery. *Nucleic Acids Research*, 40(D1):D1100–D1107, 09 2011. ISSN 0305-1048. doi: 10.1093/nar/gkr777.
- Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and Dahl, G. E. Neural message passing for quantum chemistry. In Precup, D. and Teh, Y. W. (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1263–1272, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR.
- Goh, G., Siegel, C., Vishnu, A., Hodas, N., and Baker, N. Chemception: A deep neural network with minimal chemistry knowledge matches the performance of expert-developed qsar/qspr models. *arXiv preprint arXiv:1706.06689*, 06 2017.
- Gómez-Bombarelli, R., Wei, J. N., Duvenaud, D., Hernández-Lobato, J., Sánchez-Lengeling, B., Sheberla, D., Aguilera-Iparraguirre, J., Hirzel, T. D., Adams, R. P., and Aspuru-Guzik, A. Automatic chemical design using a data-driven continuous representation of molecules. *ACS Central Science*, 4(2):268–276, 02 2018. doi: 10.1021/acscentsci.7b00572.
- Goodfellow, I., Bengio, Y., and Courville, A. *Deep Learning*. MIT Press, 2016.
- Guo, M., Zhang, Y., and Liu, T. Gaussian transformer: a lightweight approach for natural language inference. In *AAAI 2019*, 2019.
- Haghighatlari, M. and Hachmann, J. Advances of machine learning in molecular modeling and simulation. *Current Opinion in Chemical Engineering*, 23:51 – 57, 2019.
- ISSN 2211-3398. doi: <https://doi.org/10.1016/j.coche.2019.02.009>. *Frontiers of Chemical Engineering: Molecular Modeling*.
- Honda, S., Shi, S., and Ueda, H. R. Smiles transformer: Pre-trained molecular fingerprint for low data drug discovery, 2019.
- Htut, P. M., Phang, J., Bordia, S., and Bowman, S. R. Do attention heads in bert track syntactic dependencies?, 2019.
- Hu, W., Liu, B., Gomes, J., Zitnik, M., Liang, P., Pande, V. S., and Leskovec, J. Pre-training graph neural networks. *CoRR*, abs/1905.12265, 2019.
- Ishiguro, K., Maeda, S., and Koyama, M. Graph warp module: an auxiliary module for boosting the power of graph neural networks. *CoRR*, abs/1902.01020, 2019.
- Jastrzębski, S., Leśniak, D., and Czarnecki, W. M. Learning to smile (s). *arXiv preprint arXiv:1602.06289*, 2016.
- Jr, J. F. R., Florea, L., de Oliveira, M. C. F., Diamond, D., and Jr, O. N. O. A survey on big data and machine learning for chemistry, 2019.
- Karpov, P., Godin, G., and Tetko, I. V. A transformer model for retrosynthesis. In *International Conference on Artificial Neural Networks*, pp. 817–830. Springer, 2019.
- Kearnes, S., McCloskey, K., Berndl, M., Pande, V., and Riley, P. Molecular graph convolutions: Moving beyond fingerprints. *Journal of Computer-Aided Molecular Design*, 30, 03 2016. doi: 10.1007/s10822-016-9938-8.
- Kim, S., Chen, J., Cheng, T., Gindulyte, A., He, J., He, S., Li, Q., Shoemaker, B. A., Thiessen, P. A., Yu, B., et al. Pubchem 2019 update: improved access to chemical data. *Nucleic acids research*, 47(D1):D1102–D1109, 2018.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Korotcov, A., Tkachenko, V., Russo, D. P., and Ekins, S. Comparison of deep learning with multiple machine learning methods and metrics using diverse drug discovery data sets. *Molecular Pharmaceutics*, 14(12):4462–4475, 12 2017. doi: 10.1021/acs.molpharmaceut.7b00578.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In Pereira, F., Burges, C. J. C., Bottou, L., and Weinberger, K. Q. (eds.), *Advances in Neural Information Processing Systems 25*, pp. 1097–1105. Curran Associates, Inc., 2012.
- Landrum, G. Rdkit: Open-source cheminformatics software. 2016.

- Li, J., Cai, D., and He, X. Learning graph-level representation for drug discovery. *arXiv preprint arXiv:1709.03741*, 2017.
- Li, R., Wang, S., Zhu, F., and Huang, J. Adaptive graph convolutional neural networks. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Lipinski, C. A., Lombardo, F., Dominy, B. W., and Feeney, P. J. Experimental and computational approaches to estimate solubility and permeability in drug discovery and development settings. *Advanced Drug Delivery Reviews*, 23(1):3–25, 1997. ISSN 0169-409X. doi: [https://doi.org/10.1016/S0169-409X\(96\)00423-1](https://doi.org/10.1016/S0169-409X(96)00423-1). In Vitro Models for Selection of Development Candidates.
- Mayr, A., Klambauer, G., Unterthiner, T., Steijaert, M., Wegner, J. K., Ceulemans, H., Clevert, D.-A., and Hochreiter, S. Large-scale comparison of machine learning methods for drug target prediction on chembl. *Chem. Sci.*, 9: 5441–5451, 2018. doi: 10.1039/C8SC00148K.
- Podlewska, S. and Kafel, R. Metstabon—online platform for metabolic stability predictions. *International journal of molecular sciences*, 19(4):1040, 2018.
- Ramsundar, B., Eastman, P., Walters, P., Pande, V., Leswing, K., and Wu, Z. *Deep Learning for the Life Sciences*. O’Reilly Media, 2019.
- Rogers, D. and Hahn, M. Extended-connectivity fingerprints. *Journal of chemical information and modeling*, 50(5):742–754, 2010.
- Schütt, K. T., Arbabzadah, F., Chmiela, S., Müller, K. R., and Tkatchenko, A. Quantum-chemical insights from deep tensor neural networks. *Nature Communications*, 8: 13890, Jan 2017. doi: 10.1038/ncomms13890.
- Segler, M., Kogej, T., Tyrchan, C., and Waller, M. Generating focused molecule libraries for drug discovery with recurrent neural networks. *ACS Central Science*, 4, 01 2017. doi: 10.1021/acscentsci.7b00512.
- Shang, C., Liu, Q., Chen, K.-S., Sun, J., Lu, J., Yi, J., and Bi, J. Edge Attention-based Multi-Relational Graph Convolutional Networks. *arXiv e-prints*, art. arXiv:1802.04944, Feb 2018.
- Shaw, P., Uszkoreit, J., and Vaswani, A. Self-attention with relative position representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pp. 464–468, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2074.
- Shin, B., Park, S., Kang, K., and Ho, J. C. Self-attention based molecule representation for predicting drug-target interaction. In Doshi-Velez, F., Fackler, J., Jung, K., Kale, D., Ranganath, R., Wallace, B., and Wiens, J. (eds.), *Proceedings of the 4th Machine Learning for Healthcare Conference*, volume 106 of *Proceedings of Machine Learning Research*, pp. 230–248, Ann Arbor, Michigan, 09–10 Aug 2019. PMLR.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. *CoRR*, abs/1706.03762, 2017.
- Veličković, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., and Bengio, Y. Graph Attention Networks. *arXiv e-prints*, art. arXiv:1710.10903, Oct 2017.
- Wallach, I., Dzamba, M., and Heifets, A. Atomnet: A deep convolutional neural network for bioactivity prediction in structure-based drug discovery. *ArXiv*, abs/1510.02855, 2015.
- Wang, S., Guo, Y., Wang, Y., Sun, H., and Huang, J. Smilesbert: Large scale unsupervised pre-training for molecular property prediction. In *Proceedings of the 10th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics, BCB ’19*, pp. 429–436, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450366663. doi: 10.1145/3307339.3342186.
- Weininger, D. Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules. *Journal of chemical information and computer sciences*, 28(1):31–36, 1988.
- Withnall, M., Lindelöf, E., Engkvist, O., and Chen, H. Building attention and edge convolution neural networks for bioactivity and physical-chemical property prediction, Sep 2019.
- Wong, C. H., Siah, K. W., and Lo, A. W. Estimation of clinical trial success rates and related parameters. *Biostatistics*, 20(2):273–286, 01 2018. ISSN 1465-4644. doi: 10.1093/biostatistics/kxx069.
- Wu, Z., Ramsundar, B., Feinberg, E. N., Gomes, J., Geniesse, C., Pappu, A. S., Leswing, K., and Pande, V. Moleculenet: a benchmark for molecular machine learning. *Chem. Sci.*, 9:513–530, 2018. doi: 10.1039/C7SC02664A.
- Yang, K., Swanson, K., Jin, W., Coley, C., Eiden, P., Gao, H., Guzman-Perez, A., Hopper, T., Kelley, B., Mathea, M., et al. Analyzing learned molecular representations for property prediction. *Journal of chemical information and modeling*, 59(8):3370–3388, 2019.

A. Dataset details.

We include below a more detailed description of the datasets used in our benchmark.

- **FreeSolv, ESOL.** Regression tasks. Popular tasks for predicting water solubility in terms of the hydration free energy (FreeSolv) and logS (ESOL). Solubility of molecules is an important property that influences the bioavailability of drugs.
- **Blood-brain barrier permeability (BBBP).** Binary classification task. The blood-brain barrier (BBB) separates the central nervous system from the bloodstream. Predicting BBB penetration is especially relevant in drug design when the goal for the molecule is either to reach the central nervous system or the contrary – not to affect the brain.
- **MetStab_{high}, MetStab_{low}.** Binary classification tasks. The metabolic stability of a compound is a measure of the half-life time of the compound within an organism. The compounds for this task were taken from (Podlewska & Kafel, 2018), where compounds were divided into three sets: high, medium, and low stability. In this paper we concatenated these sets in order to build two classification tasks: MetStab_{high} (discriminating high against others) and MetStab_{low} (discriminating low against others).
- **Estrogen Alpha, Estrogen Beta.** Binary classification tasks. Often in drug discovery, it is important that a molecule is not potent towards a given target. Modulating of the estrogen receptors changes the genomic expression throughout the body, which in turn may lead to the development of cancer. For these tasks, the compounds with known activities towards the receptors were extracted from ChEMBL (Gaulton et al., 2011) database and divided into active and inactive sets based on their reported inhibition constant (Ki), being < 100 nM and > 1000 nM, respectively.

B. Other experimental details

In this section we include details for hyperparameters and training settings used in Section 4.2.

Molecule Attention Transformer. Table 7 shows hyperparameter ranges used in experiments for MAT. A short description of these hyperparameters is listed below:

- MODEL DIM – size of embedded atom features,
- LAYERS NUMBER – number of encoder module repeats (N in Figure 1),

- ATTENTION HEADS NUMBER – number of molecule self-attention heads,
- PFFS NUMBER – number of dense layers in the position-wise feed forward block (K in Figure 1),
- λ_{att} – self-attention weight λ_{att} ,
- λ_{dist} – distance weight λ_d ,
- DISTANCE MATRIX KERNEL – function g used to transform the distance matrix **D**,
- MODEL DROPOUT – dropout applied after the embedding layer, position-wise feed forward layers, and residual layers (before sum operation),
- WEIGHT DECAY – optimizer weight decay,
- LEARNING RATE – (see Equation 3)
- EPOCHS NUMBER – number of epochs for which the model is trained
- BATCH SIZE – batch size used during the training of the model
- WARMUP FACTOR – fraction of epochs after which we end with increasing the learning rate linearly and begin with decreasing it proportionally to the inverse square root of the step number. (see Equation 3)

Table 7. Molecule Attention Transformer hyperparameters ranges

	PARAMETERS
BATCH SIZE	8, 16, 32, 64, 128
LEARNING RATE	.01, .005, .001, .0005, .0001
EPOCHS	30, 100
MODEL DIM	32, 64, 128, 256, 512, 1024
LAYERS NUMBER	1, 2, 4, 6, 8
ATTENTION HEADS NUMBER	1, 2, 4, 8, 16
PFFS NUMBER	1
λ_{att}	0, .1, .2, .3, .4, .5, .6, .7, .8, .9, 1
$\lambda_{distance}$	0, .1, .2, .3, .4, .5, .6, .7, .8, .9, 1
DISTANCE MATRIX KERNEL	'SOFTMAX', 'EXP'
MODEL DROPOUT	.0, .1, .2
WEIGHT DECAY	.0, .00001, .0001, .001, .01
WARMUP FACTOR	.0, .1, .2, .3, .4, .5

As suggested in (Vaswani et al., 2017), for optimization of MAT we used Adam optimizer (Kingma & Ba, 2014), with learning rate scheduler given by the following formula:

$$Step_{LR} = optimizer\ factor \cdot model\ dim^{-0.5} \cdot \min(step\ num^{-0.5}, step\ num \cdot warmup\ steps^{-0.5}). \quad (3)$$

Where *optimizer factor* is given by $100 \cdot \text{LEARNING RATE}$ and *warmup steps* is given by $\text{WARMUP FACTOR} \cdot \text{total train steps number}$.

After N layers embedding of the molecule is calculated by taking the mean of returned by the network vector representations of all atoms (*Global pooling* in Figure 1). Then it is passed to the single linear layer, which returns the prediction.

SVM, RF, GCN, Weave. In our experiments, we used DeepChem (Ramsundar et al., 2019) implementation of baseline algorithms (SVM, RF, GCN, Weave). We used the same hyperparameters for tuning as were used in DeepChem, having regard to their proposed default values (we list them in Tables 8 - 11).

RF and SVM work on the vector representation of molecule given by the Extended-connectivity fingerprints (Rogers & Hahn, 2010). ECFP vectors were calculated using class CIRCULARFINGERPRINT from the DeepChem package, with default parameters (RADIUS=2, SIZE=2048).

Table 8. SVM hyperparameter ranges

	PARAMETERS
C	.25, .4375, .625, .8125, 1., 1.1875, 1.375, 1.5625, 1.75, 1.9375, 2.125, 2.3125, 2.5, 2.6875, 2.875, 3.0625, 3.25, 3.4375, 3.625, 3.8125, 4.
GAMMA	.0125, .021875, .03125, .040625, .05, .059375, .06875, .078125, .0875, .096875, .10625, .115625, .125, .134375, .14375, .153125, .1625, .171875, .18125, .190625, .2

Table 9. RF hyperparameter ranges

	PARAMETERS
N ESTIMATORS	125, 218, 312, 406, 500, 593, 687, 781, 875, 968, 1062, 1156, 1250, 1343, 1437, 1531, 1625, 1718, 1812, 1906, 2000

Table 10. GCN hyperparameter ranges

	PARAMETERS
BATCH SIZE	64, 128, 256
LEARNING RATE	0.002, 0.001, 0.0005
N FILTERS	64, 128, 192, 256
N FULLY CONNECTED NODES	128, 256, 512

EAGCN Table 12 shows hyperparameter ranges used in experiments for EAGCN. For EAGCN with *weighted* struc-

Table 11. Weave hyperparameter ranges

	PARAMETERS
BATCH SIZE	16, 32, 64, 128
NB EPOCH	20, 40, 60, 80, 100
LEARNING RATE	0.002, 0.001, 0.00075, 0.0005
N GRAPH FEAT	32, 64, 96, 128, 256
N PAIR FEAT	14

ture number of convolutional features $n_sgc = n_sgc_1 + n_sgc_2 + n_sgc_3 + n_sgc_4 + n_sgc_5$.

Table 12. EAGCN hyperparameter ranges

	PARAMETERS
BATCH SIZE	16, 32, 64, 128, 256, 512
EAGCN STRUCTURE	'CONCATE', 'WEIGHTED'
NUM EPOCHS	30, 100
LEARNING RATE	.01, .005, .001, .0005, .0001
DROPOUT	.0, .1, .3
WEIGHT DECAY	.0, .001, .01, .0001
N CONV LAYERS	1, 2, 4, 6
N DENSE LAYERS	1, 2, 3, 4
N SGC 1	30, 60
N SGC 2	5, 10, 15, 20, 30
N SGC 3	5, 10, 15, 20, 30
N SGC 4	5, 10, 15, 20, 30
N SGC 5	5, 10, 15, 20, 30
DENSE DIM	16, 32, 64, 128

C. Additional results for Sec. 4.2

Predicting internal energy We run an additional experiment on a regression task related to quantum mechanics. From the Alchemy dataset (Chen et al., 2019), which is a dataset of 12 quantum properties calculated for 200K molecules, we have chosen internal energy at 298.15 K to further test the performance of our model. We hypothesize that our molecule self-attention should perform particularly well in tasks involving atom level interactions such as energy prediction.

Table 13 presents mean absolute errors of three methods: one classical method (RF), one graph method (GCN), and our pretrained MAT. We use original train/valid/test splits of the dataset. For RF and GCN we run a random search with budget of 500 hyperparameter sets. For pretrained MAT, we tune only the learning rate, that is selected from $\{1e-3, 5e-4, 1e-4, 5e-5, 1e-5, 5e-6, 1e-6\}$.

MAT achieves a slightly lower error than GCN. As can be expected, both graph methods can learn internal energy function correctly because of the locality preserved in the graph structure. The classical method based on fingerprints gives MAE that is almost two orders of magnitude higher than MAE of the other methods in the comparison.

Table 13. Test results for internal energy prediction reported as MAE. All methods were tuned with a random search with budget of 500 hyperparameter combinations.

U (INTERNAL ENERGY)	
RF	.380
GCN	.006
MAT	.004

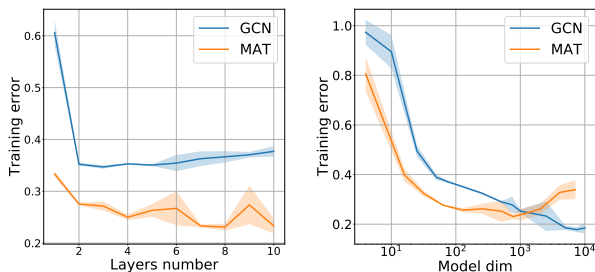


Figure 7. Training loss of MAT and GCN as a function of the number of layers (left) and model dimensionality (right).

Training error for graph-based neural networks Ishiguro et al. (2019) show that graph neural networks suffer from underfitting of the training set and their performance does not scale well with the complexity of the network. We reproduce their experiments and confirm that this problem is indeed present for both GCN and MAT. According to Figure 7, the training loss of GCN and MAT flattens at some point and stops decreasing even if we keep increasing the number of layers and model dimensionality. Despite this issue, for almost all settings, MAT achieves lower training error than GCN.

D. Additional details for Sec. 4.3

Task description. As a node-level pretraining task we chose masking from (Hu et al., 2019) which is a version of BERT masked language model adapted to graph structured data. The idea is that predicting masked nodes based on their neighbourhood will encourage model to capture domain specific relationships between atoms.

For each molecular graph we randomly replace 15% of input nodes (atom attributes) with special mask token. After forward pass we apply linear model to corresponding node embeddings to predict masked node attributes. In case of EAGCN we additionally mask attributes of edges connected to masked nodes to prevent model from learning simple value copying.

Pretraining setting. Training dataset consisted of 2 mln molecules sampled from the ZINC15 database. Models

were trained for 8 epochs with learning rate set to 0.001 and batch size 256. MAT was optimized with Noam optimizer (described in App. B), whereas for EAGCN we used Adam (Kingma & Ba, 2014). In both cases procedure minimized binary cross entropy loss.

Fine-tuning setting. All our pretrained models are fine-tuned on the target tasks for 100 epochs, with batch size equal to 32 and learning rate selected from the set of $\{1e-3, 5e-4, 1e-4, 5e-5, 1e-5, 5e-6, 1e-6\}$.

In Estrogen Alpha experiments we excluded three molecules (with the highest number of atoms) from the dataset, due to the memory issues.

Table 14. Pretrained MAT hyperparameters

PARAMETERS	
MODEL DIM	1024
LAYERS NUMBER	8
ATTENTION HEADS NUMBER	16
PFFS NUMBER	1
λ_{att}	.33
$\lambda_{distance}$.33
DISTANCE MATRIX KERNEL	'EXP'
MODEL DROPOUT	.0
WEIGHT DECAY	.0

Table 15. Pretrained EAGCN hyperparameters

PARAMETERS	
EAGCN STRUCTURE	'WEIGHTED'
DROPOUT	.0
WEIGHT DECAY	.0
N CONV LAYERS	8
N DENSE LAYERS	1
N SGC	1080

SMILES Transformer. We used pretrained weights of SMILES-Transformers conducted by Honda et al. (2019). In this setting, according to the authors, we used MLP with 1 hidden layer, with 100 hidden units, that works on the 1024-dimensional molecule embedding returned by the pretrained transformer. We trained this MLP on the target tasks for 100 epochs, with batch size equal to 32 and learning rate selected from the set of $\{1e-3, 5e-4, 1e-4, 5e-5, 1e-5, 5e-6, 1e-6\}$.

E. Additional results for Sec. 4.4

Edge features. Every bond in the molecule was embedded by the vector of edge features (we used features similar to described in (Shang et al., 2018)). Every edge feature was then passed through linear layer, followed by ReLU activation, which returned one single value for every

single edge (if there is no edge between atoms, we pass zero vector through the layer). This results in the matrix $\mathbf{E} \in \mathbb{R}^{N_{\text{atoms}} \times N_{\text{atoms}}}$ which was then used in Molecule Self-Attention layer, instead of the adjacency matrix.

Table 16. Edge Features used for experiments form Table 6

ATTRIBUTE	DESCRIPTION
BOND ORDER	VALUES FROM SET { 1, 1.5, 2, 3 }
AROMATICITY	IS AROMATIC
CONJUGATION	IS CONJUGATED
RING STATUS	IS IN A RING

F. Toy task

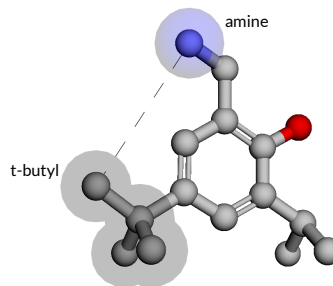
Task description. The essential feature of Molecule Attention Transformer is that it augments the self-attention module using molecule structure. Here we investigate MAT on a task heavily reliant on distances between atoms; we are primarily interested in how the performance of MAT depends on λ_a , λ_d , λ_g that are used to weight the adjacency and the distance matrices in Equation 2.

Naturally, many properties of molecules depend on their geometry. For instance, *steric effect* happens when a spatial proximity of a given group, blocks reaction from happening, due to an overlap in electronic groups. However, this type of reasoning can be difficult to learn based only on the graph information, as it does not always reflect the geometry well. Furthermore, focusing on distance information might require selecting low values for either λ_g or λ_a (see Figure 1).

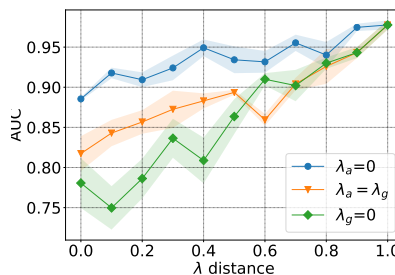
To illustrate this, we designed a toy task to predict whether or not two substructures are closer to each other in space than a predefined threshold; see also Figure 8a. We expect that MAT will work significantly better than a vanilla graph convolutional network if λ_d is tuned well.

Experimental setting. We construct the dataset by sampling 2677 molecules from PubChem (Kim et al., 2018), and use 20 Å threshold between -NH₂ fragment and *tert*-butyl group to determine the binary label. The threshold was selected so that positive and negative examples are well balanced.

Results. First, we plot Molecule Attention Transformer performance as a function of λ_d in Figure 8b for three settings of λ : $\lambda_a = 0$ (blue), $\lambda_a = \lambda_g$ (orange), and $\lambda_g = 0$ (green). In all cases we find that using distance information improves the performance significantly. Additionally, we found that GCN achieves 0.93 AUC on this task, compared to 0.98 by MAT with $\lambda_d = 1.0$. These results both motivate tuning λ , and show that MAT can efficiently use distance information if it is important for the task at hand.



(a) The toy task is to predict whether two substructures (-NH₂ fragment and *tert*-butyl group) co-occur within given distance.



(b) Molecule Attention Transformer performance on the toy task as a function of λ_d , for different settings of λ_g and λ_a .

Figure 8. MAT can efficiently use the inter-atomic distances to solve the toy task (see left). Additionally, the performance is heavily dependent on λ_d , which motivates tuning λ in the main experiments (see right).

Further details. The molecules in the toy task dataset were downloaded from PubChem. The SMARTS query used to find the compounds was C([C;H3])([C;H3])([C;H3]).[NX3H2]. All molecules were then filtered so that only those with exactly one *tert*-butyl group and one -NH₂ fragment were left. For each of them, five conformers were created with RDKit implementation of the Universal Force Field (UFF).

The task is a binary classification of the distance between two molecule fragments. If the euclidean distance between -NH₂ fragment and *tert*-butyl group is greater than a given threshold, the label is 1 (0 otherwise). As the distance we mean the distance between the closest heavy atoms in these two fragments across calculated conformers. We used 20 Å as the threshold as it leads to a balanced dataset. There are 2677 compounds total from which 1140 are in a positive class. The dataset was randomly split into training, validation, and test datasets.

In experiments the hyperparameters that yielded promising results on our datasets were used (listed in Table 17). The values of λ parameters were tuned, and their scores are shown in Figure 8b. All three λ parameters (λ_d , λ_g , λ_a) sum to 1 in all experiments.

To compare our results with a standard graph convolutional neural network, we run a grid search over hyperparameters shown in Table 18. The hyperparameters for which the best validation AUC score was reached are emboldened, and their test AUC score is 0.925 ± 0.006 .

Table 17. MAT hyperparameters used.

	PARAMETERS
BATCH SIZE	16
LEARNING RATE	0.0005
EPOCHS	100
MODEL DIM	64
MODEL N	4
MODEL H	8
MODEL N DENSE	2
MODEL DENSE OUTPUT NONLINEARITY	'TANH'
DISTANCE MATRIX KERNEL	'SOFTMAX'
MODEL DROPOUT	0.0
WEIGHT DECAY	0.001
OPTIMIZER	'ADAM_ANNEAL'
AGGREGATION TYPE	'MEAN'

Table 18. Hyperparameters used for tuning GCN.

	PARAMETERS
BATCH SIZE	16 , 32, 64
LEARNING RATE	0.0005
EPOCHS	20, 40, 60 , 80, 100
N FILTERS	64, 128
N FULLY CONNECTED NODES	128, 256

G. Interpretability analysis

Table 19. Statistics of the six attention head patterns described in the text. Each head function is defined by a SMARTS that selects atoms with high attention weights. For each atom in the dataset we calculated mean weight assigned to them by the corresponding attention head (average column value of the attention matrix). Calculated means and standard deviations show the difference between attention weights of matching atoms (μ^+ , σ^+) against the other atoms (μ^- , σ^-).

HEAD	I	II	III	IV	V	VI
SMARTS	[c:D2]	[S,S]	[N;R0]	O=*	[A:D3]	N
μ^+	.136	.330	.061	.095	.043	.228
σ^+	.080	.280	.074	.120	.032	.171
μ^-	.008	.001	.002	.006	.006	.005
σ^-	.032	.003	.016	.034	.014	.009

We found several patterns in the self-attention heads by looking at the first layer of MAT. These patterns correspond to chemical structures that can be found in molecules. For

each such pattern found in a qualitative manner, we tested quantitatively if our hypotheses are true about what these particular attention heads represent.

For each pattern found in one of the attention heads, we construct a SMARTS expression describing atoms that belong to our hypothetical molecular structures. Then, all atoms matching the pattern are extracted from the BBBP dataset, and their mean attention weights (average column value of the attention matrix) are compared against atoms that do not match the pattern. Table 19 shows the distributions of attention weights for matching and not matching atoms. Atoms which match the SMARTS expression have significantly higher attention weights ($\mu^+ > \mu^-$).