

# MassFormer: Tandem Mass Spectrum Prediction with Graph Transformers

Adamo Young<sup>1,4</sup>, Bo Wang<sup>1,2,4,5\*</sup>, and Hannes Röst<sup>1,3\*</sup>

<sup>1</sup> Department of Computer Science, University of Toronto

<sup>2</sup> Department of Laboratory Medicine and Pathobiology, University of Toronto

<sup>3</sup> Department of Molecular Genetics, University of Toronto

<sup>4</sup> Vector Institute

<sup>5</sup> Peter Munk Centre, University Health Network

\* co-senior authors

**Abstract.** Mass spectrometry is a key tool in the study of small molecules, playing an important role in metabolomics, drug discovery, and environmental chemistry. Tandem mass spectra capture fragmentation patterns that provide key structural information about a molecule and help with its identification. Practitioners often rely on spectral library searches to match unknown spectra with known compounds. However, such search-based methods are limited by availability of reference experimental data. In this work we show that graph transformers can be used to accurately predict tandem mass spectra. Our model, MassFormer, outperforms competing deep learning approaches for spectrum prediction, and includes an interpretable attention mechanism to help explain predictions. We demonstrate that our model can be used to improve reference library coverage on a synthetic molecule identification task. Through quantitative analysis and visual inspection, we verify that our model recovers prior knowledge about the effect of collision energy on the generated spectrum. We evaluate our model on different types of mass spectra from two independent MS datasets and show that its performance generalizes. Code available at [github.com/Roestlab/massformer](https://github.com/Roestlab/massformer).

**Keywords:** Mass Spectrometry · Metabolomics · Deep Learning · Chemistry · Transformers

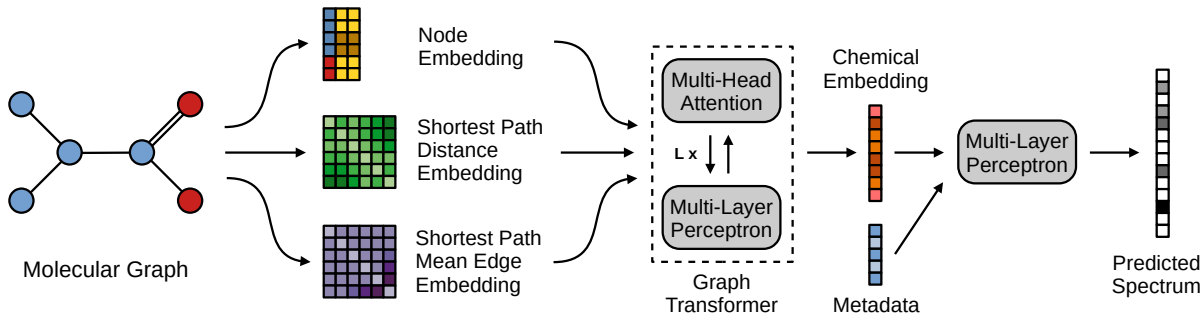
## 1 Introduction

Mass spectrometry (MS, [1]) is an analytical technique used for identifying and quantifying chemicals in a mixture. Molecules from the sample are ionized and then detected by a mass analyzer, which records information about the mass to charge ratio ( $m/z$ ) of each ion in the form of a mass spectrum. Tandem mass spectrometry (MS/MS) is a variant of MS that includes a collision step to isolate and break down charged molecules (called precursors) into smaller fragments. These charged ions subsequently appear as fragment peaks in the spectrum, and their  $m/z$  position and relative abundance can be used to make inferences about the molecular structure of the original precursor. When coupled with online liquid chromatography (LC), a technique for chemical separation, the LC-MS/MS workflow is a powerful tool for analyzing biological mixtures and is commonly employed in proteomics [2] and metabolomics [3,4] experiments. This process is also used in environmental chemistry to detect contaminants and other compounds of interest [5]. Some examples of MS/MS spectra are displayed in Figure 2.

The configuration of a mass spectrometer can have a large impact on the appearance of the generated spectra. In MS/MS, the collision energy setting determines how much fragmentation will occur: increasing collision energy encourages more fragments with lower average mass, and vice versa. MS/MS spectra are often collected at a variety of collision energies to accommodate molecules of different sizes and maximize the chance of identifying combinations of fragments that reveal useful information about the precursor’s structure. Spectra produced by different collision methods, such as collision-induced dissociation (CID) and high-energy C-trap dissociation (HCD), can result in different fragmentation patterns and are not always directly comparable. The formation of precursor adducts is also an important consideration. During the ionization step, a charged species (atom or molecule) becomes associated with the precursor molecule, changing its mass and affecting the subsequent fragmentation process. As a result, the same precursor

can produce different spectra depending on the adduct. In summary, there are a number of important experimental parameters that affect the appearance of a compound’s MS/MS spectrum. Therefore, having access to a spectrum’s “metadata” is critical for properly interpreting the information that it provides.

It is generally quite difficult to model the fragmentation process for a given precursor molecule from first principles. Expert-derived rules can be useful for certain classes of compounds with regular structure (like lipids and peptides), but are difficult to derive for smaller compounds with more structural diversity. The increasing availability of large public [6,7] and commercial [8,9,10] MS datasets makes data-driven solutions more appealing. In this work we develop a state-of-the-art molecule transformer to tackle MS/MS spectrum prediction. Our work is the first to apply transformers to this problem. We rigorously compare our approach with strong baseline models and evaluate their ability to generalize across independently collected MS/MS datasets. The model code is public and can be found [here](#).



**Fig. 1.** Overview of the method: extraction of node and topological embeddings from the molecular graph, application of the  $L$ -layer graph transformer, extraction of the chemical embedding from the readout node, concatenation of spectral metadata, and prediction of spectrum vector.

## 2 Related Work

### 2.1 Compound Identification

Compound identification from MS/MS data is a challenging problem. Most practitioners rely on database searches with large reference libraries, using spectrum similarity functions [11] or domain expertise to identify matches. However, these databases have relatively poor coverage, containing on the order of  $10^5$  spectra which only cover around  $10^4$  unique compounds. Untargeted metabolomics experiment indicate that over 95% of the human metabolome is unknown [12]. More generally, the space of small molecules is huge: enumeration analyses estimate that there are over 166 billion organic compounds with fewer than 18 heavy atoms [13]. When it comes to *de novo* identification of novel compounds from MS data, spectral searches are fundamentally limited.

The first algorithms for MS-based compound identification arose in the 1960s. DENDRAL and Meta-DENDRAL [14] were early AI projects aimed at automatically learning fragmentation rules from MS data, to help identify new molecules. Since then, algorithmic advances and increased data availability have resulted in more powerful methods. Many modern approaches rely on fragmentation trees [15,16] to represent spectra. Such tree representations are useful because they provide an interpretable description of the fragmentation process, and can be inferred directly from the peak information in the spectrum. Tree kernels can be used to compute similarity between spectra [17]. This allows for the development of kernel algorithms that infer compound identity by predicting properties from the fragmentation tree directly and searching chemical databases to find candidates that match those properties [17,18]. There are also methods for predicting substructure directly from the spectrum, without relying on fragmentation trees. [19] and [20] use neural networks to predict substructure labels and find candidate molecules that match those substructures.

## 2.2 Spectrum Prediction

A straightforward strategy to overcome the size limitation of spectral libraries is to augment them with simulated (*in silico*) mass spectra for a large number of compounds from a chemical structure database [21,22,23]. This can dramatically increase coverage for the reference library and improve the chance of finding a match, motivating the development of accurate spectrum prediction models. Such models come in a variety of forms. Rule-based fragmentation algorithms [24,25] are simple to implement and can be effective, but do not generalize to many types of compounds. Other approaches rely on combinatorial methods to identify possible fragments and their associated peak locations [26,27], but do not provide information about relative intensity. These methods can be improved by incorporating a probabilistic model to predict the relative intensity of each of the combinatorially generated fragments [28,25]. While effective, such models suffer from slow training and inference algorithms that make scaling to larger datasets more challenging. More recently, deep learning approaches [29,30] for spectrum prediction have been introduced, effectively trading model interpretability for speed and scalability.

## 2.3 Transformers

Transformers [31] are a family of neural networks characterized by their use of attention to model sequences. Originally developed for neural machine translation [32], transformer models have proven useful in a number of domains, from computer vision [33] to reinforcement learning [34], achieving state of the art performance even in problems that do not naturally lend themselves to sequence modelling. They represent an input sequence as a set of embeddings, each of which capture the semantic meaning and position of a single element. By interleaving layers of multi-head attention with small multi-layer perceptrons (MLPs), transformers iteratively process the set of embeddings and learn relationships between elements of the sequence.

A number of graph transformers have been proposed [35,36,37,38], motivated by the ability to model pairwise global interactions between all nodes in the graph. In contrast, GNNs can only model local relations in a single layer, and require large depth to model longer distance interactions. Graph transformers have generally proven useful when the input graph is small enough such that the quadratic memory footprint of the attention mechanism does not become prohibitively expensive.

# 3 Methods

## 3.1 Problem Formulation

Spectrum prediction can be viewed as a supervised learning problem, with a dataset  $\{x^i, z^i, y^i\}_{i=1}^n$  where  $x^i$  is a molecule and  $y^i$  is its spectrum under experimental conditions  $z^i$ . The goal is to learn the parameters  $\theta$  of the prediction function  $f_\theta : \mathcal{X} \times \mathcal{Z} \rightarrow \mathcal{Y}$ , where  $\mathcal{X}$  is the space of chemicals and  $\mathcal{Y}$  is the spectrum space. Mass spectra can be represented as a set of peaks, each of which have an  $m/z$  location and an intensity. By discretizing the peak locations into  $m$  fixed-width bins (as was done in [29,30]), a mass spectrum can be represented as an  $m$ -dimensional sparse vector, where each peak at location  $j$  has intensity  $y_j \geq 0$ . The problem of spectrum prediction can thus be formulated as vector regression, with  $\mathcal{Y} = \mathbb{R}^m \succeq 0$ . The spectral metadata  $z \in \mathcal{Z}$  (such as collision energy and precursor adduct) is provided as side information to the input molecule  $x$ .

## 3.2 Chemical Featurization

The choice of representation for the input molecule  $x$  is critical, as it affects the structure of the prediction function  $f_\theta$  and can have an impact on downstream performance. Three common strategies for representing molecules are fingerprints, graphs, strings. Molecular fingerprints [39] are a standard featurization that represent molecules as a vector of local graph substructure counts. Molecular graph representations capture the 2D structure of a molecule by explicitly representing atoms as nodes and bonds as edges. The node features encode various chemical properties associated with the atom (i.e. element, formal charge, number of bonded hydrogens), while the edge features encode bond information (i.e. bond type, aromaticity). Such

representations naturally lend themselves to graph neural networks (GNNs, [40]). Finally, string representations, such as SMILES [41] and SELFIES [42], encode molecules as sequences of discrete tokens, where each token describes an atom or a bond structure. This process assumes an ordering of the atoms in the molecule, often leading to multiple different string representations that describe the same underlying chemical. The benefits of using strings are compactness (relative to graphs) and compatibility with an array of sequence models, including convolutional neural networks (CNNs, [43]), recurrent neural networks (RNNs, [44]), and transformers [31].

### 3.3 Chemical Graph Transformer

A visual summary of our method is presented in Figure 1. First the embeddings (corresponding to nodes, shortest path distances, and shortest path edges) are extracted from the molecular graph. These features are passed to a graph transformer with  $L$  layers, which outputs a chemical embedding. This vector is concatenated with a vector representing the spectrum metadata and passed to an MLP, which produces a sparse positive vector that can be interpreted as the predicted mass spectrum.

Graphormer [45] is a recent graph transformer model that boasts impressive results on chemical property prediction tasks [46,45], particularly in the low-data regime (on the order of  $10^4$  samples). We adapt the Graphormer architecture as the basis for our model. The distinguishing character of this graph transformer is its unique positional encoding scheme. The model uses shortest path information between nodes, and associated edge embeddings along that path, as a form of relative positional encoding. The shortest path information is computed as a preprocessing step for each graph, using the Floyd-Warshall algorithm [47].

The attention mechanism  $a_{ij}$  is described in detail in Equations 1 and 2, where  $h_i, h_j \in \mathbb{R}^{d' \times 1}$  are representations for nodes  $i$  and  $j$  respectively.  $W_K, W_Q \in \mathbb{R}^{d \times d'}$  are the standard learnable key and query projection matrices.  $b_{ij} \in \mathbb{R}$  is a learnable scalar indexed by the shortest path distance between  $i$  and  $j$  (which is always a positive integer).  $c_{ij} \in \mathbb{R}$  is the edge embedding term, described by Equation 2. Allowing a slight abuse of notation,  $e_p \in \mathbb{R}^{d \times 1}$  is the embedding corresponding to the  $p^{th}$  edge in the shortest path between  $i$  and  $j$ , and  $w_p \in \mathbb{R}^{d \times 1}$  is a learnable weight for that position.

$$a_{ij} = \text{softmax} \left( \frac{(W_Q h_i)^T (W_K h_j)}{\sqrt{d}} + b_{ij} + c_{ij} \right) \quad (1)$$

$$c_{ij} = \frac{1}{N} \sum_p w_p^T e_p \quad (2)$$

For graph-level prediction tasks, it is useful to add a readout node to the input graph for extracting graph-level embeddings, similar to the CLS token in NLP transformers [48]. This “fake” node is initialized with a unique embedding and connected to all other nodes in the graph with a special edge type. In the final layer of the transformer, the readout node’s embedding is interpreted as a global representation of the input graph and can be used for downstream property prediction.

### 3.4 Attention Map

Attention maps are a way of visualizing a transformer’s attention on the input, and are commonly used to interpret model behaviour in NLP [49] and computer vision [33]. For MassFormer, the attention map is computed by multiplying the self-attention matrices from each multi-head attention layer in the transformer, then extracting the attention values corresponding to the readout node. The resulting vector is then normalized to form a distribution over all  $N$  elements, with higher values signifying more attention (relative to other elements). Input regions with high attention often correlate with what human practitioners consider to be important features that are relevant to the task, and can be used to help explain predictions.

### 3.5 Objective Function

Since MS peak intensities are relative <sup>6</sup>, it is advisable to use loss functions that are invariant to scaling. We choose cosine distance (Eq. 3) as the loss function, where  $\hat{y} = f_\theta(x, z)$  is the predicted spectrum and  $y$  is the

<sup>6</sup> In quantification experiments, the unnormalized value of a peak intensity can provide useful information, but in an identification context it does not.

real spectrum. This loss function, and the related cosine similarity, are commonly used to compare spectra [29,30]. Our objective also includes a weight decay term ( $l_2$  regularization on the model parameters), as is common among deep learning models.

$$CD(y, \hat{y}) = 1 - \frac{y^T \hat{y}}{\|y\|_2 \|\hat{y}\|_2} = 1 - \frac{\sum_{i=1}^m y_i \hat{y}_i}{\sqrt{\sum_{j=1}^m y_j^2 \sum_{k=1}^m \hat{y}_k^2}} \quad (3)$$

## 4 Experiments and Results

### 4.1 Datasets

We use the NIST 2020 MS/MS dataset [9] for both training and evaluation. NIST is a commercial dataset notable for its large coverage (over 1 million tandem spectra in total), standardized spectrum acquisition protocol, and high degree of manual validation. We also use the publicly available MassBank of North America (MB-NA) as a held-out evaluation set. Massbank is an important dataset for the MS community, as it is the largest publicly available online repository of MS data for small molecules. For simplicity, we only consider spectra from Fourier Transform (FT) MS instruments [1]. We split the NIST and MB-NA datasets into two subsets based on their collision cell type (HCD vs CID). The dataset statistics are summarized in Table 1.

Dataset	Collision Type	# Spectra	# Compounds
NIST	FT-HCD	717,029	21,060
NIST	FT-CID	65,753	20,388
MB-NA	FT-HCD	15,474	1,460
MB-NA	FT-CID	4,749	889

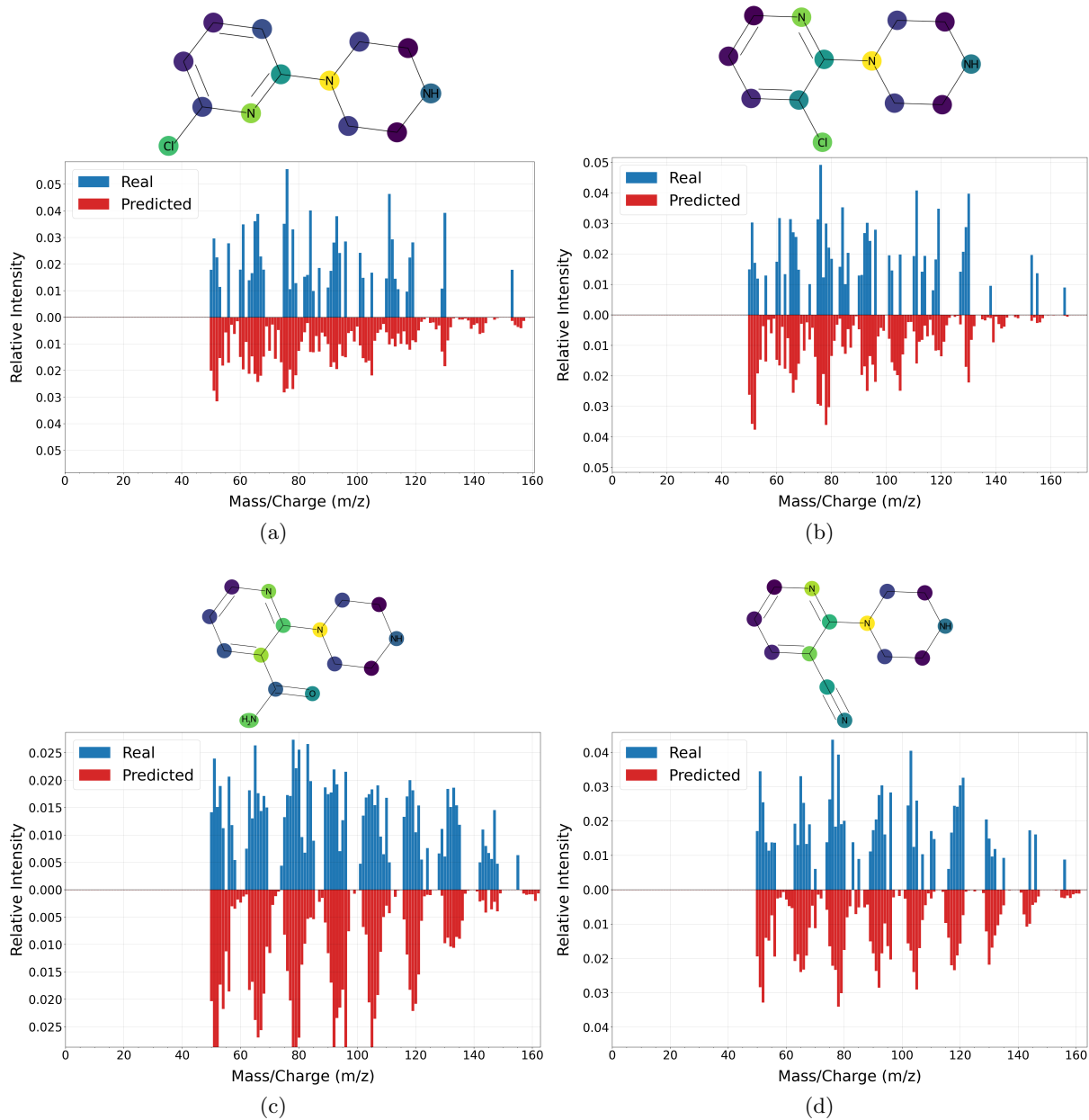
**Table 1.** Summary statistics for each MS dataset. MB-NA is MassBank of North America.

### 4.2 Models

We compare a number of models that are distinguished by the type of input chemical representation that they rely on. Each model is composed of two components: a chemical encoder network which transforms the input molecule representation into a latent representation, followed by a fully-connected neural network, conditioned on the spectral metadata, that transforms the chemical embedding into a predicted spectrum. Both parts of the network are trained jointly in an end-to-end fashion. The primary difference between models lies in the chemical encoder network architecture, which depends on the input representation; the fully-connected network is similar across models.

Our approach, MassFormer, uses a graph transformer as the encoder network (see Section 3.3 for more details). The Fingerprint model combines Morgan, MACCS, and RDKit fingerprints (all of which are available from the RDKit library [50]) and uses those directly as the chemical embedding. The WLN model uses a molecular graph representation in combination with a Weisfeiler-Lehman Network ([51]), which is a particular kind of GNN, to produce the chemical embedding. The final model, CNN, uses a SELFIES string representation [42] in combination with a CNN architecture. The Fingerprint and WLN models were chosen because they are very similar to two previously published MS prediction models ([29] and [30] respectively) that are difficult to directly compare against. The former [29] is trained on a different type of MS data (electron ionization MS), while the latter [30] does not have an available public implementation <sup>7</sup>.

<sup>7</sup> We communicated with the authors of [30] and designed our WLN model to be similar to their best model



**Fig. 2.** Real and predicted spectra for four compounds from the NIST FT-HCD test subset, along with their attention map visualization. Darker means less attention (lower magnitude). These compounds have the same Murcko scaffold and are therefore quite similar in structure.

### 4.3 Implementation Details

All models are implemented in PyTorch [52]. Our model, MassFormer, uses a modified version of the original Graphormer implementation [45]. The WLN model uses the Deep Graph Library [53,54] package for geometric deep learning. We also adapt some code from [30] for spectrum preprocessing. Before benchmarking the models, we run a bayesian hyperparameter sweep using Weights & Biases [55], with budget of 100 initializations, to find the best-performing configuration on the NIST validation set. This allows for fairer model comparisons, reducing the chance that one model outperforms another as a result of suboptimal hyperparameter settings. The list of hyperparameters that we optimize is as follows: learning rate, weight decay, dropout, minibatch size, and network-architecture specific parameters (such as hidden dimension and number of layers). For MassFormer we do not vary the transformer architecture in the hyperparameter search, sticking with “small” default settings described in [45]. All of our models use the AdamW optimizer [56,57] with a plateau learning rate scheduler [52], meaning the learning rate is linearly decayed after a certain number of epochs without improvement.

### 4.4 Spectrum Similarity Experiment

The NIST dataset is split into training (70%), validation (10%), and test (20%) subsets by Murcko scaffold [58,50], a form of molecular “backbone” which coarsely describes structure. Special care is taken to remove any molecules from the NIST subsets whose Murcko scaffold is present in the heldout MB-NA set, to prevent similar compounds from leaking across datasets. The results are summarized in Table 2. Our approach performs best across both datasets, and the performance seems to translate across the different types of MS data.

Method	Cosine Similarity			
	FT-HCD		FT-CID	
	NIST	MB-NA	NIST	MB-NA
Fingerprint	.515 $\pm$ .002	.384 $\pm$ .004	.293 $\pm$ .003	.226 $\pm$ .003
WLN	.559 $\pm$ .002	<b>.439 <math>\pm</math> .003</b>	.337 $\pm$ .003	.221 $\pm$ .007
CNN	.523 $\pm$ .002	.383 $\pm$ .002	.335 $\pm$ .002	.249 $\pm$ .002
MassFormer	<b>.565 <math>\pm</math> .001</b>	<b>.439 <math>\pm</math> .003</b>	<b>.355 <math>\pm</math> .004</b>	<b>.269 <math>\pm</math> .005</b>

**Table 2.** Cosine Similarity scores for models trained on different types of MS data (FT-HCD and FT-CID) on two heldout datasets (NIST test partition and MB-NA). Reported score is mean  $\pm$  standard deviation over 5 random seeds. Higher is better, best scores (within standard deviation) indicated with **bold**.

### 4.5 Molecule Identification Experiment

Improving coverage of spectral libraries to help identify unknown query spectra is a key application for spectrum prediction models (see Section 2.2). However, measuring a model’s ability to match true unknown queries with predicted spectra is difficult, as it requires experimental analysis to verify the identity of the compound. A simpler evaluation can be achieved with a task where each query spectrum is associated with a known compound. Inspired by [29,59], we perform the following candidate ranking experiment, where the goal is to match a query spectrum with the correct molecule from a pool of candidates. The query set consists of real spectra from one of the heldout partitions, NIST test or MB-NA. The reference set is comprised of spectra from two sources: predicted spectra for all compounds in the heldout partition, and real spectra from the NIST training and validation partitions. By computing pairwise similarity between all spectra in both query and reference sets, it is possible to establish a ranking of spectra in the reference set (from most similar to least) for each query. Since each spectrum corresponds to a single molecule, ranking the spectra induces a ranking of candidate structures. Good prediction models will produce very accurate spectra, which will result in the correct match from the reference set being similar to the query and thus receiving a high rank. The addition of experimental spectra from the training and validation partitions makes the task more difficult by introducing realistic false matches <sup>8</sup>.

<sup>8</sup> Since our data partitions are compound disjoint, we can be sure that there are no true matches from the training/validation partition



Method	Norm Rank			
	FT-HCD		FT-CID	
	NIST	MB-NA	NIST	MB-NA
Fingerprint	.095 $\pm$ .008	.138 $\pm$ .010	.054 $\pm$ .003	.081 $\pm$ .003
WLN	.043 $\pm$ .003	<b>.063 <math>\pm</math> .004</b>	.057 $\pm$ .005	.134 $\pm$ .008
CNN	<b>.042 <math>\pm</math> .002</b>	<b>.067 <math>\pm</math> .005</b>	<b>.033 <math>\pm</math> .001</b>	<b>.062 <math>\pm</math> .001</b>
MassFormer	<b>.040 <math>\pm</math> .002</b>	.071 $\pm$ .004	<b>.031 <math>\pm</math> .003</b>	<b>.062 <math>\pm</math> .006</b>
Method	Top-5%			
	FT-HCD		FT-CID	
	NIST	MB-NA	NIST	MB-NA
Fingerprint	.533 $\pm$ .030	.433 $\pm$ .033	.669 $\pm$ .012	.529 $\pm$ .013
WLN	.748 $\pm$ .013	<b>.663 <math>\pm</math> .013</b>	.709 $\pm$ .016	.428 $\pm$ .009
CNN	.744 $\pm$ .010	.632 $\pm$ .017	.791 $\pm$ .007	.641 $\pm$ .006
MassFormer	<b>.766 <math>\pm</math> .014</b>	<b>.644 <math>\pm</math> .020</b>	<b>.835 <math>\pm</math> .011</b>	<b>.691 <math>\pm</math> .010</b>

**Table 3.** Normalized Rank and Top-5% scores on the ranking task. Reported score is mean  $\pm$  standard deviation over 5 random seeds. Lower is better for Normalized Rank, while higher is better for Top-5%. Best scores (within standard deviation) indicated with **bold**.

The results of this experiment are summarized in Table 3. Since the number of candidates per query can vary<sup>9</sup>, all metrics must be normalized. We report Normalized Rank, the rank of the true matching candidate divided by the number of candidates, and Top-5%, whether the true matching candidate was found in the top 5% of all candidates. Our model is one of the best performing in 7 out of 8 data regimes, indicating consistently strong performance.

#### 4.6 Collision Energy Experiment

Collision energy is a key experimental parameter that can have a large impact on the observed spectrum (see Section 1). We want to investigate whether MassFormer can accurately model the relationship between collision energy and predicted spectra. Figure 3 shows the relationship between collision energy and mean  $m/z$  of the peaks in the spectrum, where each peak is weighted by its (normalized) relative intensity. Qualitatively, the real and predicted distributions for the FT-HCD data are difficult to distinguish. Figure 4 compares a real FT-HCD spectrum at a high collision energy (180) with three predicted spectra corresponding to the same underlying compound but at a variety of collision energies (20, 77, 180). Note that there are no ground truth spectra for for the two lower energies. In this case, predicted spectra with larger collision energies have lower mean peak  $m/z$ , which is what we would expect.

#### 4.7 Atom Attention Analysis

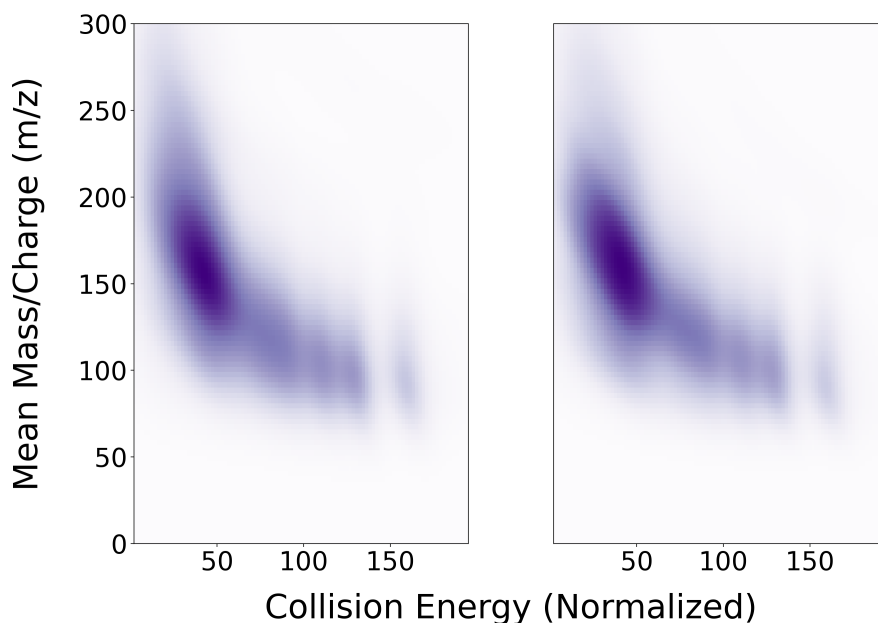
Visualizing the attention map (see Section 3.3) created by MassFormer on the atoms of the input molecules provides some insight into the model’s behaviour. Figure 2 qualitatively demonstrates that similar molecules have similar patterns of attention (and similar predictions). These molecules were chosen because they have the same underlying Murcko Scaffold [58]. In all four of these examples, the model seems to pay more attention to atoms that are not carbon (also known as heteroatoms). Looking across the entire NIST FT-HCD test set, Figure 5 shows that on average the model attends 1.2x to 2.0x more to heteroatoms than carbon atoms. This simple observation is consistent with our understanding of mass spectrometry, as heteroatoms are reactive and play a key role in an organic compound’s chemistry, and by extension, its fragmentation patterns [1].

## 5 Conclusion

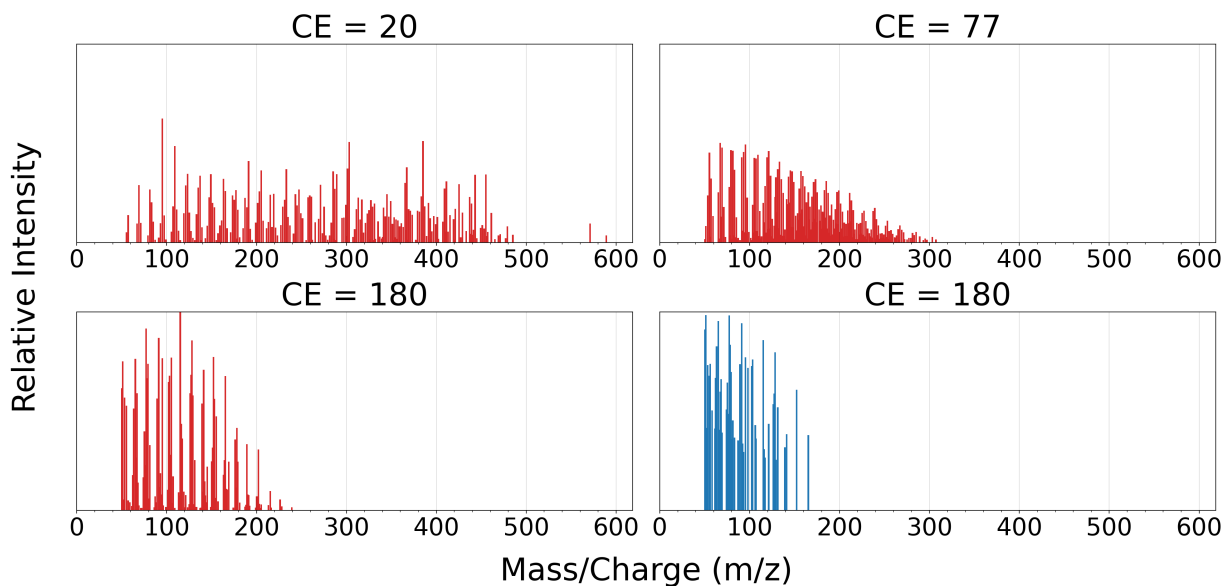
In this work we propose MassFormer, a novel method for predicting MS/MS spectra from chemical structures using a graph transformer architecture. We validate our method across two independent MS datasets, NIST

<sup>9</sup> We only choose candidates from the reference set with the same metadata as the query

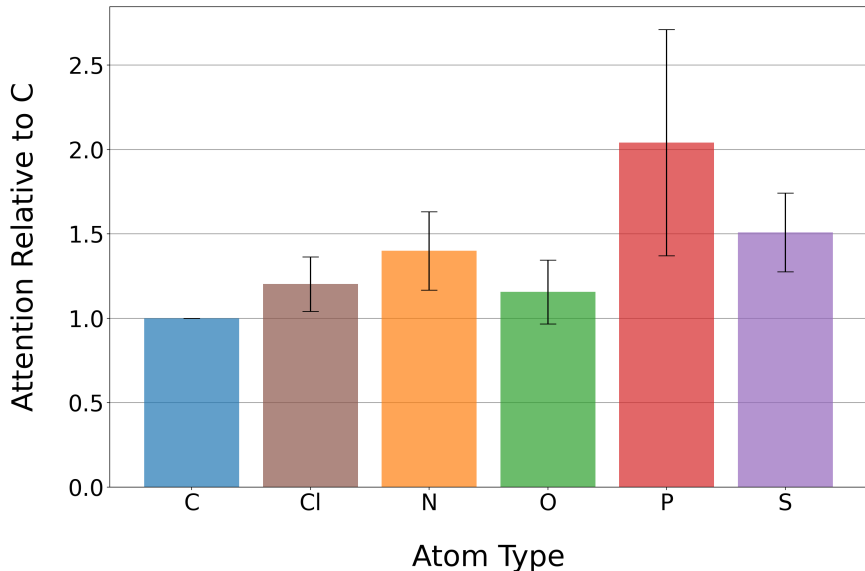




**Fig. 3.** The joint distribution of mean MS peak location (y-axis) and collision energy (x-axis). The real NIST HCD test distribution is on the left, while the predicted distribution is on the right. As collision energy increases, mean peak location decreases.



**Fig. 4.** Increasing collision energy for a given molecule results in peaks with lower  $m/z$ . The blue spectrum is a real spectrum from the NIST HCD validation subset (collision energy = 180), the red spectra are predicted ones corresponding to the same compound but at differing collision energies.



**Fig. 5.** The relative attention of 5 heteroatoms (Chlorine, Nitrogen, Oxygen, Phosphorus, and Sulfur) compared to Carbon, averaged over the NIST FT-HCD test set. Relative attention is defined as the ratio of average attention between two atom types in a molecule.

and MassBank, and show that it can produce realistic spectra. We demonstrate that the model captures prior knowledge about the fragmentation process by investigating the effect of collision energy on the predicted spectrum. We benchmark the model with a candidate ranking task and show that it can be useful for MS/MS-based molecule identification.

## 6 Discussion and Future Directions

Computer-assisted metabolite identification has long been a goal for the MS/MS community. The development of accurate spectrum prediction models is an important step in this direction. Online MS databases such as MassBank [6] and GNPS [7] host large collections of experimental and simulated spectra [25] to help practitioners match their own unlabelled query spectra. Currently, it is common to use the trained models “as-is”, either by running them on compounds of interest or simply by using available precomputed spectra. Many research labs have private repositories of MS data that they use for in-house analysis. There is great potential in developing models that are flexible enough to deal with different kinds of MS data, while being relatively inexpensive to finetune on custom datasets. Deep learning approaches scale well with increasing data size, and can be more flexible to changes in data modality than other methods that heavily rely on domain knowledge. In future work, we hope to test these ideas by training our model on a wider diversity of MS data that incorporates different types of mass spectrometers and mass analyzers. It would be interesting to see if joint training on different MS modalities could improve accuracy across the board.

Pretraining is especially useful for deep learning models when data is scarce. While the number of compounds available in spectral databases is on the order of  $10^4 - 10^5$ , there are currently millions of organic compounds in public chemical databases (like PubChem [21]) that lack spectral information. A promising direction for future work would be to pretrain the graph transformer by applying self-supervised learning strategies [60], allowing the model to learn meaningful chemical representations without requiring expensive experimental data. Additionally, performing standard supervised training on other chemical property prediction tasks with lots of data (such as HOMO-LUMO gap regression, [61]) may also help.

Our long-term goal is to apply MassFormer to a real-world molecule identification task, where the query spectrum is truly unknown, and then verify compound identity with analytical chemistry techniques.

## References

1. W. M. Niessen and D. Falck, "Introduction to Mass Spectrometry, a Tutorial," in *Analyzing Biomolecular Interactions by Mass Spectrometry*. John Wiley & Sons, Ltd, 2015, pp. 1–54. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/9783527673391.ch1>
2. R. Aebersold and M. Mann, "Mass-spectrometric exploration of proteome structure and function," *Nature*, vol. 537, no. 7620, pp. 347–355, 09 2016.
3. G. N. Gowda and D. Djukovic, "Overview of Mass Spectrometry-Based Metabolomics: Opportunities and Challenges," *Methods in molecular biology (Clifton, N.J.)*, vol. 1198, pp. 3–12, 2014. [Online]. Available: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4336784/>
4. T. De Vijlder, D. Valkenburg, F. Lemi re, E. P. Romijn, K. Laukens, and F. Cuyckens, "A tutorial in small molecule identification via electrospray ionization-mass spectrometry: The practical art of structural elucidation," *Mass Spectrometry Reviews*, vol. 37, no. 5, pp. 607–629, 2018. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/mas.21551>
5. A. T. Lebedev, "Environmental Mass Spectrometry," *Annual Review of Analytical Chemistry*, vol. 6, no. 1, pp. 163–189, 2013. [Online]. Available: <https://doi.org/10.1146/annurev-anchem-062012-092604>
6. H. Horai, M. Arita, S. Kanaya, Y. Nihei, T. Ikeda, K. Suwa, Y. Ojima, K. Tanaka, S. Tanaka, K. Aoshima, Y. Oda, Y. Kakazu, M. Kusano, T. Tohge, F. Matsuda, Y. Sawada, M. Y. Hirai, H. Nakanishi, K. Ikeda, N. Akimoto, T. Maoka, H. Takahashi, T. Ara, N. Sakurai, H. Suzuki, D. Shibata, S. Neumann, T. Iida, K. Tanaka, K. Funatsu, F. Matsuura, T. Soga, R. Taguchi, K. Saito, and T. Nishioka, "MassBank: a public repository for sharing mass spectral data for life sciences," *Journal of Mass Spectrometry*, vol. 45, no. 7, pp. 703–714, Jul. 2010. [Online]. Available: <http://doi.wiley.com/10.1002/jms.1777>
7. M. Wang, J. J. Carver, V. V. Phelan, L. M. Sanchez, N. Garg, Y. Peng, D. D. Nguyen, J. Watrous, C. A. Kapono, T. Luzzatto-Knaan, C. Porto, A. Bouslimani, A. V. Melnik, M. J. Meehan, W. T. Liu, M. Cr usemann, P. D. Boudreau, E. Esquenazi, M. Sandoval-Calder n, R. D. Kersten, L. A. Pace, R. A. Quinn, K. R. Duncan, C. C. Hsu, D. J. Floros, R. G. Gavilan, K. Kleigrew, T. Northen, R. J. Dutton, D. Parrot, E. E. Carlson, B. Aigle, C. F. Michelsen, L. Jelsbak, C. Sohlenkamp, P. Pevzner, A. Edlund, J. McLean, J. Piel, B. T. Murphy, L. Gerwick, C. C. Liaw, Y. L. Yang, H. U. Humpf, M. Maansson, R. A. Keyzers, A. C. Sims, A. R. Johnson, A. M. Sidebottom, B. E. Sedio, A. Klitgaard, C. B. Larson, C. A. B. P, D. Torres-Mendoza, D. J. Gonzalez, D. B. Silva, L. M. Marques, D. P. Demarque, E. Pociute, E. C. O'Neill, E. Briand, E. J. N. Helfrich, E. A. Granatosky, E. Glukhov, F. Ryffel, H. Houson, H. Mohimani, J. J. Kharbush, Y. Zeng, J. A. Vorholt, K. L. Kurita, P. Charusanti, K. L. McPhail, K. F. Nielsen, L. Vuong, M. Elfeki, M. F. Traxler, N. Engene, N. Koyama, O. B. Vining, R. Baric, R. R. Silva, S. J. Mascuch, S. Tomasi, S. Jenkins, V. Macherla, T. Hoffman, V. Agarwal, P. G. Williams, J. Dai, R. Neupane, J. Gurr, A. M. C. Rodr guez, A. Lamsa, C. Zhang, K. Dorrestein, B. M. Duggan, J. Almaliti, P. M. Allard, P. Phapale, L. F. Nothias, T. Alexandrov, M. Litaudon, J. L. Wolfender, J. E. Kyle, T. O. Metz, T. Peryea, D. T. Nguyen, D. VanLeer, P. Shinn, A. Jadhav, R. M ller, K. M. Waters, W. Shi, X. Liu, L. Zhang, R. Knight, P. R. Jensen, B. O. Palsson, K. Pogliano, R. G. Linington, M. Guti rrez, N. P. Lopes, W. H. Gerwick, B. S. Moore, P. C. Dorrestein, and N. Bandeira, "Sharing and community curation of mass spectrometry data with Global Natural Products Social Molecular Networking," *Nat Biotechnol*, vol. 34, no. 8, pp. 828–837, 08 2016.
8. S. Stein, "Mass Spectral Reference Libraries: An Ever-Expanding Resource for Chemical Identification," *Analytical Chemistry*, vol. 84, no. 17, pp. 7274–7282, Sep. 2012. [Online]. Available: <https://doi.org/10.1021/ac301205z>
9. X. Yang, P. Neta, and S. E. Stein, "Quality Control for Building Libraries from Electrospray Ionization Tandem Mass Spectra," *Analytical Chemistry*, vol. 86, no. 13, pp. 6393–6400, Jul. 2014. [Online]. Available: <https://pubs.acs.org/doi/10.1021/ac500711m>
10. C. Guijas, J. R. Montenegro-Burke, X. Domingo-Almenara, A. Palermo, B. Warth, G. Hermann, G. Koellensperger, T. Huan, W. Uritboonthai, A. E. Aisporna, D. W. Wolan, M. E. Spilker, H. P. Benton, and G. Siuzdak, "METLIN: A Technology Platform for Identifying Knowns and Unknowns," *Analytical Chemistry*, vol. 90, no. 5, pp. 3156–3164, Mar. 2018. [Online]. Available: <https://doi.org/10.1021/acs.analchem.7b04424>
11. S. E. Stein and D. R. Scott, "Optimization and testing of mass spectral library search algorithms for compound identification," *Journal of the American Society for Mass Spectrometry*, vol. 5, no. 9, pp. 859–866, Sep. 1994. [Online]. Available: [https://pubs.acs.org/doi/10.1016/1044-0305\(94\)00009-8](https://pubs.acs.org/doi/10.1016/1044-0305(94)00009-8)
12. R. R. da Silva, P. C. Dorrestein, and R. A. Quinn, "Illuminating the dark matter in metabolomics," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 112, no. 41, pp. 12 549–12 550, Oct. 2015.
13. L. Ruddigkeit, R. van Deursen, L. C. Blum, and J.-L. Reymond, "Enumeration of 166 Billion Organic Small Molecules in the Chemical Universe Database GDB-17," *Journal of Chemical Information and Modeling*, vol. 52, no. 11, pp. 2864–2875, Nov. 2012. [Online]. Available: <https://pubs.acs.org/doi/10.1021/ci300415d>

14. B. G. Buchanan and E. A. Feigenbaum, "DENDRAL and Meta-DENDRAL: Their Applications Dimension." STANFORD UNIV CALIF DEPT OF COMPUTER SCIENCE, Tech. Rep., Feb. 1978. [Online]. Available: <https://apps.dtic.mil/sti/citations/ADA054289>
15. K. Scheubert, F. Hufsky, F. Rasche, and S. Böcker, "Computing Fragmentation Trees from Metabolite Multiple Mass Spectrometry Data," in *Research in Computational Molecular Biology*, ser. Lecture Notes in Computer Science, V. Bafna and S. C. Sahinalp, Eds. Berlin, Heidelberg: Springer, 2011, pp. 377–391.
16. S. Böcker and K. Dührkop, "Fragmentation trees reloaded," *Journal of Cheminformatics*, vol. 8, no. 1, p. 5, Feb. 2016. [Online]. Available: <https://doi.org/10.1186/s13321-016-0116-8>
17. K. Dührkop, H. Shen, M. Meusel, J. Rousu, and S. Böcker, "Searching molecular structure databases with tandem mass spectra using csi:fingerid," *Proceedings of the National Academy of Sciences*, vol. 112, no. 41, pp. 12 580–12 585, 2015. [Online]. Available: <https://www.pnas.org/content/112/41/12580>
18. K. Dührkop, M. Fleischauer, M. Ludwig, A. A. Aksenov, A. V. Melnik, M. Meusel, P. C. Dorrestein, J. Rousu, and S. Böcker, "SIRIUS 4: a rapid tool for turning tandem mass spectra into metabolite structure information," *Nature Methods*, vol. 16, no. 4, pp. 299–302, Apr. 2019. [Online]. Available: <http://www.nature.com/articles/s41592-019-0344-8>
19. B. U. Curry and D. Rumelhart, "A Neural Network That Classifies Mass Spectra," 2001. [Online]. Available: <https://www.semanticscholar.org/paper/A-Neural-Network-That-Classifies-Mass-Spectra-Curry-Rumelhart/09cc37106e61f2fd1f1dc881de24951a6498b354>
20. J. Lim, J. Wong, M. X. Wong, L. H. E. Tan, H. L. Chieu, D. Choo, and N. K. N. Neo, "Chemical Structure Elucidation from Mass Spectrometry by Matching Substructures," *arXiv:1811.07886 [physics, stat]*, Nov. 2018. [Online]. Available: <http://arxiv.org/abs/1811.07886>
21. S. Kim, J. Chen, T. Cheng, A. Gindulyte, J. He, S. He, Q. Li, B. A. Shoemaker, P. A. Thiessen, B. Yu, L. Zaslavsky, J. Zhang, and E. E. Bolton, "PubChem 2019 update: improved access to chemical data," *Nucleic Acids Research*, vol. 47, no. Database issue, pp. D1102–D1109, Jan. 2019. [Online]. Available: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6324075/>
22. D. S. Wishart, Y. D. Feunang, A. Marcu, A. C. Guo, K. Liang, R. Vázquez-Fresno, T. Sajed, D. Johnson, C. Li, N. Karu, Z. Sayeeda, E. Lo, N. Assempour, M. Berjanskii, S. Singhal, D. Arndt, Y. Liang, H. Badran, J. Grant, A. Serra-Cayuela, Y. Liu, R. Mandal, V. Neveu, A. Pon, C. Knox, M. Wilson, C. Manach, and A. Scalbert, "HMDB 4.0: the human metabolome database for 2018," *Nucleic Acids Research*, vol. 46, no. D1, pp. D608–D617, Jan. 2018.
23. M. Kanehisa, M. Furumichi, Y. Sato, M. Ishiguro-Watanabe, and M. Tanabe, "KEGG: integrating viruses and cellular organisms," *Nucleic Acids Research*, vol. 49, no. D1, pp. D545–D551, Jan. 2021.
24. L. Ridder, J. J. J. v. d. Hooft, and S. Verhoeven, "Automatic Compound Annotation from Mass Spectrometry Data Using MAGMa," *Mass Spectrometry*, vol. 3, no. SpecialIssue.2, pp. S0033–S0033, 2014.
25. F. Wang, J. Liigand, S. Tian, D. Arndt, R. Greiner, and D. S. Wishart, "CFM-ID 4.0: More Accurate ESI-MS/MS Spectral Prediction and Compound Identification," *Analytical Chemistry*, Aug. 2021. [Online]. Available: <https://doi.org/10.1021/acs.analchem.1c01465>
26. S. Wolf, S. Schmidt, M. Müller-Hannemann, and S. Neumann, "In silico fragmentation for computer assisted identification of metabolite mass spectra," *BMC Bioinformatics*, vol. 11, no. 1, p. 148, Mar. 2010. [Online]. Available: <https://doi.org/10.1186/1471-2105-11-148>
27. C. Ruttkies, E. L. Schymanski, S. Wolf, J. Hollender, and S. Neumann, "MetFrag relaunched: incorporating strategies beyond in silico fragmentation," *Journal of Cheminformatics*, vol. 8, no. 1, p. 3, Jan. 2016. [Online]. Available: <https://doi.org/10.1186/s13321-016-0115-9>
28. F. Allen, R. Greiner, and D. Wishart, "Competitive fragmentation modeling of ESI-MS/MS spectra for putative metabolite identification," *Metabolomics*, vol. 11, no. 1, pp. 98–110, Feb. 2015. [Online]. Available: <https://doi.org/10.1007/s11306-014-0676-4>
29. J. N. Wei, D. Belanger, R. P. Adams, and D. Sculley, "Rapid prediction of electron-ionization mass spectrometry using neural networks," *ACS Central Science*, vol. 5, no. 4, pp. 700–708, Apr. 2019. [Online]. Available: <https://pubs.acs.org/doi/10.1021/acscentsci.9b00085>
30. H. Zhu, L. Liu, and S. Hassoun, "Using Graph Neural Networks for Mass Spectrometry Prediction," *arXiv:2010.04661 [cs]*, Oct. 2020. [Online]. Available: <http://arxiv.org/abs/2010.04661>
31. A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention Is All You Need," *arXiv:1706.03762 [cs]*, Dec. 2017. [Online]. Available: <http://arxiv.org/abs/1706.03762>
32. Z. Tan, S. Wang, Z. Yang, G. Chen, X. Huang, M. Sun, and Y. Liu, "Neural machine translation: A review of methods, resources, and tools," *AI Open*, vol. 1, pp. 5–21, Jan. 2020. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2666651020300024>
33. A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," *arXiv:2010.11929 [cs]*, Jun. 2021. [Online]. Available: <http://arxiv.org/abs/2010.11929>

34. M. Janner, Q. Li, and S. Levine, "Reinforcement Learning as One Big Sequence Modeling Problem," *arXiv:2106.02039 [cs]*, Jul. 2021. [Online]. Available: <http://arxiv.org/abs/2106.02039>
35. A. H. Khasahmadi, K. Hassani, P. Moradi, L. Lee, and Q. Morris, "Memory-Based Graph Networks," *arXiv:2002.09518 [cs, stat]*, Jun. 2020. [Online]. Available: <http://arxiv.org/abs/2002.09518>
36. G. Mialon, D. Chen, M. Selsos, and J. Mairal, "GraphiT: Encoding Graph Structure in Transformers," *arXiv:2106.05667 [cs]*, Jun. 2021. [Online]. Available: <http://arxiv.org/abs/2106.05667>
37. L. Maziarka, T. Danel, S. Mucha, K. Rataj, J. Tabor, and S. Jastrzębski, "Molecule Attention Transformer," *arXiv:2002.08264 [physics, stat]*, Feb. 2020. [Online]. Available: <http://arxiv.org/abs/2002.08264>
38. Y. Rong, Y. Bian, T. Xu, W. Xie, Y. Wei, W. Huang, and J. Huang, "Self-Supervised Graph Transformer on Large-Scale Molecular Data," *arXiv:2007.02835 [cs, q-bio]*, Oct. 2020. [Online]. Available: <http://arxiv.org/abs/2007.02835>
39. D. Rogers and M. Hahn, "Extended-Connectivity Fingerprints," *Journal of Chemical Information and Modeling*, vol. 50, no. 5, pp. 742–754, May 2010. [Online]. Available: <https://doi.org/10.1021/ci100050t>
40. P. W. Battaglia, J. B. Hamrick, V. Bapst, A. Sanchez-Gonzalez, V. Zambaldi, M. Malinowski, A. Tacchetti, D. Raposo, A. Santoro, R. Faulkner, C. Gulcehre, F. Song, A. Ballard, J. Gilmer, G. Dahl, A. Vaswani, K. Allen, C. Nash, V. Langston, C. Dyer, N. Heess, D. Wierstra, P. Kohli, M. Botvinick, O. Vinyals, Y. Li, and R. Pascanu, "Relational inductive biases, deep learning, and graph networks," *arXiv:1806.01261 [cs, stat]*, Oct. 2018. [Online]. Available: <http://arxiv.org/abs/1806.01261>
41. D. Weininger, "SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules," *Journal of Chemical Information and Computer Sciences*, vol. 28, no. 1, pp. 31–36, Feb. 1988, publisher: American Chemical Society. [Online]. Available: <https://pubs.acs.org/doi/abs/10.1021/ci00057a005>
42. M. Krenn, F. Häse, A. Nigam, P. Friederich, and A. Aspuru-Guzik, "Self-Referencing Embedded Strings (SELFIES): A 100% robust molecular string representation," *Machine Learning: Science and Technology*, vol. 1, no. 4, p. 045024, Nov. 2020. [Online]. Available: <http://arxiv.org/abs/1905.13741>
43. Y. LeCun, P. Haffner, L. Bottou, and Y. Bengio, "Object Recognition with Gradient-Based Learning," in *Shape, Contour and Grouping in Computer Vision*, ser. Lecture Notes in Computer Science, D. A. Forsyth, J. L. Mundy, V. di Gesù, and R. Cipolla, Eds. Berlin, Heidelberg: Springer, 1999, pp. 319–345. [Online]. Available: [https://doi.org/10.1007/3-540-46805-6\\_19](https://doi.org/10.1007/3-540-46805-6_19)
44. S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997. [Online]. Available: <https://doi.org/10.1162/neco.1997.9.8.1735>
45. C. Ying, T. Cai, S. Luo, S. Zheng, G. Ke, D. He, Y. Shen, and T.-Y. Liu, "Do transformers really perform bad for graph representation?" *Neural Information Processing Systems (NeurIPS)*, 2021.
46. W. Hu, M. Fey, M. Zitnik, Y. Dong, H. Ren, B. Liu, M. Catasta, and J. Leskovec, "Open Graph Benchmark: Datasets for Machine Learning on Graphs," *arXiv:2005.00687 [cs, stat]*, Feb. 2021. [Online]. Available: <http://arxiv.org/abs/2005.00687>
47. R. W. Floyd, "Algorithm 97: Shortest path," *Communications of the ACM*, vol. 5, no. 6, p. 345, Jun. 1962. [Online]. Available: <https://doi.org/10.1145/367766.368168>
48. J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," *arXiv:1810.04805 [cs]*, May 2019. [Online]. Available: <http://arxiv.org/abs/1810.04805>
49. J. Vig, "Visualizing Attention in Transformer-Based Language Representation Models," *arXiv:1904.02679 [cs, stat]*, Apr. 2019. [Online]. Available: <http://arxiv.org/abs/1904.02679>
50. G. Landrum, "Rdkit: Open-source cheminformatics." [Online]. Available: <http://www.rdkit.org>
51. W. Jin, C. W. Coley, R. Barzilay, and T. Jaakkola, "Predicting Organic Reaction Outcomes with Weisfeiler-Lehman Network," *arXiv:1709.04555 [cs, stat]*, Dec. 2017. [Online]. Available: <http://arxiv.org/abs/1709.04555>
52. A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "PyTorch: An Imperative Style, High-Performance Deep Learning Library," *arXiv:1912.01703 [cs, stat]*, Dec. 2019. [Online]. Available: <http://arxiv.org/abs/1912.01703>
53. M. Wang, D. Zheng, Z. Ye, Q. Gan, M. Li, X. Song, J. Zhou, C. Ma, L. Yu, Y. Gai, T. Xiao, T. He, G. Karypis, J. Li, and Z. Zhang, "Deep Graph Library: A Graph-Centric, Highly-Performant Package for Graph Neural Networks," *arXiv:1909.01315 [cs, stat]*, Aug. 2020. [Online]. Available: <http://arxiv.org/abs/1909.01315>
54. M. Li, J. Zhou, J. Hu, W. Fan, Y. Zhang, Y. Gu, and G. Karypis, "DGL-LifeSci: An Open-Source Toolkit for Deep Learning on Graphs in Life Science," *arXiv:2106.14232 [cs, q-bio]*, Jun. 2021. [Online]. Available: <http://arxiv.org/abs/2106.14232>
55. L. Biewald, "Experiment tracking with weights and biases," 2020, software available from wandb.com. [Online]. Available: <https://www.wandb.com/>
56. D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," *arXiv:1412.6980 [cs]*, Jan. 2017. [Online]. Available: <http://arxiv.org/abs/1412.6980>

57. I. Loshchilov and F. Hutter, “Decoupled Weight Decay Regularization,” *arXiv:1711.05101 [cs, math]*, Jan. 2019. [Online]. Available: <http://arxiv.org/abs/1711.05101>
58. G. W. Bemis and M. A. Murcko, “The Properties of Known Drugs. 1. Molecular Frameworks,” *Journal of Medicinal Chemistry*, vol. 39, no. 15, pp. 2887–2893, Jan. 1996. [Online]. Available: <https://doi.org/10.1021/jm9602928>
59. E. L. Schymanski, C. Ruttkies, M. Krauss, C. Brouard, T. Kind, K. Dührkop, F. Allen, A. Vaniya, D. Verdegem, S. Böcker, J. Rousu, H. Shen, H. Tsugawa, T. Sajed, O. Fiehn, B. Ghesquière, and S. Neumann, “Critical Assessment of Small Molecule Identification 2016: automated methods,” *Journal of Cheminformatics*, vol. 9, no. 1, p. 22, Mar. 2017. [Online]. Available: <https://doi.org/10.1186/s13321-017-0207-1>
60. W. Hu, B. Liu, J. Gomes, M. Zitnik, P. Liang, V. Pande, and J. Leskovec, “Strategies for Pre-training Graph Neural Networks,” *arXiv:1905.12265 [cs, stat]*, Feb. 2020, arXiv: 1905.12265. [Online]. Available: <http://arxiv.org/abs/1905.12265>
61. M. Nakata and T. Shimazaki, “PubChemQC Project: A Large-Scale First-Principles Electronic Structure Database for Data-Driven Chemistry,” *Journal of Chemical Information and Modeling*, vol. 57, no. 6, pp. 1300–1308, Jun. 2017. [Online]. Available: <https://doi.org/10.1021/acs.jcim.7b00083>