
Using Graph Neural Networks for Mass Spectrum Prediction

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Mass Spectrum (MS) is widely used to assign chemical identities to small molecules
2 in many biological and biomedical applications. In this work, we explore using
3 graph neural networks (GNNs) for MS prediction. The input to our model is
4 a molecular graph. The model is trained and tested on the NIST 17 LC-MS
5 dataset. We compare our results to NEIMS, a neural network model that utilizes
6 molecular fingerprints as inputs. Our results show that GNN-based models offer
7 higher performance than NEIMS. Importantly, we show that ranking results heavily
8 depend on the candidate set size and on the similarity of the candidates to the target
9 molecule, thus highlighting the need for consistent, well-characterized evaluation
10 protocols for this domain.

11 1 Introduction

12 Mass Spectrum (MS) techniques coupled with liquid or gas chromatography separation techniques,
13 LC-MS or GC-MS, have become a standard analytical platform for small molecules produced in
14 biological systems [1, 2]. The LC or GC step aims to separate compounds within the sample, while the
15 MS step aims to ionize, fragment and detect a fragmentation pattern. For each particular compound
16 and its fragments, this pattern forms a *spectral signature*, comprising a chromatographic retention
17 time (RT) paired with mass-to-charge ratio (m/z) and their respective relative intensities.

18 Interpreting MS measures requires *annotation*, the process of assigning putative chemical identities
19 to each spectral signature. From a computational perspective, mapping molecules to their respective
20 spectral signatures represents a “forward problem”. Mapping spectral signatures back to their respec-
21 tive molecular identity is an “inverse problem”. The inverse problem is exceptionally challenging
22 as not all fragments are measured and many isomers (same molecular formula but different atom
23 configurations) have almost indistinguishable spectra.

24 Current annotation techniques attempt to solve the forward problem, and can be broadly classified
25 into two categories. Database lookup relies on comparing measured spectra against experimentally
26 generated spectra cataloged in spectral databases [3, 4, 5, 6]. Coverage of such libraries however is
27 limited, as experiments are required to generate signatures from known, trusted chemical standards.
28 Based on a predetermined set of *candidate molecules*, *in-silico* annotation tools recommend a
29 candidate molecule that best explains the measured spectra. The candidate set is typically culled
30 from large molecular databases, such as PubChem, based on molecular mass or formula, if possible.
31 Earlier works generated fragmentation patterns of candidate metabolites using rule-based approaches
32 [7, 8, 9] followed by combinatorial enumeration methods [10, 11, 12, 13]. More recently, machine-
33 learning algorithms have been investigated. CFM-ID trains a probabilistic generative model of the
34 fragmentation process to predict patterns of fragmentation [14, 15]. CSI:FingerID [16] first predicts
35 a fragmentation tree based on a spectral signature and then uses multiple-kernel learning [17] and
36 support vector machines (SVMs) to predict molecular properties that are then searched against

37 properties of candidate molecules. All of these techniques *explicitly* model compound fragmentation.
38 A recent study proposed a model, NEIMS (Neural Electron-Ionization Mass Spectrometry), to
39 augment existing NIST 17 GC-MS libraries with synthetic spectra predicted from candidate
40 molecules [18]. The molecules are first mapped to their ECFP (Extended-Connectivity Fingerprints)
41 fingerprint that record the count of molecular subgraphs within a specified radius centered on each
42 atom in the molecule. Using a radius of 2, the fingerprint consisted of 4096 entries. The fingerprint is
43 input to a fully connected feed forward neural network (FFNN) with a gated bidirectional design to
44 improve the prediction accuracy. However, fingerprints, while common and useful, are not tailored
45 to the prediction task. Moreover, NEIMS was only evaluated on the NIST data generated via the
46 GC-MS technique. GC-MS fragmentation patterns are simpler than those obtained using LC-MS.
47 Further GC-MS is typically used to measure volatile compounds (or compounds that can be extracted
48 into an organic solvent and vaporized using GC), which typically have masses less than 500 Daltons.
49 Here, we evaluate the feasibility of using Graph Neural Networks (GNNs) on the MS prediction
50 task. GNNs have been shown powerful in terms of learning representations from structured data
51 [19, 20, 21, 22], such as social networks, knowledge graphs and molecules. Here, we represent each
52 molecule as a graph, where atoms are represented as nodes, and bonds are represented as edges. We
53 explore the use of both Graph Convolutional Networks (GCN) [20] and Graph Attention Networks
54 (GAT) [21]. We train and evaluate our technique on the NIST 17 LC-MS dataset.

55 2 Methods

56 2.1 Datasets

57 For training and evaluation, we focused our efforts on the tandem MS/MS data in the NIST 17 dataset,
58 and selected spectra obtained using HCD (higher energy collisional dissociation), which provides a
59 richer and more varied spectra than CID (collision induced dissociation). The HCD option was used
60 to measure 69.6% of the MS/MS data. We also restricted the precursor types to only include ones that
61 are common in the dataset (e.g., [M+H]⁺, [M-H]⁻, [M+H-H₂O]⁺, [M-H-H₂O]⁻, etc.). We selected
62 one spectrum with the largest collision energy under 40eV for each training and test molecule.

63 From the reduced dataset, we randomly select 1000 molecules as a test set. We then split the remaining
64 dataset into training and validation (4:1 ratio), yielding 6,188 training and 1,539 test molecules. The
65 validation set is used for model selection while the test set is used to report performance. The relevant
66 candidate sets for the test molecules were queried from PubChem using the exact molecular formula
67 of each test molecule. The average number of candidates per test molecule is 1,530.

68 2.2 Data Preprocessing

69 A molecule is represented as a graph $G = (V, E)$, where atoms correspond to the node set V and
70 bonds correspond to the edge set E . Let $\mathbf{X} = \{x_1, x_2, \dots, x_N\}$, $x_i \in \mathbb{R}^F$ denote a set of node
71 features, where N is the number of nodes and F is the number of node features. The connectivity
72 among nodes is described by an adjacency matrix $\mathbf{A} \in \{0, 1\}^{N \times N}$. Node features include standard
73 atom-level features, such as atomic weight, atom type, number of bonds, etc..

74 The spectra data is a list of paired mass-to-charge ratio (m/z) and their relative intensities. Each m/z
75 value was rounded down to the nearest integer m/z bin. If more than one m/z values are rounded to
76 the same bin, we record the highest intensity. The range of intensity values has a long tail of large
77 values, so we take either the logarithm or the square-root of these values and denote the vector as \mathbf{y} .

78 2.3 Neural Network

79 An overview of our neural network (NN) architecture is illustrated in Figure 1. The model comprises
80 multiple GNN layers, a pooling layer and a fully-connected feed forward regression model. In
81 each GNN layer, node information is propagated along graph edges. More specifically, let $\mathbf{h} =$
82 $\{h_1, h_2, \dots, h_N\}$, $h_i \in \mathbb{R}^H$ denote the input of a GNN layer, where H is the dimension of input node
83 embedding vector. Let \mathbf{h}' denote the output with H' as the dimension of the new node embedding.
84 The weight matrix of such a layer is $\mathbf{W} \in \mathbb{R}^{H' \times H}$ and the bias term is b . In the first GNN layer, \mathbf{h} is
85 set to the original node features, \mathbf{X} . By stacking several layers of GNN layers, information on each
86 node propagates along edges to a broader neighborhood within the graph.

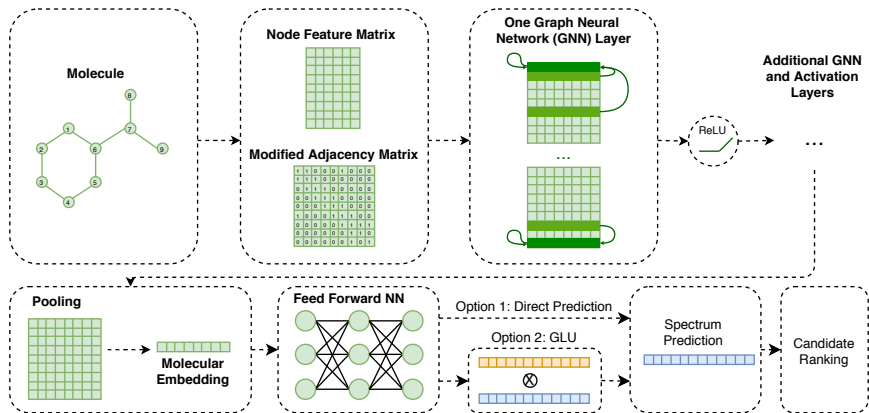


Figure 1: Illustration of the network architecture used in this study

87 We explore using two GNN implementations: GCN [20] and GAT [21]. The propagation rule of a
 88 single GCN layer is given as:

$$h'_i = \sigma \left(\sum_{j \in \mathcal{N}(i)} \frac{1}{c_{ij}} \mathbf{W} h_j + b \right)$$

89 where $\mathcal{N}(i)$ is the set of neighbors of node i and σ is an activation function. c_{ij} is a normalization
 90 term and it is set equal to $\sqrt{|\mathcal{N}(i)| |\mathcal{N}(j)|}$. The normalization penalizes nodes with too many
 91 connections to avoid extreme values. In contrast to GCN, GAT utilizes attention mechanism to
 92 propagate information as follows:

$$h'_i = \sum_{j \in \mathcal{N}(i)} \alpha_{ij} \mathbf{W} h_j$$

93 where α_{ij} is the attention term and equals to $\text{softmax}_i(\sigma(\vec{a}^T [\mathbf{W} h_i || \mathbf{W} h_j]))$. In essence, GAT
 94 introduces additional trainable attention weights, \vec{a}^T , on the concatenated $\mathbf{W} h_i$ and $\mathbf{W} h_j$
 95 vectors, to model a weighting term that controls how the message $\mathbf{W} h_j$ from each neighbor j
 96 should be propagated to node i . The softmax activation ensures that the sum of these weighting
 97 terms equals to 1 for each node. By design, GAT does not have a bias term.

98 After the learned node representation \mathbf{h} is obtained, node information is “pooled” into a graph
 99 embedding vector $\mathbf{v} \in \mathbb{R}^H$, with information about the entire graph. We compared several different
 100 pooling methods including Global Maximum, Global Average and Global Attention [23].

101 After the pooling layer, the graph embedding vector \mathbf{v} is fed into a FFNN that predicts $\hat{\mathbf{y}} \in \mathbb{R}^{1000}$.
 102 We also evaluated a gated linear unit (GLU) [24] instead of a dense layer for the outcome prediction.
 103 A GLU predicts two outputs simultaneously while one of the two outputs is activated by sigmoid and
 104 acts as a gate for the other output. The final output is activated by a ReLU function [25].

105 2.4 Training

106 The model was trained by minimizing the mean square error (MSE) between $\hat{\mathbf{y}}$ and \mathbf{y} after normal-
 107 ization. To reduce overfitting, we used L2 regularization with lambda set to 1.0 and a dropout rate at
 108 0.5. All models were trained on a Nvidia P100 for a maximum of 1,000 epochs using Adam [26]
 109 with early stopping on validation loss with a window size of 15.

110 2.5 Evaluation

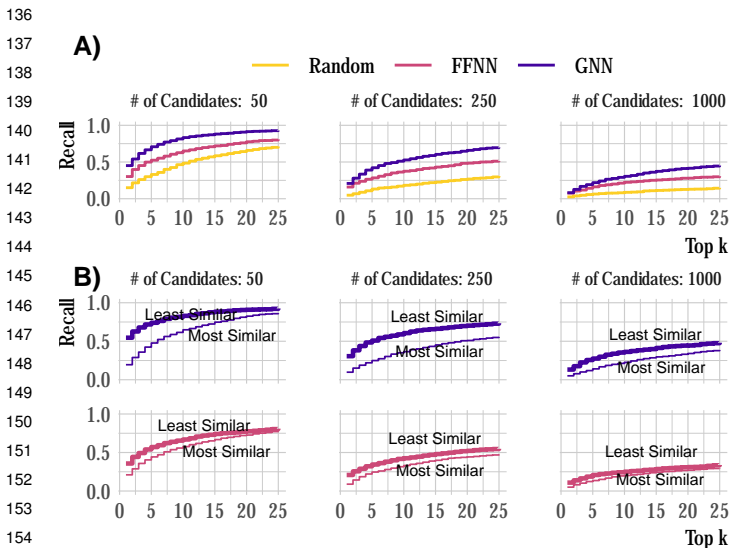
111 We evaluated our models using two metrics. First, the cosine similarity between the predicted and
 112 target spectra is used to assess the quality of the predicted spectra. Second, recall@k, a common
 113 metric for evaluating annotation tools, measures the the portion of correctly identified molecular
 114 identities in the test dataset when considering the top k ranked candidate molecules for each test
 115 spectra. We used the NEIMS model [18] as a baseline.

116 3 Results

117 Compared with the NEIMS
 118 model (Table 1), GNN-based
 119 NN models generate signif-
 120 icantly more accurate spec-
 121 tra and exhibit higher rec-
 122 all@k. In our hyper-
 123 parameter search, we found
 124 that including adjacent bond
 125 information was beneficial.
 126 Log transformation on spec-
 127 tral intensity consistently
 128 provides better performance
 129 than square-root transformation. GAT performs better than GCN, suggesting that the attention
 130 mechanism is effective for this task. We also observed the reported effect that the performance of
 131 GCN quickly dropped as the number of GCN layers increase [22]. However, by using GAT, it was
 132 possible to build deeper GNN models. The best performance happens when there are 10 layers of
 133 stacked GATs with 64 hidden units, with the output predicted by GLU. Among the three pooling
 134 methods, Global Max provided better performance than the other two. After the model was trained, it
 135 achieved ~11,000 predictions per second on an Nvidia V100 GPU card.

Table 1: Summary of results in cosine similarity and recall@k

Experiment	Similarity	Avg Rank	Recall@1	Recall@5	Recall@10
NIST 17 + Candidates from PubChem (Sampled with Average Size = 50)					
NEIMS	0.157	16.7	30.2	52.5	65.1
GCN (3 layers)	0.426	10.8	33.9	62.3	75.9
GAT (3 layers)	0.468	9.2	36.0	65.2	76.6
GAT (10 layers)	0.512	7.1	41.2	70.7	81.4
GAT + GLU	0.517	6.8	45.1	70.8	83.3



156 Figure 2: Impact of: (A) candidate size and (B) candidate similar-
 157 ity on ranking results

158 Ranking on a candidate set using
 159 recall@k is widely used for annotation evaluation. Current evaluation datasets (in terms of test
 160 molecules and their candidate sets), however, vary tremendously. Our results show that ranking
 161 results for both models heavily depend on the candidate set size and the similarity of the candidates
 162 to the target molecule, thus highlighting the need for better and consistent evaluation protocols for
 163 annotation tools.

164 4 Conclusion

165 We investigated several GNN-based models to predict the mass spectra for query molecular structure.
 166 Our model outperforms previously reported NN models. Importantly, we found that ranking results
 167 are heavily dependent on the candidate set size as well as the similarity of candidate molecules
 168 with target molecule. We encourage researchers to standardize performance evaluation for the MS
 169 prediction task, and to consider GNN-based methods for annotation.

To evaluate how the candidate set size impacts ranking performance, we calculate the ratio between the average number of candidates (1,530) and 50/250/1000 and used these ratios to sample a proportion of candidates for each test molecule appropriately. As shown in Figure 3A, the recall@k performances for both models decrease as the average number of candidates increases. With more candidates to choose from, correct identification becomes more challenging. In Figure 3B, where MACCS Fingerprints were used to identify the most similar and least similar compounds in the candidate set, recall@k performances are low when candidate molecules are similar with their target molecule.

References

- 170
- 171 [1] Arnald Alonso, Sara Marsal, and Antonio Julià. Analytical Methods in Untargeted
172 Metabolomics: State of the Art in 2015. *Frontiers in Bioengineering and Biotechnology*,
173 3, mar 2015.
- 174 [2] Naomi L. Kuehnbaum and Philip Britz-McKibbin. New Advances in Separation Science for
175 Metabolomics: Resolving Chemical Diversity in a Post-Genomic Era. *Chemical Reviews*,
176 113(4):2437–2468, apr 2013.
- 177 [3] Colin A Smith, Grace O?? Maille, Elizabeth J Want, Chuan Qin, Sunia A Trauger, Theodore R
178 Brandon, Darlene E Custodio, Ruben Abagyan, and Gary Siuzdak. METLIN: A Metabolite
179 Mass Spectral Database. *Therapeutic Drug Monitoring*, 27(6):747–751, dec 2005.
- 180 [4] David S. Wishart, Dan Tzur, Craig Knox, Roman Eisner, An Chi Guo, Nelson Young, Dean
181 Cheng, Kevin Jewell, David Arndt, Summit Sawhney, Chris Fung, Lisa Nikolai, Mike Lewis,
182 Marie-Aude Coutouly, Ian Forsythe, Peter Tang, Savita Shrivastava, Kevin Jeroncic, Paul
183 Stothard, Godwin Amegbey, David Block, David D. Hau, James Wagner, Jessica Miniaci,
184 Melisa Clements, Mulu Gebremedhin, Natalie Guo, Ying Zhang, Gavin E. Duggan, Glen D.
185 MacInnis, Alim M. Weljie, Reza Dowlatabadi, Fiona Bamforth, Derrick Clive, Russ Greiner,
186 Liang Li, Tom Marrie, Brian D. Sykes, Hans J. Vogel, and Lori Querengesser. HMDB: The
187 Human Metabolome Database. *Nucleic Acids Research*, 35(suppl_1):D521–D526, jan 2007.
- 188 [5] MassBank: A public repository for sharing mass spectral data for life sci-
189 ences - Horai - 2010 - Journal of Mass Spectrometry - Wiley Online Library.
190 <https://onlinelibrary.wiley.com/doi/full/10.1002/jms.1777>.
- 191 [6] Henry Lam, Eric W. Deutsch, James S. Eddes, Jimmy K. Eng, Nichole King, Stephen E. Stein,
192 and Ruedi Aebersold. Development and validation of a spectral library searching method for
193 peptide identification from MS/MS. *Proteomics*, 7(5):655–667, March 2007.
- 194 [7] Christoph Bueschl, Bernhard Kluger, Marc Lemmens, Gerhard Adam, Gerlinde Wiesenberger,
195 Valentina Maschietto, Adriano Marocco, Joseph Strauss, Stephan Bödi, Gerhard G. Thallinger,
196 Rudolf Krska, and Rainer Schuhmacher. A novel stable isotope labelling assisted workflow for
197 improved untargeted LC–HRMS based metabolomics research. *Metabolomics*, 10(4):754–769,
198 aug 2014.
- 199 [8] Accurately Predict Mass Spec Fragmentation | ACD/MS Fragmenter.
200 https://www.acdlabs.com/products/adh/ms/ms_frag/.
- 201 [9] Jiarui Zhou, Ralf J. M. Weber, J. William Allwood, Robert Mistrik, Zexuan Zhu, Zhen Ji, Siping
202 Chen, Warwick B. Dunn, Shan He, and Mark R. Viant. HAMMER: Automated operation
203 of mass frontier to construct in silico mass spectral fragmentation libraries. *Bioinformatics*,
204 30(4):581–583, feb 2014.
- 205 [10] Sebastian Wolf, Stephan Schmidt, Matthias Müller-Hannemann, and Steffen Neumann. In
206 silico fragmentation for computer assisted identification of metabolite mass spectra. *BMC*
207 *bioinformatics*, 11:148, mar 2010.
- 208 [11] Markus Heinonen, Ari Rantanen, Taneli Mielikäinen, Juha Kokkonen, Jari Kiuru, Raimo A.
209 Ketola, and Juho Rousu. FiD: A software for ab initio structural identification of product
210 ions from tandem mass spectrometric data. *Rapid Communications in Mass Spectrometry*,
211 22(19):3043–3052, 2008.
- 212 [12] André Wegner, Daniel Weindl, Christian Jäger, Sean C. Sapcaru, Xiangyi Dong, Gregory
213 Stephanopoulos, and Karsten Hiller. Fragment Formula Calculator (FFC): Determination of
214 Chemical Formulas for Fragment Ions in Mass Spectrometric Data. *Analytical Chemistry*,
215 86(4):2221–2228, feb 2014.
- 216 [13] Nicola Segata, Jacques Izard, Levi Waldron, Dirk Gevers, Larisa Miropolsky, Wendy S. Garrett,
217 and Curtis Huttenhower. Metagenomic biomarker discovery and explanation. *Genome Biology*,
218 12(6):R60, jun 2011.

- 219 [14] Felicity Allen, Allison Pon, Michael Wilson, Russ Greiner, and David Wishart. CFM-ID: A
220 web server for annotation, spectrum prediction and metabolite identification from tandem mass
221 spectra. *Nucleic Acids Research*, 42(Web Server issue):W94–W99, jul 2014.
- 222 [15] Felicity Allen, Russ Greiner, and David Wishart. Competitive fragmentation modeling of
223 ESI-MS/MS spectra for putative metabolite identification. *Metabolomics*, 11(1):98–110, feb
224 2015.
- 225 [16] Kai Dührkop, Huibin Shen, Marvin Meusel, Juho Rousu, and Sebastian Böcker. Searching
226 molecular structure databases with tandem mass spectra using CSI:FingerID. *Proceedings of*
227 *the National Academy of Sciences*, 00÷(41):12580–12585, oct 2015.
- 228 [17] Huibin Shen, Kai Dührkop, Sebastian Böcker, and Juho Rousu. Metabolite identification
229 through multiple kernel learning on fragmentation trees. *Bioinformatics (Oxford, England)*,
230 30(12):i157–164, jun 2014.
- 231 [18] Jennifer N. Wei, David Belanger, Ryan P. Adams, and D. Sculley. Rapid Prediction of Elec-
232 tron–Ionization Mass Spectrometry Using Neural Networks. *ACS Central Science*, 5(4):700–
233 708, apr 2019.
- 234 [19] Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl.
235 Neural Message Passing for Quantum Chemistry. *arXiv:1704.01212 [cs]*, jun 2017.
- 236 [20] Thomas N. Kipf and Max Welling. Semi-Supervised Classification with Graph Convolutional
237 Networks. *arXiv:1609.02907 [cs, stat]*, feb 2017.
- 238 [21] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua
239 Bengio. Graph Attention Networks. *arXiv:1710.10903 [cs, stat]*, feb 2018.
- 240 [22] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S. Yu. A
241 Comprehensive Survey on Graph Neural Networks. *IEEE Transactions on Neural Networks*
242 *and Learning Systems*, pages 1–21, 2020.
- 243 [23] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. Gated Graph Sequence
244 Neural Networks. *arXiv:1511.05493 [cs, stat]*, sep 2017.
- 245 [24] Yann N. Dauphin, Angela Fan, Michael Auli, and David Grangier. Language Modeling with
246 Gated Convolutional Networks. *arXiv:1612.08083 [cs]*, sep 2017.
- 247 [25] Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann
248 machines. In *Proceedings of the 27th International Conference on International Conference on*
249 *Machine Learning, ICML’10*, page 807–814, Madison, WI, USA, 2010. Omnipress.
- 250 [26] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization.
251 *arXiv:1412.6980 [cs]*, jan 2017.