
AN END-TO-END DEEP LEARNING FRAMEWORK FOR TRANSLATING MASS SPECTRA TO DE-NOVO MOLECULES

Eleni E. Litsa¹, Vijil Chenthamarakshan², Payel Das^{2*}, and Lydia E. Kavragi^{1†}

¹Department of Computer Science, Rice University, Houston, TX

²IBM Research, IBM Thomas J. Watson Research Center, Yorktown Heights, NY

ABSTRACT

1 Elucidating the structure of a chemical compound is a fundamental task in chemistry with applications
2 in multiple domains including drug discovery, precision medicine, and biomarker discovery. The
3 common practice for elucidating the structure of a compound is to obtain a mass spectrum and
4 subsequently retrieve its structure from spectral databases. However, these methods fail for novel
5 molecules that are not present in the reference database. We propose Spec2Mol, a deep learning
6 architecture for molecular structure recommendation given mass spectra alone. Spec2Mol is inspired
7 by the Speech2Text deep learning architectures for translating audio signals into text. Our approach
8 is based on an encoder-decoder architecture. The encoder learns the spectra embeddings, while the
9 decoder, pre-trained on a massive dataset of chemical structures for translating between different
10 molecular representations, reconstructs SMILES sequences of the recommended chemical structures.
11 We have evaluated Spec2Mol by assessing the molecular similarity between the recommended
12 structures and the original structure. Our analysis showed that Spec2Mol is able to identify the
13 presence of key molecular substructures from its mass spectrum, and shows on par performance,
14 when compared to existing fragmentation tree methods particularly when test structure information is
15 not available during training or present in the reference database.

16 1 Introduction

17 The identification of the chemical compounds that are present in a sample of chemical matter is a fundamental task in
18 chemical analysis with applications in multiple domains. The field of metabolomics, for example, seeks to identify the
19 chemical molecules that are present in a biological sample. In humans, the metabolome, that is the set of all chemical
20 molecules that can be found in human tissues, is a great source for biomarker discovery as it reflects changes at a
21 genetic, proteomic or environmental level [1]. Additionally, mapping the human metabolome will lead to a better
22 understanding of human physiology and disease etiology and pathology which is essential for the identification of
23 new therapeutic targets for developing new treatments. The increasing interest in mapping the metabolome extends to
24 other organisms as well, such as plants which have been a great source of bioactive compounds for multiple products
25 including drugs and supplements [2]. The identification of chemical compounds is also critical in product development
26 such as in the production of pharmaceuticals and agrochemicals. Structure elucidation practices are being used for
27 quality control and detection of impurities, as well as in safety studies for identifying potential metabolites that can be
28 formed in the human body. Finally, structure elucidation practices are being employed in forensics analysis.

29 The identification of the structure of a chemical compound is perceived as one of the most time consuming and laborious
30 task in chemical analysis. This is often performed through analytical techniques such as mass spectroscopy (MS) and
31 nuclear magnetic resonance (NMR) [3, 4, 5] with MS being used more often due to its higher sensitivity and specificity
32 [3]. In MS, the molecules that are present in a biological sample are first separated using a chromatographic technique,
33 such as liquid chromatography (LC) and gas chromatography (GC), with the latter being used more commonly [1, 6].
34 After the separation, the molecule is fragmented into positive or negative charged ions using an ionization source such
35 as electron ionization (EI), chemical ionization (CI) and electrospray ionization source (ESI) [1, 6]. What the instrument
36 records is the mass-to-charge (m/z) ratios of the generated fragment ions. The information that is collected from this

*daspa@us.ibm.com

†kavraki@rice.edu

37 process is presented in the mass spectrum which is a graph with the m/z of each recorded fragment in the horizontal
38 axis and the relative abundance in the vertical axis. In order to obtain more detailed information on the query structure,
39 a sequential fragmentation process is often used called tandem mass spectrometry [5]. Once the molecule has been
40 fragmented into ions, a set of them, called precursor ions, is selected and further fragmented to generate MS2 (also
41 called MS/MS) spectra. These second-level ions can be fragmented even further giving MS3 spectra and so on. The
42 peaks and their intensity in the resulting spectrum depend not only on the structure of the chemical molecule that is
43 being fragmented, but also on the experimental conditions, that is the instrument used, the collision energy, the selected
44 precursor ion and the ionization mode, as it is illustrated in Figure 1.

45 Once the mass spectrum is obtained, it is matched against the content of spectral databases of reference compounds
46 in order to retrieve its structure. Various chemical databases provide spectra data of metabolites [7] such as Human
47 Metabolome Database, METLIN, MassBank and mzCloud [7]. Certain databases are focused on the metabolites of
48 specific organisms, such as the Human Metabolome Database, or on specific molecular classes, such as the LIPID
49 MAPS Structure Database, while others have greater coverage such as METLIN. However, despite the intense ongoing
50 efforts to map the metabolome of various organisms, existing databases cover only a small percentage of the actual
51 metabolites that occur in organisms. Particularly for humans, it is estimated that less than 10% of metabolites have
52 experimental reference mass spectra [8], which means that the current practice cannot identify a large percentage of the
53 molecules that are found in human tissues. It is estimated that in untargeted metabolomics studies less than 2% of the
54 detected spectral features are identified [8].

55 An approach that has been developed to address the problem of limited amount of experimental spectra data is *in silico*
56 fragmentation which essentially attempts to solve the inverse problem. This approach aims at enhancing the content of
57 existing spectra databases with computed spectra of known molecular structures which have no available experimental
58 spectra. Essentially this approach seeks to close the gap between spectral and structural databases. *In silico* fragmentation
59 tools predict the fragmentation process either relying on fragmentation rules or using combinatorial/optimization-based
60 approaches or employing machine learning methodologies [6, 9, 10]. Fragment prediction methods have been especially
61 successful for predicting spectra of peptides, however, fragmentation of small molecules into ions is a more stochastic
62 process that is especially challenging to predict [6].

63 A more direct approach to the structure elucidation problem would be to reconstruct the underlying chemical structures
64 given spectra features. Such an undertaking though is computationally challenging as it requires the generation of a
65 molecular structure. Indeed, this approach is performed as a two step process to circumvent the need for generating
66 molecular structures: A machine learning model is used to map the spectrum to an intermediate vector representation
67 such as a molecular fingerprint. Once the fingerprint is obtained then it is matched against the content of structural
68 databases in order to identify candidate molecular structures with similar fingerprints [11, 12]. This method though will
69 also fail for molecules that are not present in the structural database and especially for novel molecules. A more direct
70 association of spectra features with molecular structures through a rule-based approach has also been explored [13].
71 More specifically, this approach extracts rules, that associate spectra features with substructures, from spectra databases
72 aiming at a partial structure identification.

73 An additional concept that has been introduced to facilitate the interpretation of mass spectra, and subsequently structure
74 identification, is that of fragmentation trees [6, 14]. A fragmentation tree is derived computationally from tandem
75 mass spectra using optimization algorithms such that its nodes correspond to fragments or precursor ions and the
76 edges correspond to fragmentation reactions. Fragmentation trees have various uses such as identifying the molecular
77 formula and clustering molecules by aligning fragmentation trees [15]. They have also been used for the prediction of
78 molecular fingerprints that are subsequently used to search structural databases [16, 17]. The information in a mass
79 spectrum is thought to be insufficient to explain the fragmentation process by itself while the fragmentation tree provides
80 complementary information by elucidating the dependencies between the mass peaks [6]. However, fragmentation trees
81 are expensive to compute and often approximations are preferred for practical applications.

82 A more thorough review of existing methodologies for metabolite identification, including *in silico* fragmentation tools,
83 fingerprint prediction and fragmentation trees, was recently presented by Nguyen et al. with a focus on machine learning
84 (ML) approaches [6]. It should be noted here that early ML-based approaches were built on shallow ML models, such
85 as Support Vector Machines (SVMs) and Random Forests (RFs), applied either on features extracted from the mass
86 spectra or the fragmentation trees, and also kernel-based methods to determine similarity between either spectra or
87 fragmentation trees. However, lately there is a growing interest in exploring Deep Learning (DL) architectures for the
88 development of computation tools to support structure elucidation. There have been efforts to learn spectra embeddings
89 that can be subsequently used to assess spectral similarity when searching in spectral databases [18, 12]. Additionally,
90 there are DL-based methodologies for clustering spectra, either for identifying the compound class [19, 12] or for aiding
91 medical diagnosis by differentiating between healthy and cancerous tissues [20]. Most DL-based methodologies that
92 operate directly on spectra data are based on Convolutional Neural Networks (CNNs) representing the spectrum as a

93 vector that indicates the intensities of each fragment mass [20, 21, 22]. The CNN attempts to automatically identify
94 spectra features replacing the need for manual featurization. Architectures that have adopted concepts from Natural
95 Language Processing (NLP) have also emerged representing the mass spectrum as text and the mass peaks as words
96 [18]. Due to the limited amount of mass spectra data, different workarounds have been investigated including hybrid
97 approaches [19], combining statistical ML models and DL architectures, and approaches based on transfer learning
98 [20].

99 It should be noted that, at the same time, DL-based approaches are being investigated for identifying protein sequences
100 from mass spectra in proteomics studies [23, 21, 22]. A noteworthy effort, DeepNovo, consists of an end-to-end
101 DL architecture for de novo peptide sequencing from mass spectra [22], that is a direct reconstruction of the peptide
102 sequence from the mass spectra data. Structure elucidation of small molecules though is perceived as a more challenging
103 problem due to the stochastic nature of the fragmentation process. On top of that, the structure of small molecules has a
104 graph-like representation as opposed to the linear nature of a peptide sequence. Existing approaches essentially attempt
105 to retrieve molecules from structure databases that have a spectrum similar to the query spectrum. This method though,
106 cannot identify novel molecules, that is molecules whose structure currently remains unknown and therefore they do
107 not exist in chemical databases.

108 In this paper, we present Spec2Mol, an end-to-end DL architecture for translating MS/MS spectra to molecular
109 structures. Spec2Mol is intended for recommending molecular structures that can explain observed MS/MS spectra.
110 We represent molecular structures as sequences using the SMILES notation [24] and MS/MS spectra as vectors of
111 fragment intensities. Spec2Mol consists of an encoder, that learns an embedding for the MS/MS spectrum, and a
112 decoder that generates the SMILES sequences of the recommended chemical molecules. Due to the limited amount of
113 available spectra data, our approach is based on unsupervised pre-training on a large dataset of unlabeled molecules. In
114 particular, we pre-trained the decoder as part of an auto-encoder (AE) architecture which is trained to reconstruct a
115 molecule through its SMILES sequence. The encoder is subsequently trained such that the spectra embeddings match
116 the embeddings that the AE has learnt. In the following sections, we discuss the data used to develop and evaluate the
117 model, the architecture of Spec2Mol, as well as, the evaluation of the model.

118 The main contributions of this work are as follows:

- 119 • To our knowledge, this is the first approach for generating potential molecular structures from mass spectrometry
120 data that is not based solely on database retrieval.
- 121 • Our method can facilitate database retrieval and additionally de novo molecular structure recommendation.
- 122 • Our approach takes advantage of large datasets of unlabeled molecules using unsupervised pre-training.
- 123 • We introduce metrics to assess the similarity of the generated molecules with the reference ones and we
124 perform a comparative evaluation with a widely accepted method that makes use of additional information,
125 that is fragmentation trees.

126 2 Results and Discussion

127 2.1 Reconstruction accuracy of the autoencoder

128 As a sanity check, we evaluated the ability of the pre-trained AE to reconstruct the SMILES of the molecules in
129 the testing set of the spectra dataset. This is performed by comparing the canonicalized input SMILES and the
130 canonicalized output SMILES and evaluating whether there is an exact match between the two. The autoencoder is
131 trained by minimizing the mean reconstruction error on a single-character level for each input sequence. Therefore,
132 the reconstruction accuracy is estimated on a single-character level, by comparing the correct character in the target
133 sequence with the most probable character in the decoder RNN’s output at each position. It should be noted, that the
134 reconstructed SMILES, as well as neural fingerprints derived from SMILES [25, 26, 27], has been successfully used in
135 similarity search and have been found to be more informative, when compared to molecular fingerprints.

136 The AE was able to correctly reconstruct the SMILES sequence for about 93.3% of the NIST molecules. This is
137 very close to the reconstruction rate of the AE on a held out test set which was 94.95%. This demonstrates that the
138 pre-trained model has been trained on a diverse set of molecules and therefore it is able to handle the large variability of
139 the molecules in the NIST dataset.

140 2.2 Spec2Mol performance evaluation

141 Spec2Mol generates a set of recommended molecular structures given MS/MS spectra. Our evaluation focuses on
142 assessing the similarity between the generated structures and the reference molecular structure from the NIST dataset.

143 We recall here that the information in an MS/MS spectrum may not be sufficient to fully reconstruct the molecular
144 structure. It is possible that more than one molecular structures may explain a given spectrum. For that reason our
145 analysis has been focused on assessing whether the model has learnt to identify key features in the molecular structure
146 from the mass spectra rather than identifying the exact same structure with the reference molecule from the NIST
147 dataset.

148 For the evaluation of the model, we first perform a coarse-level comparison taking into account physicochemical
149 properties and more specifically the molecular weight and the element composition of the molecule. Next, we assess
150 molecular similarity at the substructure level. In particular, we compute the fingerprint similarity as well as the
151 maximum common substructure between the generated structures and the reference structure. The specifications for
152 each metric are given below. We evaluate the overall performance in the entire test set as well as the performance of the
153 model when not all four required spectra are available as input. Additionally, we assess the contribution of each of the
154 two strategies for generating the recommended structures.

- 155 • **Physicochemical attributes:** A property of special interest is the molecular weight since it is directly reflected
156 in the mass spectrum. In particular, the spectra indicates the mass of the fragments and therefore the mass of
157 the original, non-fragmented, molecule can be approximated more easily given the mass spectra as opposed to
158 determining the composition or the structure of the molecule. We record the difference between the molecular
159 weight of the generated structures and the reference structure and we report the relative average-minimum
160 difference, that is, the average-minimum difference over all the predicted structures divided by the average
161 molecular weight of the reference structures (DMW_{min}). We also report the average-average difference over
162 all the predicted structures divided by the average molecular weight of the reference structures (DMW_{avg}).
163 Additionally, we also evaluate whether the model is able to identify the element composition of the molecule.
164 In particular, we assess whether the atom species that are present in the reference molecule have been identified
165 in the predicted structures ignoring the numbers of atoms for each atom species. More specifically, for each
166 atom species we report sensitivity and specificity for detecting the presence of this species. In order to account
167 for discrepancies in the number of atoms per atom species, we also report the difference between the molecular
168 formulas of the predicted structures and the reference structure (DMF). We define the distance between
169 two molecular formulas as the number of atoms that differ between two molecules when accounting for the
170 atom species and the number of atoms for each species (without including hydrogen atoms). We report the
171 minimum distance over all predictions divided by the average number of heavy atoms (DMF_{min}) as well as
172 the average distance over all predictions divided by the average number of heavy atoms (DMF_{avg}). The exact
173 mathematical formulas for the calculation of the DMW and DMF are provided in the supplementary material
174 (Supplementary Methods 3).
- 175 • **Fingerprint similarity:** Fingerprints are vector representations of chemical molecules, which indicate the
176 presence of certain substructures in the molecule, and are widely used as an efficient way to judge similarity
177 between molecules [28]. We extracted fingerprint representations based on the Morgan algorithm [29] using the
178 RDKit toolkit [30] and used the cosine coefficient to assess similarity ($Fngp_{cosine}$). The Morgan fingerprints
179 are computed for radius 2 and 1024 bits. We report the maximum fingerprint similarity among all model
180 predictions when compared with the reference structure as well as the average similarity of all predicted
181 structures.
- 182 • **Maximum common substructure (MCS):** We computed the MCS between two molecular structures using the
183 RDKit toolkit [30] with the following constraints: the substructure match respects the atom species, the bond
184 orders, as well as the ring bonds, that is ring bonds are only matched to ring bonds. From the computed MCS we
185 extracted the following three metrics: i) MCS ratio, ii) MCS Tanimoto, and iii) overlap coefficient, which are
186 defined as follows, respectively: $MCS_{ratio} = \frac{a_{MCS}}{a_r}$, $MCS_{tan} = \frac{a_{MCS}}{a_r + a_p - a_{MCS}}$, $MCS_{ovrlp} = \frac{a_{MCS}}{\min(a_r, a_p)}$, where
187 a_{MCS} denotes the number of atoms in the MCS, a_r the number of atoms in the reference compound, and a_p
188 the number of atoms in the predicted compound. For each metric, we report the maximum value as well as the
189 average value over all predictions.

190 Table 1 summarizes the evaluation of the effect of missing data in the predictions. More specifically, we present the
191 evaluation metrics on four different partitions of the test-set depending on the number of the available spectra. We
192 recall that the input to the model consists of four different spectra obtained through different specifications. However,
193 not all molecules in the dataset have all four spectra available. Our results indicate that missing only one spectrum
194 does not severely impact performance, but performance starts to degrade when less than three spectra are available.
195 This is expected as the number of spectral peaks that will be observed in one spectrum (or two) most likely will not be
196 adequate to reconstruct the molecular structure. It should be noted though that other factors, such as the molecular
197 size, are also potentially contributing to the variability observed among the different subsets of the test-set. The set of
198 molecules with three available spectra for example, includes molecules that on average have smaller molecular weight

199 and shorter SMILES representation. The model appears to have the highest performance on this subset of the test-set
200 since reconstructing shorter SMILES is expected to be less of a challenge for the decoder. The evaluation of the model
201 on the training set is presented in the supplementary material (Supplementary Note 1, Table S3).

202 Next, we evaluate the effect of the strategy that is used to generate the recommended molecules. The analysis is shown
203 in Table 2. We recall that the recommended structures are obtained either directly through decoding the computed
204 embeddings or indirectly by identifying the closest embeddings from the pre-trained dataset. In particular, we are
205 comparing the top-20 predictions, as ranked using the molecular weight criterion, through i) only the direct strategy, ii)
206 only the indirect strategy, and, iii) the two strategies combined. According to the results, the indirect approach, that
207 generates molecules through decoding the closest embeddings from the pre-trained dataset appears to have a larger
208 contribution on the effectiveness of the method to generate relevant structures. However, combining the two strategies
209 appears to slightly improve performance.

210 Overall, the results illustrate that the predicted structures have a molecular weight that is significantly close to the
211 molecular weight of the reference compound. This is not surprising as the generated molecules are ranked based on
212 the molecular weight. The molecular formula though seems to also be considerably close to the reference one. The
213 model was able to retrieve the exact structure for a small percentage of the test cases (7%) while it identified the exact
214 molecular formula for a considerably larger percentage (26%). The performance of the model was significantly better
215 when at least 3 out of the 4 input spectra where available.

216 Regarding the structural similarity between the predicted structures and the reference structure, the obtained values for
217 the respective metrics demonstrate that the structures share common substructures. More specifically, the metrics that
218 are based on the MCS between the reference and the predicted structures indicate that the common substructure is, on
219 average, nearly 70% of the size of the reference structure for the closest structure and more than 50% for the average
220 prediction. This result is in agreement with the high correlation between the molecular fingerprints.

221 Regarding the ability of the model to identify the presence of each atom species in the molecular structure, it varies
222 significantly and it correlates with the frequency of each atom species in the training dataset, as it is shown in Table 3.
223 More specifically, the model has very high sensitivity for nitrogen (N) and oxygen (O) which are the most common atom
224 species in the dataset (excluding carbon which is not included in this analysis as it is present in all molecules). However,
225 the specificity for oxygen is significantly lower than that of nitrogen which means that there is a significant number of
226 false positives for oxygen compared to nitrogen. Regarding the more rare atom species, the opposite phenomenon is
227 observed: specificity is significantly high while sensitivity is low. This means that for the rare species there is a very
228 small number of false positives which is expected as these atoms are under-represented in the training set. However,
229 sensitivity is at least 0.5 for all atoms, which shows that the model is able to capture the presence of rare atoms quite
230 well considering that some atom species are severely under-represented in the training set.

231 Finally, we investigated the effect of the molecular weight as well as the presence of heteroatoms on the ability of the
232 model to identify the exact structure or the exact molecular formula. More specifically, we divided the test set molecules
233 into those that have molecular weight (MW) less than 300Da and those that have molecular weight greater than or equal
234 to 300Da (the average molecular weight in the test set is 275Da). Furthermore, we created four categories based on
235 the presence of heteroatoms: 1) molecules that have only C and O, 2) molecules in which N is present, 3) molecules
236 in which S is present, and, 4) molecules in which a halogen (one of Br, Cl, F, I) is present. Table 4 summarizes this
237 analysis. The model is able to identify the atom species and atom counts for almost half of the molecules (45.4%) with
238 MW less than 300Da and for more than 60% of the molecules that contain only C and O (63.6%). Higher molecular
239 weight as well as presence of atoms that are under-represented in the training set (S and halogens) degrades the ability
240 of the model to identify the molecular structure or formula.

241 Figure 3 shows a few examples of successful cases with the model correctly identifying key substructures such as rings
242 and long chains, and the presence of rare atoms and functional groups. Given the vast space of possible molecular
243 structures, these cases demonstrate that the model has indeed learnt to associate spectra features with molecular
244 structures.

245 We also identify two general scenarios where the model has a difficulty in predicting relevant structures: (1) Molecules
246 with large rings and (2) Molecules that have poor quality spectra. An example of the first case is illustrated in Figure 4.
247 We believe this is because molecules with large rings are significantly under-represented in the dataset that was used
248 to pre-train the decoder. Also, it is hard to generate a valid SMILES sequence for molecules with very large rings.
249 Regarding the second cases of poor quality input spectra, it includes cases where there is a very small number of peaks
250 in the spectra and therefore not adequate information to reconstruct the SMILES sequence.

251 2.3 Comparative evaluation

252 In order to perform a comparative evaluation, we have used SIRIUS 4 [31], which offers multiple functions including
253 chemical formula, as well as molecular structure, identification from mass spectra. SIRIUS' structure elucidation
254 method, called CSI:FingerID, is a database retrieval method [16]. It relies on Support Vector Machines (SVMs) for
255 predicting a molecular fingerprint and subsequently compares the predicted fingerprint against those of a reference
256 database in order to identify candidate structures. The input to the SVM is the MS/MS spectrum along with the
257 corresponding computed fragmentation tree. CSI:FingerID has shown superior performance when compared to other
258 existing tools for automatic identification of molecular structures from spectra data. In particular, it was the best
259 performing method in the Critical Assessment of Small Molecule Identification (CASMI) contest for 2016 and 2017
260 [31]. However, the performance of this method degrades significantly for cases that are not covered in the training set
261 [31]. Additionally, the dependence of CSI:FingerID on fragmentation tree data adds significantly to the running time of
262 this method.

263 We run SIRIUS on the same test set we developed for evaluating Spec2Mol. As input, we provided SIRIUS with
264 the positive mode spectra (that is [M+H]⁺ at low and high energy) as they were selected for Spec2Mol. The spectra
265 from negative ions were not used since a single run for SIRIUS accepts spectra from a single precursor which may
266 be obtained through different energies. As 53 test cases out of the 1000 cases of the test set did not have any positive
267 mode spectra and therefore the test set used for the comparison consists of 947 cases. As a side note, SIRIUS performs
268 structure elucidation after identifying the molecular formula. The number of molecular formulas to be explored is one
269 of the parameters of the tool which we set to 10. An additional parameter is the reference database which we set to
270 PubChem, which is the largest available source offered by SIRIUS. Finally, SIRIUS allows the user to define the set of
271 chemical elements to be considered when performing the search which we set to: C, H, O, N, S, Cl, F, Br, P and I. It
272 should be noted that expanding the pre-defined set of atoms (C, H, N, O, P, S) to account for more rare atoms, which
273 were present in the NIST dataset, significantly increased the running time.

274 On the test set of 947 cases, SIRIUS found the correct formula for about 98% of the test cases while it found the correct
275 structure for about 67%. For 6 cases out of 947 SIRIUS did not return any structures. It should be highlighted that the
276 CSI:FingerID method from SIRIUS has been trained on the NIST dataset (NIST v17). As it is discussed in the original
277 study on the SIRIUS tool, the presence of spectra for a given test structure in the training set can significantly boost
278 performance even if the spectra that are used when testing are not the exact same spectra as the ones used in training
279 [31].

280 The comparative evaluation between SIRIUS and Spec2Mol was performed on the cases where SIRIUS failed to find
281 the exact molecular structure. Since Spec2Mol is intended for recommending potential molecular structures given mass
282 spectra, our intention here is to evaluate how relevant the recommendations are, when compared to a widely accepted
283 and state-of-the-art method like SIRIUS. By focusing our comparison on the cases where SIRIUS did not find an exact
284 match, we are essentially evaluating the relevance of the recommended structures when an exact match is not found,
285 which points to the case of novel molecules. In particular, we compared SIRIUS and Spec2Mol on the 307 cases, for
286 which SIRIUS failed to find an exact match, using the metrics based on fingerprint similarity and MCS. It should be
287 noted here that failure to identify the exact structure includes cases where SIRIUS either did not return any structure
288 as well as cases where the reference structure was not among the predicted structures. The results are summarized in
289 Table 5. The comparison on the full test set (including cases where SIRIUS found the exact structure) is provided in the
290 supplementary material (Supplementary Note 2, Table S4). According to our analysis, the structures recommended by
291 Spec2Mol are at least as relevant as the ones recommended by SIRIUS. More specifically, Spec2Mol achieved slightly
292 better cosine similarity for the closest structure, while almost all metrics based on the MCS are improved in the case of
293 Spec2Mol. This outcome is especially interesting and encouraging, given that Spec2Mol is an end-to-end approach that
294 does not take into account any prior knowledge. Spec2Mol generates potential molecular structures by solely looking at
295 raw MS/MS spectra. On the other hand, the combination of CSI:FingerID and SIRIUS attempts to retrieve the exact
296 molecular structure from a reference database taking as input the computed fragmentation tree on top of the raw mass
297 spectra. It should be stressed that a direct comparison of the two methods is not possible since they differ significantly:
298 CSI:FingerID uses predicted fingerprints from the MS/MS spectrum of an unknown compound to find the best match
299 against a chemical structure database, while Spec2Mol aims for de-novo generation of potential molecular structures
300 rather than attempting a best match retrieval from a database. Therefore, Spec2Mol is useful in situations where a
301 reference database is not available or CSI-FingerID cannot find an exact match. For that reason, the comparison is
302 performed on the cases where CSI-FingerID failed to identify the exact structure and the metrics used aim at evaluating
303 molecular similarity rather than exact matches.

304 Still the outcome of our comparative evaluation demonstrates that the molecular structures generated by Spec2Mol are
305 at least as successful as the ones obtained by state-of-the-art tools when considering novel molecules despite the fact
306 that Spec2Mol relies solely on raw MS/MS spectra.

307 **3 Conclusions**

308 Elucidating the structure of chemical compounds is a fundamental, but cumbersome, task in metabolomics studies,
309 as well as in chemical analysis in various domains including drug development and forensics analysis. The available
310 computational tools for aiding structure elucidation are based on fragment annotation and database retrieval methods.
311 This approach fails to identify molecules that are not present in the reference database which, in practice, may
312 correspond to a considerably large percentage of the query spectra. We have developed Spec2Mol, an end-to-end deep
313 learning architecture for directly generating molecular structures (as SMILES sequences) from the input MS/MS spectra.
314 Spec2Mol is based on an encoder-decoder architecture that generates molecular SMILES sequences, given mass spectra.
315 While the proposed architecture supports the retrieval of molecules from a database that best matches the input spectra,
316 it can also generate new molecules that have not been seen before in any dataset. Our analysis demonstrates that
317 the recommended molecules are structurally, and physiochemically, similar to the reference compounds, suggesting
318 that the latent space has indeed learnt informative associations between the spectra and the structural features. When
319 compared to an existing method that depends on the fragmentation tree annotation, on top of the raw spectra for
320 molecule identification, Spec2Mol performed on par for the task of recommending potential molecular structures.
321 Our results indicate that the proposed approach of recommending de-novo molecules directly from input MS spectra
322 provides critical insights on the characteristics of the underlying molecular structure, and, can complement existing
323 tools especially when the current tools fail to identify the right molecule from existing databases. We speculate that
324 incorporating prior knowledge in the model, for example in the form of fragmentation trees, can further boost the
325 performance of the proposed method. Further, even though the main focus of our work is on de-novo generation
326 of molecules given an input spectrum, the indirect method proposed by our paper can be extended to identify the
327 correct molecule from a library of a plausible set of molecules, similar to the work proposed by Lim et. al [32]. A
328 substructure-constrained similarity search or a nearest neighbor search on the embeddings of the molecule library with
329 the spectra embedding as a query can be used to identify the best candidates from a relevant library.

330 **4 Methodology**

331 Spec2Mol consists of an encoder that learns spectra embeddings and a pre-trained decoder, which has been trained as
332 part of an autoencoder architecture. The autoencoder has been trained on a large set of molecules (molecule dataset
333 discussed in section 4.1.1), while the encoder has been trained on a set of molecules for which MS/MS data are available
334 (spectral dataset discussed in section 4.1.2).

335 **4.1 Datasets**

336 **4.1.1 Molecule dataset**

337 The autoencoder, from which the Spec2Mol decoder has been derived, was pre-trained on about 135 million molecules
338 which were sourced from the PubChem [33] and ZINC-12 [34] datasets. The structures of these molecules are repre-
339 sented using the SMILES notation [24]. Stereochemistry information was not indicated in the SMILES representation.
340 The reason for not accounting for stereochemistry is that, in the subsequent task of spectra translation, recovering
341 stereochemistry information from the mass spectra is especially challenging or possibly even impossible and therefore
342 it is out of the scope of this work.

343 **4.1.2 Spectral dataset**

344 The mass spectra data for training the encoder has been derived from the NIST Tandem Mass Spectral Library
345 2020 which is a commercial dataset of more than 1M spectra obtained from more than 30K compounds [35, 36].
346 The largest percentage of the NIST dataset (60%) corresponds to metabolites (6K human metabolites and 8K plant
347 metabolites) while a significant amount of the data is drugs (20%). The rest corresponds to peptides, lipids, forensics,
348 surfactants/contaminants and sugars/glycans. The dataset contains low and high resolution MS/MS spectra, obtained
349 through different fragmentation techniques. Each molecule in the dataset may be associated with more than one spectra
350 which may be obtained through different experimental conditions, that is, different fragmentation instrument, precursor
351 ion, ionization mode, collision energy or fragmentation level (MS2, MS3 or MS4). Statistics of the dataset regarding
352 common molecular properties (e.g. molecular weight, number of atoms and number of rings), as well as the atom
353 species coverage, are presented in the supplementary material (Supplementary Methods 1, S1.2, Tables S1-S2).

354 4.2 Data processing and representation

355 In order to minimize variations in the spectra data, due to differences in the experimental conditions, we chose to keep
356 certain variables in the dataset fixed. Details on the filtering process that we followed for constructing the spectral
357 dataset are provided in the supplementary material (Supplementary Methods 1, S1.1). More importantly, we used only
358 the spectra that are obtained through the most common precursor ions, that is $[M+H]^+$ and $[M-H]^-$. For each precursor
359 ion, we used two spectra, one obtained using low collision energy (35% NCE) and one with high collision energy (130%
360 NCE). Therefore, each instance in the dataset we constructed is characterized by four MS/MS spectra derived from two
361 different precursor ions and two energy levels. The four spectra constitute the input to the spectra encoder as described
362 in paragraph 3.2. It should be highlighted though, that not all molecules in the NIST dataset have experimental data for
363 the specific precursors and energy levels. However, we have allowed cases with missing data in the dataset and the
364 missing spectra are represented as empty spectra, that is spectra with no peaks, in an attempt to develop a model that is
365 robust to missing data. Therefore, the model is being trained and evaluated on cases that may not have available all four
366 spectra.

367 4.2.1 Data representation

368 We represent each MS/MS spectrum as a vector in which each bit corresponds to a specific mass-over-charge (m/z)
369 value, representing the m/z value of the recorded fragments, while the value of each bit corresponds to the intensity,
370 or otherwise frequency, of the fragments that have been recorded with that specific mass-over-charge value. We have
371 normalized the intensity values by dividing with the maximum intensity over all the vector bits of a given spectrum.
372 More details on the representation of the MS/MS spectra are provided in the supplementary material (Supplementary
373 Methods 1, S1.3). Regarding the molecular structures, we represent them using canonical SMILES without indicating
374 stereochemistry information.

375 4.2.2 Data augmentation

376 The variability in the spectra for a given molecule opens up the possibility for data augmentation. In particular, although
377 some spectra from the same molecule may differ significantly, as shown in Figure 1, in many cases the obtained spectra
378 are closely related. One such case is when the collision energies that are being used are relatively close.

379 In order to augment the dataset, for each instance in the training set we are creating an additional training instance
380 by slightly perturbing the collision energy in all four spectra. In particular, each spectrum, out of the four spectra
381 that are used to represent an instance in the dataset, is replaced with a spectrum that has the closest collision energy
382 in the dataset while all other parameters (precursor ion, instrument) are shared. More information is provided in the
383 supplementary material (Supplementary Methods 1, S1.4).

384 4.2.3 Data partition

385 After the data filtering process, the acquired dataset consists of 23K molecules, each one of them is associated with
386 four MS/MS spectra, or more precisely, up to four MS/MS spectra given that there are cases with missing spectra.
387 This dataset was partitioned into a training, a validation and a test set with the validation and test set having about 1K
388 molecules each. For the test set specifically, we used fingerprint similarity, based on the Tanimoto coefficient [28], in
389 order to ensure that no test molecule is either in the train or in the validation set. The validation set was used to select
390 the model hyper-parameters and the test set was used to evaluate the performance of the model.

391 4.3 Spec2Mol architecture

392 Spec2Mol uses an encoder-decoder architecture for recommending molecular structures from MS/MS spectra. The
393 Spec2Mol encoder generates spectra embeddings while the decoder reconstructs the SMILES sequence from a spectra
394 embedding. The encoder and the decoder have been trained separately as it is shown in figure 2. First, the decoder is
395 trained as part of an autoencoder architecture for reconstructing the SMILES sequence from a SMILES embedding. Next,
396 the spectra encoder is trained such that the learnt spectra embeddings match the corresponding SMILES embeddings.
397 Finally, for making inference on unseen cases, Spec2Mol uses the spectra encoder to obtain the spectra embedding
398 which is subsequently used in order to decode potentially novel molecules and also to retrieve molecules from the
399 pre-training dataset.

400 The specifications for training each model are given in the following paragraphs while more details on the architectures
401 of the models, hyperparameters and training parameters are provided in the supplementary material (Supplementary
402 Methods 2).

403 **4.4 Pre-training the AE on chemical structures**

404 The autoencoder is trained on a translation task where a randomized input SMILES is translated into its corresponding
405 canonical SMILES, similar to the work of Winter et al [25]. The encoder and the decoder of the AE are both based
406 on gated recurrent units (GRU) which is a variation of the standard long short term memory (LSTM) models, that are
407 commonly used for learning sequence representations, with fewer parameters. The details regarding the autoencoder
408 architecture are in the supplementary material (Supplementary Methods 2, S2.1).

409 **4.5 Training the spectra encoder**

410 The spectra encoder is trained in a supervised manner such that the learnt spectra embeddings are the same as the
411 SMILES embeddings that the AE has learnt. More specifically, the input of the spectra encoder consists of the four
412 spectra that have been pre-selected to represent each molecule. The spectra encoder is based on 1-D CNNs and in
413 particular consists of two 1-D CNN layers and two fully connected layers. The four spectra are represented as 4 discrete
414 vectors which are fed into the 1-D CNN as data from four different channels. Each channel corresponds to a specific
415 precursor ([M+H]⁺ or [M-H]⁻) and energy level (low or high). If any of the required four spectra is not available, then
416 the input to the respective channel is an all-zeros vector. The output of the spectra encoder is a 1-D vector which is the
417 latent representation of the spectra in the embedding space. The model is trained such that the distance (root mean
418 square error) between the latent representation that is learnt by the spectra encoder and the latent representation that is
419 obtained from the pre-trained SMILES encoder is minimized. Details regarding the architecture and training of the
420 spectra encoder are provided in the supplementary material (Supplementary Methods 2, S2.2).

421 **4.6 Recommending molecular structures for unseen spectra**

422 Spec2Mol provides as output molecular structures that can potentially explain the observed spectra peaks. The
423 recommended molecules for unseen spectra are obtained using two strategies: a direct and an indirect molecule
424 generation strategy. The direct molecule generation strategy generates molecular structures using the SMILES decoder
425 from the computed MS/MS embedding. Multiple SMILES are generated for each MS/MS embedding using a pure
426 sampling strategy [37], and subsequently filtered in order to retain only the valid ones, i.e., the sequences that are in
427 accordance with the SMILES syntax. The indirect strategy retrieves molecular structures from the dataset that was used
428 for pre-training the AE based on the distance in the embedding space. More specifically, for each MS/MS embedding
429 we find the closest embeddings from the pool of molecules used to pre-train the AE and decode those embeddings into
430 SMILES sequences.

431 The predicted molecules obtained through these two strategies are combined and ranked based on their discrepancy
432 from the expected molecular weight. The molecular weight of the underlying chemical structure is easily inferred from
433 the mass spectrum and therefore in this work we consider it as known. The molecular structures that have molecular
434 weight closer to the reference weight are highly ranked. The top-20 ranked predictions are returned to the user.

435 **5 Data availability**

436 The spectra dataset used for training and evaluating the model cannot be made publicly available as it is a commercial
437 dataset.

438 **6 Code availability**

439 The trained models and code are available in <https://github.com/KavrakiLab/Spec2Mol> .

440 **7 Acknowledgements**

441 E.E.L and L.E.K have been supported in part by Rice University funds. L.E.K has also been supported by NIH
442 U01CA258512. C.V. and P.D. have been supported by IBM Research.

443 **8 Author Contributions**

444 1) Conceptualization: E.E.L., 2) Methodology: E.E.L., P.D., and L.E.K., 3) Software: E.E.L for Spectra Encoder and
445 V.C. for Pre-trained Autoencoder, 4) Data analysis: E.E.L., 5) Interpretation of results: All authors, 6) Visualization,

446 figures and tables: E.E.L. 7) Supervision: P.D. and L.E.K., 8) Manuscript—original draft: E.L. Review and editing: All
447 authors. All authors approved the manuscript.

448 9 Competing interests

449 The authors declare no competing interests.

450 References

- 451 [1] S. Nalbantoğlu. Metabolomics: Basic principles and strategies. In S. Nalbantoğlu and H. Amri, editors, *Molecular*
452 *Medicine*. IntechOpen, 2019.
- 453 [2] S. Lee, D.G. Oh, D. Singh, J.S. Lee, S. Lee, and C.H. Lee. Exploring the metabolomic diversity of plant species
454 across spatial (leaf and stem) components and phylogenic groups. *BMC Plant Biology*, 20(1), 2020. PMID:
455 31992195; PMCID: PMC6986006.
- 456 [3] A.H. Emwas. The strengths and weaknesses of NMR spectroscopy and mass spectrometry with particular focus
457 on metabolomics research. *Methods in Molecular Biology*, pages 161–193, 2015.
- 458 [4] D.S. Wishart. Computational strategies for metabolite identification in metabolomics. *Bioanalysis*, 1(9):1579–
459 1596–193, 2009.
- 460 [5] Diogo Ribeiro Demartini. A short overview of the components in mass spectrometry instrumentation for proteomics
461 analyses. In Ana Varela Coelho and Catarina de Matos Ferraz Franco, editors, *Tandem Mass Spectrometry -*
462 *Molecular Characterization*. IntechOpen, 2013.
- 463 [6] Nguyen DH, Nguyen CH, and Mamitsuka H. Recent advances and prospects of computational methods for
464 metabolite identification: a review with emphasis on machine learning approaches. *Briefings in Bioinformatics*,
465 20(6):2028–2043, 2019.
- 466 [7] Maria Vinaixa, Emma L. Schymanski, Steffen Neumann, Miriam Navarro, Reza M. Salek, and Oscar Yanes. Mass
467 spectral databases for lc/ms- and gc/ms-based metabolomics: State of the field and future prospects. *TrAC Trends*
468 *in Analytical Chemistry*, 78:23–35, 2016.
- 469 [8] D.S. Wishart, Y.D. Feunang, A. Marcu, A.C. Guo, K. Liang, R. Vázquez-Fresno, T. Sajed, D. Johnson, C. Li,
470 N. Karu, Z. Sayeeda, E. Lo, N. Assempour, M. Berjanskii, S. Singhal, D. Arndt, Y. Liang, H. Badran, J. Grant,
471 A. Serra-Cayuela, Liu Y., R. Mandal, V. Neveu, A. Pon, C. Knox, M. Wilson, C. Manach, and A. Scalbert. HMDB
472 4.0: the human metabolome database for 2018. *Nucleic Acids Research*, 46(D1):D608–D617, 2018.
- 473 [9] Yannick Djoumbou-Feunang, Allison Pon, Naama Karu, Jiamin Zheng, Carin Li, David Arndt, Maheswor Gautam,
474 Felicity Allen, and David S. Wishart. Cfm-id 3.0: Significantly improved esi-ms/ms prediction and compound
475 identification. *Metabolites*, 9(4), 2019.
- 476 [10] Jennifer N. Wei, David Belanger, Ryan P. Adams, and D. Sculley. Rapid prediction of electron–ionization mass
477 spectrometry using neural networks. *ACS Central Science*, 5(4):700–708, 2019.
- 478 [11] M. Heinonen, H. Shen, N. Zamboni, and J. Rousu. Metabolite identification and molecular fingerprint prediction
479 through machine learning. *Bioinformatics*, 28:2333–2341, 2012. PMID: 22815355.
- 480 [12] Hongchao Ji, Hanzi Deng, Hongmei Lu, and Zhimin Zhang. Predicting a molecular fingerprint from an electron
481 ionization mass spectrum with deep neural networks. *Analytical Chemistry*, 92(13):8649–8653, 2020. PMID:
482 32584545.
- 483 [13] Youzhong Liu, Aida Mrzic, Pieter Meysman, Thomas De Vijlder, Edwin P. Romijn, Dirk Valkenburg, Wout
484 Bittremieux, and Kris Laukens. Messar: Automated recommendation of metabolite substructures from tandem
485 mass spectra. *PLOS ONE*, 15(1):1–17, 01 2020.
- 486 [14] Arpana Vaniya and Oliver Fiehn. Using fragmentation trees and mass spectral trees for identifying unknown
487 compounds in metabolomics. *TrAC Trends in Analytical Chemistry*, 69:52–61, 2015.
- 488 [15] Florian Rasche, Kerstin Scheubert, Franziska Hufsky, Thomas Zichner, Marco Kai, Aleš Svatoš, and Sebastian
489 Böcker. Identifying the unknowns by aligning fragmentation trees. *Analytical Chemistry*, 84(7):3417–3426, 2012.
490 PMID: 22390817.
- 491 [16] Kai Dührkop, Huibin Shen, Marvin Meusel, Juho Rousu, and Sebastian Böcker. Searching molecular structure
492 databases with tandem mass spectra using csi:fingerid. *Proceedings of the National Academy of Sciences*,
493 112(41):12580–12585, 2015.

- 494 [17] Huibin Shen, Kai Dührkop, Sebastian Böcker, and Juho Rousu. Metabolite identification through multiple kernel
495 learning on fragmentation trees. *Bioinformatics*, 30(12):i157–i164, 06 2014.
- 496 [18] Florian Huber, Lars Ridder, Stefan Verhoeven, Jurriaan H. Spaaks, Faruk Diblen, Simon Rogers, and Justin J. J.
497 van der Hooft. Spec2vec: Improved mass spectral similarity scoring through learning of structural relationships.
498 *PLOS Computational Biology*, 17(2):1–18, 02 2021.
- 499 [19] Kai Dührkop, Louis-Félix Nothias, Markus Fleischauer, Raphael Reher, Marcus Ludwig, Martin A. Hoffmann,
500 Daniel Petras, William H. Gerwick, Juho Rousu, Pieter C. Dorrestein, and Sebastian Böcker. Systematic
501 classification of unknown metabolites using high-resolution fragmentation mass spectra. *Nature Biotechnology*,
502 2020.
- 503 [20] Khawla Seddiki, Philippe Saudemont, Frédéric Precioso, Nina Ogrinc, Maxence Wisztorski, Michel Salzet,
504 Isabelle Fournier, and Arnaud Droit. Cumulative learning enables convolutional neural network representations
505 for small mass spectrometry data classification. *Nature Communications*, 11:5595, 2020.
- 506 [21] Yang-Ming Lin, Ching-Tai Chen, and Jia-Ming Chang. Ms2cnn: predicting ms/ms spectrum based on protein
507 sequence using deep convolutional neural networks. *BMC Genomics*, 20, 2019.
- 508 [22] Ngoc Hieu Tran, Xianglilan Zhang, Lei Xin, Baozhen Shan, and Ming Li. De novo peptide sequencing by deep
509 learning. *Proceedings of the National Academy of Sciences*, 114(31):8247–8252, 2017.
- 510 [23] Fatema Tuz Zohora, M. Ziaur Rahman, Ngoc Hieu Tran, Lei Xin, Baozhen Shan, and Ming Li. DeepIso: A deep
511 learning model for peptide feature detection from lc-ms map. *Scientific Reports*, 9, 2019.
- 512 [24] David Weininger. Smiles, a chemical language and information system. 1. introduction to methodology and
513 encoding rules. *Journal of Chemical Information and Computer Sciences*, 28(1):31–36, 1988.
- 514 [25] Robin Winter, Floriane Montanari, Frank Noé, and Djork-Arné Clevert. Learning continuous and data-driven
515 molecular descriptors by translating equivalent chemical representations. *Chemical science*, 10(6):1692–1701,
516 2019.
- 517 [26] Jerret Ross, Brian Belgodere, Vijil Chenthamarakshan, Inkit Padhi, Youssef Mroueh, and Payel Das. Large-scale
518 chemical language representations capture molecular structure and properties. *Nature Machine Intelligence*,
519 4(12):1256–1264, 2022.
- 520 [27] Brian Belgodere, Vijil Chenthamarakshan, Payel Das, Pierre Dognin, Toby Kurien, Igor Melnyk, Youssef Mroueh,
521 Inkit Padhi, Mattia Rigotti, Jarret Ross, et al. Cloud-based real-time molecular screening platform with molformer.
522 In *ECML PKDD*, 2022.
- 523 [28] Gerald Maggiora, Martin Vogt, Dagmar Stumpfe, and Jürgen Bajorath. Molecular similarity in medicinal chemistry.
524 *Journal of Medicinal Chemistry*, 57(8):3186–3204, 2014. PMID: 24151987.
- 525 [29] H. L. Morgan. The generation of a unique machine description for chemical structures—a technique developed at
526 chemical abstracts service. *Journal of Chemical Documentation*, 5(2):107–113, 1965.
- 527 [30] Rdkit: Open-source cheminformatics software. <https://www.rdkit.org/>.
- 528 [31] Kai Dührkop, Markus Fleischauer, Marcus Ludwig, Alexander A. Aksenov, Alexey V. Melnik, Marvin Meusel,
529 Pieter C. Dorrestein, Juho Rousu, and Sebastian Böcker. Sirius 4: a rapid tool for turning tandem mass spectra
530 into metabolite structure information. *Nature Methods*, 16:299–302, 2019.
- 531 [32] Jing Lim, Joshua Wong, Minn Xuan Wong, Lee Han Eric Tan, Hai Leong Chieu, Davin Choo, and Neng Kai Nigel
532 Neo. Chemical structure elucidation from mass spectrometry by matching substructures, 2018.
- 533 [33] Sunghwan Kim, Jie Chen, Tiejun Cheng, Asta Gindulyte, Jia He, Siqian He, Qingliang Li, Benjamin A Shoemaker,
534 Paul A Thiessen, Bo Yu, Leonid Zaslavsky, Jian Zhang, and Evan E Bolton. PubChem in 2021: new data content
535 and improved web interfaces. *Nucleic Acids Research*, 49(D1):D1388–D1395, 11 2020.
- 536 [34] John J. Irwin and Brian K. Shoichet. Zinc a free database of commercially available compounds for virtual
537 screening. *Journal of Chemical Information and Modeling*, 45(1):177–182, 2005. PMID: 15667143.
- 538 [35] X. Yang, P. Neta, and S. Stein. Extending a tandem mass spectral library to include MS2 spectra of fragment ions
539 produced in-source and MSn spectra. *Journal of the American Society for Mass Spectrometry*, 28:2280–2287,
540 2017.
- 541 [36] NIST 20 dataset. [https://chemdata.nist.gov/dokuwiki/lib/exe/fetch.php?media=chemdata:
542 asms2020:xiaoyu_yang_asms2020_presentation.pdf](https://chemdata.nist.gov/dokuwiki/lib/exe/fetch.php?media=chemdata:asms2020:xiaoyu_yang_asms2020_presentation.pdf). Accessed: 2021-04-04.
- 543 [37] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration.
544 *arXiv preprint arXiv:1904.09751*, 2019.

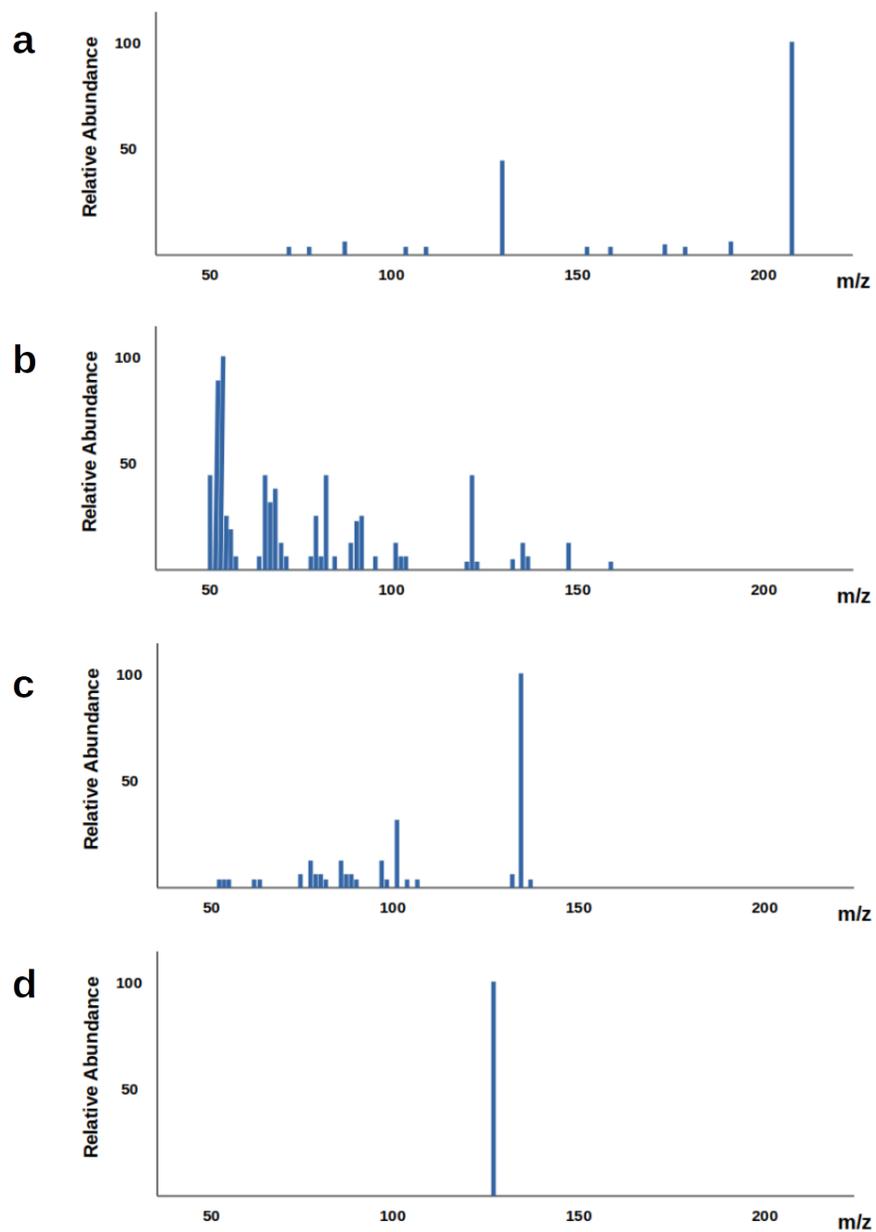


Figure 1: MS/MS spectra from different experimental conditions for the same molecule. MS/MS spectra obtained through different experimental conditions from the same molecule (approximate spectra based on data obtained from the Human Metabolome Database). (a) Precursor ion: $[M+H]^+$, NCE: 35%, Instrument: HCD. (b) Precursor ion: $[M+H]^+$, NCE: 130%, Instrument: HCD. (c) Precursor ion: $[M+H-Br]^+$, NCE: 35%, Instrument: HCD. (d) Precursor ion: $[M+H+2i]^+$, NCE: 35%, Instrument: IT-FT.

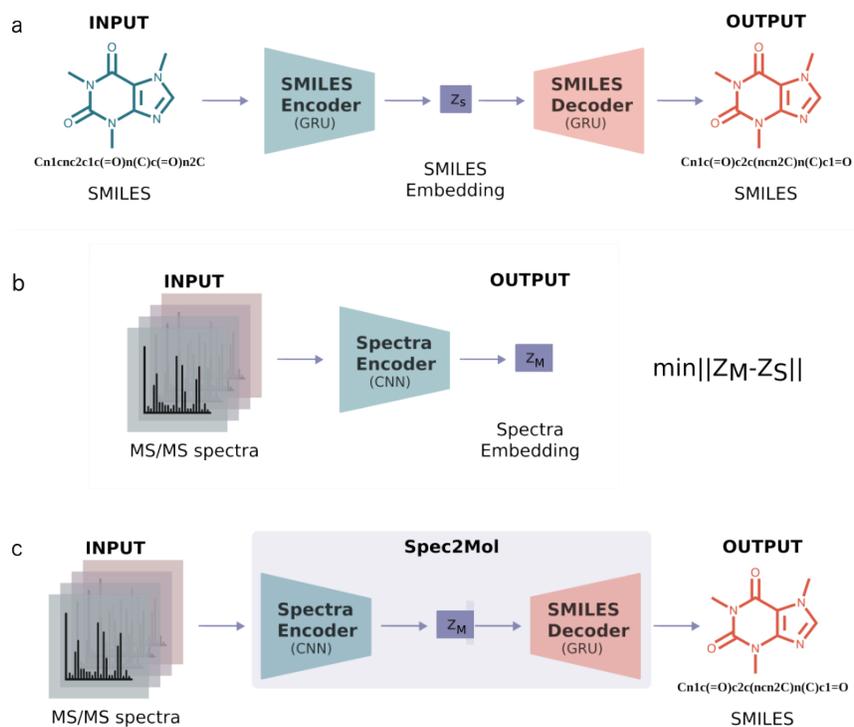


Figure 2: Spec2Mol architecture. The Spec2Mol model consists of a spectra encoder and a SMILES decoder which have been trained separately but share the same embedding space. (a) The AE is pre-trained to translate from a random SMILES to the canonical SMILES string. (b) The spectra encoder is trained to learn the same embedding as the SMILES encoder. (c) During inference, the spectra encoder and the SMILES decoder of the pre-trained model are used to translate spectra into molecular structures.

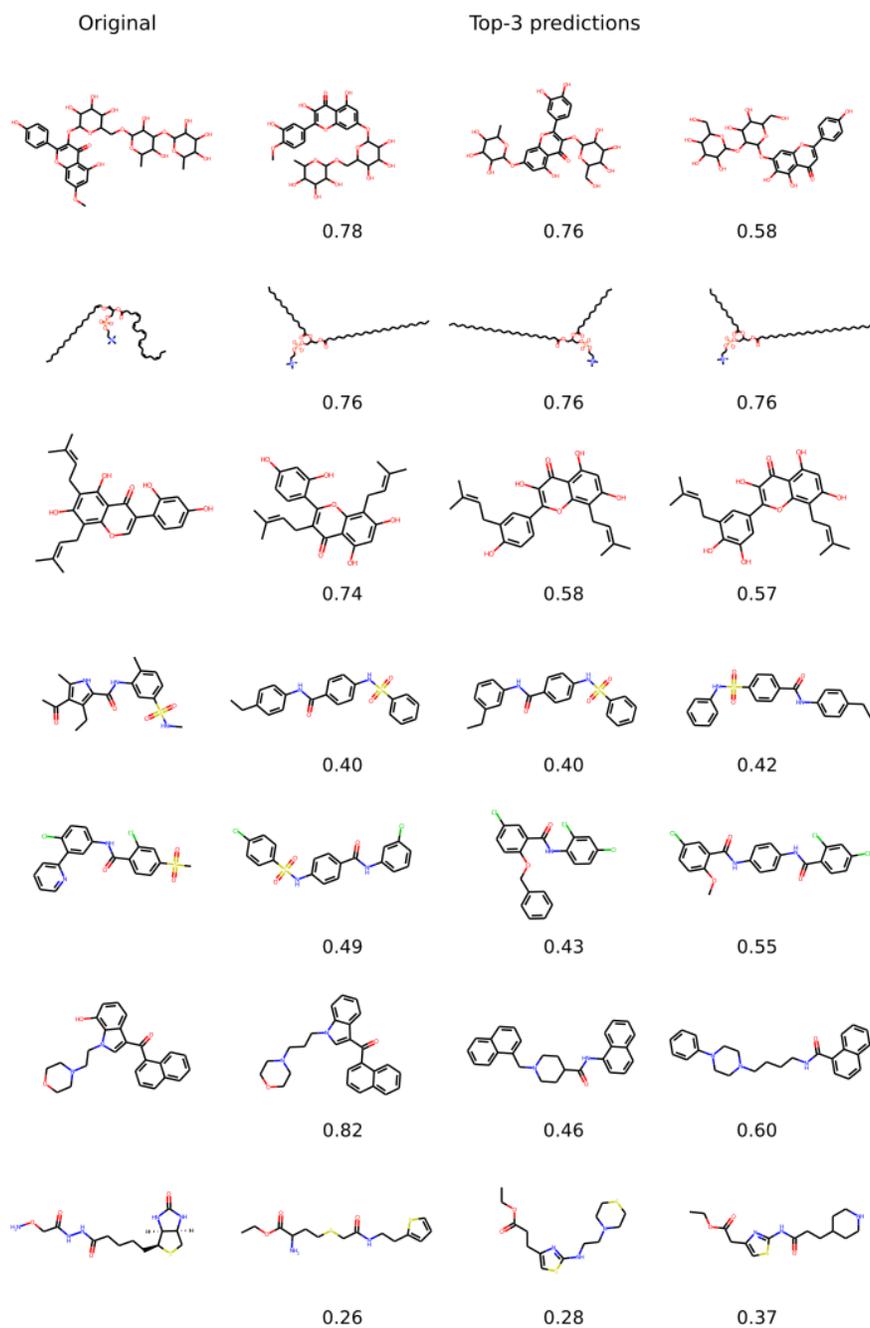


Figure 3: Examples of cases where Spec2Mol successfully identified key substructures. Examples of the most likely predicted structures from Spec2Mol along with the cosine similarity values with respect to the original reference structures.

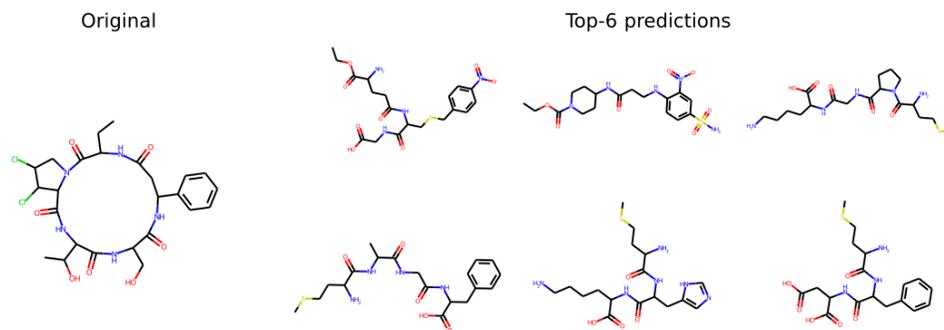


Figure 4: A case where Spec2Mol did not identify relevant structures. An example where Spec2Mol failed to identify a similar structure for a reference compound containing a large ring.

Table 1: Effect of missing spectra in the model input. Evaluation metrics when considering the entire test set and the test-data partitions that have available all 4, only 3, only 2 and only 1 spectrum. The arrows show the desired trend for each metric.

metric		full dataset	4 spectra	3 spectra	2 spectra	1 spectrum
# test cases		1000	413	65	483	39
Avg. MW		275.3	287.5	242.6	267.4	300.3
Avg. SMILES length		34.5	37.0	28.5	32.5	43.6
correct molecules (\uparrow)	(%)	7.0	9.2	15.2	4.1	5.1
correct formulas (\uparrow)	(%)	39.3	45.1	46.9	34.8	20.5
DMW% (\downarrow)	min	2.3	1.6	0.5	2.4	9.5
	avg	6.3	5.5	3.9	6.6	14.6
DMF% (\downarrow)	min	9.2	6.5	8.1	10.8	21.1
	avg	21.7	17.8	24.5	24.0	32.9
Fngp _{cosine} (\uparrow)	max	0.53	0.56	0.57	0.50	0.45
	avg	0.36	0.39	0.38	0.34	0.31
MCS _{ratio} (\uparrow)	max	0.68	0.70	0.72	0.66	0.57
	avg	0.51	0.53	0.55	0.50	0.43
MCS _{tan} (\uparrow)	max	0.55	0.58	0.60	0.53	0.44
	avg	0.38	0.39	0.41	0.36	0.30
MCS _{coef} (\uparrow)	max	0.71	0.73	0.74	0.69	0.63
	avg	0.54	0.55	0.58	0.53	0.48

Table 2: Effect of the molecule generation strategy. Comparative evaluation of the top-20 predictions using the direct strategy, the indirect strategy and the two strategies combined. The arrows show the desired trend for each metric.

metric		direct	indirect	combined
correct molecules (\uparrow)	(%)	0.8	6.9	7.0
correct formulas (\uparrow)	(%)	26.1	28.0	39.3
DMW% (\downarrow)	min	3.1	4.4	2.3
	avg	11.6	9.3	6.3
DMF% (\downarrow)	min	10.4	11.9	9.2
	avg	24.2	22.4	21.7
Fngp _{cosine} (\uparrow)	max	0.46	0.53	0.53
	avg	0.33	0.36	0.36
MCS _{ratio} (\uparrow)	max	0.65	0.66	0.68
	avg	0.50	0.51	0.51
MCS _{tan} (\uparrow)	max	0.50	0.55	0.55
	avg	0.34	0.38	0.38
MCS _{coef} (\uparrow)	max	0.68	0.71	0.71
	avg	0.53	0.56	0.54

Table 3: Sensitivity and specificity for detecting the presence of each atom species in the entire test set, having as reference the frequency of each species in the training spectra dataset.

	O	N	S	Cl	F	Br	P	I
Sensitivity	0.94	0.86	0.50	0.68	0.48	0.79	0.53	0.51
Specificity	0.50	0.76	0.96	0.91	0.92	0.98	0.99	0.99
Frequency (%)	85.4	71.5	18.4	15.2	11.5	7.5	2.5	1.4

Table 4: Effect of molecular weight and presence of heteroatoms.

	MW<300	MW \geq 300	only C and O	N present	S present	Halogen present
number of cases	668	332	184	769	199	318
exact structure (%)	8.5	3.9	9.8	6.1	5.5	5.7
exact formula (%)	45.4	27.1	63.6	34.1	23.6	25.8

Table 5: Comparative evaluation between SIRIUS and Spec2Mol, based on structural similarity between the recommended structures and the reference structure, on the subset of the test set where SIRIUS failed to identify an exact match.

Method		F_{ngp_cosine}	MCS_{ratio}	MCS_{tan}	MCS_{coef}
SIRIUS	max	0.49	0.65	0.54	0.66
	avg	0.33	0.49	0.35	0.49
Spec2Mol	max	0.49	0.66	0.53	0.69
	avg	0.34	0.50	0.36	0.53