

# DeepSPInN - multimodal Deep learning for molecular Structure Prediction from Infrared and NMR spectra

Sriram Devata,<sup>†,§</sup> Bhuvanesh Sridharan,<sup>†,§</sup> Sarvesh Mehta,<sup>†,§</sup> Yashaswi Pathak,<sup>†</sup>  
Siddhartha Laghuvarapu,<sup>‡</sup> Girish Varma,<sup>¶</sup> and Deva Priyakumar<sup>\*,†</sup>

<sup>†</sup>*Center for Computational Natural Sciences and Bioinformatics, International Institute of  
Information Technology, Hyderabad, India*

<sup>‡</sup>*IHub-Data, International Institute of Information Technology, Hyderabad, India*

<sup>¶</sup>*Center for Security, Theory and Algorithms Research, International Institute of  
Information Technology, Hyderabad, India*

<sup>§</sup>*Contributed equally to this work*

E-mail: deva@iiit.ac.in;sriram.devata@research.iiit.ac.in

## Abstract

Molecular spectroscopy studies the interaction of molecules with electromagnetic radiation, and interpreting the resultant spectra is invaluable for deducing the molecular structures. However, it is a strenuous task that requires highly specific domain knowledge. DeepSPInN predicts the molecular structure when given Infrared and <sup>13</sup>C Nuclear magnetic resonance spectra without referring to any pre-existing spectral databases or molecular fragment knowledge bases. DeepSPInN does this by formulating the molecular structure prediction problem as a Markov decision process (MDP) and employs Monte-Carlo tree search to explore and choose the actions in the formulated MDP. On the QM9 dataset, DeepSPInN is able to predict the molecular structure

for 95.98% of the input spectra in an average time of 160 seconds for molecules with less than 10 heavy atoms. This study is the first of its kind that uses multimodal data for molecular structure prediction, and is a leap forward in automated molecular spectral analysis.

## Introduction

Molecular spectroscopy is the analysis of the electronic, vibrational, and rotational excitations of the nuclei of molecules as they interact with electromagnetic radiation. It is widely used as a tool to identify and characterize molecules for quantitative and qualitative analysis of materials. The spectrum of a molecule is the measured absorption or emission of the incident electromagnetic radiation. Each molecule produces a unique spectrum for a particular spectroscopic method, allowing the spectrum to be used as a fingerprint of the molecule.

Infrared (IR) spectroscopy is a spectroscopic technique that sheds light on the vibrational modes of a molecule that changes its dipole moment.<sup>1</sup> These vibrational modes cause the molecules to absorb electromagnetic radiation in the Infrared spectral region, lying in the range of wavenumbers  $4000 - 400 \text{ cm}^{-1}$ . Functional groups have unique absorbances in the region of peaks beyond  $1500 \text{ cm}^{-1}$  called the functional group region.<sup>2</sup> Peaks with wavenumbers  $< 1500 \text{ cm}^{-1}$  are considered to be in the fingerprint region<sup>2</sup> since the elaborate patterns of peaks here are highly specific to a molecule and are often too complex to interpret.

Nuclear magnetic resonance (NMR) spectroscopy is another widely used spectroscopic technique to characterize the structure of molecules.<sup>3</sup> In NMR spectroscopy, an external magnetic field is applied to a molecule and the nuclei of some isotopes (e.g.  $^1\text{H}$ ,  $^{13}\text{C}$ ) absorb radio waves of specific frequencies to change their nuclear spin. In  $^{13}\text{C}$  NMR for example, any small changes in the local environment of the atom in the molecule cause the  $^{13}\text{C}$  nuclei to absorb radio waves of different frequencies. The relative differences of these frequencies against a reference  $^{13}\text{C}$  NMR frequency of tetramethylsilane (TMS) are measured in parts per million (ppm)<sup>4</sup> to give the chemical shifts of the nuclei. The spin-spin coupling of the

adjacent protons of the  $^{13}\text{C}$  nuclei cause the splitting of the corresponding NMR signal and allows the calculation of the multiplicity of each peak. This chemical split of each  $^{13}\text{C}$  nuclei's chemical shift is indicative of the number of directly attached hydrogen atoms. Together, the chemical shift and chemical split values of a  $^{13}\text{C}$  NMR spectrum allow the deduction of the atom type and chemical environment of each carbon atom, and subsequently the complete structure of the molecule.

For a structure to be elucidated from molecular spectra, all structural fragments are identified by interpreting the peaks in the spectra as the first step. These structural fragments are combined to list the possible molecular structures that can be made. These structures are then verified by cross-referencing the expected peaks of the functional groups in the input spectra, or by comparing their predicted spectra with the input spectra. CASE (Computer Aided Structure Elucidation) programs have evolved a lot since their introduction and have made good progress for structure elucidation from spectra, but they are still expected to have a degree of intervention from chemists and spectrometrists.<sup>5</sup> These programs also typically require 2D spectra in addition to any 1D IR, NMR, and MS spectra as the input.<sup>6</sup> Even today, most computational methods to identify a substance from its spectral data rely on matching against a database of already known spectra or by searching through knowledge bases of substructures<sup>7-15</sup>. Such methods restrict their applicability to the cases where the molecule's spectra is already stored in the database, or cases where the structural motifs are adequately represented in the dataset.

Recently, new methods have made use of Machine learning (ML) algorithms to solve problems in computational chemistry such as predicting new drug molecules<sup>16-18</sup>, performing molecular dynamics simulations<sup>19-21</sup>, protein stability and binding site prediction<sup>22,23</sup>, and predicting physical molecular properties.<sup>24-26</sup> Efforts for finding correlations between the spectral features of molecules and their structural features using ML can be dated back to the 1990s.<sup>27</sup> Interpretation of spectra to understand the complex relationship between a spectrum and the molecular structure is a difficult task. Recent developments in deep

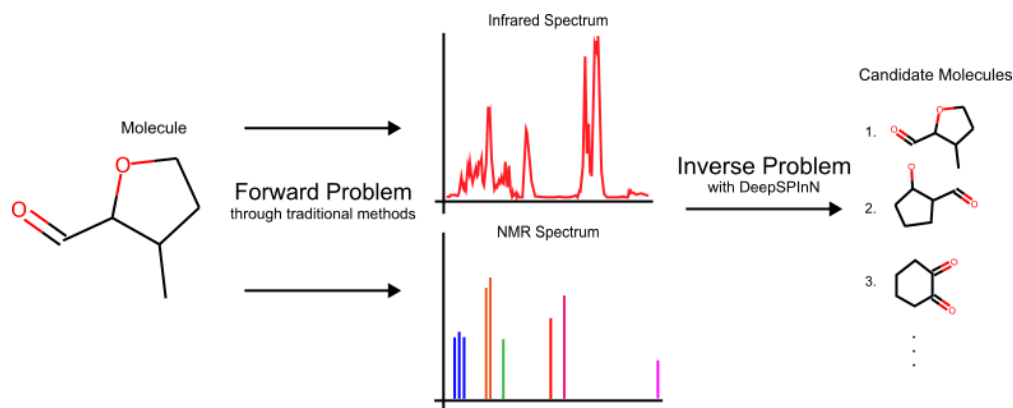


Figure 1: The IR and NMR spectra of 3-methyloxane-2-carbaldehyde to highlight the definitions of a *forward problem* and its corresponding *inverse problem*

learning open new avenues to explore the mapping between the molecular structure and the information-rich spectral data.

The *forward problem* can be defined as the prediction of the spectra of a given molecular structure, and the corresponding *inverse problem* is generating the molecular structure given the spectra (Figure 1). Although they are computationally intensive, quantum mechanical methods can be used to obtain various molecular spectra. Many recent works made progress in solving the forward problem of predicting the spectra of a molecule where they utilize ML for predicting IR<sup>28–32</sup>, NMR<sup>33–35</sup>, UV-visible<sup>36</sup>, and photoionization<sup>37,38</sup> spectra.

There have been works demonstrating how deep learning can solve inverse problems<sup>39</sup> in various domains. For the inverse problem in molecular structure elucidation, there have been works that aimed to automate the process of interpretation of IR spectra<sup>40,41</sup>. Many of them use only the functional group region of the spectra for their interpretation. Wang et al.<sup>40</sup> use a support vector machine to do multi-class classification for spectra from the OMNIC FTIR spectral library. The trained support vector machine identified 16 functional groups with a prediction accuracy of 93.3%. Fine et al.<sup>41</sup> introduce a multi-label neural network to identify functional groups present in a sample using a combination of FTIR and MS spectra. Jonas<sup>42</sup> and Howarth et al.<sup>43</sup> used a deep neural network that works with proton-coupled <sup>13</sup>C NMR to predict the molecular structure. Zhang et al.<sup>44</sup> use ChemTS<sup>45</sup> to identify

a molecule from its NMR spectrum using Monte Carlo tree search (MCTS) guided by a recurrent neural network (RNN). Huang et al.<sup>46</sup> propose an ML-based algorithm that takes  $^1H$  and  $^{13}C$  NMR as input and predicts the correct molecule as the top scoring candidate molecule with an accuracy of 67.4%. Pesek et al.<sup>47</sup> introduce a rule based combinatorial approach in which the framework uses  $^1H$  and  $^{13}C$  NMR, IR, and mass spectra to elucidate the structure of an unknown compound. Such knowledge engineering approaches would limit the capability of the solution since they inherit the biases of the rules programmed.

Elyashberg and Argyropoulos<sup>5</sup> predict that using deep learning algorithms would improve the performance and robustness of CASE systems. They also highlight AlphaZero's success in mastering games<sup>48</sup> as a testament to how deep learning can learn to perform complicated tasks. A concurrent work<sup>49</sup> proposes a transformer model that utilizes IR spectra to achieve a top-1 accuracy of  $\sim 55\%$  on molecules with less than 10 heavy atoms. Another similar concurrent work<sup>50</sup> utilizes both  $^1H$  and  $^{13}C$  NMR spectra to achieve a top-1 accuracy of  $\sim 70\%$  on molecules with less than 10 heavy atoms. It has recently been shown that a Monte-Carlo tree search (MCTS) algorithm can be used for the elucidation of molecular structure from NMR spectra, achieving a top-1 accuracy of 57.2%<sup>51</sup> for molecules with less than 10 heavy atoms on the nmrshiftdb2<sup>52</sup> dataset that contains experimentally calculated  $^{13}C$  NMR spectra of 2134 molecules.

Elucidation of spectra to molecular structure by a chemist typically involves the analysis of multiple types of spectra (UV-Vis, IR, NMR, etc.) rather than relying on any one type of spectrum. This is because the various techniques to capture these spectra capture complementary information of the molecules, which chemists with highly specific domain knowledge are able to interpret.

In this paper, our main contribution is a framework that performs multimodal analysis by utilizing both IR and  $^{13}C$  NMR spectra to accurately identify the molecular structure without any knowledge engineering or database searches. DeepSPInN formulates the molecular structure prediction problem as an MDP and employs MCTS to generate and traverse a

search tree while using a set of pre-trained Graph Convolution Networks<sup>53</sup> to guide the tree search. DeepSPInN is able to achieve an accuracy of 95.98% on molecules with less than 10 heavy atoms, outperforming previous and concurrent works on structure elucidation from molecular spectra.

## Methods

### Dataset

The QM9<sup>54,55</sup> dataset is a subset of the GDB-17<sup>56</sup> chemical universe and consists of 134k stable small organic molecules with up to nine heavy atoms (CNOF). We first identified molecules in the QM9 dataset for which IR and <sup>13</sup>C NMR spectra were calculated using the Gaussian 09<sup>57</sup> suite of programs. We were able to calculate both IR and <sup>13</sup>C NMR spectra for 119,062 molecules. We then chose molecules where the smallest ring (if any ring(s) exist(s)) in the molecule has at least 5 atoms to account for ring strain, and molecules where none of the atoms have any formal charge. This left us with about 50k molecules to use as the input data for this work. A train-test split of 80-20 was used to make the train and test dataset of molecules.

To calculate the IR absorbance spectra, the geometrical optimization and the subsequent calculation of the vibrational frequencies were done using the B3LYP density functional methods with a 6-31g(2df,p) basis set in the gas phase. The spectra from these DFT calculations for each molecule is a set of frequency-intensity pairs. These infinitely sharp stick spectra were broadened to mimic actual gas-phase spectra using a peak broadening function as described and trained by McGill et al.<sup>58</sup>. This function is a two-layer fully connected neural network followed by an exponential transform, and takes frequency-intensity pairs to give a continuous spectrum. The intensities of the resulting spectra were binned with a bin-width of 2 cm<sup>-1</sup> in the spectral range from 400 – 4000 cm<sup>-1</sup>. This results in the gas-phase IR absorbance spectrum for each molecule being represented by a 1801-length vector.

To make a dataset of proton-coupled  $^{13}\text{C}$  NMR spectra, the peak positions (chemical shift) and peak splits (C-H spin-spin coupling) were obtained from the QM9-NMR dataset.<sup>59</sup> The QM9-NMR dataset has the gas phase mPW1PW91/6-311+G(2d,p)-level atom-wise isotropic shielding for the QM9 dataset. These  $^{13}\text{C}$  isotropic shielding ( $\sigma_{\text{iso}}$ ) values were converted to  $^{13}\text{C}$  chemical shifts ( $\delta_{\text{iso}}$ ) through  $\delta_{\text{iso}} = \sigma_{\text{iso}}^{\text{reference}} - \sigma_{\text{iso}}$ ,<sup>60</sup> where  $\sigma_{\text{iso}}^{\text{reference}}$  is the reference value for tetramethylsilane (TMS), which is a standard reference compound. The mean absolute error between the  $^{13}\text{C}$  NMR spectra obtained in this way against spectra from the experimental nmrshiftdb2<sup>52</sup> database for the common molecules is 10.78 ppm per peak. As a reference,  $^{13}\text{C}$  NMR shift values are typically between 0 and 200.

## DeepSpInN Framework

The methods section is divided into five parts to explain the proposed framework:

- i. description of how molecular structure prediction can be modelled as a Markov decision process (MDP)
- ii. description of how MCTS can be used to generate a search tree of molecules and refine the policy at each state
- iii. explanation of the architecture of the prior and value model used by DeepSpInN
- iv. explanation of how  $^{13}\text{C}$  NMR split values are used to prune the MCTS search tree
- v. description of the training methodology used to train the prior and value model

### MDP formulation

The problem of molecular structure prediction can be modelled as a finite Markov decision process (MDP)<sup>61,62</sup> in a way similar to the formulation in Sridharan et al.<sup>51</sup>. An MDP is defined as a tuple  $\langle \mathcal{S}, \mathcal{A}, \{P_s\}, R \rangle$  with states  $\mathcal{S}$ , actions  $\mathcal{A}$ , policy  $\{P_s\}$ , and reward function

$R$ .<sup>63</sup> The goal is to learn the policies  $P_s$  which gives the transition probabilities over the action space  $\mathcal{A}$  at a particular state  $s \in \mathcal{S}$ .

Each state  $s \in \mathcal{S}$  consists of a molecular graph  $m$  and the target IR spectrum  $y_{\text{IR}}$ . A molecular graph represents a molecule where the atoms and bonds are mapped to nodes and edges in a graph.  $m$  also has the information about the target  $^{13}\text{C}$  NMR spectrum encoded as node-wise features. In the initial state, the molecular graph is a null graph with nodes representing each atom in the molecular formula and no edges. The molecule  $\text{mol}_s$  at a state  $s$  is the largest connected component in the molecular graph. The remaining individual nodes in  $m$  might join  $\text{mol}_s$  after taking an action  $a \in \mathcal{A}$ . In the initial state,  $\text{mol}_s$  is just a single carbon atom corresponding to any of the nodes in  $m$ .

An action  $a \in \mathcal{A}$  adds an edge between two nodes in  $m$ , which is equivalent to the addition of a bond between two atoms. Since the QM9 dataset has molecules that have a maximum of 9 atoms (number of nodes) and since there are 3 types of bonds (edges), the action space  $\mathcal{A}$  has  $9 * 9 * 3 = 243$  actions. For the molecular graphs to represent chemically valid molecules, only a subset of these actions can be considered to be valid. If a state has no valid actions that can be taken to reach any children states, it is a terminal state. In the action space for a state  $s$ , the valid actions are those that satisfy these conditions:

- Out of the two nodes that the action adds an edge between, at least one of the nodes must belong to the largest connected component ( $\text{mol}_s$ ) of the molecular graph, i.e. the current molecule of the state.
- The edge added by the action should satisfy the chemical valency rules of the two nodes. If all the edges of a node do not satisfy the octet of the corresponding atom type, it is implicitly assumed that hydrogen atoms contribute to the octet.
- The action should not create a self-loop since atoms do not form bonds with themselves.
- The action does not add an edge between two nodes that already belong to the same cycle.



- The action does not create a cycle whose length is less than 5, since rings with less than 5 atoms have high ring strain if they have double or triple bonds.

The reward function  $\mathcal{R}$  returns a non-zero reward for all terminal states and a zero reward for all non-terminal states. For the terminal states, the reward is a function of the spectral distance between the input IR spectrum and the IR spectrum of  $\text{mol}_s$  as predicted by Chemprop-IR.<sup>58</sup> Chemprop-IR is an extension to the Chemprop<sup>64</sup> architecture and uses a Directed Message Passing Neural Network<sup>65</sup>(D-MPNN) to predict the IR spectrum of an input molecular graph.  $\mathcal{R}$  is the Spectral Information Similarity<sup>58</sup> (SIS) metric which is calculated by rescaling the spectral divergence between two IR spectra found by their Spectral Information Divergence<sup>66</sup> (SID). The reward function  $\mathcal{R}$  is given by:

$$\mathcal{R} = \text{SIS}(A, B) = \frac{1}{1 + \text{SID}(A, B)} = \left( 1 + \sum_i (A_i \ln \frac{A_i}{B_i} + B_i \ln \frac{B_i}{A_i}) \right)^{-1}$$

where A and B are two IR spectra.

## Generating and exploring the search tree with MCTS

With this MDP formulation, we can use search algorithms to build a tree of state-labelled nodes.<sup>67,68</sup> We can build such a tree by repeatedly starting at the root state and reaching children states by taking any of the valid actions at each state. We use MCTS to estimate the optimal policy for the modelled reinforcement learning (RL) task.<sup>69</sup>

Starting from a root node, MCTS has 4 stages - selection, expansion, roll-out, and back-propagation (see Figure 2). In the selection stage, the algorithm chooses actions with probabilities proportional to their UCT<sup>67</sup> (Upper Confidence Bound applied to trees) values, until it reaches a leaf node. The UCT value of an action  $a$  at state  $s$  is given by

$$\text{UCT}(s, a) = Q(s, a) + c_{\text{puct}} \cdot \pi_s^a \cdot \frac{\sqrt{\sum_b N(s, b) + 1}}{N(s, a) + 1}$$

where  $Q(s, a)$  is the expected reward of taking action  $a$  from state  $s$ ,  $c_{\text{puct}}$  is a parameter to

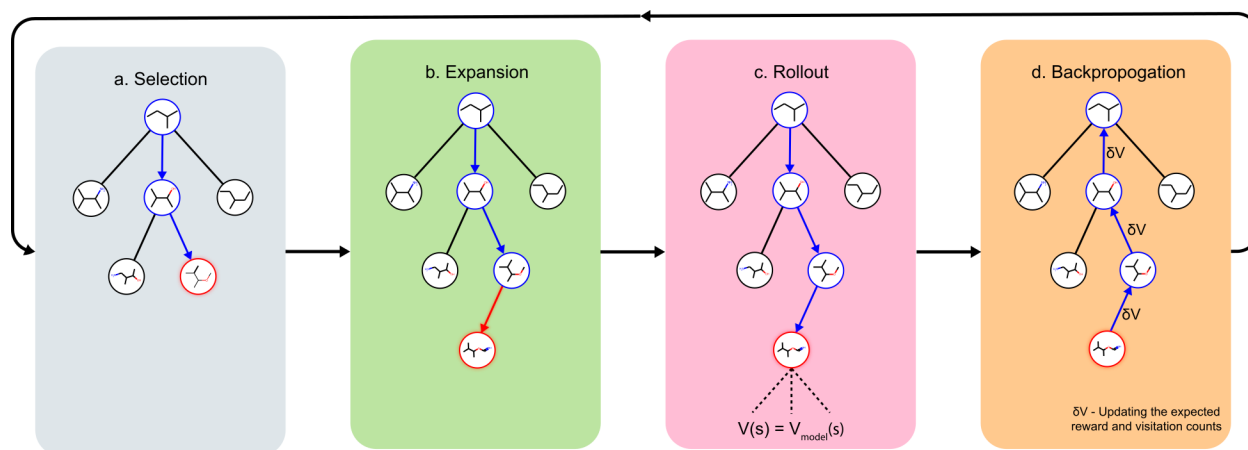


Figure 2: MCTS progresses in 4 stages to generate the search tree. a) Selection: starting from the root node of the tree, choose actions based on the UCT values b) Expansion: when the tree search reaches a leaf node, add a new child state to the tree c) Rollout: calculate the expected reward of the new child state through a series of random roll-outs d) Backpropagation: update the UCT values of all ancestors of the new child state

balance exploration and exploitation in the tree search,  $\pi_s^a$  is the probability of taking action  $a$  from state  $s$  according to the policy returned by a prior model,  $N(s, a)$  is the number of times action  $a$  has been taken from state  $s$ , and  $\sum_b N(s, b)$  is the number of times state  $s$  has been reached.

In the process of traversing the search tree according to the UCT values, the algorithm would reach a point where taking an action  $a$  from state  $s$  would lead to a state  $s'$  that does not exist in the search tree. This leads to the expansion stage of MCTS where the new state  $s'$  is added to the search tree.

Once a new child node  $s'$  is added in the expansion stage, the rollout stage is used to evaluate the value of  $s'$ . An ideal way to calculate this value is to calculate the expected reward by a series of random rollouts. Due to the computational complexity of calculating the expected reward in the ideal way, we approximate the value using an offline-trained value model.<sup>48,70</sup> The value of  $s'$  is recursively back-propagated through all its parent nodes till the root node to update the ancestors' values and visitation counts. If  $s$  is a terminal state that already exists in the tree, the reward of  $s$  is back-propagated to update the values of all ancestor nodes. A state  $s$  is considered to be terminal if it has no valid actions, or if its

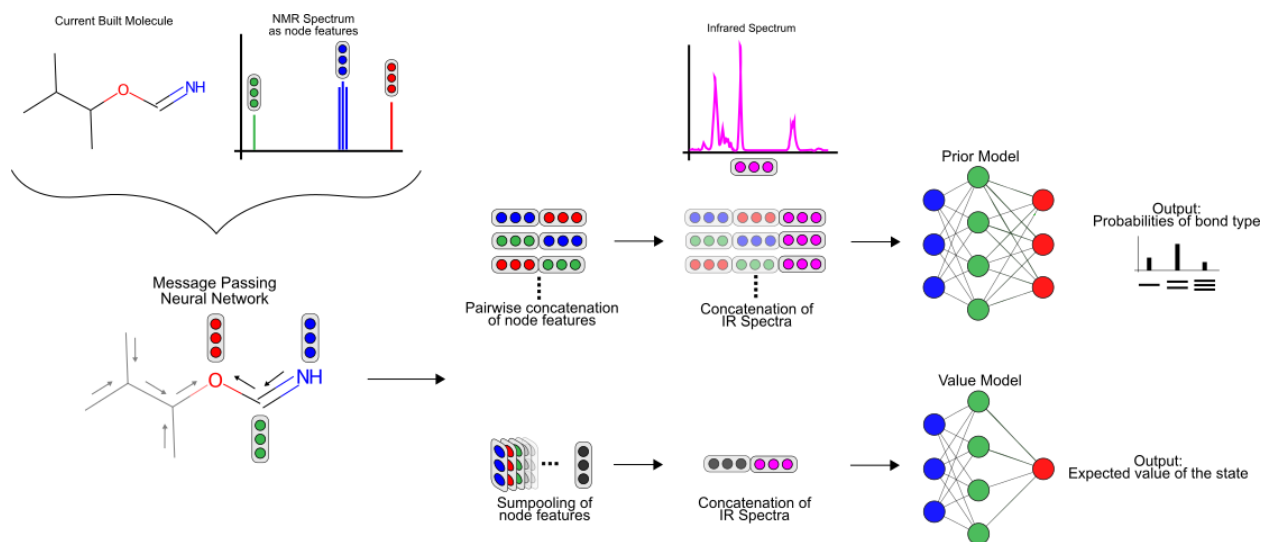


Figure 3: A prior model and a value model are used with the MCTS algorithm to get the probabilities over the action space and to predict the value of a particular state. An MPNN uses the initial node-wise features that contain the  $^{13}\text{C}$  NMR spectrum to give node-wise embeddings after three message passing steps. The prior model uses the pair-wise node embeddings and the IR spectrum to predict the probability of each pair of nodes having a single, double, or triple bond between them. The value model uses the sum-pooled node-wise embeddings and the IR spectrum to predict the value of a particular state.

reward exceeds a particular threshold (explained in the Supplementary Information). All 4 MCTS stages are repeated  $n_{\text{mcts}}$  number of times which is a hyper-parameter of DeepSPInN. After  $n_{\text{mcts}}$  repetitions of the above 4 MCTS stages, a true action is taken according to the final policy at this state.

### Description of the prior and value model

To featurize the built molecule at each state, both the prior and value model use a Message Passing Neural Network<sup>53,71</sup> (MPNN) that run for three time steps (see Figure 3). Consider a molecular graph  $G(V, E)$  where each node has initial node features  $x_v, \forall v \in V$ .  $x_v$  contain the chemical description of the atom and the NMR spectrum of the atom corresponding to node  $v$  as listed in Table 1. Each node  $v$  also has hidden features  $h_v$  that are initialized to  $x_v$ , with the MPNN updating these hidden features in each time step of the forward pass. All edges in the molecular graph have edge features  $e_{vw}, \forall v, w \in V$  as listed in Table

Table 1: Featurization of nodes and edges in the molecular graph

| Node Feature              | Description                                                                                |
|---------------------------|--------------------------------------------------------------------------------------------|
| Element Type              | one-hot of [C,N,O,F]                                                                       |
| Hybridization             | one-hot of [ $sp$ , $sp^2$ , $sp^3$ ]                                                      |
| Implicit Valency          | one-hot of [0,1,2,3,4]                                                                     |
| Radical Electrons         | one-hot of [0,1,2]                                                                         |
| Formal Charge             | one-hot of [-2, -1, 0, 1, 2]                                                               |
| $^{13}\text{C}$ NMR split | one-hot of [0,1,2,3]                                                                       |
| $^{13}\text{C}$ NMR shift | a gaussian with $\sigma = 2$ centered at the chemical shift value discretized into 64 bins |
| Edge Feature              | Description                                                                                |
| Bond Type                 | one-hot of [single, double, triple, aromatic]                                              |
| Bond Conjugation          | boolean of whether the bond is conjugated                                                  |
| Presence in a Ring        | boolean of whether the bond is in a ring                                                   |

1. The forward pass of an MPNN has  $T$  message passing time steps and a final gathering step. The message passing steps use a message function  $M_t$  to form messages from the hidden features of neighbouring nodes  $N(v)$  and the features of their corresponding edges. An update function  $U_t$  updates the hidden features of a node based on its current hidden features and the messages it received from its neighbouring nodes.

$$m_v^{t+1} = \sum_{w \in N(v)} M_t(h_v^t, h_w^t, e_{vw})$$

$$h_v^{t+1} = U_t(h_v^t, m_v^{t+1})$$

After  $T$  message passing steps, a gathering function  $G_T$  uses the initial node features  $x_v$  and the final hidden features  $h_v$  to give the node-wise features  $F_v$ .

$$F_v = G_T(x_v, h_v^t)$$

In DeepSPInN,  $M_t$  and  $U_t$  are fully connected neural neural networks, and  $G_T$  is an element-wise addition operation.

Using the node-wise features from the MPNN, the prior model generates all possible pairs of nodes and concatenates the node-wise features of all these pairs of nodes to get pair-wise features.  $y_{\text{IR}}$  is compressed by passing through a two-layer fully connected neural network to give  $y'_{\text{IR}}$  and is appended to all these pair-wise features. The product of this concatenation is passed through another two-layer fully connected neural network  $Pr_{\text{model}}$  to predict the probabilities of a bond of each of the three types (single, double, triple) existing between the pair of nodes. The prior model works as follows

$$P_{\text{bond}} = Pr_{\text{model}}([F_v, F_u, y'_{\text{IR}}]), \text{ for each pair of nodes } u, v \in V$$

where, “[ ]” represents a concatenation operation,  $Pr_{\text{model}}$  is the prior model, and  $P_{\text{bond}}$  is a 3-tuple giving the probabilities of nodes  $u$  and  $v$  having a single, double, and triple bond respectively.

The value model first performs a sum-pooling operation on the node-wise features obtained from the MPNN. It then appends the compressed IR spectrum to the sum-pooled feature vector of the molecule and passes this through a two-layer fully connected neural network  $V_{\text{model}}$  to predict the value of this state. The value model works as follows

$$V_s = V_{\text{model}}\left(\left[\sum_i F_i, y'_{\text{IR}}\right]\right)$$

where,  $\sum_i F_i$  is the result of the sum-pooling operation of all node-wise features in the molecular graph.

## RL environment

The RL environment handles the algorithm’s queries to return children states when given state-action pairs according to the transition process. This environment also checks the validity of the molecular graphs by enforcing the chemical valencies of each atom in the graph. The  $^{13}\text{C}$  NMR split values help prune the search tree by identifying molecules that are

invalid according to the input NMR spectrum and nullifies the rewards for these molecules. This discourages the tree search from exploring these molecules. The NMR split values for the NMR shift of each carbon atom are equivalent to the number of hydrogen atoms attached to it. Each carbon atom can be either a singlet (S, quaternary), doublet (D, tertiary), triplet (T, secondary), or a quartet (Q, primary) atom with each of these denoted by S, D, T, and Q respectively. Since each valid action is defined as the addition of a bond between two atoms in the MDP reformulation, each valid action can only convert the carbon atoms from  $Q \rightarrow T \rightarrow D \rightarrow S$ . If the target Q-splits are more than the Q-splits at one such state, this state is not valuable since no valid action from this state would be able to increase the count of Q-splits. If the Q-splits match, checking the T and D-splits subsequently in the same way further identify more states that get zero-rewards.

## Training Methodology

The prior and value model are trained on a set of experiences generated from a guided tree search on the molecules in the training dataset. Since the target molecule is known while training, the reward function is replaced with a binary function that returns a value depending on whether the molecule built at the current state is subgraph isomorphic to the molecular graph of the target molecule. The reward for taking an action  $a$  from state  $s$  to reach state  $s'$  is:

$$r(s, a) = \begin{cases} 1 & \text{if } S(\text{mol}_{s'}, \text{mol}_{\text{target}}) \\ 0 & \text{otherwise} \end{cases}$$

where  $\text{mol}_{s'}$  is the molecular graph of the molecule at state  $s'$ ,  $\text{mol}_{\text{target}}$  is the molecular graph of the target molecule, and  $S(\text{mol}_{s'}, \text{mol}_{\text{target}})$  is RDKit's<sup>72</sup> substructure search that does a subgraph isomorphism check and returns a boolean value.

The policies and values derived from the above training runs are stored as experiences and are used to train the prior and value model. The training took about 45 hours on a

system with a Intel(R) Xeon(R) CPU E5-2640 v4 processor and a GeForce RTX 2080 Ti GPU.

To evaluate DeepSPInN, we run 40 episodes for each set of input spectra. Each episode builds the MCTS tree from scratch by going through all four phases of MCTS  $n_{mcts}$  times and returns a final molecule. All the unique candidate molecules from these 40 episodes are then ranked using the reward function as a scoring function. We report multiple *Top N* metrics where each *Top N* metric denotes whether the target molecule was present in the top  $N$  candidate molecules.

## Results

This section describes the performance of DeepSPInN in identifying the correct molecule given a set of IR and NMR spectra. To rigorously evaluate DeepSPInN, we present the results of a few experiments in the following subsections. The first subsection compares the performance of DeepSPInN for different  $n_{mcts}$  values. The next subsection compares the final rewards for correctly and incorrectly predicted molecules. In the following subsection, the time taken to predict the molecules for different  $n_{mcts}$  values is analyzed. In the subsequent subsection, performance of the model is discussed when only one of IR or NMR spectrum is given as the input. The final subsection describes and presents the results of an experiment to check the generalizability of DeepSPInN.

### Performance of DeepSPInN for varying $n_{mcts}$ values

Table 2 compares the results for different values of  $n_{mcts}$  when given both IR and NMR spectra. For  $n_{mcts} = 800$ , DeepSPInN correctly identifies the molecule  $\sim 96\%$  of the time as the top candidate molecule. Although the previous MCTS-based structure elucidation method<sup>51</sup> uses only NMR spectra, DeepSPInN outperforms the previous method even with  $n_{mcts} = 100$  by achieving a *Top 1* (%) accuracy of  $\sim 82\%$  compared to the previous work's

Table 2: *Top N* metrics for varying  $n_{mcts}$  values

| $n_{mcts}$ | IR+NMR |        |        |               |
|------------|--------|--------|--------|---------------|
|            | 100    | 200    | 400    | 800           |
| Top 1 (%)  | 82.311 | 90.773 | 94.934 | <b>95.980</b> |
| Top 3 (%)  | 82.874 | 91.597 | 95.839 | <b>96.925</b> |
| Top 5 (%)  | 82.994 | 91.778 | 96.060 | <b>97.106</b> |
| Top 10 (%) | 83.015 | 91.798 | 96.120 | <b>97.166</b> |
| Top 40 (%) | 83.015 | 91.798 | 96.120 | <b>97.206</b> |

best *Top 1* (%) accuracy of  $\sim 60\%$  for  $n_{mcts} = 200$ . Across various  $n_{mcts}$  values, the *Top 1* (%) accuracy increases as  $n_{mcts}$  increases. There is a stark increase in the *Top 1* (%) accuracy between  $n_{mcts} = 100, 200$  and  $n_{mcts} = 200, 400$ , but there is only a marginal difference between  $n_{mcts} = 400, 800$ . This shows that increasing  $n_{mcts}$  further will result in diminishing increase in performance while taking a disproportionately greater amount of time as explained in a subsequent subsection.

Even within each  $n_{mcts}$  value, the *Top N* (%) metrics increase marginally starting from *Top 1* (%) to *Top 40* (%). The increases across the *Top N* (%) metrics are due to an imperfect scoring function being used to rank all the candidate molecules. If the correct target molecule is not ranked as the top candidate molecule, it would contribute to one of the *Top N* (%) metrics. Still, we observe that the scoring function proposed in DeepSPInN is significantly better than the one used in Sridharan et al.<sup>51</sup> since there are great differences in metrics across the *Top N* (%) metrics in the original work. DeepSPInN does not show such great differences in the *Top N* metrics, illustrating that the scoring function used here performs better in ranking the candidate molecules. In DeepSPInN, if the correct molecule is found to be one of the candidate molecules, it is almost always ranked as the top candidate.

### Comparison of rewards for correctly and incorrectly predicted molecules

Figure 4 contains the histograms of the rewards for the cases when DeepSPInN was and was not able to predict the correct molecule as the top candidate. The histogram of the rewards



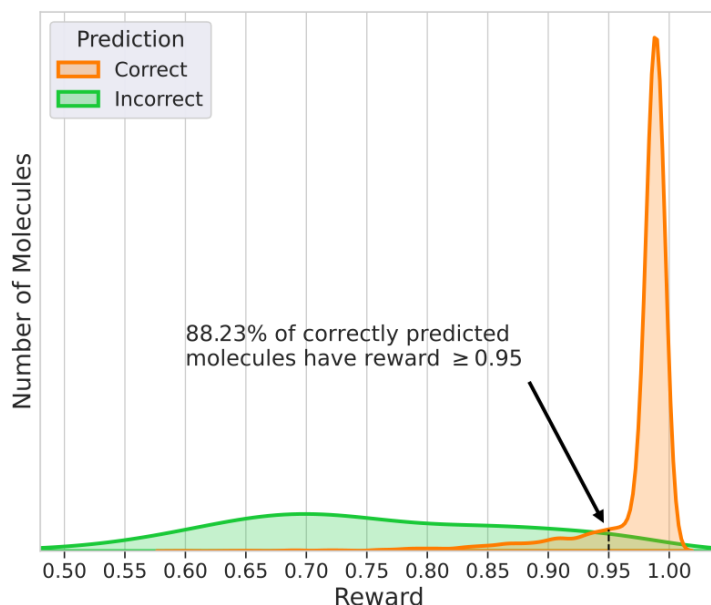
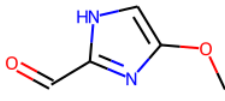
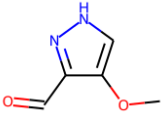
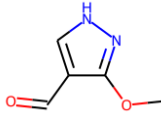
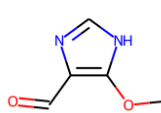
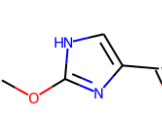
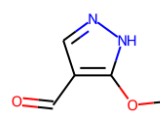
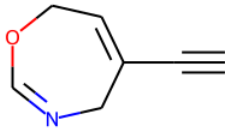
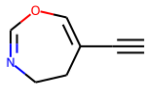
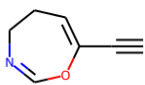
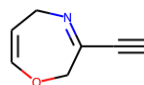
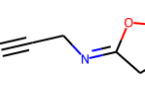
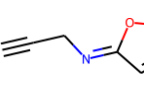
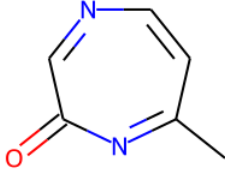
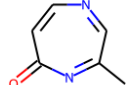
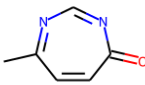
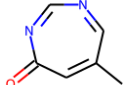
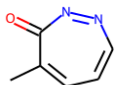
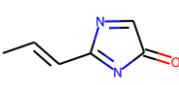
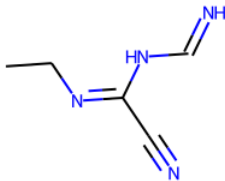
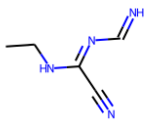
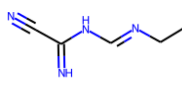
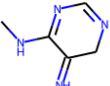
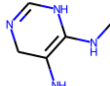
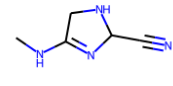
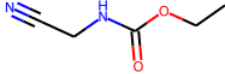
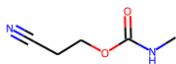
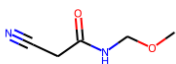
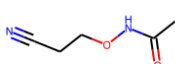
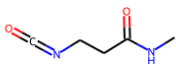
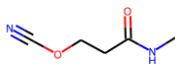
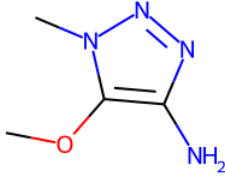
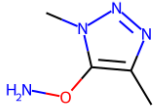
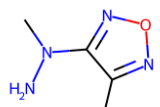
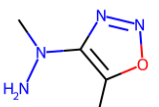
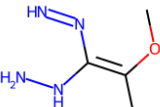
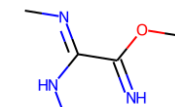


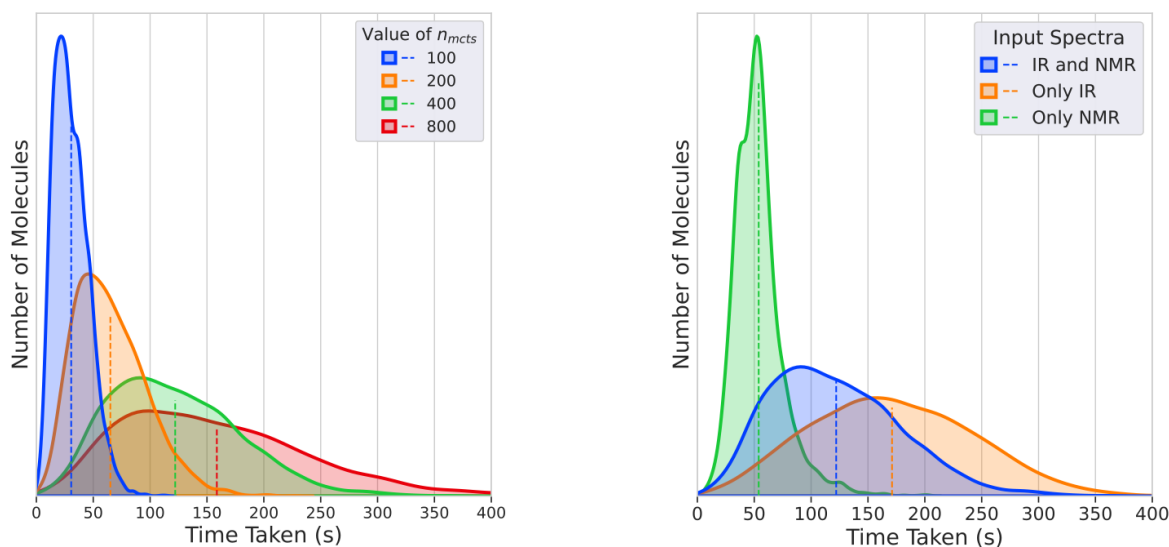
Figure 4: Histogram of the rewards of molecules that had the correct and incorrect structure as the top ranked candidate molecule for  $n_{\text{mcts}} = 800$

for the correctly predicted molecules has a very narrow distribution and has an average reward of 0.975. It is also left-skewed with most of the correctly predicted molecules receiving a higher reward when compared to the incorrectly predicted molecules. The histogram of the rewards for the incorrectly predicted molecules has a broader distribution with an average reward of 0.748. 88.23% of the correctly predicted molecules had a reward  $\geq 0.95$  while only 6.5% of the incorrectly predicted molecules had a reward  $\geq 0.95$ . DeepSPInN would allow researchers to use the final reward as a confidence measure of the correctness of the prediction. When DeepSPInN gives a final reward  $\geq 0.95$  for a set of input spectra, the top candidate is the target molecule 99.69% of the time. Table 3 shows 5 randomly selected molecules for which DeepSPInN did not have the correct molecule in any of the predicted candidate molecules. The top candidate molecules even for these incorrectly predicted molecules are structurally similar to the correct molecule, with the average Tanimoto similarity between the correct molecule and the top candidate molecule being 0.970 for the test set.

Table 3: Rewards of the top 5 ranked candidate structures for molecules that were never predicted for  $n_{\text{mcts}} = 800$

| Target Molecule                                                                     | Top 5 Candidate Molecules<br>(and their final rewards)                                       |                                                                                              |                                                                                               |                                                                                                |                                                                                                |
|-------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------|
|    | <br>0.611   | <br>0.601   | <br>0.574   | <br>0.563   | <br>0.500   |
|    | <br>0.858   | <br>0.755   | <br>0.735   | <br>0.673   | <br>0.658   |
|   | <br>0.771  | <br>0.754  | <br>0.745  | <br>0.679  | <br>0.634  |
|  | <br>0.652 | <br>0.560 | <br>0.549 | <br>0.542 | <br>0.541 |
|  | <br>0.648 | <br>0.624 | <br>0.600 | <br>0.529 | <br>0.519 |
|  | <br>0.731 | <br>0.697 | <br>0.642 | <br>0.575 | <br>0.573 |

## Analysis of the time taken for the predictions



(a) Histograms of the time taken to predict each molecule when given both IR and NMR spectra for varying  $n_{mcts}$  values

(b) Histograms of the time taken to predict each molecule when given either IR or NMR spectra for  $n_{mcts} = 400$

Figure 5: Histograms of time taken to predict each molecule when given both IR and NMR spectra or either one spectrum

Figure 5a shows the distribution of times taken for DeepSPInN to predict candidate molecules for input IR and NMR spectra for different values of  $n_{mcts}$ . For  $n_{mcts} = 800$ , the average time taken is 157 seconds with about 95% of the test molecules taking less than 300 seconds. The predictions for  $n_{mcts} = 100, 200, 400$  take 30, 65, and 122 seconds on average.

### Importance of multimodal data as input (IR and NMR spectra vs either IR or NMR)

To compare the distinguishing ability of IR and NMR and to compare the utility of having both IR and NMR spectra as the input, we performed ablation studies where we ran the model with either one of the spectra as the input for  $n_{mcts} = 400$ . Table 4 shows the *Top N* metrics for the models that received both IR and NMR, only IR, and only NMR spectra as input. The IR-and-NMR-trained model has a *Top 1* accuracy of 94.93% while the IR-trained

Table 4: Performance of IR-and-NMR-trained, IR-trained, and NMR-trained models for  $n_{mcts} = 400$

|            | IR and NMR | Only IR | Only NMR | Only NMR ranked by IR |
|------------|------------|---------|----------|-----------------------|
| Top 1 (%)  | 94.934     | 74.075  | 40.462   | 73.708                |
| Top 3 (%)  | 95.839     | 74.396  | 61.849   | 74.391                |
| Top 5 (%)  | 96.060     | 74.396  | 68.201   | 74.472                |
| Top 10 (%) | 96.120     | 74.396  | 73.105   | 74.472                |
| Top 40 (%) | 96.120     | 74.423  | 74.472   | 74.472                |

and NMR-trained models have a *Top 1* accuracy of 74.07% and 40.46% respectively. All *Top N* metrics for the IR-and-NMR-trained model are greater than the models that work with either one of the spectra. This implies that the model is able to learn complementary information from both the spectra and subsequently performs better than the models with either one of the spectra as the input. Among the models that work on either one of the spectra, the IR-trained model performed better than the NMR-trained model in most of the *Top N* metrics. Figure 5b shows the distributions of times taken by IR-and-NMR-trained, IR-trained, and NMR-trained models to predict candidate molecules. The NMR-trained model has the fastest average prediction time of 54 seconds, while the IR-trained model has the slowest average prediction time of 171 seconds.

Table 4 also compares the *Top N* metrics for two NMR-trained models that only differ in how they rank the candidate structures. An NMR-trained model was used to first obtain candidate molecules for an input NMR spectrum. These candidate molecules were then ranked by two scoring functions that each used either IR or NMR spectral distances. The first scoring function used the spectral distance between the predicted NMR spectra of the candidate molecules and the input NMR spectrum. The second scoring function used the spectral distance between the predicted IR spectra of the candidate molecules, and the actual IR spectrum of the molecule corresponding to the input NMR spectrum. There is a decreased variability in the *Top N* metrics when the candidate molecules are ranked by scoring functions that use the IR spectral distance, which illustrates that IR spectra are able

Table 5: Training on molecules with  $\leq 7$  atoms and testing on molecules with  $\geq 8$  atoms for  $n_{\text{mcts}} = 400$

|            | $\geq 8$ atom molecules | 8-atom molecules | 9-atom molecules |
|------------|-------------------------|------------------|------------------|
| Top 1 (%)  | 80.971                  | 96.024           | 77.948           |
| Top 3 (%)  | 82.046                  | 97.247           | 78.992           |
| Top 5 (%)  | 82.148                  | 97.247           | 79.115           |
| Top 10 (%) | 82.199                  | 97.247           | 79.176           |
| Top 40 (%) | 82.250                  | 97.247           | 79.238           |

to distinguish between molecules better. Even limiting the usage of IR spectra to just sort the candidate molecules vastly improves the performance of DeepSPInN by ensuring that the correct molecule is ranked highly.

### Generalizability of DeepSPInN in understanding the action space

To understand how well DeepSPInN generalizes learning about the actions, the prior and value models were first trained on all molecules with less than 8 heavy atoms. It was then tested on a subset of the remaining molecules using these prior and value models. Table 5 shows the *Top N* metrics for this subset of test molecules, and the *Top N* metric for 8-atom molecules and 9-atom molecules in this subset. DeepSPInN achieves a *Top 1* accuracy of 80.97% even when all the test molecules have more heavy atoms than the molecules that DeepSPInN was trained on. The decreased accuracy when compared to the original model might be because there were very few molecules for training the prior and value models in this experiment. More details about the training and testing for this experiment are available in the Supplementary Information.

## Conclusions

DeepSPInN predicts the molecular structure when given an input IR and  $^{13}\text{C}$  NMR spectra without searching any pre-existing spectral databases or enumerating the possible structural

motifs present in the input spectra. After formulating the molecular structure prediction problem as an MDP, DeepSPInN employs MCTS to explore and choose the actions in the MDP. After building a null molecular graph from the molecular formula, DeepSPInN builds the molecular graph by treating the addition of each edge as an action in the MDP with the help of offline-trained GCNs to featurize each state in the MDP. DeepSPInN is able to correctly predict the molecular structure for 95.98% of input IR and  $^{13}\text{C}$  NMR spectra in an average time of 160 seconds for molecules with  $< 10$  heavy atoms.

DeepSPInN currently works on molecules that have less than 10 heavy atoms and future work could extend DeepSPInN to work on bigger molecules, or perhaps introduce other approaches that can easily be extended to bigger molecules. DeepSPInN currently requires the molecular formula of the molecule to be inferred from another chemical characterization technique apart from the input spectra. Although interpreting the entire structure from the spectra is harder, enumerating the splitting of the  $^{13}\text{C}$  NMR peaks is still a tough task because of the overlap of nearby peaks and also due to the noise in experimental  $^{13}\text{C}$  NMR spectra. Removing these strict requirements on the pre-processing of the  $^{13}\text{C}$  NMR spectra and removing the dependence on the molecular formula can be directions for future avenues of research in structure elucidation. Additionally, it will be interesting to see if DeepSPInN's accuracy improves with the addition of other spectral information such as UV-Vis spectra and mass spectra. We believe that DeepSPInN is a valuable demonstration of how machine learning can contribute to molecular structure prediction, and that it would help spur further research in the application of deep learning in high-throughput synthesis to enable faster and more efficient drug discovery pipelines.

## Author contributions

S.D., B.S., and S.M. wrote the bulk of the code specific to this work. B.S., S.M., Y.P., and S.L conceptualized the original approach. S.D. prepared the draft manuscript text and

figures with input from B.S., S.M, G.V., and U.D.P.. All co-authors were involved in editing and reviewing the final manuscript. U.D.P supervised the project.

## Competing interests

The authors have been partly or entirely supported by IHub-Data, DST-SERB (CRG/2021/008036), and IIIT Hyderabad's Kohli Center on Intelligent Systems. The funders however did not have any role in the conceptualization, design, data collection, analysis, decision to publish, or preparation of the manuscript.

## Correspondence

Any correspondence requests can be addressed to U.D.P. or S.D..

## Data Availability

The data and code used in this work can be accessed through the GitHub repository.

## References

- (1) Atkins PW, d. P. J. *Elements of physical chemistry (5th ed.)*; Oxford: Oxford U.P, 2009; p 459.
- (2) Janice Gorzynski Smith, S. J. *Organic Chemistry*; McGraw-Hill., Chapter 13 Mass Spectrometry and Infrared Spectroscopy.
- (3) Brian E. Mann, B. F. T. *<sup>13</sup>C NMR data for organometallic compounds*; Academic Press, 1981.

- (4) *The Theory of NMR - Chemical Shift*; University of Colorado, Boulder, Chemistry and Biochemistry Department, 2011.
- (5) Elyashberg, M.; Argyropoulos, D. Computer Assisted Structure Elucidation (CASE): Current and future perspectives. *Magnetic Resonance in Chemistry* **2021**, *59*, 669–690.
- (6) Burns, D. C.; Mazzola, E. P.; Reynolds, W. F. The role of computer-assisted structure elucidation (CASE) programs in the structure elucidation of complex natural products. *Nat. Prod. Rep.* **2019**, *36*, 919–933.
- (7) Steinbeck, C. Recent developments in automated structure elucidation of natural products. *Nat. Prod. Rep.* **2004**, *21*, 512–518.
- (8) Kwon, Y.; Lee, D.; Choi, Y.-S.; Kang, S. Molecular search by NMR spectrum based on evaluation of matching between spectrum and molecule. *Scientific Reports* **2021**, *11*, 1–9.
- (9) Dührkop, K.; Shen, H.; Meusel, M.; Rousu, J.; Böcker, S. Searching molecular structure databases with tandem mass spectra using CSI: FingerID. *Proceedings of the National Academy of Sciences* **2015**, *112*, 12580–12585.
- (10) Elyashberg, M. E.; Williams, A.; Blinov, K. *Contemporary Computer-Assisted Approaches to Molecular Structure Elucidation*; New Developments in NMR; The Royal Society of Chemistry, 2012; pp P001–482.
- (11) Hemmer, M. C.; Gasteiger, J. Prediction of three-dimensional molecular structures using information from infrared spectra. *Analytica Chimica Acta* **2000**, *420*, 145–154.
- (12) Valli, M.; Russo, H. M.; Pilon, A. C.; Pinto, M. E. F.; Dias, N. B.; Freire, R. T.; Castro-Gamboa, I.; da Silva Bolzani, V. Computational methods for NMR and MS for structure elucidation I: software for basic NMR. *Physical Sciences Reviews* **2019**, *4*, 20180108.



- (13) Valli, M.; Russo, H. M.; Pilon, A. C.; Pinto, M. E. F.; Dias, N. B.; Freire, R. T.; Castro-Gamboa, I.; da Silva Bolzani, V. Computational methods for NMR and MS for structure elucidation II: database resources and advanced methods. *Physical Sciences Reviews* **2019**, *4*, 20180167.
- (14) Bitchagno, G. T. M.; Tanemossu, S. A. F. Computational methods for NMR and MS for structure elucidation III: More advanced approaches. *Physical Sciences Reviews* **2019**, *4*, 20180109.
- (15) Elyashberg, M.; Blinov, K.; Martirosian, E. A new approach to computer-aided molecular structure elucidation: the expert system Structure Elucidator. *Laboratory Automation & Information Management* **1999**, *34*, 15–30.
- (16) Vamathevan, J.; Clark, D.; Czodrowski, P.; Dunham, I.; Ferran, E.; Lee, G.; Li, B.; Madabhushi, A.; Shah, P.; Spitzer, M.; Zhao, S. Applications of machine learning in drug discovery and development. *Nature Reviews Drug Discovery* **2019**, *18*, 463–477.
- (17) Ekins, S.; Puhl, A. C.; Zorn, K. M.; Lane, T. R.; Russo, D. P.; Klein, J. J.; Hickey, A. J.; Clark, A. M. Exploiting machine learning for end-to-end drug discovery and development. *Nature Materials* **2019**, *18*, 435–441.
- (18) Mehta, S.; Laghuvarapu, S.; Pathak, Y.; Sethi, A.; Alvala, M.; Priyakumar, U. D. MEMES: Machine learning framework for Enhanced MolEcular Screening. *Chem. Sci.* **2021**, *12*, 11710–11721.
- (19) Manzhos, S.; Carrington, T. Neural Network Potential Energy Surfaces for Small Molecules and Reactions. *Chemical Reviews* **2021**, *121*, 10187–10217.
- (20) Pattnaik, P.; Raghunathan, S.; Kalluri, T.; Bhimalapuram, P.; Jawahar, C. V.; Priyakumar, U. D. Machine Learning for Accurate Force Calculations in Molecular Dynamics Simulations. *The Journal of Physical Chemistry A* **2020**, *124*, 6954–6967.

- (21) Noé, F.; Tkatchenko, A.; Müller, K.-R.; Clementi, C. Machine Learning for Molecular Simulation. *Annual Review of Physical Chemistry* **2020**, *71*, 361–390, PMID: 32092281.
- (22) Samaga, Y. B. L.; Raghunathan, S.; Priyakumar, U. D. SCONES: Self-Consistent Neural Network for Protein Stability Prediction Upon Mutation. *The Journal of Physical Chemistry B* **2021**, *125*, 10657–10671.
- (23) Aggarwal, R.; Gupta, A.; Chelur, V.; Jawahar, C. V.; Priyakumar, U. D. DeepPocket: Ligand Binding Site Detection and Segmentation using 3D Convolutional Neural Networks. *Journal of Chemical Information and Modeling* **2021**,
- (24) Pathak, Y.; Laghuvarapu, S.; Mehta, S.; Priyakumar, U. D. Chemically Interpretable Graph Interaction Network for Prediction of Pharmacokinetic Properties of Drug-Like Molecules. *Proceedings of the AAAI Conference on Artificial Intelligence* **2020**, *34*, 873–880.
- (25) Laghuvarapu, S.; Pathak, Y.; Priyakumar, U. D. BAND NN: A Deep Learning Framework for Energy Prediction and Geometry Optimization of Organic Small Molecules. *Journal of Computational Chemistry* **2020**, *41*, 790–799.
- (26) Goel, M.; Raghunathan, S.; Laghuvarapu, S.; Priyakumar, U. D. MoleGuLAR: Molecule Generation Using Reinforcement Learning with Alternating Rewards. *Journal of Chemical Information and Modeling* **2021**, *61*, 5815–5826.
- (27) Ricard, D.; Cachet, C.; Cabrol-Bass, D.; Forrest, T. P. Neural network approach to structural feature recognition from infrared spectra. *Journal of Chemical Information and Computer Sciences* **1993**, *33*, 202–210.
- (28) Ren, H.; Li, H.; Zhang, Q.; Liang, L.; Guo, W.; Huang, F.; Luo, Y.; Jiang, J. A machine learning vibrational spectroscopy protocol for spectrum prediction and spectrum-based structure recognition. *Fundamental Research* **2021**, *1*, 488–494.

- (29) Yao, K.; Herr, J. E.; Toth, D.; Mckintyre, R.; Parkhill, J. The TensorMol-0.1 model chemistry: a neural network augmented with long-range physics. *Chem. Sci.* **2018**, *9*, 2261–2269.
- (30) Kananenka, A. A.; Yao, K.; Corcelli, S. A.; Skinner, J. L. Machine Learning for Vibrational Spectroscopic Maps. *Journal of Chemical Theory and Computation* **2019**, *15*, 6850–6858.
- (31) Gastegger, M.; Behler, J.; Marquetand, P. Machine learning molecular dynamics for the simulation of infrared spectra. *Chem. Sci.* **2017**, *8*, 6924–6935.
- (32) Schütt, K. T.; Kessel, P.; Gastegger, M.; Nicoli, K. A.; Tkatchenko, A.; Müller, K.-R. SchNetPack: A Deep Learning Toolbox For Atomistic Systems. *Journal of Chemical Theory and Computation* **2019**, *15*, 448–455.
- (33) F.M., P.; A., H.; F., M.; S., D.; M., C.; L., E. Chemical shifts in molecular solids by machine learning. *Nature Communications* **2018**, *9*, All Open Access, Gold Open Access, Green Open Access.
- (34) Jonas, E.; Kuhn, S. Rapid prediction of NMR spectral properties with quantified uncertainty. *Journal of Cheminformatics* **2019**, *11*, 50.
- (35) Yang, Z.; Chakraborty, M.; White, A. D. Predicting chemical shifts with graph neural networks. *Chem. Sci.* **2021**, *12*, 10802–10809.
- (36) Ye, S.; Hu, W.; Li, X.; Zhang, J.; Zhong, K.; Zhang, G.; Luo, Y.; Mukamel, S.; Jiang, J. A neural network protocol for electronic excitations of  $\text{N}_i\text{-methylacetamide}$ . *Proceedings of the National Academy of Sciences* **2019**, *116*, 11612–11617.
- (37) Ghosh, K.; Stuke, A.; Todorović, M.; Jørgensen, P. B.; Schmidt, M. N.; Vehtari, A.; Rinke, P. Deep Learning Spectroscopy: Neural Networks for Molecular Excitation Spectra. *Advanced Science* **2019**, *6*, 1801367.

- (38) S., K. G.; U., S.; J.M., R. Purifying Electron Spectra from Noisy Pulses with Machine Learning Using Synthetic Hamilton Matrices. *Physical Review Letters* **2020**, *124*, Cited by: 7; All Open Access, Green Open Access, Hybrid Gold Open Access.
- (39) Ongie, G.; Jalal, A.; Metzler, C. A.; Baraniuk, R. G.; Dimakis, A. G.; Willett, R. Deep Learning Techniques for Inverse Problems in Imaging. 2020; <https://arxiv.org/abs/2005.06001>.
- (40) Wang, Z.; Feng, X.; Liu, J.; Lu, M.; Li, M. Functional groups prediction from infrared spectra based on computer-assist approaches. *Microchemical Journal* **2020**, *159*, 105395.
- (41) Fine, J. A.; Rajasekar, A. A.; Jethava, K. P.; Chopra, G. Spectral deep learning for prediction and prospective validation of functional groups. *Chem. Sci.* **2020**, *11*, 4618–4630.
- (42) Jonas, E. Deep imitation learning for molecular inverse problems. *Advances in Neural Information Processing Systems*. 2019.
- (43) Howarth, A.; Ermanis, K.; Goodman, J. M. DP4-AI automated NMR data analysis: straight from spectrometer to structure. *Chem. Sci.* **2020**, *11*, 4351–4359.
- (44) Zhang, J.; Terayama, K.; Sumita, M.; Yoshizoe, K.; Ito, K.; Kikuchi, J.; Tsuda, K. NMR-TS: de novo molecule identification from NMR spectra. *Science and Technology of Advanced Materials* **2020**, *21*, 552–561.
- (45) Yang, X.; Zhang, J.; Yoshizoe, K.; Terayama, K.; Tsuda, K. ChemTS: an efficient python library for de novo molecular generation. *Science and technology of advanced materials* **2017**, *18*, 972–976.
- (46) Huang, Z.; Chen, M. S.; Woroch, C. P.; Markland, T. E.; Kanan, M. W. A framework

- for automated structure elucidation from routine NMR spectra. *Chem. Sci.* **2021**, *12*, 15329–15338.
- (47) Pesek, M.; Juvan, A.; Jakos, J.; Kosmrlj, J.; Marolt, M.; Gazvoda, M. Database Independent Automated Structure Elucidation of Organic Molecules Based on IR, <sup>1</sup>H NMR, <sup>13</sup>C NMR, and MS Data. *Journal of chemical information and modeling* **2020**, *61*, 756–763.
- (48) Silver, D.; Hubert, T.; Schrittwieser, J.; Antonoglou, I.; Lai, M.; Guez, A.; Lanctot, M.; Sifre, L.; Kumaran, D.; Graepel, T.; Lillicrap, T.; Simonyan, K.; Hassabis, D. Mastering Chess and Shogi by Self-Play with a General Reinforcement Learning Algorithm. 2017; <https://arxiv.org/abs/1712.01815>.
- (49) Alberts, M.; Laino, T.; Vaucher, A. C. Leveraging Infrared Spectroscopy for Automated Structure Elucidation. 2023.
- (50) Alberts, M.; Zipoli, F.; Vaucher, A. C. Learning the Language of NMR: Structure Elucidation from NMR spectra using Transformer Models. 2023.
- (51) Sridharan, B.; Mehta, S.; Pathak, Y.; Priyakumar, U. D. Deep Reinforcement Learning for Molecular Inverse Problem of Nuclear Magnetic Resonance Spectra to Molecular Structure. *The Journal of Physical Chemistry Letters* **2022**, 4924–4933.
- (52) Kuhn, S.; Schlörer, N. E. Facilitating quality control for spectra assignments of small organic molecules: nmrshiftdb2 – a free in-house NMR database with integrated LIMS for academic service laboratories. *Magnetic Resonance in Chemistry* **2015**, *53*, 582–589.
- (53) Gilmer, J.; Schoenholz, S. S.; Riley, P. F.; Vinyals, O.; Dahl, G. E. Neural Message Passing for Quantum Chemistry. *CoRR* **2017**, *abs/1704.01212*.
- (54) Ruddigkeit, L.; van Deursen, R.; Blum, L. C.; Reymond, J.-L. Enumeration of 166

- Billion Organic Small Molecules in the Chemical Universe Database GDB-17. *Journal of Chemical Information and Modeling* **2012**, *52*, 2864–2875, PMID: 23088335.
- (55) Ramakrishnan, R.; Dral, P. O.; Rupp, M.; von Lilienfeld, O. A. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific Data* **2014**, *1*, 140022.
- (56) Ruddigkeit, L.; van Deursen, R.; Blum, L. C.; Reymond, J.-L. Enumeration of 166 Billion Organic Small Molecules in the Chemical Universe Database GDB-17. *Journal of Chemical Information and Modeling* **2012**, *52*, 2864–2875.
- (57) Frisch, M. J. et al. Gaussian 09, Revision A.1. 2016; Gaussian Inc. Wallingford CT.
- (58) McGill, C.; Forsuelo, M.; Guan, Y.; Green, W. H. Predicting Infrared Spectra with Message Passing Neural Networks. *Journal of Chemical Information and Modeling* **2021**, *61*, 2594–2609.
- (59) Gupta, A.; Chakraborty, S.; Ramakrishnan, R. Revving up <sup>13</sup>C NMR shielding predictions across chemical space: benchmarks for atoms-in-molecules kernel machine learning with new data for 134 kilo molecules. *Machine Learning: Science and Technology* **2021**, *2*, 035010.
- (60) Mehring, M. *High resolution NMR spectroscopy in solids*; Springer Science & Business Media, 2012; Vol. 11.
- (61) Silver, D. et al. Mastering the game of Go with deep neural networks and tree search. *Nature* **2016**, *529*, 484–489.
- (62) James, S.; Konidaris, G.; Rosman, B. An Analysis of Monte Carlo Tree Search. Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence. 2017; p 3576–3582.
- (63) Sutton, R. S.; Barto, A. G. *Reinforcement Learning: An Introduction*, 2nd ed.; The MIT Press, 2018.

- (64) Yang, K.; Swanson, K.; Jin, W.; Coley, C.; Eiden, P.; Gao, H.; Guzman-Perez, A.; Hopper, T.; Kelley, B.; Mathea, M.; Palmer, A.; Settels, V.; Jaakkola, T.; Jensen, K.; Barzilay, R. Analyzing Learned Molecular Representations for Property Prediction. *Journal of Chemical Information and Modeling* **2019**, *59*, 3370–3388.
- (65) Dai, H.; Dai, B.; Song, L. Discriminative Embeddings of Latent Variable Models for Structured Data. *CoRR* **2016**, *abs/1603.05629*.
- (66) Chang, C.-I. An information-theoretic approach to spectral variability, similarity, and discrimination for hyperspectral image analysis. *IEEE Transactions on Information Theory* **2000**, *46*, 1927–1932.
- (67) Kocsis, L.; Szepesvári, C. Bandit Based Monte-Carlo Planning. Machine Learning: ECML 2006. Berlin, Heidelberg, 2006; pp 282–293.
- (68) Segler, M. H. S.; Preuss, M.; Waller, M. P. Planning chemical syntheses with deep neural networks and symbolic AI. *Nature* **2018**, *555*, 604–610.
- (69) Moerland, T. M.; Broekens, J.; Plaat, A.; Jonker, C. M. Monte Carlo Tree Search for Asymmetric Trees. *arXiv preprint arXiv:1805.09218* **2018**,
- (70) Silver, D. et al. Mastering the game of Go without human knowledge. *Nature* **2017**, *550*, 354–359.
- (71) Duvenaud, D. K.; Maclaurin, D.; Iparraguirre, J.; Bombarell, R.; Hirzel, T.; Aspuru-Guzik, A.; Adams, R. P. Convolutional Networks on Graphs for Learning Molecular Fingerprints. Advances in Neural Information Processing Systems. 2015.
- (72) Landrum, G. RDKit: Open-Source Cheminformatics Software. **2016**,