
Can LLMs Solve Molecule Puzzles? A Multimodal Benchmark for Molecular Structure Elucidation

Kehan Guo^{1*}, Bozhao Nan^{2*}, Yujun Zhou¹, Taicheng Guo¹, Zhichun Guo¹,
Mihir Surve², Zhenwen Liang¹, Nitesh V. Chawla¹, Olaf Wiest², Xiangliang Zhang¹

¹Department of Computer Science and Engineering, ²Department of Chemistry and Biochemistry,
University of Notre Dame

{kguo2, bnan, yzhou25, tguo2, zguo5, msurve, zliang6, nchawla, owiest, xzhang33}@nd.edu

Abstract

1 Large Language Models (LLMs) have shown significant problem-solving capabilities
2 across predictive and generative tasks in chemistry. However, their proficiency
3 in multi-step chemical reasoning remains underexplored. We introduce a new
4 challenge: molecular structure elucidation, which involves deducing a molecule’s
5 structure from various types of spectral data. Solving such a molecular puzzle, akin
6 to solving crossword puzzles, poses reasoning challenges that require integrating
7 clues from diverse sources and engaging in iterative hypothesis testing. To address
8 this challenging problem with LLMs, we present **MolPuzzle**, a benchmark comprising
9 234 instances of structure elucidation, which feature over 18,000 QA samples
10 presented in a sequential puzzle-solving process, involving three interlinked sub-
11 tasks: molecule understanding, spectrum interpretation, and molecule construction.
12 Our evaluation of more than 10 LLMs reveals that the best-performing LLM, GPT-
13 4o, performs significantly worse than humans, with only a small portion (1.4%)
14 of its answers exactly matching the ground truth. However, it performs nearly
15 perfectly in the first subtask of molecule understanding, achieving accuracy close
16 to 100%. This discrepancy highlights the potential of developing advanced LLMs
17 with improved chemical reasoning capabilities in the other two sub-tasks. Our
18 MolPuzzle dataset and evaluation code are available at this link.

19 1 Introduction

20 Artificial intelligence (AI) is revolutionizing the field of chemistry, influencing diverse sectors such as
21 industrial chemical engineering [1, 2], drug discovery [3], and chemistry education [4]. In particular,
22 recent studies have highlighted the success of large language models (LLMs) in addressing predictive
23 challenges in chemistry, including molecular property prediction [5], reaction prediction [6], and
24 experiment automation [7]. These advancements suggest significant potential for AI to enhance
25 efficiency and innovation across these critical areas.

26 We introduce a new chemical challenge to AI, **molecular structure elucidation**. **While this critical**
27 **task has been explored in other contexts, it remains unexplored for large language models (LLMs),**
28 **extending beyond familiar predictive and generative domains such as property or reaction prediction,**
29 **and representing a shift toward complex problem-solving.** Analogous to solving a detailed cross-
30 word puzzle, **molecular structure elucidation** can be seen as a **molecular puzzle**. It requires the

*Both authors contributed equally to this work, supported by the NSF Center for Computer-Assisted Synthesis (C-CAS). <https://ccas.nd.edu>

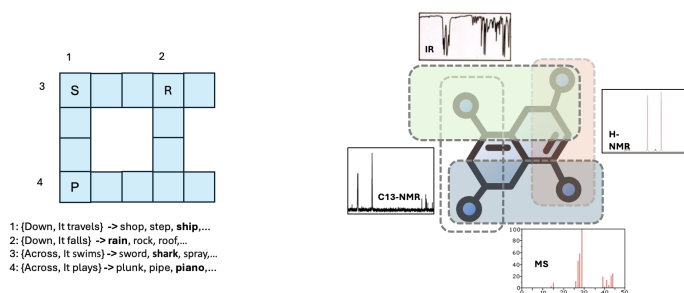


Figure 1: A crossword puzzle (left), and a molecular structure elucidation puzzle (right)

31 integration of multifaceted data, iterative hypothesis testing, and a deep understanding of chemical
 32 cues, much like piecing together clues across a crossword grid to form a coherent solution. Fig. 1
 33 illustrates the problem of molecular structure elucidation alongside its analogical counterpart, the
 34 crossword puzzle, highlighting the parallels in strategy and complexity between these two intellectual
 35 challenges.

36 Just as a crossword puzzle requires figuring out words based on given clues and fitting them together
 37 in a grid, molecular structure elucidation involves deducing a molecule's structure from various types
 38 of data such as nuclear magnetic resonance (NMR), infrared spectroscopy (IR), mass spectrometry,
 39 and others. Each type of data provides clues about different aspects of the molecular structure. In
 40 a crossword, we integrate clues from across different directions and hints to form words that fit
 41 together correctly. Similarly, in molecular structure elucidation, we need to integrate information
 42 from different spectroscopic methods to form a consistent picture of the molecule. For example,
 43 IR spectra reveal molecular vibrations and functional groups, NMR provides information about
 44 the framework of hydrogen and carbon atoms, while mass spectrometry can offer insights into the
 45 molecular weight and possible fragmentations.

46 Nevertheless, molecular structure elucidation is a challenging and time-consuming task. Training
 47 undergraduate students in chemistry to solve these puzzles has been a part of the curriculum be-
 48 cause determining the structure of molecules is a fundamental skill in the field. Typically, even a
 49 single molecule puzzle question on a final exam can take 10 to 15 minutes to solve[8], demanding
 50 considerable memory and processing skills from the students. In the domain of complex molecule
 51 research, the process of molecular deduction can become even more complex and time-consuming.
 52 Therefore, fully automating this process is highly beneficial for accelerating the design of new
 53 materials and drugs, as well as enhancing the efficiency of chemical research[9, 10]. However, it
 54 remains a challenging task due to the complexities involved in interpreting spectral data and solving
 55 intricate reasoning problems associated with molecular structures [11].

56 In this work, we aim to present molecular structure elucidation in formats that LLMs can effectively
 57 process. By adapting this complex task to be compatible with LLMs, we explore their potential as
 58 promising tools in chemical research. If successful, LLMs could significantly accelerate scientific
 59 discovery in chemistry, transforming how we approach and solve intricate molecular puzzles.

60 To achieve our objectives, we first introduce a novel dataset named **MolPuzzle**, which includes
 61 234 instances of structure elucidation challenges inspired by common chemistry tasks. Unlike
 62 datasets used in predictive or generative tasks, which typically consist of a collection of independent
 63 samples and are relatively straightforward to construct, each instance in the MolPuzzle dataset is
 64 uniquely complex. It is structured as a sequential process involving three interlinked sub-tasks:
 65 **molecule understanding**, **spectrum interpretation**, and **molecule construction**. These instances
 66 are accompanied by multimodal data, including images of IR, MASS, H-NMR, and C-NMR spectra,
 67 alongside their corresponding molecular formulas. Presenting such a complex, multimodal problem in
 68 a format that LLMs can effectively process presents a unique challenge. We, a team of AI researchers
 69 and chemists, are dedicated to formulating the molecule puzzle instances in descriptive languages
 70 that are accessible to LLMs. Our focus is on ensuring the utility of these instances, as well as their
 71 comprehensive coverage over various scenarios and challenges that mimic real-world conditions. By
 72 doing so, **MolPuzzle** opens the door for LLMs to contribute meaningfully to the field of chemistry,
 73 potentially accelerating scientific discoveries and innovations.

74 Second, we present our effort to automate the solving of molecular structure elucidation using LLMs.
75 While certain sub-tasks, such as translating an IR spectrum into a molecular formula, may be solvable
76 by encoder-decoder models [12], the comprehensive resolution of the entire molecular puzzle likely
77 requires the advanced planning and reasoning capabilities of LLMs. We tested 11 state-of-the-art
78 LLMs including GPT-4o, Gemini-pro, and Claude-3-opus. We also conducted a human baseline to
79 compare the performance of humans and LLMs in solving the same puzzles. The **key findings** are:
80 1) GPT-4o significantly outperforms other LLMs; 2) The best-performing LLM, GPT-4o, performs
81 significantly worse than humans, with only a small portion (1.4%) of its answers exactly matching
82 the ground truth; and 3) GPT-4o’s performance primarily collapses in the Stage-2 of spectrum
83 interpretation and gets worse in the Stage-3 of molecule construction, although it performs nearly
84 perfectly in Stage-1 of molecule understanding (with accuracy close to 100%).

85 To summarize, our key contributions in this work are the presentation of:

- 86 • **A new reasoning problem for AI community.** As the focus of AI development has evolved
87 from solving predictive tasks and generative tasks to engaging in complex reasoning tasks—akin
88 to system 2 level thinking—we introduce a reasoning task centered around molecular structure
89 elucidation. This crucial problem from the field of chemistry sets a high benchmark for AI models
90 to reach. Solving this task requires AI models to possess the ability to interpret spectral images,
91 engage in complex reasoning, and plan effectively across extended workflows. This not only
92 challenges the current capabilities of AI but also pushes the boundaries of what AI can achieve in
93 scientific domains, particularly in understanding and manipulating molecular structures.
- 94 • **A new light of AI solutions for chemistry community.** By proposing the **MolPuzzle dataset**,
95 we establish another bridge between the fields of AI and chemistry. This initiative leverages the
96 important capabilities of multimodal LLMs, providing the chemistry community with innovative
97 solutions to accelerate the process of structure elucidation. Our initial exploration serves as a
98 demonstration of the potential for these technologies. It sets the stage for further collaborative
99 efforts, inspiring researchers from both domains to collaboratively explore new frontiers in scientific
100 discovery.

101 The paper is organized as follows. Section 2 presents the related work. In Section 3, we elaborate
102 on the curation of the MolPuzzle dataset. In Section 4, we report the usage of multimodal LLMs in
103 solving MolPuzzle. In Section 5, we discuss the main findings and directions opened by this work. In
104 section 6, we discuss the broader impact of our work. Last, we summarize the study in Section 7 and
105 offer our conclusions.

106 2 Related Work

107 **Molecular Structure Elucidation.** Historically, chemists used basic methods such as crystalliza-
108 tion, melting point determination, and simple reactivity tests to hypothesize about a molecule’s
109 structure. As technology advanced, tools like infrared spectroscopy (IR), nuclear magnetic resonance
110 (NMR), and mass spectrometry transformed the process, enabling precise molecular insights and
111 revolutionizing chemical analysis. Recently, Alberts et al. [12] utilized a transformer-based model
112 to predict SMILES strings from IR spectra, later extending this architecture to NMR data analysis
113 [13]. However, much of the existing research focuses on molecule elucidation using single-type
114 spectrum data, which may suffice for simple molecules. In practice, complex molecules cannot be
115 fully elucidated from a single spectrum since each type of spectrum provides only partial structural
116 information. In our study, we aim to leverage the reasoning and planning capabilities of multimodal
117 large language models (MLLMs) to integrate diverse spectral data, addressing the challenges of
118 complex real-world chemistry tasks. Our focus is on solving the entire puzzle using multiple clues,
119 rather than merely deciphering one word from a single clue.

120 **Multimodal Benchmarks for LLMs.** With the advancements in developing multimodal LLMs
121 [14, 15, 16, 17], a number of multimodal benchmarks have been curated. These benchmarks are
122 crucial for evaluating and refining the capabilities of MLLMs to process and integrate diverse data
123 types, such as text, images, and audio, for a cohesive understanding. Notably, a benchmark proposed
124 by Yue et al. [18] assesses the reasoning abilities of MLLMs in various college-level subjects.

125 Similarly, MathVista [19] explores MLLMs’ multimodal reasoning capabilities in mathematics, while
126 Yin et al. [20] introduced LAMM, a dataset focusing on multimodal instruction tuning. Our research
127 shifts the focus to the chemistry domain [6]. To our knowledge, this study is the first to adopt a
128 realistic chemistry task for MLLM processing and to conduct a thorough evaluation of these models’
129 proficiency in chemistry-related reasoning and image analysis. This specialized focus will enhance
130 our understanding of MLLMs’ capabilities within a specific scientific domain.

131 3 The MolPuzzle Dataset

132 Existing benchmarks of chemical tasks primarily focused on predictive or generative tasks involving
133 collections of independent samples that were relatively straightforward to construct. In contrast,
134 our dataset, MolPuzzle, aims to characterize an intertwined assessment of chemistry reasoning and
135 visual understanding, testing the application of AI-assisted technology towards broader scientific
136 discovery. Our data collection process is rigorously designed and implemented by a team uniquely
137 qualified for this task, consisting of esteemed researchers in chemistry and experienced AI specialists
138 who have previously tackled complex chemistry problems. This collaboration ensures that the
139 MolPuzzle dataset not only accurately reflects real-world chemical phenomena and challenges but is
140 also structured in a way that optimally facilitates access and usability for LLMs.

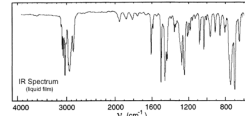
141 The basic principles guiding our data curation for the MolPuzzle dataset are: 1) ensuring compre-
142 hensive coverage by including a wide range of tasks that synthesize visual context with chemical
143 knowledge, facilitating thorough evaluations; 2) varying levels of difficulty to challenge LLMs
144 and highlight their potential limitations; 3) ensuring robust assessment outcomes, i.e., the results
145 are definitive and reliable; and 4) incorporating human expert analysis to identify strengths and
146 weaknesses in model performance, significantly enhancing our understanding of LLMs capabilities.

147 In this section, we outlined the construction process for the MolPuzzle dataset. We detailed the
148 creation of puzzle tasks in three stages (3.1), as well as the QA pairs involved in these tasks (3.2).
149 Examples are presented in Fig. 2.

150 3.1 Task Construction

151 Just like a word puzzle where each clue progressively reveals the final answer, the solution to a
152 molecule puzzle is a SMILES string that captures the interconnected substructures of a molecule. We
153 design our molecule puzzles so that solving one requires the accurate identification and integration of
154 each substructural clue, gradually unveiling the complete SMILES representation of the molecule.
155 This approach is inspired by the analytical strategies employed by chemists in the real world, who
156 interpret spectral data and chemical properties to deduce the structures of unknown molecules. Our
157 puzzle-building process mirrors this scientific exploration, arranging clues in a sequence from simple
158 to complex, where each clue builds upon the insights gained from the previous one, requiring precision
159 and careful thought at every stage. We next provide more details on our clue design methodology.

160 **The Initial Stage (Molecule Understanding).** In designing a molecule puzzle, the first stage involves
161 determining how many building blocks, or substructures, are available. This foundational step is
162 crucial as it sets the stage for constructing the molecule’s complete structure, akin to identifying the
163 key pieces in a complex jigsaw puzzle. Starting with the initial hint: [A molecular formula, derived
164 from a mass spectrum, indicates the exact types and numbers of atoms in a molecule](#) (e.g., $C_{15}H_{22}O_2$,
165 representing carbon, hydrogen, and oxygen), chemists can begin to [deduce possible structures from
166 the degree of saturation which is calculated based on the number of rings and multiple bonds
167 present in the molecule](#), the potential for forming aromatic rings, or the presence of functional
168 groups. The initial information provides a preliminary range of [building blocks](#), which can later be
169 [selected and assembled](#) to solve the molecular puzzle. To benchmark the capability of LLMs in this
170 stage, we developed 26 unique templates (see Appendix A.2 for details), targeting key analytical tasks
171 such as saturation identification, aromatic ring identification, functional group identification, and
172 saturation degree calculation. This initiative produced 6,318 QA-format pairs, effectively evaluating
173 the models’ capacity to understand and process molecular data. Details of these samples are reported
174 in Appendix A.3.

<p>1. Identify molecule substructures based on molecule formula</p> <p>Prompt: As an expert organic chemist, your task is to analyze the chemical formula C₆H₁₀O₆ and determine the potential molecular structures and the degree of unsaturation. Utilize your knowledge to systematically explore and identify plausible molecular substructure.</p>	<p>2. Refine the substructure pools based on Spectrum images.</p>  <p>Prompt: As an expert in organic chemistry, you are tasked with analyzing potential molecular structures derived from IR spectral data. Given the molecular formula and an initial set of potential fragment SMILES identified, your objective is to explore and systematically determine plausible molecular substructure that are consistent with the IR spectral data.</p>	<p>3. Select fragments from the pools and assemble molecule iteratively</p> <p>Initial selection: Prompt: Selected one fragment from the list of SMILES for the Initial structure for molecular construction: Identify one specific fragment from the [pool of fragments] provided: ensuring it's consistent with both [C13-NMR] and [H-NMR].</p> <p>Iteration: Prompt: Select one fragment from the provided list of SMILES to add to the current molecule. Identify a specific fragment from the [pool of fragments]; ensuring it is consistent with both the [C13-NMR] and [H-NMR] spectra.</p> <p>End: when run out of heavy atoms.</p>
<p>Answer: Carboxylic Acid (Yes) degree of unsaturation = 2</p>	<p>Answer: ["C(=O)O", "C(=O)OC", "C=O", "CO", "C1CO1"]</p>	<p>Answer: C1C(C(C(C(O)O)O)O)C(=O)O</p>

(a). The Initial Stage

(b). The Second Stage

(c). The Final Stage

Figure 2: Examples of QA pairs in the 3 stages of MolPuzzle

175 **The Second Stage (Spectrum Interpretation).** With the initial building blocks of the molecule
 176 identified from the molecular formula, the next critical step involves refining these components
 177 through detailed spectral analysis. Spectrum images such as IR, MASS, ¹H-NMR, and ¹³C-NMR
 178 serve as new hints, each adding layers of information akin to clues in a complex puzzle. These
 179 spectral images are pivotal in confirming or revising the initial hypotheses about the molecule's
 180 structure. For example, IR spectroscopy can verify the presence of specific functional groups, MASS
 181 spectrometry can provide the molecular MASS, molecule mass and fragmentation patterns, and NMR
 182 techniques detail the arrangement of hydrogen and carbon within the molecule. By integrating these
 183 new hints, researchers can construct a more robust and experimentally accurate model of the molecule.
 184 This process not only theoretically validates each building block but also ensures they align perfectly
 185 with empirical data, leading to a comprehensive understanding of the molecular structure. Given the
 186 importance of spectral images in this analysis, we have developed specialized question templates to
 187 evaluate the proficiency of LLMs in interpreting these images. For instance, we created 17 templates
 188 for IR and 12 for each of H-NMR, and C-NMR. Each template, such as 'Analyze the IR spectrum'
 189 includes specific queries designed to extract detailed insights, such as 'What does the absorption in
 190 3200-3600 suggest?' This structure enables us to format the questions for Visual Question Answering
 191 (VQA), facilitating a systematic approach to query handling. Our method has successfully generated
 192 a significant repository of VQA format examples, comprising 3,978 for IR and 2,808 for each of
 193 MASS, H-NMR, and C-NMR. A detailed analysis of these tasks is available in Appendix A.4.

194 **The Final Stage (Molecule Construction).** After completing the first two stages, we can assert that
 195 we have gathered the necessary building blocks to assemble the molecule. The assembly process will
 196 be guided by insights derived from NMR data. Specifically, ¹H-NMR provides information about
 197 the hydrogen environment in the molecule, such as the number of hydrogen atoms, their types (e.g.,
 198 aliphatic, aromatic), and their connectivity. On the other hand, C-NMR offers detailed insights into
 199 the carbon framework, revealing how carbon atoms are distributed and linked within the molecule.
 200 The approach to assembling the final molecular structure is iterative. Starting with initial building
 201 blocks selected from the identified fragment pool, LLMs are prompted to select one structure from
 202 the pool step by step, based on the NMR guidance, until the maximum number of iterations is reached
 203 or the fragment pool is exhausted. This systematic addition ensures that each step in the assembly
 204 process not only fits with the previous structure but also aligns perfectly with the latest spectral data,
 205 driving us closer to the accurate molecular configuration. This approach results in a total of 1,171
 206 QA samples.

207 3.2 QA Sample Derivation

208 The QA samples for Stage 1 and Stage 2 are automatically generated using their respective question
209 templates (see Appendix A.2) and RDKit [21]. RDKit is an open-source cheminformatics toolkit
210 widely employed for handling chemical informatics data, including molecular structures and finger-
211 prints. This toolkit plays a role in ensuring that the responses, based on the SMILES strings from
212 each molecule puzzle, are accurate and chemically valid. The distribution of these QA samples across
213 different categories is illustrated in Fig. 4. They form a diverse collection of samples for evaluating
214 LLMs’ ability to understand molecular formulas and spectra.

215 The fragment of each QA pair at Stage 3 is initially generated by LLMs, i.e., responding to the
216 prompt ‘select one fragment...’. To validate the reliability of these automated generations of QA
217 pairs, experts—two Ph.D. candidates from the chemistry department—manually and independently
218 verified 50 samples, labeling the generated fragments as ‘correct’ or ‘wrong’. Their verification
219 was consistent and demonstrated that 67.4% of examples have correct fragment pools in automated
220 generation. To ensure the quality of derived QA pairs in Stage 3, these chemists manually corrected
221 the fragments pool for each instance in the benchmark.

222 Fig. 3 reports the statistical distribution for the MolPuzzle dataset, which includes 234 puzzle
223 instances (the reasoning of 234 different molecules). Since one puzzle can be solved by different
224 paths, different numbers of QA samples are derived in three stages. We will next evaluate LLMs’
225 performance in solving each puzzle, as well as their capability to solve individual questions.

Statistic	Number
Total MolPuzzle Instances	234
Stage-1 QA samples	6,318
- Num. of molecule formula	176
- Max question length	128
- Average question length	94
Stage-2 QA samples	12,402
- Num. of spectrum images	944
- Max question length	340
- Average question length	264
Stage-3 QA samples	1,171
- Maximum Iteration	7
- Max question length	356
- Average question length	238

Figure 3: Statistic of the MolPuzzle dataset

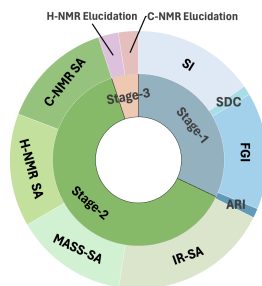


Figure 4: Inner ring: sample distribution in 3 stages. Outer ring: sample distribution across categories in each stage. SI: saturation identification, SDC: saturation degree calculation, FGI: functional group identification, ARI: aromatic ring identification, SA: spectrum analysis.

226 4 Solving MolPuzzle by Multimodal Large Language Models

227 The reasoning capabilities of foundation models in the chemistry domain remain underexplored.
228 Thus, our aim is to perform both qualitative and quantitative evaluations to systematically assess the
229 reasoning and planning abilities of these models in visual chemistry contexts, using the MolPuzzle
230 benchmark. We first conducted evaluation of a variety of LLMs for completing the individual tasks
231 in each stage, including GPT-4o [22], GPT-3.5-turbo [23], Claude-3-opus [24], Gemini-pro [25],
232 Llama-3-8B-Instruct [26], Vicuna-13B-v1.5 [27], Mistral-7B-Instruct-v0.3 [28], and in particular
233 multimodal LLMs such as Gemini-pro-vision [25], Llava-Llama-3-8B [29], Qwen-VL-Chat [30],
234 and InstructBlip-Vicuna-7B/13B [14]. Due to space limits, we present only selected results in Table 1
235 and report the complete list of results in Appendix B. We then assess LLMs’ capability to solve the
236 entire puzzles, specifically focusing on how effectively these models can derive the final molecular
237 structure from provided hints (the questions in QA samples). The results are reported in Table 2.

238 All tasks are evaluated in a zero-shot setting to determine the problem-solving capabilities of LLMs
239 without prior fine-tuning on specific task data. The evaluation process consists of three steps:
240 response generation, answer extraction, and score calculation. More details of the experimental
241 settings including prompts and hyperparameters are presented in Appendix B.1.

Table 1: F1 scores (\uparrow) of individual QA tasks in three stages. The best LLMs results are in bold font. Tasks in stage 1 are SI-Saturation Identification, ARI-Aromatic Ring Identification, FGI-Functional Group Identification, and SDC-Saturation Degree Calculation.

Stage-1 (Molecule Understanding) Tasks				
Method	SI	ARI	FGI	SDC
GPT-4o	1.00±0.000	0.943±0.016	0.934±0.005	0.667±0.003
GPT-3.5-turbo	0.451±0.025	0.816±0.017	0.826±0.075	0.5±0.099
Claude-3-opus	0.361±0.009	0.988±0.015	0.934±0.001	0.856±0.016
Llama3	0.228±0.043	0.696±0.051	0.521±0.003	0.000±0.000
Human	1.00±0.000	1.000±0.000	0.890±0.259	0.851±0.342
Stage-2 (Spectrum Interpretation) Tasks				
Method	IR Interpretation	MASS Interpretation	H-NMR Interpretation	C-NMR Interpretation
GPT-4o	0.656±0.052	0.609±0.042	0.618±0.026	0.639±0.010
LLava	0.256±0.026	0.101±0.021	0.118±0.008	0.254±0.015
Human	0.753±0.221	0.730±0.11	0.764±0.169	0.769±0.101
Stage-3 (Molecule Construction) Tasks				
Method	H-NMR Elucidation		C-NMR Elucidation	
GPT-4o	0.433±0.013		0.408±0.034	
Llama3	0.211±0.012		0.342±0.007	
Human	0.867±0.230		0.730±0.220	

242 To gain an in-depth understanding of the performance of LLMs in comparison with human experts,
 243 particularly their failed cases, we invited six Ph.D. candidates in chemistry [add acknowledgment](#)
 244 [later](#) to solve the puzzles in MolPuzzle, and also assess LLMs’ results. More comprehensive details
 245 of this **human baseline** and evaluation process are presented in Appendix B.2. The reported
 246 performance, including human baselines, is presented as an average with standard deviation over all
 247 samples.

248 4.1 LLMs’ Performance on Solving Molecule Puzzles

249 4.1.1 Addressing individual QA tasks in three stages

250 In Table 1, we report the performance of selected LLMs on conducting individual QA tasks in the three
 251 stages, including GPT-4o, GPT-3.5-turbo, Claude-3-opus (three top-performing proprietary models),
 252 Llama-3-8B-Instruct (the best performing open-source model), and the reference human baseline
 253 performance. In stage 2, the variant of Llama3 for a multimodal setting, LLava-Llama-3-8B, is used
 254 for handling spectrum image analysis. Since each task involves performing a question-answering
 255 task, we evaluate the performance using F1 and accuracy by comparing the LLMs’ answers with the
 256 ground truth. F1 scores are reported in Table 1, while the accuracy and performance of more LLMs
 257 can be found in Appendix B.

258 The results of Stage-1 (in Table 1 and Appendix Table 3) show that the GPT-4o model excels in these
 259 tasks (achieving near-perfect F1 score in 3 out of 4 tasks). The high scores in SI, AI, and FI suggest
 260 that LLMs are able to succeed in relatively straightforward chemistry analysis tasks, performing
 261 comparably to human experts. However, open-sourced models like LLama3 have limitations in
 262 addressing these tasks, possibly due to their limited reasoning abilities in chemistry text-reasoning
 263 tasks. In addition, GPT-4o’s comparative performance to humans indicates significant advancements
 264 in the use of LLMs for complex scientific tasks, suggesting a promising future for leveraging advanced
 265 LLMs to improve the efficiency of scientific analysis and discovery.

266 For the multimodal tasks of Stage-2, GPT-4o remains the top performer, though it exhibits intermedi-
 267 ate performance in spectrum interpretation. The F1 scores for the four types of spectra average around
 268 0.6, indicating a moderate level of accuracy in this complex aspect of the challenge. This performance
 269 is notably less competitive compared to human baselines, which succeed in approximately 73-77% of
 270 the tasks across the four types of spectrum interpretation. This indicates that spectrum interpretation
 271 is inherently challenging. While GPT-4o has made significant strides in automated spectrum analysis,

Table 2: The performance of LLMs and human baseline in solving MolPuzzle. The best LLM results are in bold font. Acc. stands for the Accuracy of Exact Match.

Method	Acc. (\uparrow)	Levenshtein (\downarrow)	Validity (\uparrow)	MACCS FTS (\uparrow)	RDKit FTS (\uparrow)	Morgan FTS (\uparrow)
GPT-4o	0.014\pm0.004	11.653\pm0.013	1.000\pm0.000	0.431\pm0.009	0.293\pm0.013	0.232 \pm 0.007
Claude-3-opus	0.013 \pm 0.008	12.680 \pm 0.086	1.000\pm0.000	0.383 \pm 0.050	0.264 \pm 0.040	0.241\pm0.037
Gemini-pro	0.000 \pm 0.000	12.711 \pm 0.196	1.000\pm0.000	0.340 \pm 0.017	0.208 \pm 0.002	0.171 \pm 0.007
Human	0.667 \pm 0.447	1.332 \pm 2.111	1.000 \pm 0.000	0.985 \pm 0.022	0.795 \pm 0.317	0.810 \pm 0.135

272 there remains considerable room for improvement to bridge the gap between its capabilities and
 273 human expertise. More details are presented in Appendix B.4.

274 The results for Stage-3 indicate that the most advanced LLM, GPT-4o, significantly underperforms
 275 compared to the human baseline, with nearly a 40% difference. This might be caused by the fact that
 276 the reasoning ability required for these tasks is complex and multifaceted. When information con-
 277 verges, such as identifying equivalent hydrogen or ring arrangements, a comprehensive understanding
 278 of the NMR peaks and their corresponding structures is essential. See more details in Appendix B.5.

279 4.1.2 Addressing entire molecule puzzles

280 For solving the entire molecule puzzles, the evaluation is limited to the three most advanced mul-
 281 timodal LLMs: GPT-4o [22], Claude-3-opus [24], and Gemini-pro [25], due to the involvement
 282 of spectrum image analysis in Stage-2. The results of these models are reported in Table 2, along
 283 with those from the human baseline. To comprehensively evaluate the performance, we employ two
 284 different types of metrics. The first type of metric measures the chemical similarity between the
 285 ground-truth molecules and the generated molecules, assessed using FTS (Fingerprint Tanimoto Simi-
 286 larity) [31] in terms of MACCS [32], RDKit [21], and Morgan [33]. Since the generated molecules are
 287 in SMILES string format, we also employ natural language processing metrics including the Accuracy
 288 of Exact Match [34], and Levenshtein distance [35] (the minimum number of single-character editing
 289 required to transform one string into another). Finally, to evaluate whether constructed molecules are
 290 valid, we use RDKit [21] to check the validity of constructed molecules and report the percentage of
 291 molecules that are confirmed as valid.

292 The results in Table 2 show that the best-performed LLM, GPT-4o, is performing much worse than
 293 humans, indicating a huge gap between LLMs and humans in solving the molecule puzzles. It is
 294 worth noting that all the constructed molecules are valid, even though only a small portion of them
 295 (1.4%) exactly match the ground truth. Considering that the accuracy of the exact match is too strict,
 296 we use FTS to analyze more about the chemical closeness of LLMs’ answer to the ground truth. A
 297 MACCS FTS of 0.431 suggests that the generated molecules maintain a significant level of structural
 298 similarity. This indicates that even if the answers are not perfect replicas of the ground truth, they
 299 can still be chemically valid and potentially useful as structured hypotheses that could be relived by
 300 human scientists.

301 4.2 Success and Failure Analysis

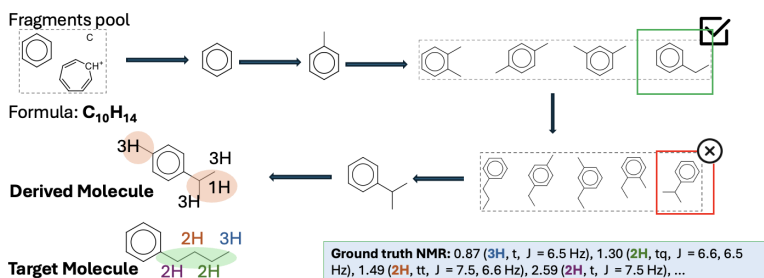


Figure 5: Errors in solving the molecule puzzle

302 The above analysis indicates that the most capable model, GPT-4o, performs **nearly perfectly**
303 in Stage-1 of molecule understanding. However, its performance **drops** in Stage-2 for spectrum
304 interpretation, and **worsens further in Stage-3** for molecule construction. We investigate in-depth
305 how GPT-4o eventually fails on most of the puzzles after progressing through the tasks of these three
306 stages. With the help of human evaluators, we gathered all the intermediate steps involved in solving a
307 molecule puzzle and engaged them to scrutinize these steps. Fig. 5 presents case studies that illustrate
308 the iterative steps involved in Stage-3, showcasing the most common errors made by GPT-4o: **the**
309 **accumulation of errors in iterative steps, which can lead to catastrophic failures**. Note that
310 this stage focuses on selecting the correct fragments and assembling them step by step to form the
311 final molecular structure. We find that GPT-4o can initially succeed in picking the correct fragment
312 when the structure is comparatively simple. However, as the process progresses, it does not select
313 structures that satisfy all the requirements indicated by the NMR data. This difficulty arises because
314 the reasoning requirements expand dramatically as more information and additional constraints need
315 to be incorporated. More qualitative examples can be found in Appendix B.6.

316 5 Findings and Open Directions

317 Our evaluation has revealed specific limitations of state-of-the-art LLMs in automating molecular
318 structure elucidation. We urge further collaborative efforts from the AI and chemistry communities to
319 design more effective solutions, especially for the tasks in Stage 2 and Stage 3. Based on our findings,
320 we next present the open directions for future research and development.

321 **Development of Specialized Multimodal LLMs Spectrum Interpretation in Stage 2.** As indi-
322 cated in our results, the performance of LLMs notably declines beginning in Stage 2, where they
323 struggle with the visual interpretation of ^1H and ^{13}C NMR spectra. This difficulty arises because
324 NMR spectra feature sharp, unlabeled peaks that also display multiplicities with very small chemical
325 shift differences, making them challenging to discern for visual models. [These multiplicities contain](#)
326 [important information on the chemical connectivity of the fragments](#). Similarly, [closely spaced IR](#)
327 [absorptions to identify key function groups](#). To address this, there is a significant opportunity to
328 develop specialized multimodal LLMs that can more effectively interpret these subtle and complex
329 spectral details.

330 **Development of New Strategies for Leveraging LLMs in Chemical-related Planning and Reason-**
331 **ing.** The failure analysis from Stage 3 has inspired us to explore more effective ways to capitalize
332 on LLMs' capabilities in planning and reasoning for fragment selection and assembly. The first imme-
333 diate solution is to employ the chain-of-thought approach [36] to provide more effective instructions
334 for solving the puzzle. However, despite our efforts to implement this method, the results were not
335 satisfying and actually performed worse than those in the zero-shot setting we reported in the paper.
336 We will continue the study and try different implementations. The second solution is to leverage
337 LLMs as agents in a more dynamic and interactive manner. This approach involves incorporating
338 feedback loops where LLMs can iteratively refine their responses based on new information or
339 corrections. In this way, there is a hope to mitigate the accumulation of errors in iterative steps and
340 prevent catastrophic failures.

341 6 Broader Impact

342 Our work has broad impacts across multiple dimensions. First, it offers valuable insights and
343 recommendations for both AI researchers and chemists in academia and industry. These perspectives
344 enhance the effective utilization of LLMs and guide future advancements in the field. Second,
345 our approach to benchmarking and improving LLMs through real-world tasks like the MolPuzzle
346 can also foster greater collaboration between computational scientists and chemists. By aligning
347 AI technologies with traditional chemical research, these interdisciplinary efforts can accelerate
348 the discovery of new materials, drugs, and chemical processes, potentially leading to significant
349 advancements in healthcare and industry.

350 **7 Conclusion**

351 In this paper, we introduced MolPuzzle, a new benchmark challenge to advance our capabilities in
352 molecular structure elucidation. We evaluated state-of-the-art LLMs on this task, revealing their
353 strengths and limitations in handling complex chemical reasoning. Our analysis highlights significant
354 performance gaps, particularly in spectrum interpretation and molecule construction. These findings
355 not only suggest ways to improve LLM performance but also set the stage for transforming approaches
356 to chemical research. MolPuzzle serves as a critical step toward harnessing the potential of LLMs
357 in chemistry, fostering innovation and collaboration within the AI and chemistry communities to
358 enhance scientific inquiry and application.

References

- [1] Venkat Venkatasubramanian. The promise of artificial intelligence in chemical engineering: Is it here, finally? *AIChE Journal*, 65(2):466–478, 2019.
- [2] Zachary J Baum, Xiang Yu, Philippe Y Ayala, Yanan Zhao, Steven P Watkins, and Qiongqiong Zhou. Artificial intelligence in chemistry: current trends and future directions. *Journal of Chemical Information and Modeling*, 61(7):3197–3212, 2021.
- [3] Alexandre Blanco-Gonzalez, Alfonso Cabezon, Alejandro Seco-Gonzalez, Daniel Conde-Torres, Paula Antelo-Riveiro, Angel Pineiro, and Rebeca Garcia-Fandino. The role of ai in drug discovery: challenges, opportunities, and strategies. *Pharmaceuticals*, 16(6):891, 2023.
- [4] Xuan-Quy Dao, Ngoc-Bich Le, Bac-Bien Ngo, and Xuan-Dung Phan. LLMs’ capabilities at the high school level in chemistry: Cases of chatgpt and microsoft bing ai chat. 2023.
- [5] Suryanarayanan Balaji, Rishikesh Magar, Yayati Jadhav, et al. GPT-MolBERTa: GPT Molecular Features Language Model for molecular property prediction. *arXiv preprint arXiv:2310.03030*, 2023.
- [6] Taicheng Guo, Bozhao Nan, Zhenwen Liang, Zhichun Guo, Nitesh Chawla, Olaf Wiest, Xi-angliang Zhang, et al. What can large language models do in chemistry? a comprehensive benchmark on eight tasks. *Advances in Neural Information Processing Systems*, 36:59662–59688, 2023.
- [7] Andres M Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D White, and Philippe Schwaller. Chemcrow: Augmenting large-language models with chemistry tools. *arXiv preprint arXiv:2304.05376*, 2023.
- [8] Alan M Rosan. Organic structures from spectra, (field, ld; sternhell, s.; kalman, jr), 2002.
- [9] Roman M Balabin, Ekaterina I Lomakina, and Ravilya Z Safieva. Neural network (ANN) approach to biodiesel analysis: analysis of biodiesel density, kinematic viscosity, methanol and water contents using near infrared (NIR) spectroscopy. *Fuel*, 90(5):2007–2015, 2011.
- [10] Liu Cao, Mustafa Guler, Azat Tagirdzhanov, Yi-Yuan Lee, Alexey Gurevich, and Hosein Mohimani. Moldiscovery: Learning mass spectrometry fragmentation of small molecules. *Nature communications*, 12(1):3718, 2021.
- [11] Xi Xue, Hanyu Sun, Minjian Yang, Xue Liu, Hai-Yu Hu, Yafeng Deng, and Xiaojian Wang. Advances in the application of artificial intelligence-based spectral data interpretation: A perspective. *Analytical Chemistry*, 95(37):13733–13745, 2023.
- [12] Marvin Alberts, Teodoro Laino, and Alain C Vaucher. Leveraging infrared spectroscopy for automated structure elucidation. 2023.
- [13] Marvin Alberts, Federico Zipoli, and Alain C Vaucher. Learning the Language of NMR: Structure Elucidation from NMR spectra using Transformer Models. 2023.
- [14] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *Advances in Neural Information Processing Systems*, 36, 2024.
- [15] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.

- 404 [16] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances*
405 *in neural information processing systems*, 36, 2024.
- 406 [17] Zhenwen Liang, Kehan Guo, Gang Liu, Taicheng Guo, Yujun Zhou, Tianyu Yang, Jiajun Jiao,
407 Renjie Pi, Jipeng Zhang, and Xiangliang Zhang. Scemqa: A scientific college entrance level
408 multimodal question answering benchmark. *arXiv preprint arXiv:2402.05138*, 2024.
- 409 [18] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens,
410 Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal
411 understanding and reasoning benchmark for expert agi. *arXiv preprint arXiv:2311.16502*, 2023.
- 412 [19] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao
413 Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical
414 reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.
- 415 [20] Zhenfei Yin, Jiong Wang, Jianjian Cao, Zhelun Shi, Dingning Liu, Mukai Li, Xiaoshui Huang,
416 Zhiyong Wang, Lu Sheng, Lei Bai, et al. Lamm: Language-assisted multi-modal instruction-
417 tuning dataset, framework, and benchmark. *Advances in Neural Information Processing Systems*,
418 36, 2024.
- 419 [21] G. A. Landrum. Rdkit: Open-source cheminformatics software. <http://www.rdkit.org>, 2020.
- 420 [22] OpenAI. Hello GPT-4o. <https://openai.com/index/hello-gpt-4o/>, 2023.
- 421 [23] OpenAI. GPT-3.5-Turbo: Enhancements and Applications. <https://openai.com/models/gpt-3.5-turbo>, 2023.
- 422 [24] Anthropic. Introducing the Claude-3 Family. <https://www.anthropic.com/news/claude-3-family>, 2023.
- 423 [25] Google. Introducing gemini: our largest and most capable ai model, 2023.
- 424 [26] Meta. Introducing Meta Llama 3. <https://llama.meta.com/llama3/>, 2023.
- 427 [27] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,
428 Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot
429 impressing gpt-4 with 90%* chatgpt quality, march 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna>, 3(5), 2023.
- 430 [28] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh
431 Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile
432 Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 433 [29] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual
434 instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023.
- 435 [30] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,
436 Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- 437 [31] Taffee T Tanimoto. Elementary mathematical theory of classification and prediction. *Journal of*
438 *Biomedical Science and Engineering*, 1958.
- 439 [32] David Ratcliff, John W.; Metzener. Pattern matching: The gestalt approach, 1988.
- 440 [33] Debadutta Dash, Rahul Thapa, Juan M Banda, Akshay Swaminathan, Morgan Cheatham,
441 Mehr Kashyap, Nikesh Kotecha, Jonathan H Chen, Saurabh Gombhar, Lance Downing, et al.
442 Evaluation of GPT-3.5 and GPT-4 for supporting real-world information needs in healthcare
443 delivery. *arXiv preprint arXiv:2304.13714*, 2023.
- 444 [34] Carl Edwards, Tuan Lai, Kevin Ros, Garrett Honke, and Heng Ji. Translation between molecules
445 and natural language. *arXiv preprint arXiv:2204.11817*, 2022.

- 447 [35] Frederic P Miller, Agnes F Vandome, and John McBrewster. Levenshtein distance: Information
448 theory, computer science, string (computer science), string metric, damerau? Levenshtein
449 distance, spell checker, hamming distance, 2009.
- 450 [36] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le,
451 Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models.
452 *Advances in neural information processing systems*, 35:24824–24837, 2022.
- 453 [37] Sunghwan Kim, Jie Chen, Tiejun Cheng, Asta Gindulyte, Jia He, Siqian He, Qingliang Li,
454 Benjamin A Shoemaker, Paul A Thiessen, Bo Yu, et al. Pubchem 2023 update. *Nucleic acids
455 research*, 51(D1):D1373–D1380, 2023.
- 456 [38] Kevin Maik Jablonka, Luc Patiny, and Berend Smit. Making molecules vibrate: Interactive web
457 environment for the teaching of infrared spectroscopy, 2022.
- 458 [39] Fei Wang, Dana Allen, Siyang Tian, Eponine Oler, Vasuk Gautam, Russell Greiner, Thomas O
459 Metz, and David S Wishart. Cfm-id 4.0—a web server for accurate ms-based metabolite
460 identification. *Nucleic acids research*, 50(W1):W165–W174, 2022.

461 **A MolPuzzle Benchmark Details**

462 This section complements Section 3 with a fine-grained summary of the dataset collection, results
463 validation, and evaluation procedure, along with a fuller characterization of the task instances and the
464 corresponding prompts.

465 **A.1 Data Collection**

466 The initial molecules were selected by referencing the textbook *Organic Structures from Spectra, 4th*
467 *Edition*, available as an online PDF on ResearchGate. We chose 234 molecules based on spectrum
468 tasks involving IR, MS, $^1\text{H-NMR}$, and $^{13}\text{C-NMR}$ to reflect a difficulty level suitable for graduate
469 students[8].

470 To address copyright concerns, we excluded molecules with publicly available mass spectrometry
471 (MS) spectra in open-source databases from our study. The remaining spectra were sourced from
472 public resources, notably the PubChem database[37]. For additional spectra that were not available,
473 we used simulation methods[38][39] and provided a Jupyter notebook to generate these data, ensuring
474 high-quality spectra for analysis. Our final dataset comprised 200 molecules.

475 Given the challenges associated with NMR spectrum images, some spectra were obtained from
476 simulated data in text format for $^1\text{H-NMR}$ and $^{13}\text{C-NMR}$. This approach ensured clarity and accuracy
477 in the evaluation of molecular structures.

478 To assess the multiple-stage abilities of LLMs, we designed a unique question-and-answer evaluation.
479 This framework tested the LLMs' capabilities in interpreting and integrating data from different types
480 of spectra, simulating real-world challenges. Details of this evaluation framework are provided in the
481 next section.

482 **A.2 Template design**

483 Each template was crafted to target specific skills within molecular understanding. For instance,
484 saturation identification challenges the models' ability to discern the degree of saturation in a molecule,
485 which is crucial for understanding its chemical reactivity and stability. Aromatic ring identification
486 tests the models' ability to recognize benzene-like structures, which are fundamental in organic
487 chemistry due to their common occurrence and unique properties. Saturation degree calculation
488 pushes the models to apply quantitative analysis, requiring not just recognition but also computation
489 based on molecular structures.

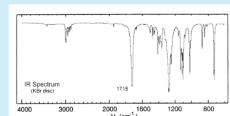
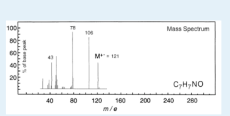
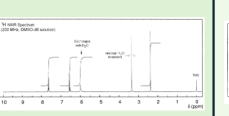
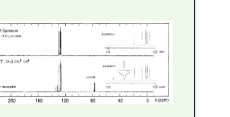
490 By diving deeper into the rationale behind each template and the kind of chemical knowledge they
491 are designed to test, we can better appreciate how these tasks simulate real-world applications in
492 chemistry. This approach not only tests the models' basic recognition abilities but also their capacity
493 to perform complex reasoning and apply theoretical knowledge practically. The template examples
494 are in A.3.

495 **A.3 Stage1 QA Samples**

Table 3: QA samples for the molecule understanding task

Task	Prompt
Saturation Identification	Question: Could the molecule with the formula C ₈ H ₁₀ O potentially be Saturated? Answer: No Model response: No.
Aromatic Ring Identification	Question: Could the molecule with the formula C ₈ H ₁₀ O have aromatic rings? Answer: Yes Model response: Yes.
Functional Group Identification	Question: Could the molecule with the formula C ₆ H ₁₄ O ₂ potentially contain a Amine group, given the Degree of Unsaturation is 0.0? Answer: No Model response: No, the molecule doesn't contain Amine group
Saturation Degree Calculation	Question: Calculate the Degree of Unsaturation of the molecule with the formula C ₈ H ₁₀ O? Answer: 4.0 Model response: 2

496 **A.4 Stage2 QA Samples**

IR Interpretation	MASS Interpretation	H-NMR Interpretation	C-NMR Interpretation
 <p>Question: Does the IR spectrum contains broad absorption peak of N-H stretching around 3200-3600 cm⁻¹?</p>	 <p>Question: Examine the MASS spectrum to determine if the molecule could potentially contain specific fragments: Ether.</p>	 <p>Question: Examine the H-NMR spectrum to determine if the molecule could potentially contain specific functional groups: Phenol?</p>	 <p>Question: Examine the C-NMR spectrum to determine if the molecule could potentially contain specific fragments: Ester.</p>
Answer: No Model response: No	Answer: No Model response: Yes	Answer: No Model response: No	Answer: No Model response: Yes

497 **A.5 Stage3 QA Samples**

Table 4: QA samples for the molecule construction task

Task	Prompt
H-NMR Elucidation	Question: Calculate the number of different types of hydrogen atoms present in the molecule, based on the given H-NMR: 4.51-4.61 (4H, 4.56 (s), 4.56 (s)), 7.06-7.32 (10H, 7.13 (dddd, J = 7.9, 7.7, 1.8, 0.6 Hz), 7.13 (dddd, J = 7.9, 7.7, 1.8, 0.6 Hz), 7.25 (dddd, J = 7.9, 1.5, 1.3, 0.6 Hz), 7.25 (dddd, J = 7.9, 1.5, 1.3, 0.6 Hz), 7.26 (tt, J = 7.7, 1.5 Hz), 7.26 (tt, J = 7.7, 1.5 Hz)) Answer: 4 Model response: 3.
C-NMR Elucidation	Question: Analyze the given C-NMR data and determine the number of different types of carbon atoms present in the molecule, based on given C-NMR: 39.3 (1C, s), 63.4 (1C, s), 127.8 (1C, s), 128.4 (2C, s), 128.8 (2C, s), 134.2 (1C, s). Only output the number. Answer: 6 Model response: 8

498 **B Evaluation Experiments**499 **B.1 Experimental Setting**

500 During our testing phase, we selected 100 questions and employed two distinct prompting strategies
501 with the large language model (LLM). Initially, the LLM was tasked with directly answering the

502 questions. In a subsequent approach, the same queries were presented, but the model was prompted to
503 execute a chain-of-thought reasoning process before responding. Each question in our dataset begins
504 with a comprehensive description of the chemical context, along with specified answer formats and
505 detailed guiding rules. To ensure a balanced representation of each task category, for tasks in Stage 1,
506 the distribution ratio for Saturation Identification (SI), Functional Group Identification (FI), Aromatic
507 Ring Identification (AI), and Saturation Degree Calculation (SC) is set at 2:3:3:2. In Stage 2, we
508 have randomly selected 100 questions from each category of the spectrum. For Stage 3, we randomly
509 selected 100 questions focused on H-NMR and C-NMR analyses.

510 We carried out this evaluation over three rounds, analyzing responses using both accuracy and the
511 F1 score for tasks involving Saturation Identification (SI), Functional Group Identification (FI), and
512 Aromatic Ring Identification (AI). For Saturation Degree Calculation (SDC), which yields numerical
513 results, we assessed accuracy by comparing the count of correct matches to the ground truth data.
514 The detailed results are reported in Table A.3. To ensure that all results are presented in a way that
515 facilitates direct comparison, only those using similar evaluation metrics(AI, FI, AI) are included
516 in the main table. For the SI, AI, and FI tasks, we use the F1 score and Accuracy to evaluate their
517 performance since they are classification tasks. For the SDC task, the answer is a numerical number,
518 so we only use the accuracy score to measure the performance of the LLMs. This approach helps to
519 keep the evaluation coherent and focused on comparable data points.

520 **B.2 Human Evaluation**

521 To evaluate the performance of large language models (LLMs) on specialized tasks against expert
522 humans, we recruited six graduate students from chemistry department to solve the MolPuzzle
523 benchmark. These students, having recently completed a graduate-level course in Molecular Structural
524 Elucidation, represented a highly skilled group of human participants.

525 For the experiment, we randomly selected six questions from the MolPuzzle dataset for each stage of
526 the study. These questions were presented to the students in different formats according to the stage:
527 In Stages 1 and 2, the questions were simple Yes/No or required short answers. In Stage 3, to align
528 with the conventional methods chemists use to express chemical structures, students were asked to
529 upload images of their hand-drawn structures instead of using SMILES strings. These images were
530 manually compared to the ground truth to calculate scores.

531 We also imposed self-regulated time constraints to mirror the challenging nature of molecular
532 structural elucidation. Beyond individual stage evaluations, we presented each participant with a
533 complete molecule puzzle, consisting of a formula and four spectral images. The students were tasked
534 with solving these puzzles within a 20-minute time frame. Impressively, all participants successfully
535 submitted their solutions within the allotted period.

536 Our study included a component where human evaluators were involved to assess the performance
537 of the AI models. To ensure the protection and ethical treatment of all participants, we conducted a
538 thorough risk assessment. Potential risks identified included privacy concerns and the mental strain
539 of repetitive tasks. Mitigation strategies, such as ensuring anonymity and providing breaks, were
540 implemented to protect our evaluators.

541 The study was submitted for review and received approval from our Institutional Review Board (IRB).
542 The IRB approval number is [insert approval number], which verifies that our protocols met all ethical
543 guidelines for research involving human subjects. Throughout the project, we adhered strictly to
544 these protocols to ensure ongoing compliance with ethical standards.

545 **B.3 Stage1**

546 Molecule understanding requires comprehensive analysis and interpretation of molecular structures,
547 with a focus on chemical properties and spectroscopic data. In our study, we created a dataset of
548 234 molecules and developed eight distinct question templates across four categories: **Saturation**
549 **Identification(SI), Functional Group Identification(FI), Aromatic Ring Identification(AI), and**

550 **Saturation Degree Calculation(SC)**. These templates assess the ability to identify substructures,
 551 compute saturation levels, and infer structural presence, incorporating concepts in the chemistry
 552 reasoning process. Each question also necessitates a deep understanding of molecular bonding,
 553 stereochemistry, and functional group identification. Responses were generated using the RDKit
 554 library, ensuring precise and reliable answers grounded in established chemical informatics.

Table 3: The accuracy(\uparrow), F1 score(\uparrow) in 4 different molecule understanding categories, the best LLMs are in bold font.

Model	CoT	SI		AI		FI		SC
		F1	Acc	F1	Acc	F1	Acc	Acc
GPT-4o	-	1±0.0	1±0.0	0.943±0.016	0.944±0.015	0.934±0.005	0.966±0.0	0.667±0.003
GPT-4o	✓	1±0.0	1±0.0	0.911±0.031	0.911±0.031	0.689±0.025	0.766±0.027	0.816±0.062
GPT-3.5	-	0.451±0.025	0.825±0.075	0.816±0.017	0.816±0.075	0.826±0.075	0.683±0.016	0.5±0.099
GPT-3.5	✓	0.448±0.026	0.816±0.008	0.798±0.025	0.800±0.027	0.526±0.053	0.622±0.031	0.533±0.131
Claude-3-opus	-	0.361±0.009	0.556±0.023	0.988±0.015	0.988±0.015	0.934±0.001	0.966±0.001	0.856±0.016
Claude-3	✓	0.760±0.189	0.903±0.046	0.878±0.025	0.867±0.001	0.547±0.112	0.843±0.081	0.900±0.025
Gemini-pro	-	0.285±0.020	0.399±0.040	0.775±0.093	0.788±0.083	0.646±0.052	0.748±0.051	0.200±0.004
Gemini-pro	✓	0.391±0.045	0.651±0.108	0.685±0.088	0.688±0.087	0.562±0.018	0.629±0.023	0.283±0.062
LLama3	-	0.367±0.018	0.583±0.047	0.490±0.030	0.533±0.027	0.472±0.133	0.588±0.0	0.0±0.0
LLama3	✓	0.473±0.011	0.899±0.040	0.384±0.026	0.533±0.0	0.570±0.035	0.799±0.047	0.017±0.001
Vicuna-13b	-	0.031±0.022	0.033±0.025	0.500±0.087	0.522±0.083	0.308±0.038	0.311±0.041	0.0±0.0
Vicuna-13b	✓	0.380±0.023	0.616±0.062	0.342±0.006	0.522±0.157	0.516±0.080	0.855±0.016	0.0±0.0
Mistral-7b	-	0.221±0.014	0.283±0.025	0.384±0.005	0.500±0.0	0.319±0.014	0.322±0.157	0.0±0.0
Mistral-7b	✓	0.433±0.007	0.766±0.023	0.342±0.006	0.522±0.016	0.601±0.102	0.877±0.031	0.0±0.0

555 B.4 Stage2

556 The Spectrum interpretation tasks mainly measure the capability of LLMs in analyzing images
 557 related to identifying key substructures indicated by the spectrum plot. In this study, we utilize
 558 four distinct types of spectral images: nuclear magnetic resonance (NMR), infrared spectroscopy
 559 (IR), mass spectrometry, and others. Each type of data offers insights into various aspects of the
 560 molecular structure. We’ve created specific question templates for each spectrum, targeting peak
 561 and substructure identification factors. These templates are designed manually and emphasize the
 562 intricate connection between the spikes or troughs in the figures and the structures of the molecules.
 563 Responses were generated using the RDKit library to ensure correctness.

564 The findings from Stage 2 are presented in Table 4. We exclusively focus on the zero-shot learning
 565 outcomes, as our observations indicate that implementing chain-of-thought prompting leads to a
 566 deterioration in model performance. To address this limitation, we offer qualitative insights in B.6.

Table 4: The accuracy(\uparrow), F1 score(\uparrow) for IR, MASS spectrum, H-NMR, and C-NMR interpretation tasks. "-" means the results are not interoperable

Model	Stage-2 Tasks							
	IR Interpretation		MASS Interpretation		H-NMR Interpretation		C-NMR Interpretation	
	F1	Acc	F1	Acc	F1	Acc	F1	Acc
GPT-4o	0.656±0.052	0.713±0.06	0.609±0.042	0.767±0.042	0.618±0.026	0.864±0.007	0.639±0.107	0.892±0.049
Claude-3-opus	0.440±0.006	0.476±0.055	0.398±0.032	0.466±0.019	0.572±0.190	0.842±0.017	0.554±0.075	0.716±0.042
Gemini-3-pro-vision	0.194±0.002	0.119±0.016	0.116±0.036	0.124±0.038	0.545±0.048	0.851±0.062	0.492±0.016	0.619±0.044
LLava1.5-8b	0.256±0.026	0.414±0.044	0.101±0.021	0.104±0.26	0.118±0.008	0.186±0.011	0.254±0.015	0.472±0.023
Qwen-VL-Chat	0.243±0.027	0.392±0.043	0.125±0.006	0.116±0.021	0.255±0.007	0.611±0.031	-	-
InstructBLIP-7b	0.239±0.020	0.263±0.014	0.101±0.021	0.104±0.26	-	-	0.044±0.006	0.064±0.023
InstructBLIP-13b	0.239±0.020	0.263±0.014	0.101±0.021	0.104±0.26	-	-	0.047±0.014	0.067±0.025

567 B.5 Stage-3

568 Constructing a molecule involves a detailed analysis of NMR data, which is critical for understanding
 569 its structure. H-NMR data are essential as they provide information about the hydrogen environments
 570 within the molecule, including the number and types of hydrogen atoms (such as aliphatic or

571 aromatic), as well as their connectivity. Conversely, C-NMR data offer in-depth insights into the
 572 carbon framework, illustrating the distribution and linkage of carbon atoms within the molecule.
 573 In our study, to evaluate the ability of large language models (LLMs) to interpret NMR data, we
 574 generated 1,171 question-and-answer (QA) pairs. These pairs focus on key NMR interpretation tasks,
 575 such as counting hydrogen atom types and identifying substructures, which are critical for accurate
 576 analysis.

577 Despite observing moderate accuracy from the LLMs in Stage 2 of our testing, we enhanced the
 578 quality of the QA pairs in Stage 3 by providing the LLMs with verified NMR data, generated by using
 579 nmrd[?]. This approach ensures that the data used is reliable and helps maintain the integrity of
 580 our results. The findings from Stage 2 are presented in Table. We exclusively focus on the zero-shot
 581 learning outcomes, as our observations indicate that implementing chain-of-thought prompting leads
 582 to a deterioration in model performance. To address this limitation, we offer qualitative insights in

Table 5: The F1 score(\uparrow) for H-NMR, and C-NMR Structure Elucidation

Method	H-NMR Elucidation	C-NMR Elucidation
GPT-4o	0.433±0.013	0.411±0.034
Claude-3-opus	0.395±0.008	0.313±0.029
Gemini-pro	0.333±0.012	0.308±0.031
Llama3	0.211±0.012	0.342±0.007
Vicuna-13b	0.181±0.013	0.244±0.001
Mistral-7b	0.131±0.032	0.122±0.027

583 B.6 Qualitative Results

584 In this section, we present several examples using GPT-4’s chain-of-thought (CoT) reasoning to
 585 facilitate a clearer understanding of the results. We have enlisted two Ph.D. candidates from the
 586 chemistry department to evaluate these CoT outcomes. The analysis uses color coding to indicate the
 587 accuracy of the generated text: green signifies correct responses, red indicates incorrect ones, and
 588 yellow denotes responses that are partially correct.

589 B.6.1 Stage 2 examples

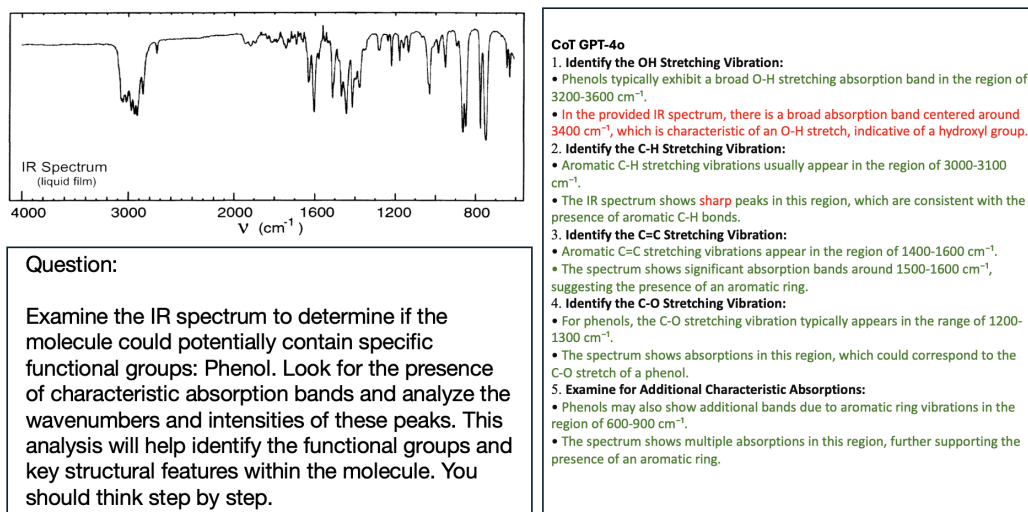


Figure 6: Human annotated IR spectrum interpretation

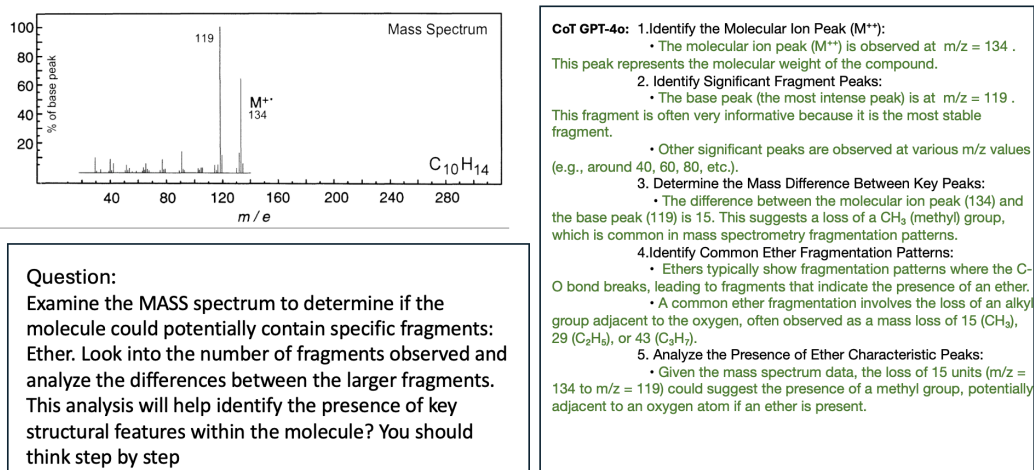


Figure 7: Human annotated MASS spectrum interpretation

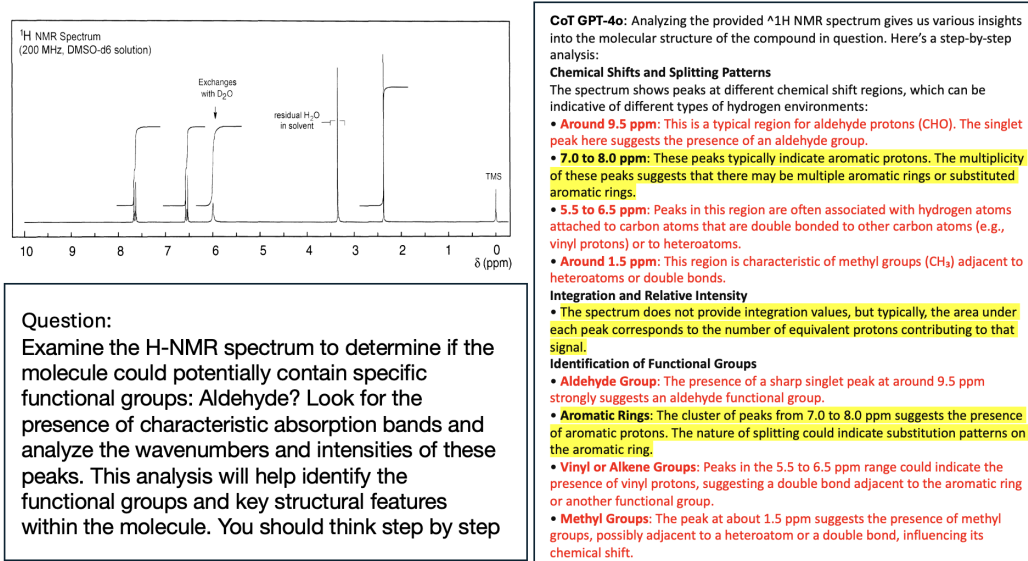


Figure 8: Human annotated H-NMR spectrum interpretation

590 B.6.2 Stage 3 examples

591 Based on the observation in spectrum interpretation, it is safe to say that the GPT-4o model has a
 592 basic grasp of the concept and the task at hand, but their deficiency mainly resides in obtaining fully
 593 correct information from the spectrum images and also they lack in-depth understanding. As for the
 594 problems in Stage 3, although the GPT-4o outputs are mostly correct, they are too general which
 595 shows the poor understanding of GPT-4o models in solving this intricate task. The stage 3 task can
 596 be roughly broken down into 3 subtasks: obtaining the correct information from the spectrum image,
 597 deducing the correct structural information from the spectral information, and finally translating this
 598 structural information into a correct molecular structure. GPT-4o models seem to perform well in the
 599 second subtask, and moderately for simple structures in the third subtask but seem to be especially
 600 struggling with the first subtask in the case of NMR spectra. This indicates the gap in current LLMs
 601 in fully interpreting data therefore more advanced models and approaches should be developed to
 602 tackle the problem.

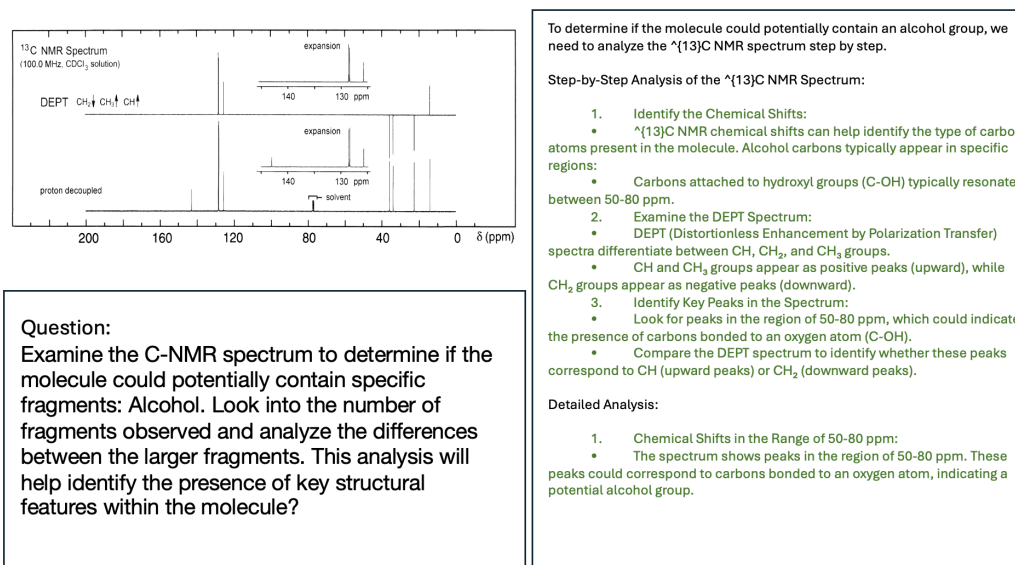


Figure 9: Human annotated C-NMR spectrum interpretation

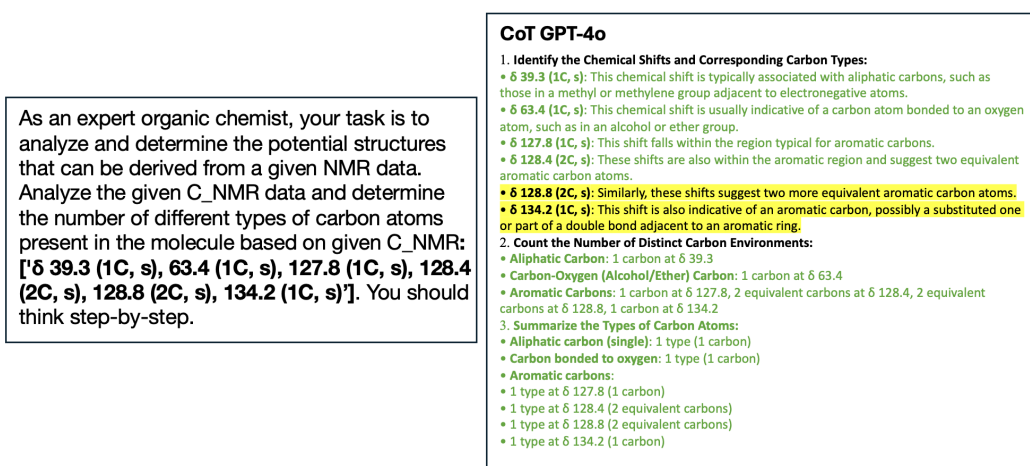


Figure 10: Human annotated H-NMR Elucidation

603 B.6.3 Complex Molecules

604 In addition to presenting molecules extracted from textbooks, we also demonstrate how the large
 605 language model (LLM) handles complex molecular structures. As illustrated in Figure 12, complex
 606 molecules typically have a larger pool of fragments. This expansion results in a greater number
 607 of valid elucidation paths, complicating the selection process for an appropriate starting point.
 608 Successfully navigating this enlarged pool necessitates an in-depth understanding of each fragment's
 609 properties and the associated, more intricate NMR data. In this context, LLMs may struggle because
 610 they often lack the nuanced chemical intuition and detailed analytical capabilities that human experts
 611 possess. Such limitations can lead to inaccuracies in interpreting complex interactions within NMR
 612 spectra, making LLMs less reliable.

As an expert organic chemist, your task is to analyze and determine the potential structures that can be derived from a given NMR data. Generate the number of different types of hydrogen atoms present in the molecule based on given H_NMR: [δ 2.23 (3H, s), 3.79 (3H, s), 7.04 (2H, ddd, J = 8.3, 1.3, 0.4 Hz), 7.97 (2H, ddd, J = 8.3, 1.6, 0.4 Hz)]. You should think step by step.

CoT GPT-4o

1. Identify the Chemical Shifts and Corresponding Hydrogen Types:

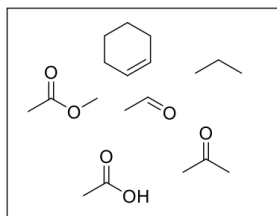
- δ 2.23 (3H, s): This singlet indicates 3 equivalent hydrogen atoms. These are likely part of a methyl group (CH₃) that is not split by neighboring hydrogens.
- δ 3.79 (3H, s): This singlet also indicates 3 equivalent hydrogen atoms, suggesting another methyl group (CH₃), possibly adjacent to an electronegative atom like oxygen (e.g., methoxy group -OCH₃).
- δ 7.04 (2H, ddd, J = 8.3, 1.3, 0.4 Hz): This complex splitting pattern (doublet of doublet of doublets) indicates 2 equivalent hydrogen atoms in an aromatic environment, with coupling constants suggesting interactions with neighboring hydrogens.
- δ 7.97 (2H, ddd, J = 8.3, 1.6, 0.4 Hz): Another complex splitting pattern, also indicating 2 equivalent hydrogen atoms in an aromatic environment, with similar coupling constants indicating interactions with neighboring hydrogens.

2. Determine the Number of Different Types of Hydrogen Environments:

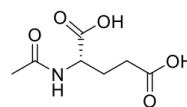
- **Aliphatic Hydrogens:**
 - 3 equivalent hydrogens at δ 2.23 (likely CH₃ group)
 - 3 equivalent hydrogens at δ 3.79 (likely CH₃ group, possibly part of a methoxy group)
- **Aromatic Hydrogens:**
 - 2 equivalent hydrogens at δ 7.04
 - 2 equivalent hydrogens at δ 7.97
 - The molecule has a total of 4 distinct types of hydrogen atoms:
- **Methyl group hydrogens:**
 - 3 equivalent hydrogens at δ 2.23 (CH₃)
 - 3 equivalent hydrogens at δ 3.79 (likely part of a methoxy group, -OCH₃)
- **Aromatic hydrogens:**
 - 2 equivalent hydrogens at δ 7.04
 - 2 equivalent hydrogens at δ 7.97

Figure 11: Human annotated C-NMR Elucidation

Fragment pool:



Target Molecule



H_NMR: 1.85-1.97 (5H, 1.91 (q, J = 7.3 Hz), 1.91 (q, J = 7.3 Hz), 1.91 (s)), 2.39-2.51 (2H, 2.45 (t, J = 7.4 Hz), 2.45 (t, J = 7.4 Hz)), 4.29 (1H, t, J = 7.3 Hz)

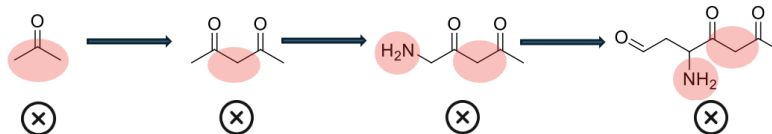


Figure 12: Complex molecule Structure Elucidation

613 C Compute Resources

614 For the execution of various models in our experiments, distinct compute resources were utilized
 615 based on the model's accessibility and computational requirements. Specifically, for models like
 616 Claude 3, GPT, and Gemini, we employed API calls to facilitate their operation, leveraging the
 617 existing infrastructure provided by their respective platforms. This approach allowed us to access
 618 these models without the need for local computational resources, thereby streamlining the process.
 619 Conversely, for all other open-sourced models employed in our study, we conducted the experiments
 620 locally using an NVIDIA A100 GPU. This high-performance computing unit was chosen due to its
 621 advanced capabilities in handling extensive computations and large model requirements efficiently.

622 Checklist

- 623 1. For all authors...
- 624 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
625 contributions and scope? [Yes]
- 626 (b) Did you describe the limitations of your work? [Yes], see Section 6
- 627 (c) Did you discuss any potential negative societal impacts of your work? [Yes], we have
628 discussed the broader impact in session 6
- 629 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
630 them? [Yes]
- 631 2. If you are including theoretical results...
- 632 (a) Did you state the full set of assumptions of all theoretical results? [No]
- 633 (b) Did you include complete proofs of all theoretical results? [N/A]
- 634 3. If you ran experiments (e.g. for benchmarks)...
- 635 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
636 mental results (either in the supplemental material or as a URL)? [Yes], the code is
637 available at <https://github.com/KehanGuo2/MolPuzzle>.
- 638 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
639 were chosen)? [Yes]
- 640 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
641 ments multiple times)? [Yes], we report the standard deviation for our result.
- 642 (d) Did you include the total amount of computing and the type of resources used (e.g.,
643 type of GPUs, internal cluster, or cloud provider)? [Yes], the total GPU usage is
644 reported in Appendix C.
- 645 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 646 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 647 (b) Did you mention the license of the assets? [Yes]
- 648 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 649 (d) Did you discuss whether and how consent was obtained from people whose data you're
650 using/curating? [Yes]
- 651 (e) Did you discuss whether the data you are using/curating contains personally identifiable
652 information or offensive content? [Yes]
- 653 5. If you used crowdsourcing or conducted research with human subjects...
- 654 (a) Did you include the full text of instructions given to participants and screenshots, if
655 applicable? [Yes], see Appendix section B.2
- 656 (b) Did you describe any potential participant risks, with links to Institutional Review
657 Board (IRB) approvals, if applicable? [Yes], see Appendix section B.2.
- 658 (c) Did you include the estimated hourly wage paid to participants and the total amount
659 spent on participant compensation? [Yes], see Appendix section B.2.