## A  MolPuzzle Benchmark Details

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

### A.1  Data Collection

The initial molecules were selected by referencing the textbook *Organic Structures from Spectra, 4th Edition*, available as an online PDF on ResearchGate. We chose 234 molecules based on spectrum tasks involving IR, MS, [1]H-NMR, and [13]C-NMR to reflect a difficulty level suitable for graduate students[36].

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

Given the challenges associated with NMR spectrum images, some spectra were obtained from simulated data in text format for [1]H-NMR and [13]C-NMR. This approach ensured clarity and accuracy in the evaluation of molecular structures.

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

### A.2  Template design

Each template was crafted to target specific skills within molecular understanding. For instance, saturation identification challenges the models' ability to discern the degree of saturation in a molecule, which is crucial for understanding its chemical reactivity and stability. Aromatic ring identification tests the models' ability to recognize benzene-like structures, which are fundamental in organic chemistry due to their common occurrence and unique properties. Saturation degree calculation pushes the models to apply quantitative analysis, requiring not just recognition but also computation based on molecular structures.
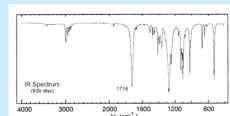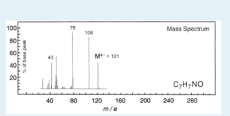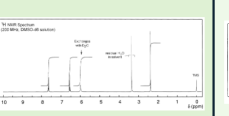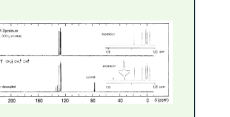
By diving deeper into the rationale behind each template and the kind of chemical knowledge they are designed to test, we can better appreciate how these tasks simulate real-world applications in chemistry. This approach not only tests the models' basic recognition abilities but also their capacity to perform complex reasoning and apply theoretical knowledge practically. The template examples are in A.3.

## 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 C8H10O potentially be Saturated?<br><br>Answer: No<br>Model response: No. |
| Aromatic Ring Identification | Question: Could the molecule with the formula C8H10O have aromatic rings?<br><br>Answer: Yes<br>Model response: Yes. |
| Functional Group Identification | Question:Could the molecule with the formula C6H14O2 potentially contain a Amine group, given the Degree of Unsaturation is 0.0?<br><br>Answer: No<br>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 C8H10O?<br><br>Answer: 4.0<br>Model response: 2 |

## A.4  Stage2 QA Samples

## 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))<br><br>Answer: 4<br>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.<br><br>Answer: 6<br>Model response: 8 |

# B  Evaluation Experiments

## B.1  Experimental Setting

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

14

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

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

## B.2  Human Evaluation

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

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

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

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

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

## B.3  Stage1

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

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

Table 3: The accuracy(↑), F1 score(↑)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 |

## B.4 Stage2

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

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

Table 4: The accuracy(↑), F1 score(↑) 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 |

## B.5 Stage-3

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

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

Despite observing moderate accuracy from the LLMs in Stage 2 of our testing, we enhanced the quality of the QA pairs in Stage 3 by providing the LLMs with verified NMR data, generated by using nmrdb[40]. This approach ensures that the data used is reliable and helps maintain the integrity of our results. The findings from Stage 2 are presented in Table. We exclusively focus on the zero-shot learning outcomes, as our observations indicate that implementing chain-of-thought prompting leads to a deterioration in model performance. To address this limitation, we offer qualitative insights in

Table 5: The F1 score(↑) 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 |

## B.6 Qualitative Results

In this section, we present several examples using GPT-4's chain-of-thought (CoT) reasoning to facilitate a clearer understanding of the results. We have enlisted two Ph.D. candidates from the chemistry department to evaluate these CoT outcomes. The analysis uses color coding to indicate the accuracy of the generated text: green signifies correct responses, red indicates incorrect ones, and yellow denotes responses that are partially correct.
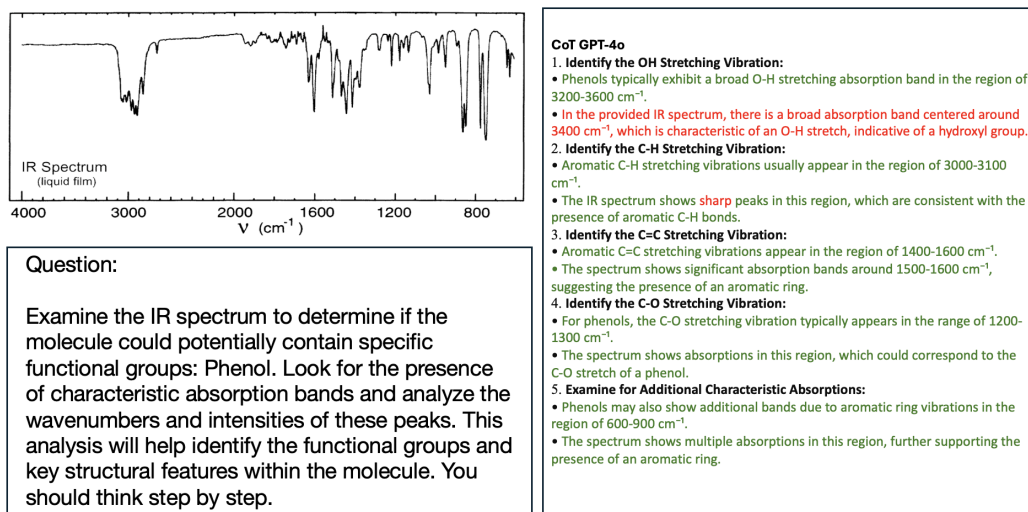
### B.6.1 Stage 2 examples
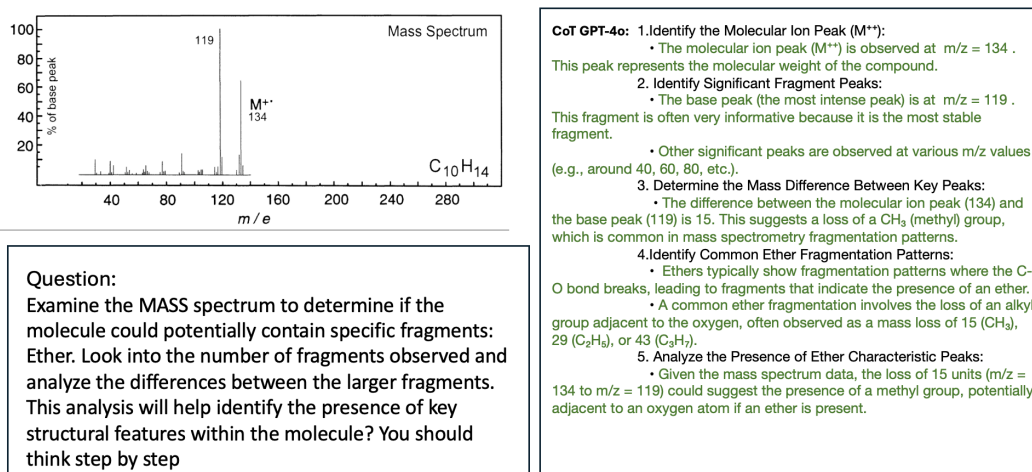


Figure 6: Human annotated IR spectrum interpretation

17

**Figure 7: Human annotated MASS spectrum interpretation**



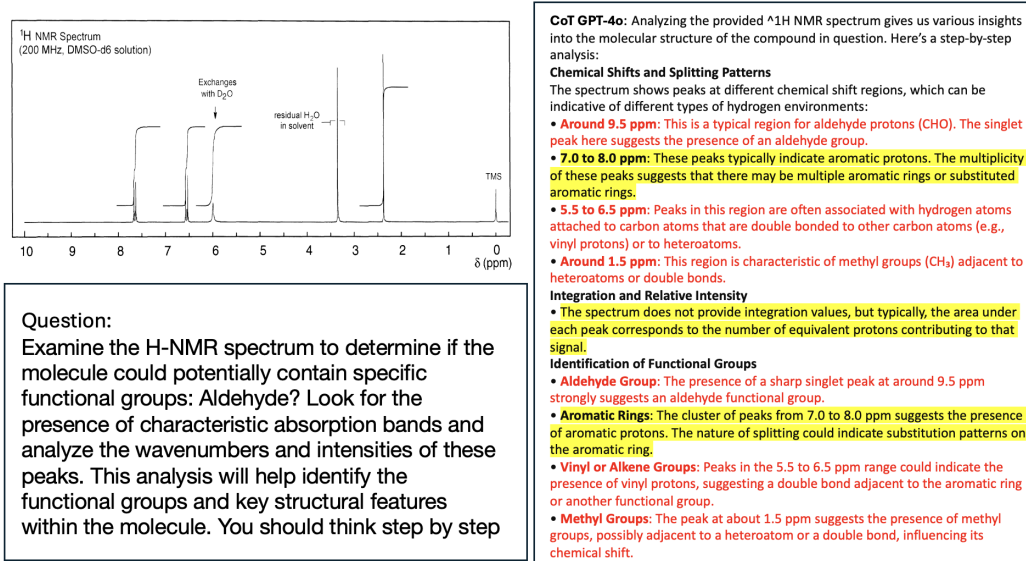**Figure 8: Human annotated H-NMR spectrum interpretation**

### B.6.2  Stage 3 examples

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

18

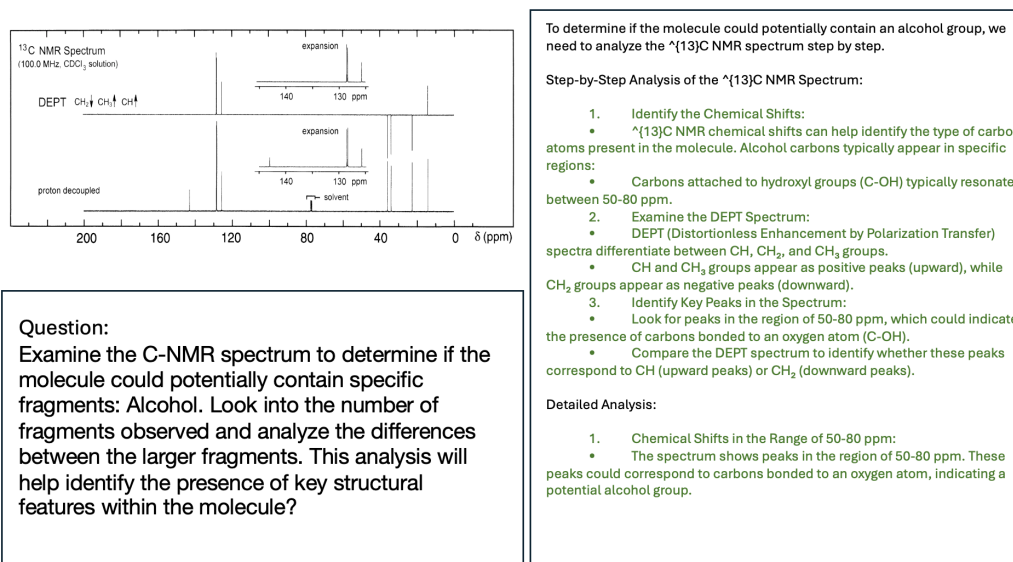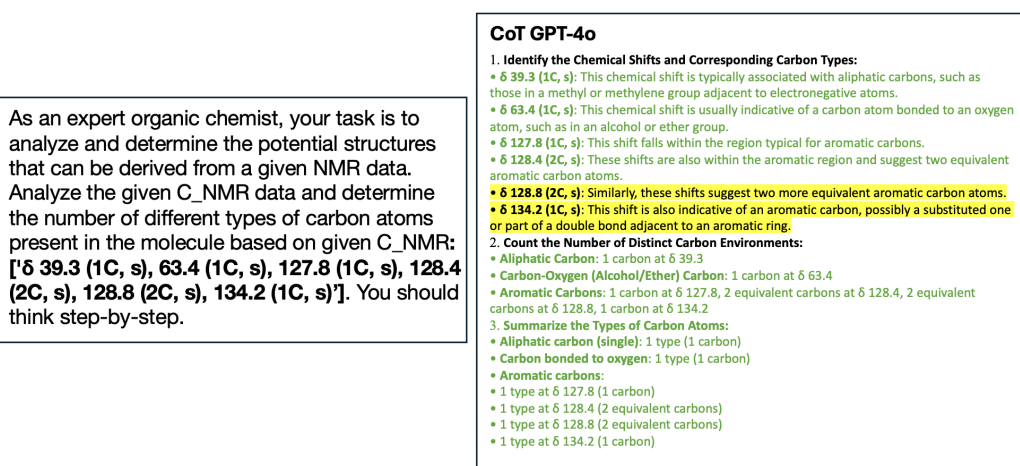Figure 9: Human annotated C-NMR spectrum interpretation



Figure 10: Human annotated H-NMR Elucidation

### B.6.3 Complex Molecules

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

**CoT GPT-4o**

1. **Identify the Chemical Shifts and Corresponding Hydrogen Types:**
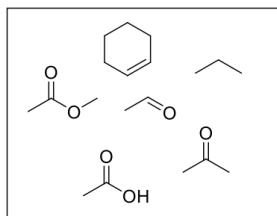• **δ 2.23 (3H, s):** This singlet indicates 3 equivalent hydrogen atoms. These are likely part of a methyl group ($CH_3$) that is not split by neighboring hydrogens.
• **δ 3.79 (3H, s):** This singlet also indicates 3 equivalent hydrogen atoms, suggesting another methyl group ($CH_3$), possibly adjacent to an electronegative atom like oxygen (e.g., methoxy group $-OCH_3$).
• **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.

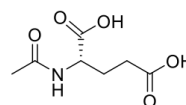2. **Determine the Number of Different Types of Hydrogen Environments:**
• **Aliphatic Hydrogens:**
• 3 equivalent hydrogens at δ 2.23 (likely $CH_3$ group)
• 3 equivalent hydrogens at δ 3.79 (likely $CH_3$ group, possibly part of a methoxy group)
• **Aromatic Hydrogens:**
• 2 equivalent hydrogens at 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 ($CH_3$)
• 3 equivalent hydrogens at δ 3.79 (likely part of a methoxy group, $-OCH_3$)
• **Aromatic hydrogens:**
• 2 equivalent hydrogens at δ 7.04
• 2 equivalent hydrogens at δ 7.97

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.
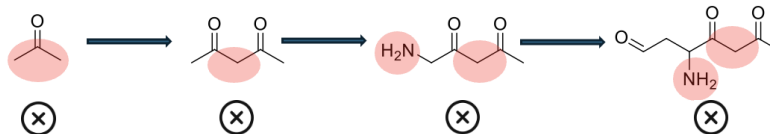
Figure 11: Human annotated C-NMR Elucidation



Figure 12: Complex molecule Structure Elucidation

## C  Compute Resources

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

20