

MolTC: Towards Molecular Relational Modeling In Language Models

Anonymous ACL submission

Abstract

Molecular Relational Learning (MRL), aiming to understand interactions between molecular pairs, plays a pivotal role in advancing biochemical research. Recently, the adoption of large language models (LLMs), known for their vast knowledge repositories and advanced logical inference capabilities, has emerged as a promising way for efficient and effective MRL. Despite their potential, these methods predominantly rely on textual data, thus not fully harnessing the wealth of structural information inherent in molecular graphs. Moreover, the absence of a unified framework exacerbates the issue of insufficient data exploitation, as it hinders the sharing of interaction mechanism learned across various datasets. To address these challenges, this work proposes a novel LLM-based multi-modal framework for **Molecular inTeration** modeling following Chain-of-Thought (CoT) theory, termed **MolTC**, which effectively integrate graphical information of two molecules in pair. For achieving a unified training paradigm, MolTC innovatively develops a *Dynamic Parameter-sharing Strategy* for cross-dataset information exchange. Moreover, to train this integrated framework efficiently, we introduce a *Multi-hierarchical CoT* theory to refine its training paradigm, and conduct a comprehensive *Molecular Interactive Instructions* dataset for the development of biochemical LLMs involving MRL. Our experiments, conducted across twelve datasets involving over 4,000,000 molecular pairs, exhibit the superiority of our method over current GNN and LLM-based baselines. Code is available at <https://anonymous.4open.science/r/MolTC-F>.

1 Introduction

Molecular Relational Learning (MRL) (Lee et al., 2023a), aiming to understand interactions between molecular *pairs*, has gained significant interest due to its wide range of applications (Roden et al., 2020). For example, Drug-Drug Interactions

(DDIs) are critical in pharmacology and drug development (Lin et al., 2020), while solute-solvent interactions (SSIs) are fundamental in solution chemistry and the design of chemical processes (Varghese and Mushrif, 2019; Chung et al., 2022). However, the exhaustive experimental validation of these interactions is notoriously time-consuming and costly. In response, adopting large language models (LLMs) (Brown et al., 2020; Taylor et al., 2022), known for their vast knowledge repositories and advanced logical inference capabilities, has emerged as an efficient and effective alternative for MRL (Park et al., 2022; Jha et al., 2022a).

Despite their promise, a primary concern of current LLM-based paradigm is the *insufficient data exploitation*. Specifically, they predominantly rely on the textual data such as SMILES (Simplified Molecular Input Line Entry System) and property descriptions, thus not fully harnessing the wealth of structural information inherent in molecular graphs (Sagawa and Kojima, 2023), as indicated in Figure 1 (a). Current studies have indicated that it is challenging for LLMs to fully understand the complex graphs based solely on textual data, hence, it’s crucial to explicitly model these structures given their significance in MRL (Park et al., 2022).

Compounding this concern is the absence of a unified framework for LLM-based MRL (Livne et al., 2023; Pei et al., 2023). Concretely, this absence impedes the sharing and integration of interaction mechanisms learned across various datasets, leading to a fragmentation in collective insights. Especially, it poses a catastrophic challenge for tasks with a limited number of labeled pairs (Chung et al., 2022), where LLMs often struggle with due to the high risk of overfitting, as illustrated in Figure 1 (b). Worse still, such limited datasets are prevalent in MRL since the experimental acquisition is often constrained by high costs (Lee et al., 2023a).

To overcome these limitations, in this work, we propose **MolTC**, a unified multi-modal frame-

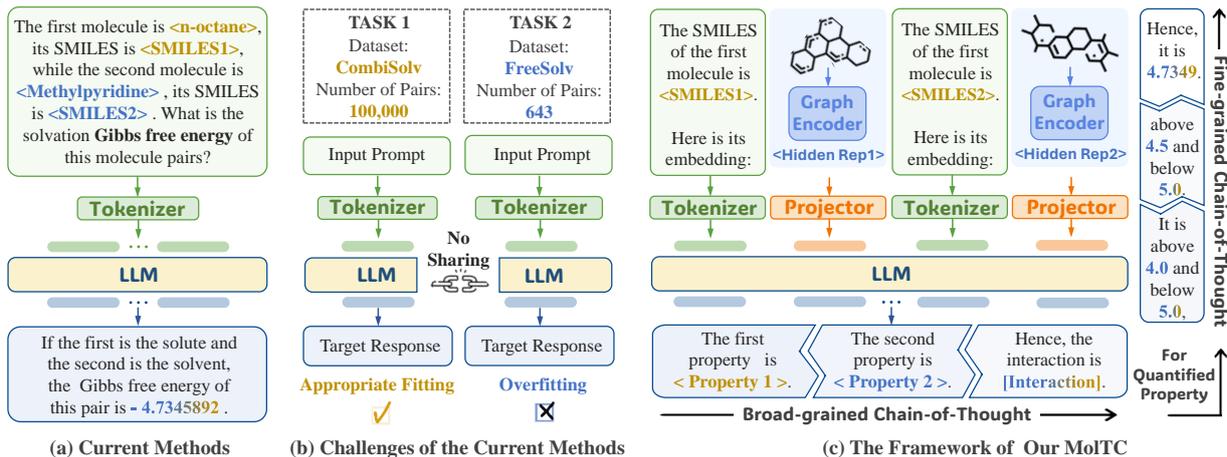


Figure 1: Comparison between the current methods leveraging LLMs to model molecule interactions and our MolTC. (a) The prevailing paradigm of current methods. (b) The challenge of applying the current paradigm to the tasks involving datasets with a small number of samples. (c) The framework of our proposed MolTC, which is enhanced by the principle of CoT. Best viewed in color.

work for **Molecular inTeration** modeling following the **Chain-of-thought** theory (Wei et al., 2022). As depicted in Figure 1 (c), MolTC employs the Graph Neural Networks (GNNs) (Kipf and Welling, 2017), known for their proficiency in graph modeling, to explicitly gather graphical information of molecular pairs, and integrates them into the input space of LLMs by two meticulously crafted projectors. In response to empirical findings that LLMs may confuse two input molecules in pair, MolTC incorporates the molecules’ SMILES information to reinforce the concept of molecular order. More importantly, to achieve a unified learning paradigm, MolTC develops a *Dynamic Parameter-sharing* strategy for bolstering cross-dataset information exchange, which can boost the efficiency and effectiveness simultaneously.

Based on these, a two-pronged approach is developed to train this integrated framework efficiently:

(1) Training Paradigm Refinement: As shown in Figure 1 (c), we introduce a *Multi-hierarchical CoT* theory to guide the training paradigm of MolTC. Concretely, the broad-grained CoT guides the pre-training stage to identify individual molecular properties before predicting interactions, ensuring an acute awareness of each molecule’s unique attribute. For quantitative interaction tasks, which are challenging for LLMs, a fine-grained CoT enables the fine-tuning stage to initially predict a range, and then progressively refining it to a precise value.

(2) Dataset Foundation Construction: In sight of the absence of a comprehensive MRL datasets

for biochemical LLMs, we construct a **Molecular inTera**ctive instructions dataset, termed **MoT-instruction**. Specifically, we first conduct twelve well-established MRL datasets across various domain, and source their detailed molecular properties from authoritative biochemical databases. Based on this, we meticulously compile these properties and empirically determine their optimal instructions. These process ensures that MoT-instructions can not only enhance the performance of our MolTC, but also contribute to the development of other biochemical LLMs involving MRL.

Our contributions can be summarized as follows:

- We identify the issue of insufficient data exploitation in current LLM-based MRL, and take the first attempt to develop a unified multi-modal framework for LLM-based MRL, named MolTC.
- We introduce the multi-hierarchical CoT theory to enhance the MolTC’s training process, especially for quantitative interaction tasks.
- We construct MoT-instructions, the first comprehensive instruction dataset in MRL domain, to enhance the development of biochemical LLMs involving MRL.
- Our experiments, across over 4,000,000 molecular pairs in various domains such as DDI and SSI, demonstrate the superiority of our method over current GNN and LLM-based baselines.

2 Methodology

In this section, we detail our MolTC, which harnesses the power of LLMs for comprehending

molecular interactions. We begin with the introduction of model framework in Section 2.1. Taking a step further, the training paradigm guided by the principle of Multi-hierarchical CoT is outlined in Section 2.2. Moreover, the dynamic parameter sharing strategy tailored for MolTC and our developed datasets, MoT-instructions, are elaborated in Section 2.3 and 2.4, respectively.

2.1 Framework of MolTC

Here we introduce four key components of MolTC’s framework: Graph Encoder, Representation Projector, SIMLES Injector, and the backbone LLM. The specific instantiation details of each module can be found in the experimental section and the appendix.

Graph Encoder. The first step of extracting interactions is to precisely encode the molecular graphs. In sight of this, we utilize two GNN-based encoders to capture the embedding of the given molecular pairs, leveraging the GNN’s robust capability in aggregating structural information. More formally, let $\mathcal{G}_a = \{\mathcal{V}_a, \mathcal{E}_a\}$ and $\mathcal{G}_b = \{\mathcal{V}_b, \mathcal{E}_b\}$ denote the input pair, where \mathcal{V}, \mathcal{E} represent atomic nodes and the chemical bonds, respectively. The two graph encoders $f_{\text{enc}1}$ and $f_{\text{enc}2}$ perform aggregating to obtain the atomic embedding:

$$\begin{aligned} \mathbf{H}_a &= [h_a^1, h_a^2, \dots, h_a^{|\mathcal{V}_a|}] = f_{\text{enc}1}(\mathcal{G}_a), \\ \mathbf{H}_b &= [h_b^1, h_b^2, \dots, h_b^{|\mathcal{V}_b|}] = f_{\text{enc}2}(\mathcal{G}_b), \end{aligned} \quad (1)$$

where h_a^i and h_b^i denote to the embedding of the i -th atom in molecule \mathcal{G}_a and \mathcal{G}_b ; \mathcal{V}_a and \mathcal{V}_b represent the number of nodes.

Representation Projector. After acquiring molecular pair representations \mathbf{H}_a and \mathbf{H}_b , the next step is to map them into the backbone LLM’s hidden space using Projectors $f_{\text{pro}1}$ and $f_{\text{pro}2}$. These projectors serve as pivotal connectors, translating \mathbf{H}_a and \mathbf{H}_b into LLM-comprehensible encodings \mathbf{M}_a and \mathbf{M}_b . Drawing inspiration from the state-of-the-art vision-language models, we instantiate $f_{\text{pro}1}$ and $f_{\text{pro}2}$ by Querying Transformers (Q-Formers) (Li et al., 2023a; Dai et al.). More formally,

$$\begin{aligned} \mathbf{M}_a &= [m_a^1, m_a^2, \dots, m_a^q] = f_{\text{pro}1}(\mathbf{H}_a), \\ \mathbf{M}_b &= [m_b^1, m_b^2, \dots, m_b^q] = f_{\text{pro}2}(\mathbf{H}_b), \end{aligned} \quad (2)$$

where q denotes the number of learnable query tokens of Q-Former’s transformer.

In detail, our Projectors, based on the BERT architecture, incorporate an additional cross-attention

module positioned between the self-attention and feed-forward modules. This instantiation offers two key benefits. Firstly, it supports seamless integration with conventional BERT-based text encoders, allowing $f_{\text{pro}1}$ and $f_{\text{pro}2}$ pre-training with extensive molecular graph-text pairs. Secondly, it maintains compatibility with various input dimensions d , and allows adjustments in the size of learnable query tokens to align with the LLM’s token embedding size. These advantages lay a solid foundation for the thorough interaction of two molecules during the LLM’s inference process. Future work will also explore more projector designs, such as streamlining it through specially tailored MLPs (Yang et al., 2023).

SMILES Tokenization. When directly analyzing the representations \mathbf{M}_a and \mathbf{M}_b with LLMs, our experiments suggest a potential confusion by LLMs in distinguishing the properties of each molecule in a pair. This observation naturally inspires us to integrate textual information of the molecules to strengthen the concept of their sequential order. Here MolTC employs SMILES due to its ubiquity and specificity. Additionally, SMILES serves as a conduit, linking the task-specific prompts with the corresponding biochemical knowledge stored within the LLM. Therefore, we directly input the SMILES of both molecules into the backbone LLM, utilizing the inherent encoder to acquire their tokens \mathbf{S}_a and \mathbf{S}_b .

Backbone LLM. MolTC leverages Galactica, a decoder-only transformer built on the OPT framework, as its backbone LLM. Pretrained on an extensive collection of scientific literature, Galactica demonstrates exceptional proficiency in biochemistry knowledge. This expertise, particularly in parsing molecular sequences such as SMILES and SELFIES strings, enables Galactica to adeptly capture the properties crucial for molecular interactions. Specifically, the goal of MolTC is to harness Galactica’s advanced inferential skills to interpret the contextual interactions between two molecular sets of token collections, $\{\mathbf{M}_a, \mathbf{S}_a\}$ and $\{\mathbf{M}_b, \mathbf{S}_b\}$. More formally, we denote an integrated prompt sequence as follows:

$$\begin{aligned} \mathbf{X} &= \{\mathbf{P}, \mathbf{M}_a, \mathbf{S}_a, \mathbf{M}_b, \mathbf{S}_b\} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_l] \\ &\text{s.t. } \mathbf{P} \sim \mathcal{P}_{\mathbf{r}}, \end{aligned} \quad (3)$$

where l is the integrated input length, \mathbf{P} denotes the task-specific prompt, and $\mathcal{P}_{\mathbf{r}}$ represents a collection of various manually designed prompts, each

tailored for the molecular interaction task r . The generation process adopts a causal mask to generate a response encapsulating key interactive properties with length T :

$$\hat{\mathbf{X}} = [\hat{x}_1, \hat{x}_2, \dots, \hat{x}_T]. \quad (4)$$

Utilizing Galactica’s autoregressive framework, the training objective involves regressing the target response based on the input prompt \mathbf{X} . Specifically, the output for i -th token \hat{x}_i , is computed based on its preceding tokens as follows for $t \in (1, T)$:

$$p(\hat{\mathbf{X}}_{[1:t]}|\mathbf{X}) = \prod_{i=1}^t p(\hat{x}_i|\mathbf{X}, \hat{\mathbf{X}}_{[1:i-1]}). \quad (5)$$

2.2 Training Paradigm of MolTC

In this part, we elaborate the training paradigm of MolTC, including pretraining and fine-tuning processes, which is guided by the principle of Multi-hierarchical CoT, as shown in Figure 2.

2.2.1 Broad-grained CoT Guided Pretraining

Given the challenge of directly understanding complex interactions between two input molecules in pair, the broad-grained CoT guides MolTC to initially identify individual molecular properties. By thoroughly understanding each molecule’s characteristics, MolTC establishes a solid foundation for accurately predicting their interactions. Specifically, in the pretraining stage, the prompt is uniformly designed as follows:

Prompt for Pretraining Stage	
Input Prompt	<SMILES1>, <GraEmb1>, the front is the first molecule, followed by the second molecule: <SMILES2>. <GraEmb2>. Please provide the biochemical properties of the two molecules one by one.
Target Response	The properties of the first molecule are [Property1], and the properties of the second molecule are [Property2].

This prompt design enable MolTC to delineate key properties of two molecules sequentially. Based on it, MolTC utilize the generation loss of the backbone LLM to train Graph Encoders, f_{enc1} and f_{enc2} , as well as the Representation Projectors, f_{pro1} and f_{pro2} . Notably, during this phase, the backbone LLM remains frozen.

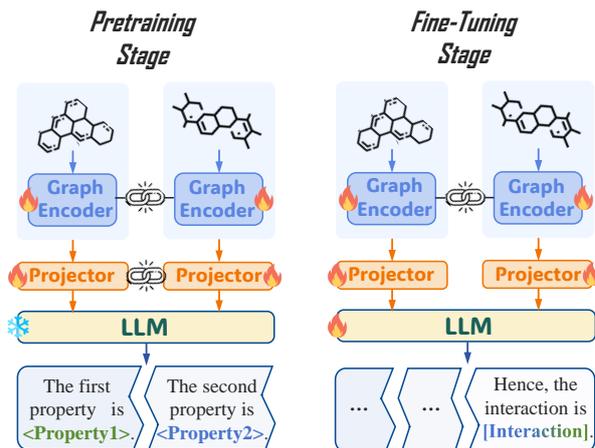


Figure 2: The training process of our MolTC. The flame symbol denotes the parameter update, the snowflake symbol indicates the parameter freezing, and the chain symbol depicts the parameter sharing between two modules. Best viewed in color.

Dataset Construction for Pretraining. To ensure backbone LLM can understand the individual characteristics of each molecule, it is pivotal to prepare a comprehensive dataset comprising molecule pairs and their corresponding biochemical properties. To this end, (1) we first conduct an extensive survey of various authoritative biochemical database such as PubChem¹ and Drugbank (Kim et al., 2023), and collect a large amount of molecule-textual properties pairs; (2) then, recognizing the variability in annotation quality within this dataset, we augment and enrich molecular descriptions that were less extensively annotated; (3) subsequently, to simulate diverse molecular interactions, we generated molecular pairs by randomly grouping two distinct molecules from the above database. This random pairing facilitates a broad spectrum of molecular combinations, exposing the pretraining stage to diverse interaction scenarios, thus naturally enhancing the generalizability of our MolTC.

2.2.2 Fine-grained CoT Guided Fine-tuning

During the fine-tuning phase, MolTC is trained to enable the backbone LLM to generate interaction properties based on the properties of individual molecules it initially identifies. To this end, prompts in the fine-tuning stage should be crafted for specific downstream task. For example, in DDI tasks, we construct the following prompt:

¹<https://pubchem.ncbi.nlm.nih.gov>

Prompt for DDI Tasks (Fine-tuning)	
Input Prompt	<SMILES1>, <GraEmb1>, the front is the first molecule, followed by the second molecule: <SMILES2>. <GraEmb2>. What are the side effects of these two drugs?
Target Response	The property of the first molecule is [Property1], while the property of the second molecule is [Property2]. Hence, the first drug molecule may increase the photosensitizing activities of the second drug molecule.

Despite the effectiveness of this prompt design, LLMs face notable challenges in quantitative analysis, especially in complex molecular interaction contexts such as SSI and chromophore-solvent interaction (CSI). Our experiments in Section 3 highlight this difficulty, demonstrating that LLMs tend to exhibit indecision regarding the quantitative values in their outputs. To address this, a fine-grained CoT concept is introduced to refine the training paradigm. Specifically, the backbone LLM is guided to initially suggest a range for the target numerical value, then progressively refining it to a precise value. Take a meticulously prompt for SSI tasks as an example:

Prompt for SSI Tasks (Fine-tuning)	
Input Prompt	<SMILES1>, <GraEmb1>, the front is the first molecule, followed by the second molecule: <SMILES2>. <GraEmb2>. What is the solvation Gibbs free energy of this pair of molecules?
Target Response	The property of the first molecule is [Property1], while the property of the second molecule is [Property2]. Hence, the solvation Gibbs free energy of these two molecules is above 3.0 and below 3.5, so the accurate value is 3.24791.

This step-wise refinement process fosters a more accurate and reliable resolution of numerically-intensive challenges. Based on these prompts, in the fine-tuning stage, the parameters in backbone

LLM are updated through Low-Rank Adaptation (LoRA) (Hu et al., 2021) strategy, known for its efficiency in tailoring the LLM to the requirements of downstream tasks and minimal memory demands in storing gradients. Meanwhile, to ensure that other modules are optimally adjusted to suit the specifics of the downstream tasks, Graph Encoders f_{enc1} and f_{enc2} , as well as Representation Projectors f_{pro1} and f_{pro2} are trained following the generation loss of the backbone LLM.

2.3 Dynamic Parameter Sharing Strategy

To implement the above training paradigm effectively, we introduce a novel parameter-sharing strategy, inspired by key biochemical insights:

(1) The Importance of **Role-Playing**: A molecule’s role in an interaction crucially influences the outcome. For example, in SSI scenario like the water-ethanol pair, utilizing water and ethanol as solvents, respectively, yields different energy releases (Reichardt, 2021). Sometimes, a reversal of roles can even result in the absence of interaction.

(2) The Importance of **Input Order**: In certain molecular pairs, the sequence of introducing molecules significantly impacts the interactions. For instance, the order of drug introduction can lead to varying therapeutic effects.

(3) The Importance of **Role and Order-Specific Feature Extraction**: The role and input order of molecules determine the relevance of their structural features. For example, a chemical group in a solute-solvent pair may be crucial for the release of Gibbs free energy when in the solute, but less so in the solvent (Reichardt, 2021; J et al., 2022).

These insights inspire MolTC to adaptively prioritize distinct key information, creating unique tokens for the same molecule based on its role and order. To enable this nuanced learning while also capitalizing on the shared aspects of molecular learning, we introduce the following parameter-sharing strategy, as shown in Figure 2:

(1) The GNN-based **Encoders** f_{enc1} and f_{enc2} , which focus on extracting molecular graph structures, share parameters during both pretraining and fine-tuning stages to enhance learning efficiency.

(2) The Qformer-based **Projectors** f_{pro1} and f_{pro2} , tasked with aligning molecular structures to semantic information, share parameters during pretraining stage to promote generalization and robustness. However, in the fine-tuning stage, we cease sharing

373 to allow customized semantic mappings tailored to
374 the varying roles and orders.

375 In summary, this strategy is tailored to balance
376 the need for role and order-based distinctively learn-
377 ing with the efficiency gained from commonalities
378 across molecular pairs.

379 2.4 Construction of MoT-instructions

380 Given the absence of a comprehensive instruction
381 datasets tailored for LLM-based MRL, we aim to
382 develop a molecular interactive instructions dataset,
383 termed MoT-instructions. This dataset is designed
384 to fulfill several key criteria: (1) it should include
385 extensive molecular pairs capable of interaction,
386 covering a broad spectrum of domains, (2) it should
387 detail important biochemical properties of each
388 molecule within these pairs, and (3) it should elab-
389 orate the resultant properties from molecular in-
390 teractions. Specifically, MoT-instructions are con-
391 structed through a three-step process as follows.

392 (1) We begin by aggregating twelve representa-
393 tive molecular interaction datasets across various
394 widely recognized biochemical tasks, such as DDI,
395 SSI, and CSI. Following this, we engage in a sys-
396 tematic search for textual descriptions of the bio-
397 chemical properties of each molecule involved in
398 these interactions. Specifically, we source this in-
399 formation from authoritative biochemical databases
400 such as DrugBank and PubChem.

401 (2) The next critical step is the **experimental de-**
402 **termination of the optimal instructions**. Specif-
403 ically, for all molecular pairs in step (1), we first
404 deconstruct the lengthy molecular properties into
405 a series of questions and answers, a format more
406 comprehensible to LLMs (Taylor et al., 2022). The
407 granularity of this deconstruction is decided based
408 on the performance of our MolTC. For more chal-
409 lenging quantitative tasks, instructions guided by
410 fine-grained CoT are required to provide a numeri-
411 cal range before specifying a concrete value. Given
412 the vast number of possible correct ranges, exhaus-
413 tive testing is impractical. Therefore, we initially
414 determine the optimal range for a small subset of
415 datasets using a grid search, guided by the pre-
416 dictive performance of MolTC. Subsequently, we
417 derive statistics, such as mean and standard devia-
418 tion, from these datasets to establish a relationship
419 between statistics and optimal ranges. Finally, for
420 other datasets, we determine their optimal range
421 based on this established rule.

422 (3) The final step in our dataset construction in-

423 volved filtering out pairs that lacked sufficient infor-
424 mation on molecule properties or interaction data.
425 Specifically, partial properties of a molecular pair
426 are often missing in some datasets. To maximize
427 the utilization of information from these datasets,
428 we consider extracting each property within them
429 as a separate dataset. This approach allows us to
430 naturally omit missing values without wasting other
431 information present in the molecular pair.

432 3 Experiment

433 In this section, we aim to answer the following
434 research questions:

- 435 • **RQ1:** Is MolTC capable of generating the inter-
436 active property, involving the *qualitative* knowl-
437 edge, of the given molecular pair?
- 438 • **RQ2:** Does MolTC have the ability to generate
439 the interactive property, involving the *quantita-*
440 *tive* property, for a given molecular pair?
- 441 • **RQ3:** What is the impact of the proposed strate-
442 gies, such as the CoT enhancement strategy and
443 SMILES injection strategy, on the inference pro-
444 cess of our MolTC?

445 3.1 Experimental Setting

446 We evaluate MolTC on twelve well-established
447 downstream molecule interaction tasks involving
448 qualitative and quantitative analysis. Here we pro-
449 vide a brief overview of our experimental setup.
450 Detailed descriptions are presented in the appendix.

451 **Datasets.** We employ 12 datasets across various
452 domains such as DDI, SSI, and CSI. Specifically,
453 we collect *Drugbank* (Version 5.0.3), *ZhangDDI*
454 (Zhang et al., 2017), *ChChMiner* (Zitnik et al.,
455 2018), *DeepDDI* (Ryu et al., 2018), *TWOSIDES*
456 (Tatonetti et al., 2012), *Chromophore* (Joung et al.,
457 2020), *MNSol* (Marenich et al., 2020), *CompSol*
458 (Moine et al., 2017), *Abraham* (Grubbs et al., 2010),
459 *CombiSolv* (Vermeire and Green, 2021), *FreeSolv*
460 (Mobley and Guthrie, 2014) and *CombiSolv-QM*
461 (Vermeire and Green, 2021).

462 **Baselines.** For a comprehensive evaluation, we
463 conduct various baseline methods encompassing
464 distinct categories such as methods based on:
465 GNNs, DL models other than GNN, and LLMs.
466 Specifically, For DDI task, we employ *GoGNN*
467 (Wang et al., 2020), *MHCADDI* (Deac et al., 2019),
468 *DeepDDI* (Ryu et al., 2018), *SSI-DDI*, *CGIB* (Lee
469 et al., 2023a), *CMRL* (Lee et al., 2023b), *MDF-SA-*
470 *DDI* (Lin et al., 2022), *DSN-DDI* (Li et al., 2023c)

Table 1: Comparative performance of various methods in qualitative interactive tasks. The best-performing methods are highlighted with a gray background, while the second-best methods are underscored for emphasis.

Baseline Model		Drugbank		ZhangDDI		ChChMiner		DeepDDI	
		Accuracy	AUC-ROC	Accuracy	AUC-ROC	Accuracy	AUC-ROC	Accuracy	AUC-ROC
GNN Based	GoGNN	84.78 \pm 0.57	91.63 \pm 0.66	84.10 \pm 0.46	92.35 \pm 0.48	91.17 \pm 0.46	96.64 \pm 0.40	93.54 \pm 0.35	92.71 \pm 0.27
	SSI-DDI	94.12 \pm 0.33	98.38 \pm 0.31	86.97 \pm 0.37	93.76 \pm 0.34	93.26 \pm 0.31	97.81 \pm 0.22	95.27 \pm 0.25	98.42 \pm 0.31
	DSN-DDI	<u>94.93</u> \pm 0.14	<u>99.01</u> \pm 0.12	87.65 \pm 0.13	94.63 \pm 0.18	84.30 \pm 0.17	94.25 \pm 0.26	95.64 \pm 0.18	98.01 \pm 0.16
	CMRL	94.83 \pm 0.12	98.76 \pm 0.10	<u>87.78</u> \pm 0.36	<u>94.68</u> \pm 0.23	94.23 \pm 0.26	98.37 \pm 0.12	<u>96.37</u> \pm 0.34	<u>98.98</u> \pm 0.31
	CGIB	94.68 \pm 0.34	98.60 \pm 0.25	87.32 \pm 0.71	94.18 \pm 0.60	<u>94.25</u> \pm 0.39	<u>98.45</u> \pm 0.31	96.23 \pm 0.52	98.45 \pm 0.64
ML Based	DeepDDI	93.15 \pm 0.25	98.06 \pm 0.54	83.35 \pm 0.49	91.13 \pm 0.58	90.34 \pm 0.62	95.73 \pm 0.37	92.39 \pm 0.38	98.11 \pm 0.42
	MHCADDI	78.50 \pm 0.80	86.33 \pm 0.35	77.86 \pm 0.59	86.94 \pm 0.68	84.26 \pm 0.54	89.33 \pm 0.82	87.01 \pm 0.77	88.64 \pm 0.83
	MDF-SA-DDI	93.86 \pm 0.31	97.65 \pm 0.29	86.89 \pm 0.25	94.03 \pm 0.23	93.64 \pm 0.20	98.10 \pm 0.19	95.12 \pm 0.30	97.84 \pm 0.36
LLM Based	Galactica	79.16 \pm 0.35	86.23 \pm 0.33	67.20 \pm 0.46	78.74 \pm 0.58	74.61 \pm 0.44	83.51 \pm 0.63	71.50 \pm 0.41	79.07 \pm 0.41
	Chem T5	85.83 \pm 0.31	91.97 \pm 0.38	72.34 \pm 0.42	89.31 \pm 0.30	80.79 \pm 0.52	85.65 \pm 0.46	75.58 \pm 0.66	84.42 \pm 0.43
	MolCA	87.95 \pm 0.52	94.00 \pm 0.37	68.21 \pm 0.59	88.53 \pm 0.62	90.15 \pm 0.43	92.92 \pm 0.60	82.95 \pm 0.58	88.52 \pm 0.77
	MolT5	89.49 \pm 0.47	93.08 \pm 0.26	76.46 \pm 0.30	89.06 \pm 0.33	84.70 \pm 0.25	91.18 \pm 0.32	86.82 \pm 0.46	90.08 \pm 0.57
MolTC (Ours)		<u>95.98</u> \pm 0.15	<u>99.12</u> \pm 0.31	89.40 \pm 0.12	95.48 \pm 0.18	95.59 \pm 0.20	98.66 \pm 0.09	96.70 \pm 0.26	99.05 \pm 0.32

Table 2: Comparative performance of various methods in quantitative interactive tasks. The best-performing methods are highlighted with a gray background, while the second-best methods are underscored for emphasis.

Baseline Model		FreeSolv		Abraham		CompSol		CombiSolv	
		MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE
GNN Based	CIGIN	0.589 \pm 0.053	0.931 \pm 0.066	0.314 \pm 0.004	0.607 \pm 0.011	0.197 \pm 0.003	0.349 \pm 0.005	0.288 \pm 0.005	0.664 \pm 0.012
	D-MPNN	0.702 \pm 0.014	1.231 \pm 0.029	0.484 \pm 0.012	0.705 \pm 0.025	0.205 \pm 0.006	0.373 \pm 0.007	0.482 \pm 0.013	0.895 \pm 0.055
	GEM	0.598 \pm 0.018	1.188 \pm 0.049	<u>0.254</u> \pm 0.004	0.531 \pm 0.005	0.203 \pm 0.006	0.337 \pm 0.007	0.290 \pm 0.009	0.783 \pm 0.020
	CGIB	<u>0.541</u> \pm 0.009	<u>0.917</u> \pm 0.055	0.258 \pm 0.008	<u>0.530</u> \pm 0.009	<u>0.178</u> \pm 0.004	<u>0.301</u> \pm 0.003	<u>0.230</u> \pm 0.004	<u>0.394</u> \pm 0.009
ML Based	GOVER	0.636 \pm 0.026	1.074 \pm 0.049	0.347 \pm 0.005	0.625 \pm 0.016	0.184 \pm 0.005	0.371 \pm 0.014	0.412 \pm 0.016	0.728 \pm 0.034
	SolvBert	0.602 \pm 0.029	1.034 \pm 0.044	0.496 \pm 0.007	0.693 \pm 0.014	0.192 \pm 0.008	0.353 \pm 0.008	0.418 \pm 0.018	0.711 \pm 0.020
	Uni-Mol	0.575 \pm 0.060	1.012 \pm 0.070	0.355 \pm 0.007	0.602 \pm 0.024	0.198 \pm 0.002	0.344 \pm 0.003	0.267 \pm 0.005	0.669 \pm 0.017
	SMD	0.599 \pm 0.037	1.202 \pm 0.036	0.400 \pm 0.022	0.646 \pm 0.037	0.199 \pm 0.006	0.348 \pm 0.007	0.657 \pm 0.011	1.023 \pm 0.029
LLM Based	Galactica	0.882 \pm 0.010	1.438 \pm 0.066	0.645 \pm 0.008	1.064 \pm 0.016	0.594 \pm 0.006	0.854 \pm 0.008	0.831 \pm 0.018	1.486 \pm 0.035
	Chem T5	0.802 \pm 0.036	1.377 \pm 0.057	0.629 \pm 0.010	0.910 \pm 0.017	0.445 \pm 0.008	0.734 \pm 0.010	0.882 \pm 0.015	1.297 \pm 0.024
	MolCA	0.760 \pm 0.033	1.271 \pm 0.039	0.581 \pm 0.007	0.897 \pm 0.008	0.467 \pm 0.006	0.716 \pm 0.022	0.648 \pm 0.033	1.125 \pm 0.035
	MolT5	0.705 \pm 0.047	1.135 \pm 0.069	0.549 \pm 0.008	0.832 \pm 0.006	0.476 \pm 0.003	0.695 \pm 0.013	0.652 \pm 0.023	1.124 \pm 0.027
MolTC (Ours)		<u>0.502</u> \pm 0.011	<u>0.684</u> \pm 0.042	<u>0.194</u> \pm 0.009	<u>0.388</u> \pm 0.010	<u>0.171</u> \pm 0.006	<u>0.295</u> \pm 0.004	<u>0.172</u> \pm 0.004	<u>0.465</u> \pm 0.008

as the backbone. For SSI and CSI tasks, we utilize *D-MPNN* (Vermeire and Green, 2021), *SolvBert* (Yu et al., 2023), *SMD* (Meng et al., 2023), *CGIB* (Lee et al., 2023a), *CIGIN* (Pathak et al., 2020), *GEM* (Fang et al., 2022), *GOVER* (Rong et al., 2020), *Uni-Mol* (Zhou et al., 2023) as the backbone. Furthermore, all downstream tasks adopt LLM-based methods, such as Galactica (Taylor et al., 2022), Chem T5 (Christofidellis et al., 2023), MolT5 (Edwards et al., 2022) and MolCA (Liu et al., 2023) as the backbone.

Metrics. For qualitative tasks, we employ prediction *Accuracy* and *AUC-ROC* (Area Under the Receiver Operating Characteristic curve) as comparative metrics, while for quantitative tasks, *MAE* (Mean Absolute Error) and *RMSE* (Root Mean Square Error) are utilized as the standards.

3.2 Qualitative Prediction Performance (RQ1)

Table 1 exhibits the performance in qualitative interactive tasks. Due to page width limitations, only a subset of the results is presented, with additional results detailed in the appendix. From Table 1, we deduce the following observations:

Obs.1: MolTC consistently outshines its counterparts in qualitative interaction predictions, While GNN-based methods demonstrate commendable performance, maintaining over 90% accuracy across numerous datasets, MolTC transcends these figures in every evaluated scenario. For instance, it marks a notable 1.05% improvement in accuracy on the drugback dataset, a feat attributable to the synergy between the LLMs’ reasoning faculties and the GNNs’ proficiency in graph modeling.

Table 3: Performance comparison of various models on different datasets.

Dataset	Metric	w/o SMILES	w/o CoT	
			Broad	Fine
DDI	Accuracy Rate (\downarrow)	6.42 \pm 0.13 7.08 %	2.01 \pm 0.05 2.13 %	—
	ACC-AUC Rate (\downarrow)	7.87 \pm 0.32 8.22 %	2.98 \pm 0.08 3.10 %	—
SSI	MAE Rate (\uparrow)	0.025 \pm 0.004 11.32 %	0.010 \pm 0.002 4.56 %	0.036 \pm 0.007 16.40 %
	RMSE Rate (\uparrow)	0.045 \pm 0.007 9.47 %	0.014 \pm 0.003 2.95 %	0.054 \pm 0.009 11.37 %
CSI Abs.	MAE Rate (\uparrow)	2.06 \pm 0.11 15.03 %	0.51 \pm 0.03 3.72 %	2.65 \pm 0.16 19.34 %
	RMSE Rate (\uparrow)	3.37 \pm 0.20 15.18 %	1.18 \pm 0.12 5.31 %	4.84 \pm 0.29 21.80 %
CSI Emis.	MAE Rate (\uparrow)	3.10 \pm 0.17 16.23 %	0.85 \pm 0.04 5.23 %	4.42 \pm 0.36 23.14 %
	RMSE Rate (\uparrow)	4.99 \pm 0.28 18.34 %	1.47 \pm 0.12 5.40 %	7.29 \pm 0.44 26.80 %
CSI Life.	MAE Rate (\uparrow)	0.085 \pm 0.003 13.70 %	0.026 \pm 0.002 4.19 %	0.072 \pm 0.004 11.61 %
	RMSE Rate (\uparrow)	0.101 \pm 0.010 12.16 %	0.034 \pm 0.008 4.09 %	0.093 \pm 0.010 11.20 %

Obs.2: The variability of MolTC’s outcomes, as indicated by the standard deviation, is consistently minimal in comparison to other models. On average, the standard deviation for MolTC is 35.41% lower than GNN-based models and 46.86% lower than LLM-based models. The precision in MolTC’s performance is largely attributed to the training paradigm enhanced by the multi-hierarchical CoT, which ensures a meticulous and accurate inference process.

3.3 Quantitative Prediction Performance (RQ2)

Table 2 shows the performance in a subset of quantitative tasks, with an exhaustive set of results detailed in the appendix. The datasets offer four-dimensional molecular information, comprising atom type, chirality tag, bond type, and bond direction. Key observations from Table 2 include:

Obs.3: MolTC continues to lead in quantitative analysis tasks, an area typically challenging for LLMs. Despite the strong baseline set by CGIB, characterized by low MAE and RMSE across datasets, MolTC outperforms it in every metric. For instance, it achieves a 23.98% reduction in RMSE on the CombiSolv dataset relative to CGIB. This underscores the advantage of adeptly leveraging the interaction between SMILE representations

and molecular graph structures.

Obs.4: LLM-based models, in general, exhibit sub-par performance in quantitative tasks compared to traditional DL-based models, attributed to their inadequacy in sharing and transferring learned molecular interaction insights across datasets and the absence of CoT-guided inference.

3.4 Ablation Study (RQ3)

Table 3 presents an ablation study aimed at dissecting the influence of SMILE auxiliary analysis and the optimized training paradigms based on Broad-grained and Fine-grained CoT. For the CSI dataset, properties such as the maximum absorption wavelength (Absorption), maximum emission wavelength (Emission), and excited state lifetime (Lifetime) are denoted as Abs., Emis., and Life., respectively. Key observations are as follows:

Obs.5: The three studied ablations exhibit significant influence on the results. For example, the collective impact of these three ablations registers an average drop of 12.77%, affirming the substantial enhancement imparted by the proposed strategies.

Obs.6: The most pronounced effect is observed with the ablation of the Fine-grained CoT paradigm, which incurs an average accuracy decrement of 18.82%. This underscores the pivotal role of guiding the LLM to deduce a numerical range, a strategy particularly beneficial for quantitative analysis tasks, typically a challenging domain for LLMs.

Obs.7: The least pronounced, yet significant, impact stems from the optimization of the Broad-grained CoT training paradigm, with an average accuracy reduction of 4.35%. Its importance is particularly underscored for molecular pairs involving larger and more complex molecules, where directly predicting interactive property by LLMs is arduous.

4 Conclusion

This work focuses on molecule rationale learning, which plays a pivotal role in predicting molecular interactions. Specifically, we introduce a novel, unified LLM-based framework for predicting molecular interactive properties, termed MolTC. To efficiently train it, we propose a multi-tiered CoT principle to guide the training paradigm. Experiments conducted across twelve varied datasets demonstrate the superiority of our method over the current GNN and LLM-based baselines. This breakthrough sets a new standard for integrating multimodal data in LLM-based MRL.

580 Limitations

581 While this research has undergone extensive testing
582 across a diverse array of datasets covering various
583 domains, it does have certain limitations. Specifi-
584 cally, the study has not been subjected to datasets
585 comprising exceptionally large molecules, which
586 represent extreme cases. Furthermore, the method-
587 ologies employed in this research have not yet been
588 adapted or evaluated in contexts requiring few-shot
589 or zero-shot learning scenarios. Future endeavors
590 will focus on expanding the scope of this study to
591 encompass these areas.

592 Ethics Statement

593 This work is primarily foundational in molecular
594 relational learning, focusing on the development
595 of a unified LLM-based paradigm. Its primary
596 aim is to contribute to the academic community by
597 enhancing the understanding and implementation
598 of the molecular relational modeling process. We
599 do not foresee any direct, immediate, or negative
600 societal impacts stemming from the outcomes of
601 our research.

602 References

603 Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. Scibert:
604 A pretrained language model for scientific text. *arXiv*
605 *preprint arXiv:1903.10676*.

606 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
607 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
608 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
609 Askell, et al. 2020. Language models are few-shot
610 learners. *Advances in neural information processing*
611 *systems*, 33:1877–1901.

612 Kexin Chen, Junyou Li, Kunyi Wang, Yuyang Du,
613 Jiahui Yu, Jiamin Lu, Guangyong Chen, Lanqing
614 Li, Jiezhong Qiu, Qun Fang, et al. 2023. Towards
615 an automatic ai agent for reaction condition recom-
616 mendation in chemical synthesis. *arXiv preprint*
617 *arXiv:2311.10776*.

618 Dimitrios Christofidellis, Giorgio Giannone, Jannis
619 Born, Ole Winther, Teodoro Laino, and Matteo Man-
620 ica. 2023. Unifying molecular and textual represen-
621 tations via multi-task language modelling. *arXiv*
622 *preprint arXiv:2301.12586*.

623 Yunsie Chung, Florence H Vermeire, Haoyang Wu,
624 Pierre J Walker, Michael H Abraham, and William H
625 Green. 2022. Group contribution and machine learn-
626 ing approaches to predict abraham solute parameters,
627 solvation free energy, and solvation enthalpy. *Journal*
628 *of Chemical Information and Modeling*, 62(3):433–
629 446.

W Dai, J Li, D Li, AMH Tiong, J Zhao, W Wang, B Li,
P Fung, and S Hoi. Instructblip: Towards general-
purpose vision-language models with instruction tun-
ing. *arxiv* 2023. *arXiv preprint arXiv:2305.06500*.

Andreea Deac, Yu-Hsiang Huang, Petar Veličković,
Pietro Liò, and Jian Tang. 2019. Drug-drug ad-
verse effect prediction with graph co-attention. *arXiv*
preprint arXiv:1905.00534.

Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zong-
han Yang, Yusheng Su, Shengding Hu, Yulin Chen,
Chi-Min Chan, Weize Chen, et al. 2022. Delta tuning:
A comprehensive study of parameter efficient meth-
ods for pre-trained language models. *arXiv preprint*
arXiv:2203.06904.

Carl Edwards, Tuan Lai, Kevin Ros, Garrett Honke,
Kyunghyun Cho, and Heng Ji. 2022. Translation
between molecules and natural language. *arXiv*
preprint arXiv:2204.11817.

Junfeng Fang, Xinglin Li, Yongduo Sui, Yuan Gao,
Guibin Zhang, Kun Wang, Xiang Wang, and Xiang-
nan He. 2024. Exgc: Bridging efficiency and explain-
ability in graph condensation. In *WWW*. ACM.

Junfeng Fang, Wei Liu, Yuan Gao, Zemin Liu,
An Zhang, Xiang Wang, and Xiangnan He. 2023a.
[Evaluating post-hoc explanations for graph neural
networks via robustness analysis](#). In *Thirty-seventh
Conference on Neural Information Processing Sys-
tems*.

Junfeng Fang, Xiang Wang, An Zhang, Zemin Liu, Xi-
angnan He, and Tat-Seng Chua. 2023b. Cooperative
explanations of graph neural networks. In *WSDM*,
pages 616–624. ACM.

Xiaomin Fang, Lihang Liu, Jieqiong Lei, Donglong
He, Shanzhuo Zhang, Jingbo Zhou, Fan Wang, Hua
Wu, and Haifeng Wang. 2022. Geometry-enhanced
molecular representation learning for property pre-
diction. *Nature Machine Intelligence*, 4(2):127–134.

Laura M Grubbs, Mariam Saifullah, E Nohelli, Shulin
Ye, Sai S Achi, William E Acree Jr, and Michael H
Abraham. 2010. Mathematical correlations for de-
scribing solute transfer into functionalized alkane
solvents containing hydroxyl, ether, ester or ketone
solvents. *Fluid phase equilibria*, 298(1):48–53.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan
Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
and Weizhu Chen. 2021. Lora: Low-rank adap-
tation of large language models. *arXiv preprint*
arXiv:2106.09685.

Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zit-
nik, Percy Liang, Vijay Pande, and Jure Leskovec.
2019. Strategies for pre-training graph neural net-
works. *arXiv preprint arXiv:1905.12265*.

Burrows C J, Harper J B, and et al. Sander W. 2022.
Solvation effects in organic chemistry. *The Journal*
of Organic Chemistry, 87(3):1599–1601.

685	Kanchan Jha, Sriparna Saha, and Hiteshi Singh. 2022a.	Zhiyuan Liu, Sihang Li, Yanchen Luo, Hao Fei, Yixin	739
686	Prediction of protein–protein interaction using graph	Cao, Kenji Kawaguchi, Xiang Wang, and Tat-Seng	740
687	neural networks. <i>Scientific Reports</i> , 12(1):8360.	Chua. 2023. Molca: Molecular graph-language	741
688	Kanchan Jha, Sriparna Saha, and Hiteshi Singh. 2022b.	modeling with cross-modal projector and uni-modal	742
689	Prediction of protein–protein interaction using graph	adapter. <i>arXiv preprint arXiv:2310.12798</i> .	743
690	neural networks. <i>Scientific Reports</i> , 12(1):8360.		
691	Joonyoung F Joung, Minhi Han, Minseok Jeong, and	Micha Livne, Zulfat Miftahutdinov, Elena Tutubalina,	744
692	Sungnam Park. 2020. Experimental database of opti-	Maksim Kuznetsov, Daniil Polykovskiy, Annika	745
693	cal properties of organic compounds. <i>Scientific data</i> ,	Brundyn, Aastha Jhunjhunwala, Anthony Costa,	746
694	7(1):295.	Alex Aliper, and Alex Zhavoronkov. 2023. nach0:	747
695	Sunghwan Kim, Jie Chen, Tiejun Cheng, Asta Gindu-	Multimodal natural and chemical languages founda-	748
696	lyte, Jia He, Siqian He, Qingliang Li, Benjamin A.	tion model. <i>arXiv preprint arXiv:2311.12410</i> .	749
697	Shoemaker, Paul A. Thiessen, Bo Yu, Leonid Za-		
698	slavsky, Jian Zhang, and Evan E. Bolton. 2023. Pub-	Ilya Loshchilov and Frank Hutter. 2017. Decou-	750
699	chem 2023 update. <i>Nucleic Acids Res.</i> , 51(D1):1373–	pled weight decay regularization. <i>arXiv preprint</i>	751
700	1380.	<i>arXiv:1711.05101</i> .	752
701	Thomas N. Kipf and Max Welling. 2017. Semi-	Sourab Mangrulkar, Sylvain Gugger, Lysandre De-	753
702	supervised classification with graph convolutional	but, Younes Belkada, Sayak Paul, and Benjamin	754
703	networks. In <i>ICLR (Poster)</i> .	Bossan. 2022. Peft: State-of-the-art parameter-	755
704	Namkyeong Lee, Dongmin Hyun, Gyoung S. Na, Sung-	efficient fine-tuning methods. https://github.com/huggingface/peft .	756
705	won Kim, Junseok Lee, and Chanyoung Park. 2023a.		757
706	Conditional graph information bottleneck for molec-	Aleksandr V Marenich, Casey P Kelly, Jason D Thomp-	758
707	ular relational learning. In <i>ICML</i> , volume 202 of	son, Gregory D Hawkins, Candee C Chambers,	759
708	<i>Proceedings of Machine Learning Research</i> , pages	David J Giesen, Paul Winget, Christopher J Cramer,	760
709	18852–18871. PMLR.	and Donald G Truhlar. 2020. Minnesota solvation	761
710	Namkyeong Lee, Kanghoon Yoon, Gyoung S Na, Sein	database (mnsol) version 2012.	762
711	Kim, and Chanyoung Park. 2023b. Shift-robust	Fanwang Meng, Hanwen Zhang, Juan Samuel Collins-	763
712	molecular relational learning with causal substruc-	Ramirez, and Paul W. Ayers. 2023. Something for	764
713	ture. <i>arXiv preprint arXiv:2305.18451</i> .	nothing: Improved solvation free energy prediction	765
714	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi.	with learning.	766
715	2023a. Blip-2: Bootstrapping language-image pre-	David L Mobley and J Peter Guthrie. 2014. Freesolv:	767
716	training with frozen image encoders and large lan-	a database of experimental and calculated hydration	768
717	guage models. <i>arXiv preprint arXiv:2301.12597</i> .	free energies, with input files. <i>Journal of computer-</i>	769
718	Tianhao Li, Sandesh Shetty, Advait Kamath, Ajay	<i>aided molecular design</i> , 28:711–720.	770
719	Jaiswal, Xiaoqian Jiang, Ying Ding, and Yejin Kim.	Edouard Moine, Romain Privat, Baptiste Sirjean, and	771
720	2023b. Cancergpt: Few-shot drug pair synergy	Jean-Noël Jaubert. 2017. Estimation of solvation	772
721	prediction using large pre-trained language models.	quantities from experimental thermodynamic data:	773
722	<i>ArXiv</i> .	Development of the comprehensive compsol data-	774
723	Zimeng Li, Shichao Zhu, Bin Shao, Xiangxiang Zeng,	bank for pure and mixed solutes. <i>Journal of Physical</i>	775
724	Tong Wang, and Tie-Yan Liu. 2023c. Dsn-ddi: an	<i>and Chemical Reference Data</i> , 46(3).	776
725	accurate and generalized framework for drug–drug	Arnold K Nyamabo, Hui Yu, and Jian-Yu Shi. 2021.	777
726	interaction prediction by dual-view representation	Ssi-ddi: substructure–substructure interactions for	778
727	learning. <i>Briefings in Bioinformatics</i> , 24(1):bbac597.	drug–drug interaction prediction. <i>Briefings in Bioin-</i>	779
728	Shenggeng Lin, Yanjing Wang, Lingfeng Zhang, Yanyi	<i>formatics</i> , 22(6):bbab133.	780
729	Chu, Yatong Liu, Yitian Fang, Mingming Jiang,	Gilchan Park, Sean McCorkle, Carlos Soto, Ian Blaby,	781
730	Qiankun Wang, Bowen Zhao, Yi Xiong, et al. 2022.	and Shinjae Yoo. 2022. Extracting protein-protein	782
731	Mdf-sa-ddi: predicting drug–drug interaction events	interactions (ppis) from biomedical literature using	783
732	based on multi-source drug fusion, multi-source fea-	attention-based relational context information. In	784
733	ture fusion and transformer self-attention mechanism.	<i>2022 IEEE International Conference on Big Data</i>	785
734	<i>Briefings in Bioinformatics</i> , 23(1):bbab421.	(<i>Big Data</i>), pages 2052–2061. IEEE.	786
735	Xuan Lin, Zhe Quan, Zhi-Jie Wang, Tengfei Ma, and	Yashaswi Pathak, Siddhartha Laghuvarapu, Sarvesh	787
736	Xiangxiang Zeng. 2020. Kgnn: Knowledge graph	Mehta, and U Deva Priyakumar. 2020. Chemically	788
737	neural network for drug–drug interaction prediction.	interpretable graph interaction network for prediction	789
738	In <i>IJCAI</i> , volume 380, pages 2739–2745.	of pharmacokinetic properties of drug-like molecules.	790
		In <i>Proceedings of the AAAI Conference on Artificial</i>	791
		<i>Intelligence</i> , volume 34, pages 873–880.	792

793	Qizhi Pei, Wei Zhang, Jinhua Zhu, Kehan Wu, Kaiyuan Gao, Lijun Wu, Yingce Xia, and Rui Yan. 2023. Biot5: Enriching cross-modal integration in biology with chemical knowledge and natural language associations. <i>arXiv preprint arXiv:2310.07276</i> .		
794			
795			
796			
797			
798	C. Reichardt. 2021. Solvation effects in organic chemistry: A short historical overview. <i>The Journal of Organic Chemistry</i> , 87(3):1616–1629.		
799			
800			
801	Dan M Roden, Robert A Harrington, Athena Poppas, and Andrea M Russo. 2020. Considerations for drug interactions on qtc in exploratory covid-19 treatment. <i>Circulation</i> , 141(24):e906–e907.		
802			
803			
804			
805	Yu Rong, Yatao Bian, Tingyang Xu, Weiyang Xie, Ying Wei, Wenbing Huang, and Junzhou Huang. 2020. Self-supervised graph transformer on large-scale molecular data. In <i>NeurIPS</i> .		
806			
807			
808			
809	Jae Yong Ryu, Hyun Uk Kim, and Sang Yup Lee. 2018. Deep learning improves prediction of drug–drug and drug–food interactions. <i>Proceedings of the national academy of sciences</i> , 115(18):E4304–E4311.		
810			
811			
812			
813	Tatsuya Sagawa and Ryosuke Kojima. 2023. Reaction5: a large-scale pre-trained model towards application of limited reaction data. <i>arXiv preprint arXiv:2311.06708</i> .		
814			
815			
816			
817	Yaorui Shi, An Zhang, Enzhi Zhang, Zhiyuan Liu, and Xiang Wang. 2023. Relm: Leveraging language models for enhanced chemical reaction prediction. In <i>EMNLP (Findings)</i> , pages 5506–5520. Association for Computational Linguistics.		
818			
819			
820			
821			
822	Teague Sterling and John J Irwin. 2015. Zinc 15–ligand discovery for everyone. <i>Journal of chemical information and modeling</i> , 55(11):2324–2337.		
823			
824			
825	Nicholas P Tatonetti, Patrick P Ye, Roxana Daneshjou, and Russ B Altman. 2012. Data-driven prediction of drug effects and interactions. <i>Science translational medicine</i> , 4(125):125ra31–125ra31.		
826			
827			
828			
829	Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. 2022. Galactica: A large language model for science. <i>arXiv preprint arXiv:2211.09085</i> .		
830			
831			
832			
833			
834	Jithin John Varghese and Samir H Mushrif. 2019. Origins of complex solvent effects on chemical reactivity and computational tools to investigate them: a review. <i>Reaction Chemistry & Engineering</i> , 4(2):165–206.		
835			
836			
837			
838	Florence H Vermeire and William H Green. 2021. Transfer learning for solvation free energies: From quantum chemistry to experiments. <i>Chemical Engineering Journal</i> , 418:129307.		
839			
840			
841			
842	Hanchen Wang, Defu Lian, Ying Zhang, Lu Qin, and Xuemin Lin. 2020. Gognn: Graph of graphs neural network for predicting structured entity interactions. <i>arXiv preprint arXiv:2005.05537</i> .		
843			
844			
845			
	Kun Wang, Yuxuan Liang, Xinglin Li, Guohao Li, Bernard Ghanem, Roger Zimmermann, Zhengyang Zhou, Huahui Yi, Yudong Zhang, and Yang Wang. 2023. Brave the wind and the waves: Discovering robust and generalizable graph lottery tickets. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , pages 1–17.		846 847 848 849 850 851 852
	Kun Wang, Yuxuan Liang, Pengkun Wang, Xu Wang, Pengfei Gu, Junfeng Fang, and Yang Wang. 2022. Searching lottery tickets in graph neural networks: A dual perspective. In <i>The Eleventh International Conference on Learning Representations</i> .		853 854 855 856 857
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in Neural Information Processing Systems</i> , 35:24824–24837.		858 859 860 861 862
	Hao Wu, Shilong Wang, Yuxuan Liang, Zhengyang Zhou, Wei Huang, Wei Xiong, and Kun Wang. 2023. Earthfarseer: Versatile spatio-temporal dynamical systems modeling in one model.		863 864 865 866
	Yutong Xia, Yuxuan Liang, Haomin Wen, Xu Liu, Kun Wang, Zhengyang Zhou, and Roger Zimmermann. 2023. Deciphering spatio-temporal graph forecasting: A causal lens and treatment.		867 868 869 870
	Zhengyi Yang, Jiancan Wu, Yanchen Luo, Jizhi Zhang, Yancheng Yuan, An Zhang, Xiang Wang, and Xiangnan He. 2023. Large language model can interpret latent space of sequential recommender. <i>arXiv preprint arXiv:2310.20487</i> .		871 872 873 874 875
	Jiahui Yu, Chengwei Zhang, Yingying Cheng, Yunfang Yang, Yuan-Bin She, Fengfan Liu, Weike Su, and An Su. 2023. Solvbert for solvation free energy and solubility prediction: a demonstration of an nlp model for predicting the properties of molecular complexes. <i>Digital Discovery</i> , 2(2):409–421.		876 877 878 879 880 881
	Wen Zhang, Yanlin Chen, Feng Liu, Fei Luo, Gang Tian, and Xiaohong Li. 2017. Predicting potential drug-drug interactions by integrating chemical, biological, phenotypic and network data. <i>BMC bioinformatics</i> , 18:1–12.		882 883 884 885 886
	Gengmo Zhou, Zhifeng Gao, Qiankun Ding, Hang Zheng, Hongteng Xu, Zhewei Wei, Linfeng Zhang, and Guolin Ke. 2023. Uni-mol: a universal 3d molecular representation learning framework.		887 888 889 890
	M Zitnik, R Sosi, S Maheshwari, and J Leskovec. 2018. Stanford biomedical network dataset collection. <i>Biosn. Datasets Stanford Biomed. Netw. Dataset Collect.</i>		891 892 893 894
	A Related Work		895
	Since exhaustive experimental validation of the molecule interactions is notoriously time-consuming and costly (Lee et al., 2023a), more		896 897 898

recently, adopting LLM has emerged as a promising alternative for efficient and effective molecular relational learning, which are known for their vast knowledge repositories and advanced logical inference capabilities. Specifically,

- (Park et al., 2022; Jha et al., 2022a,b) focus on employ LLM to optimize **protein-protein interactions (PPI)** tasks. In this context, proteins are represented as residue contact graphs, also known as amino acid graphs, where each node is a residue (Jha et al., 2022b). Notably, (Jha et al., 2022b) leverages the superior encoding capabilities of the biochemical LLMs, where the input to the LLM is the protein sequence, and the output is a feature vector for each amino acid in the sequence. This output is then used as node features in the residue contact graph to enhance the prediction of PPI tasks.
- (Sagawa and Kojima, 2023; Chen et al., 2023; Livne et al., 2023; Shi et al., 2023) focus using LLMs to optimize **chemical reactions**. Specifically, (Shi et al., 2023) selects in-context reaction examples with varying confidence scores closest to the target reaction query, encouraging large models to understand the relationships between these reactions. (Sagawa and Kojima, 2023) focuses on optimizing low-sample organic chemical applications by pretraining them with extensive compound libraries and fine-tuning with smaller in-house datasets for specific tasks. (Livne et al., 2023) introduces a new foundational model, nach0, capable of solving various chemical and biological tasks, including molecular synthesis.
- (Li et al., 2023b; Pei et al., 2023) focus on using LLMs to optimize tasks related to **drug molecules**. Specifically, (Pei et al., 2023) enriches cross-modal integration in biology with chemical knowledge and natural language associations, achieving significant results in multiple drug-target interaction prediction tasks. Meanwhile, (Li et al., 2023b) concentrates on few-shot drug pair synergy prediction.

B Experiments

Here, we provide a detailed experimental setup along with additional results. It is important to note that for aspects such as dataset division and hyperparameter configurations in baselines, we followed the settings established by CGIB (Lee et al.,

2023a). Moreover, all settings can be found in our code <https://anonymous.4open.science/r/MolTC-F>.

B.1 Datasets

We employ 12 datasets across various domains such as DDI, SSI and CSI.

Drugbank (version 5.0.3). this dataset consists of 1704 drugs, 191400 drug pairs, and defines 86 distinct DDI event types. Essential drug information, including DrugBank ID, drug name, molecular SMILES, and target, provided.

ZhangDDI. (Zhang et al., 2017) it contains 548 drugs and 48,548 pairwise interaction data and multiple types of similarity information about these drug pairs.

ChChMiner. (Zitnik et al., 2018) it contains 1,322 drugs and 48,514 labeled DDIs, obtained through drug labels and scientific publications.

DeepDDI. (Ryu et al., 2018) contains 192,284 labeled DDIs and their detailed side-effect information, which is extracted from Drugbank.

TWOSIDES. (Tatonetti et al., 2012) it collected 555 drugs and their 3,576,513 pairwise interactions involving 1318 interaction types from TWOSIDES.

Chromophore. (Joung et al., 2020) contains 20,236 combinations of 7,016 chromophores and 365 solvents which are given in the SMILES string format. All optical properties are based on scientific publications and unreliable experimental results are excluded after examination of absorption and emission spectra. In this dataset, we measure our model performance on predicting maximum absorption wavelength (Absorption), maximum emission wavelength (Emission) and excited state lifetime (Lifetime) properties which are important parameters for the design of chromophores for specific applications. We delete the NaN values to create each dataset which is not reported in the original scientific publications. Moreover, for Lifetime data, we use log normalized target value since the target value of the dataset is highly skewed inducing training instability.

MNSol. (Marenich et al., 2020) contains 3,037 experimental free energies of solvation or transfer energies of 790 unique solutes and 92 solvents. In this work, we consider 2,275 combinations of 372 unique solutes and 86 solvents following previous work.

FreeSolv. (Mobley and Guthrie, 2014) provides 643 experimental and calculated hydration free energy of small molecules in water. In this work, we

998 consider 560 experimental results following previ- 1048
999 ous work. 1049

1000 **CompSol.** (Moine et al., 2017) dataset is proposed 1050
1001 to show how solvation energies are influenced by 1051
1002 hydrogen-bonding association effects. We consider 1052
1003 3,548 combinations of 442 unique solutes and 259 1053
1004 solvents in the dataset following previous work. 1054

1005 **Abraham.** (Grubbs et al., 2010) dataset is a col- 1055
1006 lection of data published by the Abraham research 1056
1007 group at College London. We consider 6,091 com- 1057
1008 binations of 1,038 unique solutes and 122 solvents 1058
1009 following previous work. 1059

1010 **CombiSolv.** (Vermeire and Green, 2021) con- 1060
1011 tains all the data of MNSolv, FreeSolv, CompSol, 1061
1012 and Abraham, resulting in 10,145 combinations of 1062
1013 1,368 solutes and 291 solvents. 1063

1014 **CombiSolv-QM.** (Vermeire and Green, 2021) 1064
1015 is generated with 1 million combinations of 1065
1016 284 commonly used solvents and 11,029 so- 1066
1017 lutes. Those 1 million data points are randomly 1067
1018 selected from all possible solvent-solute com- 1068
1019 binations. Solvents and solutes with elements 1069
1020 *H, B, C, N, O, F, P, S, Cl, Br* and *I* are included 1070
1021 with a solute molar mass ranging from 2.02 g/mol 1071
1022 to 1776.89 g/mol. 1072

1023 B.2 Baselines 1073

1024 We use both specific task conventional deep learn- 1074
1025 ing models and current biochemical LLMs as the 1075
1026 baselines. Specifically, for qualitative tasks: 1076

1027 **GoGNN.** (Wang et al., 2020) It extracts features 1077
1028 from structured entity graphs and entity interaction 1078
1029 graphs in a hierarchical manner. We also propose a 1079
1030 dual attention mechanism that enables the model to 1080
1031 preserve the importance of neighbors in both levels 1081
1032 of the graph. 1082

1033 **MHCADDI.** (Deac et al., 2019) A gated informa- 1083
1034 tion transfer neural network is used to control the 1084
1035 extraction of substructures and then interact based 1085
1036 on an attention mechanism. 1086

1037 **DeepDDI.** (Ryu et al., 2018) First, the structural 1087
1038 similarity profile is calculated between the two in- 1088
1039 put drugs and other drugs, and then prediction is 1089
1040 completed based on the deep neural network. 1090

1041 **SSI-DDI.** (Nyamabo et al., 2021) it use a 4-layer 1091
1042 GAT network to extract substructures at different 1092
1043 levels, and finally complete the final prediction 1093
1044 based on the co-attention mechanism 1094

1045 **CGIB.** (Lee et al., 2023a) Based on the graph con- 1095
1046 ditional information bottleneck theory, conditional 1096
1047 subgraphs are extracted to complete the interaction 1097

between molecules. 1048

CMRL. (Lee et al., 2023b) it detects the core sub- 1049
1050 structure that is causally related to chemical reac- 1051
1052 tions. we introduce a novel conditional intervention 1052
1053 framework whose intervention is conditioned on 1053
1054 the paired molecule. With the conditional interven- 1054
1055 tion framework. 1055

MDF-SA-DDI. (Lin et al., 2022) it predicts inter- 1056
1057 action (DDI) events based on multi-source drug 1056
1058 fusion, multi-source feature fusion and transformer 1057
1059 self-attention mechanism. 1058

DSN-DDI. (Li et al., 2023c) it employs local and 1059
1060 global representation learning modules iteratively 1060
1061 and learns drug substructures from the single drug 1061
1062 ‘intra-view’) and the drug pair (‘inter-view’) simul- 1062
1063 taneously. 1063

For quantitative task, we employ the following 1064
1065 baselines: 1065

D-MPNN (Vermeire and Green, 2021) it employes 1066
1067 a transfer learning approach to predict solvation 1067
1068 free energies, integrating quantum calculation fun- 1068
1069 damentals with the heightened accuracy of experi- 1069
1070 mental measurements through two new databases, 1070
1071 CombiSolv-QM and CombiSolv-Exp. 1071

SolvBert. (Yu et al., 2023) it interprets solute 1072
1073 and solvent interactions through their combined 1073
1074 SMILES representation. Pre-trained using unsu- 1074
1075 pervised learning with a substantial computational 1075
1076 solvation free energy database, SolvBERT is adapt- 1076
1077 able to predict experimental solvation free energy 1077
1078 or solubility by fine-tuning on specific databases. 1078

SMD. (Meng et al., 2023) utilizes the quantum 1079
1080 charge density of a solute and a continuum repre- 1080
1081 sentation of the solvent. It breaks down solvation 1081
1082 free energy into two components: bulk electrostatic 1082
1083 contribution, treated through a self-consistent reac- 1083
1084 tion field using IEF-PCM, and a cavity-dispersion- 1084
1085 solvent-structure term, accounting for short-range 1085
1086 interactions in the solvation shell based on atomic 1086
1087 surface areas with geometry-dependent constants. 1087

CIGIN. (Pathak et al., 2020) is a method based 1088
1089 on graph neural networks. The proposed model 1089
1090 adopts an end-to-end framework consisting of three 1090
1091 essential phases: message passing, interaction, and 1091
1092 prediction. In the final phase, these stages are lever- 1092
1093 aged to predict solvation free energies. 1093

GEM. (Fang et al., 2022) exhibits a uniquely de- 1094
1095 signed geometry-based graph neural network archi- 1095
1096 tecture, complemented by several dedicated self- 1096
1097 supervised learning strategies at the geometry level. 1097
1098 That aims to acquire comprehensive molecular ge- 1098

Table 4: Comparative performance of various methods in qualitative and quantitative interactive tasks. The best-performing methods are highlighted with a gray background, while the second-best methods are underscored for emphasis.

Domains	Datasets	Metrics	Baselines				Ours MolTC
			Galactica	Chem T5	MolCA	MolT5	
DDI	TWO SIDES	ACC	82.01 \pm 1.76	84.43 \pm 2.58	90.07 \pm 1.86	<u>92.73</u> \pm 1.65	98.42 \pm 0.72
		AUCROC	87.99 \pm 2.41	89.52 \pm 1.64	93.68 \pm 0.83	<u>94.00</u> \pm 0.61	99.02 \pm 0.14
SSI	MNSol	MAE	0.584 \pm 0.095	0.504 \pm 0.038	0.491 \pm 0.053	<u>0.449</u> \pm 0.081	0.324 \pm 0.019
		RMSE	1.002 \pm 0.101	0.973 \pm 0.079	0.930 \pm 0.062	<u>0.858</u> \pm 0.069	0.585 \pm 0.023
CSI	Absorption	RMSE	43.16 \pm 1.38	38.70 \pm 1.84	<u>36.53</u> \pm 2.03	38.01 \pm 2.27	28.28 \pm 2.20
	Emission	RMSE	49.85 \pm 2.47	46.18 \pm 2.28	<u>43.35</u> \pm 1.94	46.06 \pm 1.65	35.43 \pm 1.88
	Lifetime	RMSE	1.951 \pm 0.115	1.633 \pm 0.069	1.480 \pm 0.092	<u>1.394</u> \pm 0.145	1.198 \pm 0.073

Table 5: Comparative performance of various methods in CombiSolv-QM. The best-performing methods are highlighted with a gray background, while the second-best methods are underscored for emphasis.

Baseline Model	CombiSolv-QM	
	MAE	RMSE
CIGIN	0.077 \pm 0.002	0.176 \pm 0.004
GNN D-MPNN	0.116 \pm 0.006	0.208 \pm 0.005
Based GEM	0.079 \pm 0.003	0.162 \pm 0.002
CGIB	<u>0.074</u> \pm 0.004	<u>0.150</u> \pm 0.005
GOVER	0.094 \pm 0.003	0.277 \pm 0.005
ML SolvBert	0.102 \pm 0.005	0.318 \pm 0.006
Based Uni-Mol	0.089 \pm 0.006	0.214 \pm 0.005
SMD	0.107 \pm 0.004	0.341 \pm 0.003
Galactica	0.303 \pm 0.004	0.601 \pm 0.008
LLM Chem T5	0.321 \pm 0.006	0.555 \pm 0.008
Based MolCA	0.298 \pm 0.004	0.545 \pm 0.007
MolT5	0.214 \pm 0.004	0.339 \pm 0.009
MolTC (Ours)	<u>0.072</u> \pm 0.002	<u>0.140</u> \pm 0.003

ometry knowledge for accurate prediction of molecular properties.

GOVER. (Rong et al., 2020) captures rich structural information from extensive unlabeled molecular data through self-supervised tasks, employing a flexible Transformer-style architecture integrated with Message Passing Networks. This allows GROVER to be trained efficiently on large-scale datasets without supervision, addressing data scarcity and bias challenges.

Uni-Mol. (Zhou et al., 2023) incorporates two pre-trained models featuring the SE(3) Transformer architecture: a molecular model pre-trained on 209 million molecular conformations and a pocket model pre-trained on 3 million candidate protein pocket data. Additionally, Uni-Mol integrates various fine-tuning strategies to effectively apply these pre-trained models across diverse downstream tasks.

B.3 Modules

In our experiments, the two graph encoder are instantiated by the five-layer GINE (Hu et al., 2019). We conduct 2 million molecules from the ZINC15 (Sterling and Irwin, 2015) dataset to pretrain them by contrastive learning following (Liu et al., 2023). Similarly, two projector are initialized with the encoder-only transformer, Sci-BERT, which is pre-trained on scientific publications (Beltagy et al., 2019), while its cross-attention modules are randomly initialized. More detailed pretraining process of our Q-Formers follows the training process in (Liu et al., 2023), such as there are 8 query tokens in Q-Formers ($N_q = 8$). Note that for LLM-based baselines, we fine-tune the backbone LLMs on task-specific datasets for fair comparison. Their prediction is considered accurate only if the outputs include words or numbers that correctly depict the interaction in question, without presenting any that describe alternative interactions.

B.4 Training Epochs

During the fine-tuning phase, the number of epochs varies for different tasks. For example, for the DDI task, we typically fine-tune for 100 epochs. For SSI datasets with more than 3000 molecular pairs, we initially fine-tune on the CombiSolv-QM (Vermeire and Green, 2021) dataset for 100 epochs, followed by an additional 30 epochs on their respective datasets. For SSI datasets with fewer than 3000 molecular pairs, this number is adjusted to 20. Furthermore, both the fine-tuning and pre-training phases employ the same configuration for the optimizer and learning rate scheduler, as detailed in the following section.

B.5 Training Strategy

We employ the AdamW optimizer (Loshchilov and Hutter, 2017) with a weight decay set at 0.05. Our learning rate strategy utilizes a combination of linear warm-up and cosine decay, optimizing the training process by initially increasing the learning rate to promote faster convergence, and then gradually decreasing it according to a cosine curve to fine-tune the model parameters. LoRA is implemented using the Open Delta library (Ding et al., 2022), and the PEFT library (Mangrulkar et al., 2022). LoRA’s rank r is set to 16, while LoRA is applied to Galactica’s modules of [q_proj, v_proj, out_proj, fc1, fc2] following (Liu et al., 2023). This configuration yields a LoRA adapter with 12M parameters which constitutes merely 0.94% of the parameters in the Galactica_{1.3B}.

B.6 More Experimental Results

Table 4 presents the experimental results not shown in the main text due to length constraints. Note that the three datasets in the CSI domain are all derived by splitting the Chromophore dataset. As discussed in Section 3.3, for a fair comparison, we limited the input features to four-dimensional molecular information, comprising atom type, chirality tag, bond type, and bond direction. Given the difficulty of convergence for some DL-based baselines under this setting, we only showcased the performance of the LLM-based baselines. Meanwhile, considering that our SSI tasks are firstly fine-tuned on the CombiSolv-QM dataset, we present the comprehensive results of this dataset, as shown in Table 5. Observations from Table 4 and 5 are largely consistent with those in the main experimental section. That is, across all tasks, our MolTC outperforms the LLM-based baseline methods in a large margin.

C Future Work

In this paper, we introduce a novel unified framework, leveraging LLM technology to predict molecular interactive properties. The future development directions of this project are twofold. First, there is an emphasis on expanding its application scope, for instance, applying it to downstream tasks such as few-shot learning. Second, we aim to enhance its capabilities by incorporating technologies like graph explainability (Fang et al., 2023a,b), graph sampling (Wang et al., 2022, 2023; Fang et al., 2024), and spatio-temporal modeling (Xia et al.,

2023; Wu et al., 2023), making it more comprehensive or enabling it to process multiple inputs simultaneously, instead of just two.

1201
1202
1203