ENTROPY-GUIDED DYNAMIC TOKENS FOR GRAPH-LLM ALIGNMENT IN MOLECULAR UNDERSTANDING

Anonymous authors

Paper under double-blind review

ABSTRACT

Molecular understanding is central to advancing areas such as scientific and drug discovery, yet Large Language Models (LLMs) struggle to understand molecular graphs effectively. Existing graph—LLM bridges often adapt the Q-Formerstyle connector with fixed-length static tokens, which is originally designed for vision tasks. These designs overlook stereochemistry and substructural context and typically require costly LLM-backbone fine-tuning, limiting efficiency and generalization. We introduce **EDT-Former**, an Entropy-guided **D**ynamic Token Trans**former** that generates tokens aligned with informative molecular patches, thereby preserving both local and global structural features for molecular graph understanding. Beyond prior approaches, EDT-Former enables alignment between frozen graph encoders and LLMs without tuning the LLM backbone, resulting in computationally efficient finetuning, and achieves state-of-the-art results on MoleculeQA and Mol-Instructions benchmarks, underscoring its effectiveness for scalable and generalizable multimodal molecular understanding.

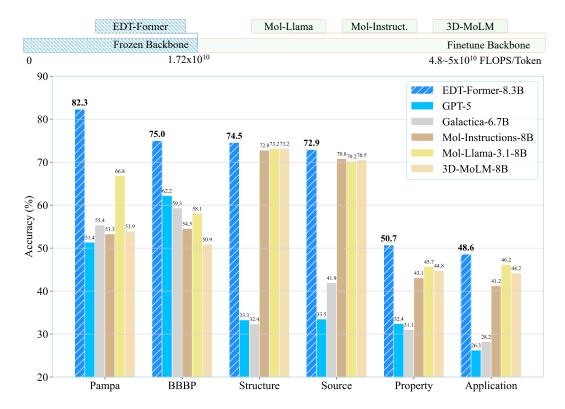


Figure 1: LLM joint fine-tuning efficiency and benchmark performance preview of EDT-Former.

1 Introduction

Background. Large language models (LLMs) such as Llama (Touvron et al., 2023a; MetaAI, 2024) and GPT (OpenAI, 2023; 2025) series demonstrate powerful multimodal reasoning across domains. This rapid progress highlights the potential adaptation to specialized scientific domains, and NatureLM (Xia et al., 2025) exemplifies this direction in scientific knowledge. Within molecular science, recent foundation models have extended large-scale pretraining to molecular graphs and geometry structures. Uni-Mol-v2 (Ji et al., 2024) leverages extensive conformer datasets to jointly model atomic, graph, and geometric representations, while UniMolT (Zhang et al., 2024b) adapts instruction-tuned LLMs to molecular reasoning and understanding tasks.

To bridge molecular graphs with natural language, multimodal connectors have been introduced. Most approaches adopt the Q-Former mechanism (Li et al., 2023), where a fixed number of modality-anchor tokens are learned to query structural encoders. Representative examples include 3D-MolT5 (Pei et al., 2024b), Mol-LLM (Lee et al., 2025), and Mol-LLaMA (Kim et al., 2025), which demonstrate the feasibility of aligning molecular structures (graph and geometry) with language models and enabling cross-modal molecular understanding.

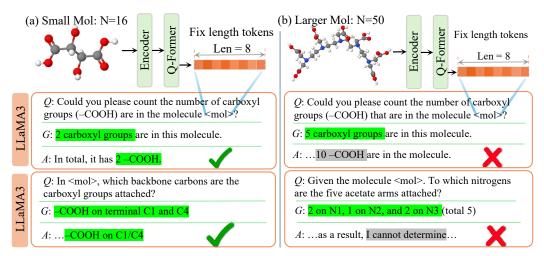


Figure 2: Illustration of motivation—Loss of structure. Comparison of molecules of different sizes (atom counts N=16 and N=50) encoded by the same fixed query length Q-Former bridge (8 query tokens) to Llama-3.1-8B backbone, with example prompts and generated responses.

Challenges. (1) Loss of structure. Current graph–LLM bridges typically fuse modalities by introducing a fixed number of learnable query tokens, which interact with molecular encoder outputs and are then fed into the language model. (Lee et al., 2025; Kim et al., 2025). However, compressing length-heterogeneous molecules into the fixed set collapses critical features such as stereochemistry and functional groups. This loss carries over to reasoning, where even chain-of-thought prompting provides little gain and can reduce accuracy on property prediction (Kim et al., 2025). As shown in Fig. 2, fixed-length fusion captures key groups in small molecules (N=16 atoms) and supports structure-aware tasks; however, in larger molecules (N=50 atoms), it results in incomplete substructure coverage, information loss, and consequently brittle, chemically unfaithful predictions.

(2) **Heavy fine-tuning.** Most prior systems train the bridge jointly with the LLM, requiring gradient updates to the backbone itself. Such updates hinder generalization: models overfit to narrow datasets, lose robustness across modalities and structural variations, and fail to transfer

Table 1: Computing costs comparison between frozen/unfrozen backbone settings of Mol-Llama.

Connector	r LLM	Trainable	FLOPs/Token	Time/Step
Finetune	Finetune	8.1B	4.9e10	0.93
Finetune	Frozen	84M	1.7e10	0.23

alignment when scaled to larger frozen backbones. Clearly, they are also computationally inefficient: as shown in Tab. 1, jointly tuning the LLM and connector requires $\times 96$ more trainable parameters than tuning the connector alone (details are in App. D.3). Consequently, a connector-only bridge with frozen encoders and a frozen LLM is urgently needed to achieve scalable, low-cost deployment.

Our approach. We propose EDT-Former, an Entropy-guided Dynamic Token Transformer that aligns chemical graphs with frozen LLMs while preserving structural fidelity. EDT-Former introduces an **Entropy-Guided Patching** strategy that segments molecules into informative sub-groups, generating dynamic query tokens that retain stereochemistry. It further incorporates the **Dynamic Query Transformer** module that integrates these dynamic tokens with learned modality anchors to form a stable cross-modal interface before mapping into the LLM embedding space. Together, these designs enable frozen-backbone alignment, yielding efficient training, robust generalization, and chemically faithful understanding (Fig. 1). Our contributions are summarized as follows:

- This work presents **EDT-Former**, the first connector-only method that aligns chemical graphs with frozen LLMs via dynamic, substructure-aware query tokens.
- We present Entropy-Guided Patching and Dynamic Query Transformer, which together enable efficient cross-modal alignment without updating backbone parameters.
- EDT-Former achieves state-of-the-art results on molecular understanding and property prediction benchmarks, demonstrating both scalability and robust generalization.

2 Related Work

Multimodal Fusion for Molecular Generative Modeling. Beyond SMILES-based models (e.g., KV-PLM (Zeng et al., 2022a) and Mol-Instructions (Fang et al., 2023)), recent work has attempted to fuse SMILES, topological graphs, and 3D conformers with large language models. GIM-LET (Zhao et al., 2023) showed that instruction tuning over paired graph-text corpora can provide LLMs with basic molecular understanding. Larger frameworks extended modality coverage: MoleculeSTM (Liu et al., 2022), MolFM (Luo et al., 2023a), BioT5+ (Pei et al., 2024a), and Pro-LLAMA (Lv et al., 2024) combine discrete structure tokens, biomedical knowledge graphs, and protein sequences, while 3D-MolT5 (Zhao et al., 2024) and UniText-3D (Zhang et al., 2024a) incorporate coarse 3D geometry for conditional generation. Most recently, UniMolT (Zhang et al., 2024b), Mol-LLM (Lee et al., 2025), and Mol-Llama (Kim et al., 2025) adopt the Q-Former mechanism from BLIP-2 (Li et al., 2023) to bridge molecular graphs and LLMs, relying on a fixed-length learnable modality anchor tokens. Such fixed-length fusion compresses stereochemistry and graph substructures, leading to structural information loss. More flexible alignment strategies are needed to retain sub-structural fidelity as molecular complexity and conformer diversity increase.

Molecular Understanding with Multimodal LLMs. Molecular understanding has been explored by combining structural encoders with instruction-tuned LLMs. MolReasoner Zhao et al. (2025), CoT-Mol (Liu et al., 2024) showed that coupling molecular tokens with prompts enables step-wise reaction explanation and property inference. Subsequent efforts introduced structural bias, such as Mol-LLM (Lee et al., 2025) with topology-conditioned attention and ProLLaMA (Lv et al., 2024) with protein–ligand co-representations. However, most approaches still depend on fixed-length anchor token representations or extensive fine-tuning of the LLM, which weakens structural fidelity and makes training inefficient. Recent chemical agents such as ChemCrow (Bran et al., 2023) and RetroGPT (Huang et al., 2023) attempt to add planner feedback but remain constrained by conformer sensitivity and limited scalability. Progress, therefore, requires a multimodal alignment method that preserves modality-specific cues without relying on heavy LLM tuning and adaptation. More related works and relationships to current works are discussed in App. B.

3 EDT-FORMER: ENTROPY-GUIDED DYNAMIC GRAPH-LLM ALIGNMENT

EDT-Former provides efficient and substructure-aware alignment between molecule and LLM. Sec. 3.1 outlines the overall architecture, after which Sec. 3.2 introduces the Entropy-Guided Patching strategy for dynamic token generation from graph node-embeddings, while Sec. 3.3 presents the Dynamic Query Transformer module. Finally, Sec. 3.4 describes the training process under the frozen-backbone regime. The implementation details are discussed in App. C.

163

164

166

167

169

170

171 172

173

174 175 176

177

178

179

181

183

185

186

187

188 189

190

191

192

193

194

195

196

197

199

200201

202

203

204

205

206

207

208209

210

211

212213

214

215

Figure 3: The architecture of EDT-Former. (a) Entropy-based Patching segments node embeddings into patches to produce dynamic query tokens. (b) EDT-Former integrates anchors and dynamic queries through Dynamic Query Transformer to align the molecular graph with the LLM.

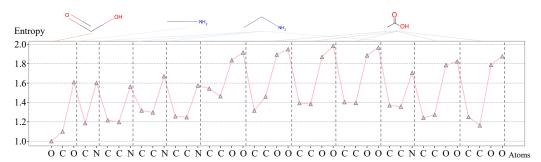


Figure 4: Illustration of entropy-guided patching on an example molecule. Atom-level entropy is plotted along the molecular sequence, and a new patch is initiated after each local maximum.

3.1 STRUCTURE OVERVIEW

EDT-Former establishes alignment between molecular encoders and frozen LLMs through two complementary components: Entropy-Guided Patching and Dynamic Query Transformer. As shown in Fig. 3, node embeddings from frozen graph encoders are processed by entropy-guided patching, which measures node-level uncertainty for language models, segments molecules into informative sub-groups, and pools them into dynamic tokens that preserve local graph features. These tokens are combined with static modality anchor tokens in a shared query bank, where self-attention propagates context and cross-attention retrieves structural evidence from molecular embeddings. The resulting representations are projected into the LLM embedding space, forming a cross-modal alignment interface that balances global consistency from modality anchors with local fidelity from dynamic query tokens, enabling modality alignment with high efficiency and structural fidelity.

3.2 Entropy-Guided Sub-Graph Patching

Setup. As shown in Fig. 3(a), let a molecule be represented by a SMILES-ordered atom sequence (a_1,\ldots,a_T) . We pre-train a lightweight Next-Atom Predictor f_θ (NAP, a small Transformer) on large SMILES corpora to model $p(a_{t+1}\mid a_{1:t})$. At inference time, given logits $L_t\in\mathbb{R}^V$ over the atom vocabulary $\mathcal V$, the probability assigned to the ground-truth next atom is

$$p_t = \text{softmax}(L_t)[a_{t+1}], \qquad t = 1, \dots, T - 1.$$
 (1)

The information content is defined (negative log-likelihood) as

$$e_t = -\log p_t$$
 (nats; use \log_2 for bits if desired). (2)

Peak-based segmentation. Instead of thresholding, we identify local maxima of the entropy signal $\{e_t\}$ and place split points at peaks (shown in Fig. 4):

$$\operatorname{peak}(t) = \left[(e_{t-1} < e_t) \land (e_t > e_{t+1}) \right],$$
 (3)

optionally enforcing a minimal separation Δ (non-maximum suppression) and a small prominence γ to remove spurious bumps:

$$NMS(t) = \left[\min_{|s-t| \le \Delta} e_t \ge e_s \right], \quad prom(t) = \left[e_t - \frac{1}{2} (e_{t-1} + e_{t+1}) \ge \gamma \right]. \quad (4)$$

```
216
           Algorithm 1: Dynamic Query Transformer (anchors + dynamic tokens)
217
           Input: Graph G; frozen graph encoder \mathcal{E}; frozen LLM \mathcal{L}; dynamic tokens Z = [z_1, \dots, z_M]
218
                       from Sec. 3.2; #anchors k; layers L; projector W_{\text{proj}}.
219
           Output: U \in \mathbb{R}^{(k+M) \times d_{\text{LLM}}} (LLM-conditioning sequence).
220
           Node embeddings: X \leftarrow \mathcal{E}(G);
                                                                                  // frozen molecular encoder
221
           Anchors: Q_{\text{fix}} \in \mathbb{R}^{k \times d}; // learnable, modality-stable Query bank: Q_{bank} \leftarrow [Q_{\text{fix}}; Z] \in \mathbb{R}^{(k+M) \times d}; // construct query bank
222
224
           Multimodal Alignment (L layers): ;
                                                                            // Stack of Transformer Blocks
225
           for \ell \leftarrow 1 to L do
                \begin{array}{l} Q_{bank} \leftarrow Q_{bank} + \text{Self-Attn}(Q_{bank}) \; ; \qquad // \; \text{mix anchors} \; \leftrightarrow \; \text{dynamics} \\ Q_{bank} \leftarrow Q_{bank} + \text{Cross-Attn}(Q_{bank}, K\!=\!X, V\!=\!X) \; ; \quad // \; \text{retrieve node-level} \end{array}
226
227
228
             Q_{bank} \leftarrow Q_{bank} + \text{FFN}(Q_{bank});
                                                               // shared FFN for all query tokens
229
230
           Projection to LLM space: U \leftarrow W_{\text{proj}} Q_{bank};
                                             231
```

Then, the valid split indices are collected as $\mathcal{T}_{\star} = \{t \in \{2, \dots, T-2\} \mid \text{peak}(t) \land \text{NMS}(t) \land \text{prom}(t)\}$ and form dynamic segments by cutting after each peak:

$$1 = \tau_0 < \tau_1 < \dots < \tau_M = T, \quad \{\tau_1, \dots, \tau_{M-1}\} = \mathcal{T}_{\star}, \quad \mathcal{S}_k = \{t \mid \tau_{k-1} \le t < \tau_k\}.$$
 (5)

Graph mapping and pooling. Let $\pi: \{1, \dots, T\} \to \{1, \dots, N\}$ map SMILES positions to graph node indices. Each segment \mathcal{S}_k induces a substructure (node set)

$$\widehat{\mathcal{S}}_k = \{ \pi(t) : t \in \mathcal{S}_k \} \subseteq \{1, \dots, N\}.$$
(6)

Given frozen graph-encoder node embeddings $X \in \mathbb{R}^{N \times d}$, we construct one query token for each substructure via average pooling:

$$z_k = \operatorname{AvgPool}(\{X_i : i \in \widehat{\mathcal{S}}_k\}) \in \mathbb{R}^d, \quad k = 1, \dots, M.$$
 (7)

The collection $\{z_k\}_{k=1}^M$ constitutes dynamic query tokens for Dynamic Query Transformer.

Rationale. The entropy e_t quantifies the predictive uncertainty of extending the current SMILES fragment with a_{t+1} . Local peaks concentrate at structural transitions (e.g., branch entries/exits, ring closures, junctions of functional groups) where multiple chemically plausible continuations exist. Peak-based cuts thus yield data-driven substructures without hand-crafted rules, while averaging the corresponding node embeddings preserves locality (bonding patterns, stereochemical context) for structure-aware alignment. This purely entropy-driven segmentation is simple and reproducible, matching dynamic-token count to molecular complexity (variable-sized molecules \Rightarrow dynamic tokens) and improving fidelity under a frozen-backbone regime. The theoretical analysis of the rationale of Entropy-Guided Patching is discussed in App. A.1.

3.3 Dynamic Query Transformer

Dynamic Query Transformer integrates fixed modality anchors with dynamic substructure tokens to form a compact, structure-aware interface to a frozen LLM, as shown in Fig. 3(b). Let $X \in \mathbb{R}^{N \times d}$ be node embeddings from the frozen graph encoder and let $Z = [z_1, \ldots, z_M] \in \mathbb{R}^{M \times d}$ be dynamic tokens from Sec. 3.2. As shown in Algorithm 1, we initialize k learnable anchors $Q_{\text{fix}} \in \mathbb{R}^{k \times d}$ and build a query bank $Q_{bank} = [Q_{\text{fix}}; Z] \in \mathbb{R}^{(k+M) \times d}$. A lightweight transformer with L layers refines Q_{bank} by (i) self-attention over queries to mix global and local context, (ii) cross-attention from queries to node embeddings X to retrieve substructure evidence, and (iii) a shared feed-forward network. The enriched queries are projected to the LLM embedding space via $U = W_{\text{proj}} Q_{bank} \in \mathbb{R}^{(k+M) \times d_{\text{LLM}}}$ and then consumed by the frozen LLM as conditioning. During training, only the bridge parameters (anchors, attention/FFN, and W_{proj}) are updated; both the molecular encoder and the LLM remain frozen, yielding an efficient graph–LLM alignment.

Table 2: Accuracy (%) on zero-shot molecular property prediction benchmarks (Pampa and BBBP). Models are evaluated with three prompting strategies: Direct, Reasoning, and Rich Instructions (see App. D.3). The best (pink) and second-best (lightpink) results are highlighted.

Models	Size		Pampa				BBBP			
Models	SILC	Direct	Reasoning	RichInst.	Avg.	Direct	Reasoning	RichInst.	Avg.	
General LLMs										
GPT-4o	-	48.65	58.23	47.17	51.35	60.82	61.34	64.43	62.20	
Llama2	7B	57.14	57.53	84.52	66.40	37.37	51.56	53.09	47.34	
Llama3.1	8B	56.51	46.19	63.64	55.45	57.07	51.03	55.15	54.42	
Molecular LLMs										
Mol-InstLlama2	7B	49.63	31.16	38.18	39.66	52.58	52.58	51.34	52.17	
Mol-LLaMA-2	7.2B	75.68	79.61	67.90	74.40	53.37	52.58	52.58	52.84	
Mol-InstLlama3.1	8B	55.91	33.50	70.47	53.29	53.44	55.31	54.91	54.55	
3D-MoLM	8B	46.93	50.00	64.86	53.93	49.14	51.65	51.91	50.90	
LLaMo	8B	49.25	64.37	48.51	54.04	55.44	55.45	56.91	55.93	
Mol-LLaMA3.1	8.2B	63.55	64.37	72.48	66.80	59.54	55.56	59.08	58.06	
Ours										
EDT-Former	8.3B	81.57	81.57	83.78	82.31	74.44	74.69	75.86	75.00	

3.4 Frozen-Backbone Alignment Training

EDT-Former is trained via a two-stage regimen (pretraining, then alignment tuning) on the Mol-LLaMA-Instruct dataset (Kim et al., 2025). In pretraining, only the Dynamic Query Transformer is optimized together with frozen graph encoders, without attaching the LLM. Objectives follow standard adapter-style practice (cf. Q-Former): a cross-modal contrastive loss for global cohesion, an anchor-modality matching loss to stabilize the fixed interface, and a masked substructure reconstruction loss to inject local chemical semantics. Details are discussed in App. C.2.

In the alignment tuning stage, the frozen LLM is attached while the graph encoders and LLM remain frozen; only the parameters of Dynamic Query Transformer are updated to align structure-aware queries with instruction prompts. This protocol preserves backbone stability, enables efficient adaptation, and yields a chemically faithful graph–LLM interface. Detail are in App. C.3

4 EXPERIMENTS

In this section, EDT-Former is evaluated on benchmarks—MoleculeQA (Lu et al., 2024), Mol-Instructions (Fang et al., 2023), Pampa, and BBBP from TDC Huang et al. (2021), addressing the following questions: Q1) Does EDT-Former outperform existing multimodal molecular LLM baselines? Q2) Can Entropy-Guided Patching capture subgraph features? Q3) Does the Dynamic Query Transformer improve LLM understanding compared to fixed-length query connectors? Q4) Does freezing the LLM reduce compute while retaining strong accuracy? See App. D for baseline and dataset descriptions, comprehensive experiment implementation, benchmark and evaluation details.

4.1 PROPERTY PREDICTION (TEXT-BASED)

As shown in Tab. 2, we evaluate on the Pampa and BBBP benchmarks (Huang et al., 2021) using official splits and zero-shot inference. Comparisons cover general LLMs and molecular LLMs (including multimodal variants). To reduce prompt sensitivity, the average accuracy over three regimes (Direct, Reasoning, Rich Instructions is reported, details are in App. D.3). Under identical conditions, EDT-Former attains the best results on both tasks, with >20% relative gains over the strongest baseline and pushing average accuracy above 75% on each dataset, indicating practical viability for LLM-based property prediction (Q1). Potential dataset imbalance is discussed in App. D.6.

4.2 MOLECULEQA (REASONING AND UNDERSTANDING)

Evaluation is conducted on the large-scale MoleculeQA benchmark (Lu et al., 2024) (Tab. 3), a multiple-choice suite covering four tasks—structure, source, property, and application. Large GPT models are assessed in a 10-shot setting, while all other models are fine-tuned on the official splits for the same number of epochs under matched hyperparameters. Under these conditions, our SFT model consistently achieves the best performance across all four tasks, demonstrating strong molec-

Table 3: Performance on the MoleculeQA benchmark. Models are compared across four tasks (Structure, Source, Property, Application) with accuracy (%) reported. The best (pink) and second-best (lightpink) results are highlighted. Modality T/G/3D represents text, graph, and 3D geometry.

Model	Modality	Imp.	Size	Strct.	Src.	Prop.	App.	Avg.	Total
Random									
Random	-	-	-	24.41	22.30	23.04	24.57	23.58	24.03
Molecular LLMs									
BioMedGPT-LM	Text	SFT	7B	54.19	60.01	38.85	40.90	48.49	52.23
Mol-InstLlama3.1	T/G	SFT	8B	72.79	70.82	43.08	41.22	56.98	65.31
3D-MoLM	T/G/3D	SFT	8B	73.17	70.50	44.79	44.19	58.16	65.96
LLaMo	Text	SFT	8B	70.56	66.63	44.60	45.18	56.74	63.74
Mol-Llama3.1	Text/Graph/3D	SFT	8.2B	73.16	70.22	45.70	46.18	58.82	66.21
General LLMs									
Galactica	Text	SFT	6.7B	32.35	41.92	31.05	28.21	33.38	33.96
BLOOM	Text	SFT	7.1B	35.01	47.51	31.46	33.56	36.89	37.31
Pythia	Text	SFT	6.9B	42.79	58.90	38.58	39.07	44.84	45.61
Llama-2-chat	Text	SFT	7B	28.75	39.84	31.33	27.71	31.91	31.54
Vicuna-v1.5	Text	SFT	13B	37.01	43.19	30.64	31.55	35.60	37.07
Large Scale General I	LLMs								
GPT-3.5	Text	10-Shot	-	25.60	37.60	28.04	32.22	30.87	29.29
GPT-4	Text	10-Shot	-	60.94	50.19	35.57	43.91	47.65	53.47
GPT-5	Text	10-Shot	-	62.78	53.22	36.42	46.91	49.83	56.12
Ours									
EDT-Former	T/G/3D	10-Shot	8.3B	66.46	60.98	40.35	36.40	51.05	58.78
EDT-Former	T/G/3D	SFT	8.3B	74.55	72.39	50.71	48.58	61.56	68.34

Table 4: Results on the Mol-Instructions dataset. Models are finetuned and evaluated on two tasks: molecular description generation (BLEU, ROUGE, and METEOR scores) and molecular property prediction (MAE). The best (pink) and second-best (lightpink) results are highlighted.

Models	Molecular Description						Property (MAE)
Wideis	BLUE-2	BLUE-4	ROUGE-1	ROUGE-2	ROUGE-L	METEOR	Troperty (WILL)
Alpaca-7B	0.068	0.014	0.178	0.041	0.136	0.107	322.109
Baize-7B	0.064	0.015	0.189	0.053	0.148	0.106	261.343
LLaMA-2-7B	0.059	0.014	0.164	0.066	0.148	0.184	5.553
Vicuna-v1.5-13B	0.052	0.011	0.151	0.055	0.130	0.168	860.051
Galatica-6.7B	0.024	0.008	0.074	0.015	0.063	0.065	0.568
Qwen3-8B	0.098	0.029	0.207	0.050	0.157	0.173	4.737
Mol-InstLlama2-7B	0.217	0.143	0.337	0.196	0.291	0.254	0.013
Mol-InstLlama3.1-8B	0.419	0.361	0.719	0.646	0.709	0.637	15.059
Mol-LLaMA2-7.2B	0.433	0.385	0.711	0.649	0.601	0.601	0.0087
Mol-LLaMA3.1-8.2B	0.445	0.398	0.717	0.656	0.709	0.617	0.0079
MolReasoner-7B	0.438	0.322	0.553	0.366	0.482	0.475	10.323
EDT-Former-8.3B	0.424	0.402	0.726	0.652	0.717	0.631	0.0062

ular reasoning. Notably, the 10-shot variant of EDT-Former outperforms the newest GPT-5 model, indicating an efficient trade-off between scale and domain alignment. These results answer **Q1**.

4.3 Mol-Instructions (Understanding)

For benchmark Mol-Instructions (Fang et al., 2023), we evaluate the molecular captioning and property tasks. For fair comparison, all baselines and EDT-Former are fine-tuned under identical settings on the official splits. As shown in Tab. 4, EDT-Former attains stronger captioning ability on most metrics (defined in App. D.3) and lower property Mean Absolute Error (MAE), indicating improved performance relative to both general and molecular LLM baselines. These results address **Q1**.

4.4 ABLATION STUDIES

Comprehensive ablation experiments are conducted for question **Q2–Q4**, ranging from multimodal fusion to graph patching strategies for Entropy-Guided Patching and connector design for Dynamic Query Transformer. All experiments use official splits with matched token budgets and hyperparameter settings, see App. D. Additional ablation studies are in App. E.

None

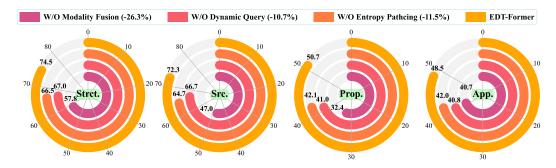


Figure 5: Ablation study of components on the MoleculeQA dataset. Accuracy is reported across four task types (Structure, Source, Property, and Application) when removing each component (modality fusion, Entropy-Guided Patching, or Dynamic Query Transformer).

Table 5: Ablations on patching strategies. Four methods are tested on the BBBP and Pampa. Accuracy (%) and F1 (subscript) are reported, with the top and second best highlighted.

39.6749.58

1				
Methods	BBBP	Pampa	Avg. Drop	L
Entropy	75.0675.06	84.5291.61	0%	M
BRICS	73.5984.74	71.9083.09	3.96%	Q
Fix-Len	69.0080.93	67.57 79.82	8.87%	Ll

78.6287.90

21.60%

Table 6: Effect of EDT-Former on different LLM backbones. Four novel LLMs are tested on the MoleculeQA benchmark. Total accuracy and average gain for each model are reported.

LLMs	LLM Only	+EDT-Former	Avg. Gain
Mistral-8B	28.00	55.63	+98.71%
Qwen3-8B	27.18	54.79	+101.58%
Llama2-7B	38.37	65.89	+71.72%
Llama3.1-8B	49.31	68.34	+38.59%

Component Effects. Multimodal fusion and each core component are ablated on the four MoleculeQA tasks (Fig. 5). The full EDT-Former performs best across Structure, Source, Property, and Application. Removing multimodal fusion produces the largest degradation (26% avg.), confirming that combining text, graph, and 3D is critical for LLM comprehension. Disabling either Dynamic Query Transformer or Entropy-Guided Patching yields more than 10% average drops (10.7% and 11.5%), indicating that dynamic queries are necessary to handle variable-sized graphs and capture more evidence than fixed-length fusion, while entropy-driven patches expose LLM-salient subgraphs for stronger structure-aware understanding (Q2, Q3).

Patching Strategies. We compare four segmentation methods on PAMPA/BBBP with matched splits, prompts, and query budgets: Entropy-Guided Patching (entropy peaks), BRICS-based fragments (Jinsong et al., 2024), fixed-length patches (k=8), and no patching. As shown in Tab. 5, entropy-guided patching yields the strongest results, BRICS trails slightly (chemically sound segmentation) while fixed-length patching degrades more (ignores true subgraph boundaries), and removing patching produces the largest drop—confirming that learned, data-driven boundaries are essential. Entropy peaks from a small next-atom transformer mark LLM-salient transitions, exposing true substructure boundaries to the query interface; this data-driven segmentation likely generalizes beyond molecules and is necessary to realize multimodal fusion gains ($\mathbf{Q2}$).

Backbone Sensitivity. On MoleculeQA with official splits, attaching our connector to Mistral, Qwen3, and Llama series models yields consistent improvements over the text-only backbones (Tab. 6). Gains are larger for weaker backbones but remain substantial across all four.

Efficiency Analysis. As shown in Tab. 7, relative to LoRA finetuning, our strategy halves GPU memory and achieves $\sim 3.5 \times$ faster training per step under identical settings; full backbone tuning is practically infeasible in our setting. These support that efficient alignment is achievable without updating the LLM (Q4).

Table 7: Estimated memory usage and training time per step for EDT-Former with Llama3.1-8B.

EDT-Former	Llama3.1-8B	Mem (GB)	Time/Step (s)
Train	Frozen	37	0.26
Train	LoRA	77	0.93
Train	Train	>200	_
	<u> </u>		<u> </u>

Fusion Benefits. In Tab. 8, four modality settings are evaluated on MoleculeQA. Text-only (SMILES) performs worst, underscoring the need for explicit structural features. Removing 3D

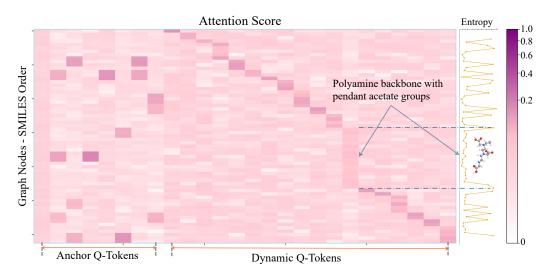


Figure 6: Attention visualization for EDT-Former. Last-layer attention map for a single molecule (N=63 atoms) with atoms in SMILES order along the vertical axis and query tokens (fixed-length anchors followed by dynamic tokens) along the horizontal axis. The right panel shows the corresponding entropy signal for the SMILES sequence. Values indicate normalized attention scores.

Table 8: Modality ablation on MoleculeQA. Four modality combinations are tested on four tasks, with the accuracies (%) reported. The top (pink) and second-best (lightpink) results are highlighted.

Modality	Structure	Source	Property	Application	Average	Total	Avg. Drop
Text/Graph/3D	74.55	72.39	50.71	48.58	61.56	68.34	0.00%
Text/Graph	70.84	68.72	49.22	46.51	58.82	65.11	4.73%
Text/3D	66.50	64.77	42.17	42.02	53.87	60.48	11.50%
Text Only	66.46	60.98	40.35	36.40	51.05	58.77	14.00%

causes a modest decline relative to full fusion, whereas removing 2D (graph) degrades substantially, indicating that graph information is essential while 3D provides complementary gains.

Attention Insights. The attention visualization in Fig. 6 reveals the complementary roles for the two query types. Fixed-length tokens behave as modality anchors: their attention is diffuse and often spans multiple, unrelated patches, which is useful for global alignment but unreliable for isolating substructures. In contrast, dynamic tokens align with the SMILES-ordered entropy segments and attend sharply to contiguous node ranges that correspond to chemically meaningful subgraphs (e.g., the polyamine backbone with pendant acetate groups highlighted). These patterns substantiate **Q2** and **Q3**: multimodal fusion benefits when local structure is exposed, and the proposed novelties (Entropy-Guided Patching plus the Dynamic Query Transformer) provide a substructure-aware interface that preserves local fidelity while anchors maintain global consistency.

5 CONCLUSION

EDT-Former is introduced as a method for aligning molecular graphs with large language models. Entropy-Guided Patching segments molecules at entropy peaks to yield substructure-aware dynamic tokens, and Dynamic Query Transformer integrates these tokens with modality anchors to project into the LLM space, providing a stable interface without updating backbone parameters. Across standard benchmarks, EDT-Former delivers consistent gains over text-only and prior multimodal baselines while markedly reducing training cost relative to LLM tuning. Overall, EDT-Former offers chemically faithful, efficient graph—LLM alignment and a general recipe for multimodal fusion that can extend to broader graph domains; code and processed data are released for reproducibility, and for current constraints and future directions, see Apps. F.1 and F.2.

6 ETHICS STATEMENT

This work does not involve human subjects, personally identifiable information, or non-public/proprietary datasets. All data are publicly available under their respective licenses, and no wet-lab or animal experiments were conducted. Although molecular AI can present dual-use risks (e.g., aiding the design of harmful compounds), our contribution focuses on multimodal representation learning and evaluation on standard benchmarks; we neither train nor release models intended to generate or optimize novel toxic molecules. We will release code and models with terms of use that prohibit harmful or unlawful applications. All authors complied with institutional policies and responsible research practices.

7 REPRODUCIBILITY STATEMENT

This work is fully open-source under the MIT license. We release end-to-end resources to reproduce every result in the paper, including: (i) complete training and evaluation code (with all ablation scripts), (ii) baseline and benchmark runners with fixed seeds, (iii) processed pretraining and fine-tuning datasets with scripts to rebuild them from the original sources, (iv) model checkpoints, and (v) Dockerfile for bitwise-reproducible environments. A single-entry script reproduces all main tables and even the figures; expected variances from stochastic training are reported in the repository. No external services are required at evaluation time (offline scoring only). The repository is actively maintained and will serve as the foundation for subsequent extensions. All materials are available via the anonymous GitHub link: https://anonymous.4open.science/r/EDT-Former-844D.

As part of reproducibility, App. C.1 and App. C.2 detail the implementation of our two novelties and the pre-training settings. The fine-tuning settings and parameters are provided in App. C.3. As well, App. D.2 documents the data processing for each dataset. Finally, App. D.3 reports the hyperparameter choices and model configurations for each benchmark.

REFERENCES

Walid Ahmad, Elana Simon, Seyone Chithrananda, Gabriel Grand, and Bharath Ramsundar. Chemberta-2: Towards chemical foundation models. In *arXiv.org*, 2022. doi: 10.48550/arXiv. 2209.01712.

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, A. Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, O. Vinyals, Andrew Zisserman, and K. Simonyan. Flamingo: a visual language model for few-shot learning. In *Neural Information Processing Systems*, 2022.

Stella Biderman, Hailey Schoelkopf, Quentin G. Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, 2023. doi: 10.48550/arXiv.2304.01373.

Andrés M Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D. White, and P. Schwaller. Augmenting large language models with chemistry tools. In *Nat. Mac. Intell.*, 2023. doi: 10. 1038/s42256-024-00832-8.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, J. Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, T. Henighan, R. Child, A. Ramesh, Daniel M. Ziegler, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Ma teusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, I. Sutskever, and Dario Amodei. Language models are few-shot learners. In Neural Information Processing Systems, 2020.

Wenliang Dai, Junnan Li, Dongxu Li, A. M. H. Tiong, Junqi Zhao, Weisheng Wang, Boyang Albert Li, Pascale Fung, and Steven C. H. Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. In *Neural Information Processing Systems*, 2023. doi: 10.48550/arXiv.2305.06500.

Yin Fang, Xiaozhuan Liang, Ningyu Zhang, Kangwei Liu, Rui Huang, Zhuo Chen, Xiaohui Fan, and Huajun Chen. Mol-instructions: A large-scale biomolecular instruction dataset for large language models. In *International Conference on Learning Representations*, 2023. doi: 10.48550/arXiv. 2306.08018.

- Thomas A Halgren. Merck molecular force field. i. basis, form, scope, parameterization, and performance of mmff94. *Journal of computational chemistry*, 17(5-6):490–519, 1996.
- Kexin Huang, Tianfan Fu, Wenhao Gao, Yue Zhao, Yusuf Roohani, Jure Leskovec, Connor W Coley, Cao Xiao, Jimeng Sun, and Marinka Zitnik. Therapeutics data commons: Machine learning datasets and tasks for drug discovery and development. *arXiv* preprint arXiv:2102.09548, 2021.
- Yuyang Huang, Wenbo Hu, Boxuan Wang, Yibo Zhang, Hong Lin, Yi Zeng, and Weinan Wang. Retrogpt: Self-correcting retrosynthesis via guided chain-of-thought. *arXiv* preprint *arXiv*:2311.16401, 2023.
- Yunhui Jang, Jaehyung Kim, and Sungsoo Ahn. Chain-of-thoughts for molecular understanding. In *arXiv.org*, 2024. doi: 10.48550/arXiv.2410.05610.
- Xiaohong Ji, Zhen Wang, Zhifeng Gao, Hang Zheng, Linfeng Zhang, Guolin Ke, and E. Weinan. Uni-mol2: Exploring molecular pretraining model at scale. In *Neural Information Processing Systems*, 2024. doi: 10.48550/arXiv.2406.14969.
- Shao Jinsong, Jia Qifeng, Chen Xing, Yajie Hao, and Li Wang. Molecular fragmentation as a crucial step in the ai-based drug development pathway. *Communications Chemistry*, 7(1):20, 2024.
- Dongki Kim, Wonbin Lee, and Sung Ju Hwang. Mol-llama: Towards general understanding of molecules in large molecular language model. In *arXiv.org*, 2025. doi: 10.48550/arXiv.2502. 13449.
- Greg Landrum. Rdkit documentation. Release, 1(1-79):4, 2013.
- Chanhui Lee, Yuheon Song, Yongjun Jeong, Hanbum Ko, Rodrigo Hormazabal, Sehui Han, Kyunghoon Bae, Sungbin Lim, and Sungwoong Kim. Mol-llm: Generalist molecular llm with improved graph utilization. *ArXiv*, abs/2502.02810, 2025.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip-2: Bootstrapping languageimage pre-training with frozen image encoders and large language models. arXiv preprint arXiv:2301.12597, 2023.
- Sihang Li, Zhiyuan Liu, Yancheng Luo, Xiang Wang, Xiangnan He, Kenji Kawaguchi, Tat-Seng Chua, and Qi Tian. Towards 3d molecule-text interpretation in language models. In *International Conference on Learning Representations*, 2024. doi: 10.48550/arXiv.2401.13923.
- Shengchao Liu, Weili Nie, Chengpeng Wang, Jiarui Lu, Zhuoran Qiao, Ling Liu, Jian Tang, Chaowei Xiao, and Anima Anandkumar. Multi-modal molecule structure-text model for text-based retrieval and editing. In *Nat. Mac. Intell.*, 2022. doi: 10.48550/arXiv.2212.10789.
- Shengchao Liu, Weili Nie, Chengpeng Wang, Jiarui Lu, Zhuoran Qiao, Ling Liu, Jian Tang, Chaowei Xiao, and Animashree Anandkumar. Multi-modal molecule structure–text model for text-based retrieval and editing. *Nature Machine Intelligence*, 5(12):1447–1457, 2023.
- Yujie Liu, Tianlong Wu, Haotian Jiang, Changchang Lu, Pengfei Li, Jiarui Liu, Tong Zhang, Mingli Song, and Philip S Yu. Cot-mol: Chain-of-thought prompting for molecular property prediction. *arXiv* preprint arXiv:2402.01498, 2024.
- Xingyu Lu, He Cao, Zijing Liu, Shengyuan Bai, Leqing Chen, Yuan Yao, Hai-Tao Zheng, and Yu Li. Moleculeqa: A dataset to evaluate factual accuracy in molecular comprehension. In *Conference on Empirical Methods in Natural Language Processing*, 2024. doi: 10.48550/arXiv.2403.08192.

- Yi Luo, Kai Yang, Massimo Hong, Xingyi Liu, and Zaiqing Nie. Molfm: A multimodal molecular foundation model. In *arXiv.org*, 2023a.
 - Yi Luo, Jiahuan Zhang, Siqi Fan, Kai Yang, Yushuai Wu, Mu Qiao, and Zaiqing Nie. Biomedgpt: Open multimodal generative pre-trained transformer for biomedicine. In *arXiv.org*, 2023b. doi: 10.48550/arXiv.2308.09442.
 - Liuzhenghao Lv, Zongying Lin, Hao Li, Yuyang Liu, Jiaxi Cui, Calvin Yu-Chian Chen, Li Yuan, and Yonghong Tian. Prollama: A protein language model for multi-task protein language processing. In *arXiv.org*, 2024.
 - Liuzhenghao Lv, Zongying Lin, Hao Li, Yuyang Liu, Jiaxi Cui, Calvin Yu-Chian Chen, Li Yuan, and Yonghong Tian. Prollama: A protein large language model for multi-task protein language processing. *IEEE Transactions on Artificial Intelligence*, 2025.
 - MetaAI. The llama 3 herd of models. In arXiv.org, 2024. doi: 10.48550/arXiv.2407.21783.
 - Mistral AI. Ministral-8b-instruct-2410. https://huggingface.co/mistralai/Ministral-8B-Instruct-2410, 2024. Model card, accessed 2025-09-23.
 - OpenAI. Gpt-4 technical report, 2023. https://openai.com/research/gpt-4.
 - OpenAI. Introducing gpt-5. OpenAI Blog, 2025. URL https://openai.com/gpt-5/. Accessed August 7, 2025.
 - Artidoro Pagnoni, Ram Pasunuru, Pedro Rodriguez, John Nguyen, Benjamin Muller, Margaret Li, Chunting Zhou, Lili Yu, Jason Weston, Luke Zettlemoyer, Gargi Ghosh, Mike Lewis, Ari Holtzman, and Srinivasan Iyer. Byte latent transformer: Patches scale better than tokens. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (ACL)*, 2025. URL https://aclanthology.org/2025.acl-long.453/.
 - Jinyoung Park, Minseong Bae, Dohwan Ko, and Hyunwoo J. Kim. Llamo: Large language model-based molecular graph assistant. In *Neural Information Processing Systems*, 2024. doi: 10.48550/arXiv.2411.00871.
 - Qizhi Pei, Lijun Wu, Kaiyuan Gao, Xiaozhuan Liang, Yin Fang, Jinhua Zhu, Shufang Xie, Tao Qin, and Rui Yan. Biot5+: Towards generalized biological understanding with iupac integration and multi-task tuning. In *Annual Meeting of the Association for Computational Linguistics*, 2024a. doi: 10.48550/arXiv.2402.17810.
 - Qizhi Pei, Lijun Wu, Kaiyuan Gao, Jinhua Zhu, and Rui Yan. 3d-molt5: Leveraging discrete structural information for molecule-text modeling. In *ICLR* 2025, 2024b.
 - Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. Alpaca: A strong, replicable instruction-following model. *Stanford Center for Research on Foundation Models. https://crfm. stanford. edu/2023/03/13/alpaca. html*, 3(6):7, 2023.
 - Ross Taylor, Vitalii Borisov, Ronan Le Bras, et al. Galactica: A large language model for science. *arXiv preprint arXiv:2211.09085*, 2022.
 - Hugo Touvron, Thibaut Lavril, Gautier Izacard, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
 - Hugo Touvron et al. Llama 2: Open foundation and fine-tuned chat models. In arXiv.org, 2023b.
- BigScience Workshop. Bloom: A 176b-parameter open-access multilingual language model. In *arXiv.org*, 2022. doi: 10.48550/arXiv.2211.05100.
 - Fang Wu, Dragomir Radev, and Stan Z Li. Molformer: Motif-based transformer on 3d heterogeneous molecular graphs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 5312–5320, 2023.

- Yingce Xia, Peiran Jin, Shufang Xie, Liang He, Chuan Cao, Renqian Luo, Guoqing Liu, Yue Wang, Zequn Liu, Yuan-Jyue Chen, et al. Naturelm: Deciphering the language of nature for scientific discovery. *arXiv e-prints*, pp. arXiv–2502, 2025.
- Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. Baize: An open-source chat model with parameter-efficient tuning on self-chat data. *arXiv preprint arXiv:2304.01196*, 2023.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- Zheni Zeng, Yuan Yao, Zhiyuan Liu, and Maosong Sun. A deep-learning system bridging molecule structure and biomedical text with comprehension comparable to human professionals. In *Nature Communications*, 2022a. doi: 10.1038/s41467-022-28494-3.
- Zheni Zeng, Yuan Yao, Zhiyuan Liu, and Maosong Sun. A deep-learning system bridging molecule structure and biomedical text with comprehension comparable to human professionals. *Nature communications*, 13(1):862, 2022b.
- Juzheng Zhang, Yatao Bian, Yongqiang Chen, and Quanming Yao. Unimot: Unified molecule-text language model with discrete token representation. In *arXiv.org*, 2024a. doi: 10.48550/arXiv. 2408.00863.
- Juzheng Zhang, Yatao Bian, Yongqiang Chen, and Quanming Yao. Unimot: Unified molecule-text language model with discrete token representation. In *arXiv.org*, 2024b. doi: 10.48550/arXiv. 2408.00863.
- Guojiang Zhao, Sihang Li, Zixiang Lu, Zheng Cheng, Haitao Lin, Lirong Wu, Hanchen Xia, Hengxing Cai, Wentao Guo, Hongshuai Wang, Mingjun Xu, Siyu Zhu, Guolin Ke, Linfeng Zhang, and Zhifeng Gao. Molreasoner: Toward effective and interpretable reasoning for molecular llms. In *arXiv.org*, 2025.
- Haiteng Zhao, Shengchao Liu, Chang Ma, Hannan Xu, Jie Fu, Zhihong Deng, Lingpeng Kong, and Qi Liu. Gimlet: A unified graph-text model for instruction-based molecule zero-shot learning. In *bioRxiv*, 2023. doi: 10.1101/2023.05.30.542904.
- Zihan Zhao, Da Ma, Lu Chen, Liangtai Sun, Zihao Li, Hongshen Xu, Zichen Zhu, Su Zhu, Shuai Fan, Guodong Shen, Xin Chen, and Kai Yu. Chemdfm: A large language foundation model for chemistry. In *arXiv.org*, 2024.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, E. Xing, Haotong Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena. *ArXiv*, abs/2306.05685, 2023.
- G. Zhou, Zhifeng Gao, Qiankun Ding, Hang Zheng, Hongteng Xu, Zhewei Wei, Linfeng Zhang, and Guolin Ke. Uni-mol: A universal 3d molecular representation learning framework. In *International Conference on Learning Representations*, 2023.

APPENDIX TABLE OF CONTENTS

· A. Theoretical Analysis and Proofs

- A.1 Entropy Patching Analysis
- A.2 Dynamic Query Analysis

• B. Relationship with Previous Methods

- B.1 From Representation to Multimodal Understanding
- B.2 Relationship to Related Works

• C. Implementation Details

- C.1 Entropy Guided Patching
- C.2 Dynamic Query Transformer
- C.3 EDT-Former Settings and Finetuning

• D. Experiments Details

- D.1 Benchmarks and Baselines
- D.2 Data Process
- D.3 Evaluation Settings
- D.4 Data Contamination Analysis
- D.5 Reproducibility

• E. Extended Ablations

- E.1 Impact of Query Token Length
- E.2 Impact of Model Size
- E.3 Analysis of Hyperparameter
- E.4 Impact of Graph Encoders

• F. Limitations and Future Work

- F.1 Limitations
- F.2 Future Work
- G LLM Usage

A THEORETICAL ANALYSIS AND PROOFS

Scope and assumptions. The analyses in App. A.1 and App. A.2 are intended as proof sketches under standard modeling assumptions (calibrated next-token predictors, piecewise-stationary generation along SMILES order, and locally Lipschitz Transformer blocks on bounded inputs). We work with token-level surprisal $e_t = -\log p(a_{t+1} \mid a_{1:t})$ (and note that entropy refers to its expectation). The results formalize why (i) segmentation at surprisal peaks aligns with change points and minimizes a budgeted upper bound on pooling/modeling loss, and (ii) anchor-augmented projection yields a stable, low-rank-expressive interface to a frozen LLM. Constants and high-probability terms are suppressed for clarity.

A.1 ENTROPY PATCHING ANALYSIS

Setup. Let a molecule be represented by a SMILES-ordered atom sequence (a_1,\ldots,a_T) generated by a piecewise-stationary conditional distribution $\mathcal{P}^\star(a_{t+1}\mid a_{1:t})$, with unknown change points $1=\kappa_0<\kappa_1<\cdots<\kappa_M=T$ such that for $t\in[\kappa_{m-1},\kappa_m-1]$ the law of $a_{t+1}\mid a_{1:t}$ belongs to a stationary family \mathcal{P}_m . A trained next-atom predictor f_θ outputs a calibrated estimate $p_\theta(a_{t+1}\mid a_{1:t})$. Define the surprisal (negative log-likelihood)

$$e_t = -\log p_{\theta}(a_{t+1} \mid a_{1:t}), \qquad t = 1, \dots, T - 1.$$

Entropy-guided patching places a cut after indices in $T^* = \{t : t \text{ is a local maximum of } e_t\}$, refined by non-maximum suppression (window Δ) and a prominence threshold γ , yielding segments $S_k = \{t : \tau_{k-1} \le t < \tau_k\}$ where $1 = \tau_0 < \cdots < \tau_M = T$ are the retained peaks.

¹Calibrated means $\mathbb{E}[\mathbf{1}\{A=a\} \mid p_{\theta}(A=a \mid \cdot) = u] = u$; consistency is not required, but strengthens the results.

Objective. Given a token budget (implicit in Δ , γ) we seek a segmentation that (i) maximizes the information that segment tokens carry about the underlying structure, and (ii) minimizes within-segment heterogeneity that would be averaged out by pooling into a single token.

Let $X \in \mathbb{R}^{N \times d}$ be frozen node embeddings and let π map SMILES indices to graph nodes. Each segment induces a node set $\widehat{S}_k = \{\pi(t) : t \in S_k\}$ and a token $z_k = \text{AvgPool}(\{X_i : i \in \widehat{S}_k\})$.

Assumption 1 (Piecewise stationarity). Within each true region $[\kappa_{m-1}, \kappa_m)$ the conditional $a_{t+1} | a_{1:t}$ has constant entropy $H_m = \mathbb{E}[e_t | t \in [\kappa_{m-1}, \kappa_m)]$ and adjacent regions satisfy a separation in total variation: $\mathrm{TV}(\mathcal{P}_m, \mathcal{P}_{m+1}) \geq \delta > 0$.

Assumption 2 (Calibration / bounded error). $|\mathbb{E}[e_t] - H^*(a_{t+1} | a_{1:t})| \le \varepsilon$ uniformly in t.

Key quantity. For a candidate segmentation $S = \{S_k\}_{k=1}^M$, define the within-segment entropy

$$\mathcal{E}(\mathcal{S}) \ = \ \sum_{k=1}^{M} \ \frac{1}{|S_k|} \sum_{t \in S_k} \mathbb{E}[e_t].$$

This upper-bounds both (a) the Bayes error of predicting local structural events within segments via Fano-type inequalities, and (b) the information lost by averaging embeddings inside a segment (see Lemma 2).

Lemma 1 (Surprisal peaks mark change points). Under Assumptions 1–2, for any interior index t that lies η away from the nearest true change point, the expected discrete second difference satisfies $\mathbb{E}[e_{t+1} - 2e_t + e_{t-1}] \leq c_1 \varepsilon - c_2 \delta$ near a change point and $\geq -c_1 \varepsilon$ away from it, for constants $c_1, c_2 > 0$. Hence the probability that a local maximum of e_t occurs within $O(\eta)$ of each change point tends to 1 as $\varepsilon \to 0$.

Sketch. At a change point, the one-step law jumps from \mathcal{P}_m to \mathcal{P}_{m+1} , inducing a kink in the log-likelihood path with magnitude controlled by $\mathrm{TV}(\mathcal{P}_m,\mathcal{P}_{m+1})$. Calibrated predictors track this kink up to ε , yielding an excess in e_t locally. Away from changes, stationarity implies e_t is a martingale difference with bounded variance, suppressing spurious peaks in expectation. The discrete curvature argument yields the stated bounds.

Lemma 2 (Pooling loss is entropy-controlled). Let g be any L-Lipschitz readout on sets of node embeddings (the self/cross-attention layer consuming z_k is L-Lipschitz under standard assumptions). Then for each segment S_k ,

$$\mathbb{E} \Big[\big\| g(\{X_i : i \in \widehat{S}_k\}) - g(\{\bar{X}_k\}^{|S_k|}) \big\| \Big] \ \leq \ C \sqrt{\operatorname{Var}(X \mid S_k)} \ \leq \ C' \sqrt{\mathcal{E}(S_k)},$$

where \bar{X}_k is the mean in the segment and the final inequality follows because variation of embeddings within a segment is controlled by variation of the local next-atom law, which is upper-bounded (up to constants) by the average surprisal in that segment.

Proposition 1 (Peak cutting minimizes an upper bound on loss). Among all segmentations with a given minimum spacing Δ (token budget), selecting cuts at the most prominent local maxima of e_t (with NMS window Δ) minimizes $\mathcal{E}(\mathcal{S})$ up to $O(\varepsilon)$, hence minimizes the stated upper bounds on both Bayes error within segments and pooling-induced representation loss.

Sketch. Because H_m is (approximately) constant inside each stationary region, any cut placed within a region increases $\mathcal{E}(\mathcal{S})$ only by splitting a constant-entropy block, while a cut placed at a change point prevents averaging across different entropies $H_m \neq H_{m+1}$ and thus strictly decreases $\mathcal{E}(\mathcal{S})$ by an amount $\Omega(\delta)$ (Assumption 1). Lemma 1 guarantees that the largest local maxima of e_t concentrate at change points, so NMS over a spacing Δ selects (approximately) one cut per change region, which is optimal under a fixed budget. Lemma 2 transfers the bound from \mathcal{E} to the representation loss after pooling.

Corollary 1 (Adaptive token count matches structural complexity). Let K_{γ} be the number of retained peaks after prominence γ and window Δ . Then $\mathbb{E}[K_{\gamma}]$ is non-decreasing in the total variation budget $\sum_{m} \mathrm{TV}(\mathcal{P}_{m}, \mathcal{P}_{m+1})$ and non-increasing in (Δ, γ) , yielding a token count that adapts to molecular complexity while remaining budget-controlled.

Robustness and determinism. The algorithm is deterministic given (Δ, γ) , linear in T for entropy and segmentation, and stable to predictor noise: if $\varepsilon < c\delta$ (Assumption 2), the selected peaks coincide with true change neighborhoods with high probability; otherwise NMS and prominence thresholding suppress spurious bumps so that $\mathcal{E}(\mathcal{S})$ remains within $O(\varepsilon)$ of the optimal budgeted value.

Takeaway. Entropy-guided patching is a principled change-point segmentation that (i) places boundaries where the conditional law changes (surprisal peaks), (ii) provably minimizes an entropy-based upper bound on within-segment error and pooling loss under a token budget, and (iii) adapts the number of dynamic tokens to structural complexity while remaining compute-efficient.

A.2 DYNAMIC QUERY ANALYSIS

Setup. Let $Z=\{z_k\}_{k=1}^K$ be the variable-length substructure tokens from Entropy-Guided Patching, $A=\{a_j\}_{j=1}^Q$ be Q modality anchors (fixed, learnable, or lightly tuned), and let $P\in\mathbb{R}^{d_{\text{in}}\times d_{\text{LLM}}}$ be a connector that linearly projects tokens into the frozen LLM embedding space. The input to the LLM is the concatenated sequence $\widetilde{X}=[PA\ ;\ PZ]\in\mathbb{R}^{(Q+K)\times d_{\text{LLM}}}$. Denote the first Transformer block (self-attention + MLP) of the frozen LLM by $f:\mathbb{R}^{(Q+K)\times d_{\text{LLM}}}\to\mathbb{R}^{(Q+K)\times d_{\text{LLM}}}$, with attention maps parameterized by frozen matrices (W_Q,W_K,W_V) and output projection W_O . Subsequent blocks are identical copies with frozen parameters.

Assumption A1 (Lipschitz block). Each LLM block is L-Lipschitz on a compact input set \mathcal{X} : $||f(X) - f(Y)|| \le L||X - Y||$ (this holds under standard bounded-activation/softmax assumptions).

Assumption A2 (Anchor coverage). Let $S_A = \text{span}\{W_K PA\} \subset \mathbb{R}^{d_k}$ be the key subspace induced by anchors. There exists $\sigma_{\min} > 0$ such that the minimum singular value of $[W_K PA]$ is $\geq \sigma_{\min}$, and $\text{span}\{W_O PA\} = S_A$ as well (anchors expose a stable key/query frame).

Assumption A3 (Bounded projection error). Let $Z^\star \in \mathbb{R}^{K \times d_{\text{LLM}}}$ be an (unknown) ideal embedding that would produce the desired frozen-LLM behavior. The learned connector satisfies $\|PZ - Z^\star\| \leq \varepsilon$ on the data manifold.

Stability via anchors. Consider attention logits in the first block: $S = XW_Q(W_K^\top X^\top)/\sqrt{d_k}$. Anchors contribute structured rows/columns S_{AA} and S_{AZ} that do not depend on Z's length or order. Because A2 ensures a well-conditioned anchor subspace S_A , the Jacobian of attention w.r.t. PZ is bounded, which yields an input-output Lipschitz bound uniform in K.

Lemma 3 (Anchor-conditioned Lipschitz stability). Under A1–A2, there exists $C = C(L, \sigma_{\min}, \|W_Q\|, \|W_K\|, \|W_V\|)$ such that for any two token sets Z_1, Z_2 (possibly different lengths) and the same anchors A,

$$||f([PA; PZ_1]) - f([PA; PZ_2])|| \le C ||PZ_1 - PZ_2||.$$

Hence the per-block perturbation is controlled by the connector-space distance, independent of K.

Sketch. Write the attention output as a composition of (i) bilinear logits $(\cdot)W_Q(W_K^{\top}, \cdot^{\top})$ and (ii) softmax-weighted value mixing. The anchor block PA produces a key/query frame whose spectrum is bounded below by σ_{\min} (A2), preventing ill-conditioning in the softmax sensitivity. Combine with A1 and standard Jacobian bounds for softmax to obtain the stated Lipschitz constant C.

Expressivity vs. LoRA at the input. Linearizing the first block at an operating point $\widetilde{X}_0 = [PA; PZ_0]$ gives $f(\widetilde{X}) \approx f(\widetilde{X}_0) + J_0(\widetilde{X} - \widetilde{X}_0)$, where J_0 is the Jacobian. A rank-r LoRA update at the input embedding $(E \mapsto E + \Delta E, \operatorname{rank}(\Delta E) \leq r)$ induces an output shift $J_0\Delta E$ of rank at most r. If P is factored as $P = UV^{\top}$ with $\operatorname{rank}(P) \leq r$, then PZ can realize an equivalent family of rank-r shifts in the first-block input, matching what a LoRA adapter would add at that boundary. Proposition 2 (Connector-only can emulate low-rank input adaptation). Fix anchors A and linearization point \widetilde{X}_0 . For any rank-r update ΔE to the input embeddings, there exists a rank-r connector $P = UV^{\top}$ and a perturbation ΔZ such that $J_0([0; \Delta E]) = J_0([0; P\Delta Z])$. Hence, in the first-order regime, connector-only training with rank-r P is as expressive as a rank-r LoRA at the input boundary.

Sketch. Let J_0 define a linear map on inputs; match images by choosing $P\Delta Z = \Delta E$ in the column space that J_0 uses. Factor P with rank $\leq r$ and select ΔZ in the corresponding r-dimensional coordinate system to reproduce $J_0\Delta E$.

Information preservation. Let ϕ be the (frozen) LLM map from the first block to the layer used by the cross-attention interface (or to logits for analysis). Define the task-relevant mutual information $I(Z;\phi([PA;PZ]))$. Because f is L-Lipschitz and the connector error is ε (A3), data processing and stability imply only small information loss.

Lemma 4 (Lower bound on retained information). *Under A1–A3 and mild regularity (bounded support and densities), there exists* $\alpha > 0$ *such that*

$$I(Z; \phi([PA; PZ])) \geq I(Z; \phi([PA; Z^*])) - \alpha \varepsilon.$$

Thus, if P approximates the ideal embedding within ε , downstream information loss is $O(\varepsilon)$.

Sketch. Combine Lemma 3 with standard continuity bounds for mutual information under Lipschitz perturbations (via Wasserstein stability or Pinsker-type inequalities). The α constant absorbs Lipschitz and measure-regularity terms.

Budget awareness and interference control. Let Π_A be the orthogonal projector onto span $\{PA\}$ in the LLM embedding space. Adding an orthogonality penalty $\|\Pi_A(PZ)\|^2$ ensures that content tokens PZ occupy a subspace complementary to anchors, reducing cross-attention interference and concentrating gradients on task-relevant directions.

Proposition 3 (Conditioning and interference). If training enforces (approximately) $\Pi_A(PZ) = 0$ and normalizes PA, then the attention Gram matrix on \widetilde{X} has condition number bounded by a constant depending only on anchors. Consequently, gradient norms for P are well-conditioned and do not deteriorate with token count K.

Sketch. Under the orthogonality constraint, the mixed Gram blocks between anchors and content vanish at first order, yielding a block-structured Gram with bounded spectrum on the anchor block (A2) and controlled spectrum on the content block via normalization and Lipschitzness, hence a uniform bound on the condition number.

Takeaways. (i) Stability: anchors provide a well-conditioned key/query frame that makes the frozen LLM locally Lipschitz in the connector space, uniformly in token count (Lemma 3). (ii) Expressivity: a low-rank connector can emulate first-order effects of a low-rank (LoRA) input adapter (Proposition 2), explaining strong performance without tuning the backbone. (iii) Information retention: if P approximates an ideal embedding, the downstream information loss is $O(\varepsilon)$ (Lemma 4). (iv) Good conditioning: orthogonalizing content against anchors controls attention conditioning and prevents interference (Proposition 3).

Overall, Dynamic Query Transformer yields a stable, expressive, and budget-aware interface to frozen LLMs, complementing Entropy-Guided Patching's adaptive tokenization.

B EXTENDED RELATED WORK

In addition to the discussion in Sec. 2, this section provides further analysis of how our approach relates to prior efforts. Specifically, we highlight the differences between our Dynamic Query Transformer and fixed-length Q-Former-style bridges, as well as connections to entropy-related works and recent multimodal molecular language models. This extended review situates our model within the broader landscape and clarifies its novel contributions.

B.1 From Representation to Multimodal Understanding

Representation models vs. language models. Early molecular representation learning focused on embedding models that map graphs or sequences to vector spaces for downstream prediction, e.g., MolFormer (Wu et al., 2023), Uni-Mol/Uni-Mol-v2 (Zhou et al., 2023; Ji et al., 2024), and

ChemBERTa-2 (Ahmad et al., 2022). These encoders excel at property classification and regression as they optimize contrastive and masked-token training over large chemical corpora, yielding stable, geometry-aware embeddings for downstream heads. However, they do not natively support natural-language generation (e.g. UniMoT (Zhang et al., 2024b)) or multi-hop reasoning (Liu et al., 2024; Jang et al., 2024); in these settings, evidence remains implicit in continuous vectors, and exact numeric targets are better served by calibrated regressors than autoregressive tokens.

From generation to reasoning with LLMs. LLM-based molecular systems provide a natural-language interface for explanation, instruction following, and tool use, enabling generative descriptions and stepwise reasoning (Fang et al., 2023; Zhao et al., 2023; Zhang et al., 2024b). Yet vanilla LLMs struggle with exact property prediction due to the mismatch between discrete token likelihoods and calibrated real-valued targets, unit/scale sensitivity, and exposure bias during numeric generation; recent work addresses these issues with instruction data, tool-augmented agents, and chain-of-thought methods tailored to molecular structure (Kim et al., 2025; Bran et al., 2023; Jang et al., 2024). Overall effectiveness hinges on how molecular structure is exposed to the language model—whether via tokenized geometry, graph-aware connectors, or substructure tokens.

Multimodal fusion for molecular generative modeling. Beyond SMILES-only models (e.g., KV-PLM (Zeng et al., 2022b)), recent work fuses SMILES, molecular graphs, and 3D conformers with LLMs. Instruction-tuned systems such as GIMLET (Zhao et al., 2023) show that paired graph-text supervision imparts basic molecular understanding; larger frameworks (e.g., MoleculeSTM (Liu et al., 2023), MolFM (Luo et al., 2023a), BioT5+ (Pei et al., 2024a), ProLLaMA (Lv et al., 2025)) broaden modality and biomedical knowledge coverage; and 3D-aware models (e.g., 3D-MolT5 (Pei et al., 2024b)) introduce coarse geometry for conditional generation. More recent graph—LLM connectors (e.g., Mol-LLM (Lee et al., 2025), Mol-LLaMA (Kim et al., 2025)), and UniMolT (Zhang et al., 2024b) adapt the BLIP-2/Q-Former (Li et al., 2023) paradigm with fixed-length, learnable query tokens. Although effective for alignment, fixed-length fusion compresses stereochemistry and substructures, resulting in structural information loss that worsens with increasing molecular size and conformer diversity; this motivates the development of dynamic and substructure-aware interfaces that present structures more faithfully to the LLM.

B.2 RELATIONSHIP TO RELATED WORKS

Relationship to Q-Former. BLIP-2 introduces a lightweight Querying Transformer (Q-Former) that uses a fixed set of learnable query tokens to extract visual evidence from a frozen image encoder via cross-attention, then projects the resulting tokens to a frozen LLM—enabling efficient vision—language alignment without updating the backbones (Li et al., 2023). The Q-Former is pretrained with vision—language objectives (e.g., contrastive/alignment and generative stages) to ensure its queries attend to informative regions and produce LLM-consumable embeddings. Subsequent systems (e.g., InstructBLIP (Dai et al., 2023)) keep the same mechanism but condition the Q-Former on textual instructions, improving instruction following while retaining the fixed query budget. Related adapters such as Flamingo's (Alayrac et al., 2022) Perceiver-Resampler also distill high-dimensional visual features to a small, fixed-length token set for a frozen LM.

As shown in Fig. 7, Dynamic Query Transformer differs in two aspects. 1) Modality. Molecules are variable-sized graphs (and conformers), where structural fidelity hinges on preserving substructures (stereochemistry, functional groups), not just salient image regions. 2) Interface. Instead of a fixed query budget, EDT-Former upgrades the Q-Former's fixed queries with **entropy-guided dynamic queries**. Concretely, Entropy-Guided Patching discovers entropy-peak—based segments and pools node embeddings into dynamic tokens, and Dynamic Query Transformer refines these tokens (with a small set of anchors) via self-/cross-attention before projection to the LLM space. This design preserves Q-Former—style efficiency (frozen backbones, a lightweight adapter) while improving graph-specific fidelity by matching the number and placement of queries to molecular structure, thereby enabling multimodal molecular models without fine-tuning the LLM backbone. In short, we adopt the Q-Former rationale—frozen encoders bridged by a compact adapter—but move beyond fixed-length queries to a dynamic structure-aligned interface tailored to graphs.

Relationship to BLT (Byte Latent Transformer). BLT groups raw bytes into dynamic variable-length patches using the entropy of the next byte, so that higher-entropy regions receive more com-

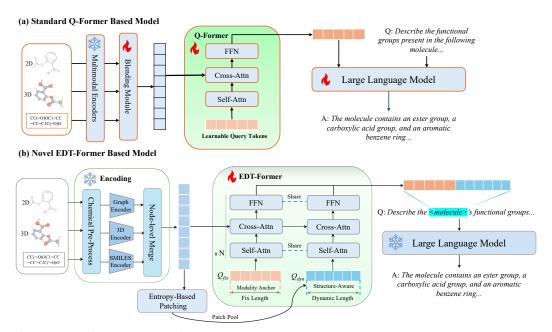


Figure 7: Architecture comparison between Q-Former-based model and EDT-Former. (a) Standard Q-Former connector that uses fixed-length learnable query tokens as modality anchors. (b) Novel EDT-Former connector that combines fixed-length tokens and entropy-guided structure-aware dynamic query tokens.

pute; a latent transformer then operates over these patches and projects to a language model, yielding tokenizer-free efficiency and robustness (Pagnoni et al., 2025). Our approach draws inspiration from the entropy-driven segmentation principle, but targets a different modality and objective: molecules are variable-sized graphs where fidelity depends on preserving substructures and stereochemistry. Concretely, we estimate next-atom uncertainty over SMILES to place peak-based cut points, map segments to graph nodes, and pool node embeddings into dynamic substructure tokens. These tokens are integrated (with a small set of anchors) via self-/cross-attention and projected to a frozen LLM, aligning query number and placement to molecular structure. Thus, while BLT's entropy patching motivates adaptive granularity, our entropy definition, boundary rule, and graph-aware interface are tailored to structural chemistry rather than byte-level text.

Comparison to other molecular LLMs. (1) Mol-Instructions (Fang et al., 2023). Instructions establishes large-scale instruction data for molecules and fine-tunes general LLMs (e.g., Llama-2/3 (Touvron et al., 2023b; MetaAI, 2024)) from SMILES inputs to demonstrate the value of instruction supervision for captioning and property-related responses. We adopt its benchmarks and reported baselines for our evaluations on molecular description and property prediction. While effective as a single-modality (text/SMILES) approach, it does not expose explicit structural signals (graphs/3D), which limits structure-aware reasoning beyond what can be inferred from tokenized strings. (2) Mol-LLaMA (Kim et al., 2025). Mol-LLaMA bridges molecular structure and language using a Q-Former-style connector: a fixed set of learnable query tokens cross-attends to strong encoders and conditions an LLM. This design proves the feasibility of graph-LLM alignment but inherits the fixed-length bottleneck, which compresses substructures and stereochemistry as molecular size and conformer diversity grow. In contrast, our method retains the lightweight adapter philosophy while replacing fixed queries with dynamic entropy-guided substructure tokens refined by a dynamic query transformer. For completeness, we employ recognized 2D/3D encoders (e.g., MoleculeSTM (Liu et al., 2023), Uni-Mol (Zhou et al., 2023)) and comparable processed training data when reproducing baselines; these choices standardize comparisons but are orthogonal to the core novelty.

Table 9: Configurations and Training settings for the Next-Atom Prediction (NAP) model, architecture, hyperparameters, together with the resulting compute footprint and parameter count.

Configurations	Pre-Training		
Item	Value	Item	Value
Vocabulary size	39	Learning rate	1×10^{-4}
Model architecture	GPT-2LMHead	Batch size	64
Maximum position embeddings	1024	Training epochs	1
Context window	512	Training steps	4,929
Transformer layers	2	Weight decay	0.01
Attention heads per layer	2	Total FLOPs	7.5395×10^{13}
Hidden size (embedding dimension)	128	Model size	0.538M
Activation function	GELU	Dataset	Pubchem
Dropout - residual connections	0.10	Max Length	128
Dropout - token/position embeddings	0.10	Concat Data	True
Dropout - attention	0.10	Precision	FP16
LayerNorm epsilon	1×10^{-5}	Objective	NAP

C IMPLEMENTATION DETAILS

This section details the implementation for reproducibility. It begins with the entropy-guided patching algorithm, including entropy computation, peak detection, the next-token predictor, and a brief rationale recap. Attention computations for integrating dynamic tokens with anchors are then formalized. Configurations of EDT-Former, covering both graph encoders and LLM backbones, are described next. The section concludes with pretraining settings (objectives, prompts, and hyperparameters) and finetuning protocols (setups, prompts, and representative inference examples).

C.1 Entropy-Guided Patching

Algorithm of Entropy-Guided Patching. Algorithm 2 segments a SMILES sequence into dynamic substructures using an entropy signal from a lightweight next-atom predictor (NAP). For each position, the predictor provides the probability of the ground-truth next atom; its negative log-probability defines the entropy trace. Local maxima are detected with a prominence threshold and non-maximum suppression, and cuts are placed after the retained peaks. Segments are mapped from SMILES indices to graph nodes and average-pooled over frozen node embeddings to form dynamic substructure tokens. The procedure is simple, deterministic given (Δ, γ) , linear in sequence length for entropy and segmentation, and produces structure-aware tokens aligned with molecular complexity.

The NAP model configuration. As summarized in Tab. 9, the NAP model adopts a GPT-2 architecture with a maximum input context of 512 tokens. This context window bounds the effective sequence length used during training and inference. The SMILES atom vocabulary shows in Tab. 10) comprises 39 tokens and spans major chemical groups—halogens, noble gases, alkali/alkaline earth metals, transition/post-transition metals, metalloids, and other nonmetals—providing broad elemental coverage for next-atom prediction.

Pre-training of NAP. We train the NAP model with standard autoregressive next-token prediction over SMILES. Given a tokenized sequence $x_{1:T}$ (prepended with $\langle b \circ s \rangle$ and optionally terminated by $\langle e \circ s \rangle$), a causal Transformer parameterized by θ defines the left-to-right factorization

$$p_{\theta}(x_{1:T}) = \prod_{t=1}^{T} p_{\theta}(x_t \mid x_{< t}).$$
 (8)

We optimize the negative log-likelihood (cross-entropy) under teacher forcing, averaged over non-padding tokens:

$$\mathcal{L}(\theta) = -\frac{1}{\sum_{t=1}^{T} m_t} \sum_{t=1}^{T} m_t \log p_{\theta}(x_t \,|\, x_{< t}), \qquad (9)$$

Table 10: Molecular atom vocabulary from SMILES grouped by category for the NAP (Next-Atom Prediction) model, with per-category counts (total tokens = 39, including 3 special tokens).

Category	Symbols	Count
Special tokens	<pad>, <bos>, <eos></eos></bos></pad>	3
Halogens	F, Cl, Br, I	4
Noble gases	He, Ne, Ar, Kr, Xe	5
Alkali metals	Li, Na	2
Alkaline earth metals	Mg, Ca	2
Transition metals	Fe, Co, Ni, Cu, Zn, Ag, Au, Pd, Pt, Mn, Hg	11
Post-transition metals	Al, Sn	2
Metalloids	B, Si, Sb, Te	4
Other nonmetals	C, N, O, P, S, Se	6
Total		39

where $m_t \in \{0,1\}$ masks out $\langle \texttt{pad} \rangle$ (and any positions beyond the context window). In practice, logits are computed over the 39-token vocabulary; optional label smoothing can be applied by replacing the one-hot targets with a smoothed distribution. During inference, we sample (or decode greedily) from $p_{\theta}(\cdot \mid x_{< t})$ until $\langle \texttt{eos} \rangle$ or the 512-token context limit is reached.

Computing costs analysis of NAP model. Using the settings in Tab. 9 on a single RTX 3090, end-to-end training lands in the "single-digit minutes" range (well below 0.1 GPU-hours for a full pass), so the cost is effectively negligible. The model is tiny (≈ 0.54 M parameters) and can also be trained comfortably on a modern multi-core CPU within an hour, avoiding the need for a GPU altogether. Compared with a typical 8B-parameter model, this NAP setup is $\approx 150 \times$ smaller and $\approx 1000 \times$ faster for both training and inference, so it does not increase the training or inference budget.

Entropy peaks from the NAP predictor identify structural transition points in SMILES; cutting at these peaks yields data-driven substructures and dynamic tokens that preserve locality and align cleanly with the language stream-simple, rule-free, and robust under a frozen backbone. The pipeline adds negligible overhead and does not materially affect training or inference costs.

C.2 Dynamic Query Transformer

Modeling Details with Anchors and Dynamic Tokens. We keep standard Transformer machinery and only specify the routed, two-stream updates. Let $Q_{\text{fix}}^{(\ell)} \in \mathbb{R}^{k \times d}$ (anchors) and $Z^{(\ell)} \in \mathbb{R}^{M \times d}$ (dynamic tokens).

$$Q_{\text{fix,sa}}^{(\ell)} = Q_{\text{fix}}^{(\ell)} + \text{SelfAttn}_{\text{fix}}(Q_{\text{fix}}^{(\ell)}, Z^{(\ell)}), \tag{10}$$

$$Z_{\rm sa}^{(\ell)} = Z^{(\ell)} + \operatorname{SelfAttn}_{Z}(Q_{\rm fix}^{(\ell)}, Z^{(\ell)}). \tag{11}$$

With graph node embeddings $X \in \mathbb{R}^{N \times d}$, each stream retrieves evidence independently:

$$\widehat{Q}_{\text{fix}}^{(\ell)} = Q_{\text{fix,sa}}^{(\ell)} + \text{CrossAttn}_{\text{fix}}(Q_{\text{fix,sa}}^{(\ell)}, X), \tag{12}$$

$$\widehat{Z}^{(\ell)} = Z_{\text{sa}}^{(\ell)} + \text{CrossAttn}_{Z}(Z_{\text{sa}}^{(\ell)}, X).$$
(13)

A shared FFN is applied to each stream separately:

$$Q_{\text{fix}}^{(\ell+1)} = \widehat{Q}_{\text{fix}}^{(\ell)} + \text{FFN}(\widehat{Q}_{\text{fix}}^{(\ell)}), \qquad Z^{(\ell+1)} = \widehat{Z}^{(\ell)} + \text{FFN}(\widehat{Z}^{(\ell)}).$$
 (14)

After L layers, project to the LLM space and concatenate to form the conditioning sequence:

$$U = \left[Q_{\text{fix}}^{(L)} W_{\text{proj}}; Z^{(L)} W_{\text{proj}} \right] \in \mathbb{R}^{(k+M) \times d_{\text{LLM}}}. \tag{15}$$

Configurations and pre-training settings. We pretrain the Dynamic Query Transformer bridge with a fixed anchor/dynamic query budget, moderate-depth Transformer, and a standard mixed-precision schedule with warmup, cosine decay, and weight decay; full hyperparameters and decoding settings are reported in Tab. 11.

1176

1177

1178

1179

1180

1181

1182 1183

1184 1185 1186

1187

```
1134
          Algorithm 2: Entropy-Guided Patching
1135
          Input: SMILES atom sequence (a_1, \ldots, a_T); trained next-atom predictor f_\theta; node
1136
                      embeddings X \in \mathbb{R}^{N \times d} (frozen graph encoder); map \pi (SMILES\rightarrowgraph); NMS
1137
                      window \Delta; prominence \gamma.
1138
          Output: Dynamic tokens \{z_k\}_{k=1}^M; SMILES segments \{S_k\}_{k=1}^M.
1139
          Entropy computation:;
1140
          \textbf{for } t \leftarrow 1 \textbf{ to } T - 1 \textbf{ do}
1141
                                                                              // logits over atom vocab {\cal V}
               L_t \leftarrow f_{\theta}(a_{1:t});
1142
               p_t \leftarrow \operatorname{softmax}(L_t)[a_{t+1}];
                                                                 // prob. of ground-truth next atom
1143
               e_t \leftarrow -\log p_t;
                                                     // information content (nats); O(T) total
1144
1145
          Peak detection (local maxima + prominence):;
          \mathcal{T}_{\star} \leftarrow \emptyset;
                                                                                // candidate split indices
1146
          for t \leftarrow 2 to T-2 do
1147
               isLocalMax \leftarrow (e_{t-1} < e_t) \land (e_t > e_{t+1});
                                                                                                      // strict peak
1148
               prominent \leftarrow \left(e_t - \frac{1}{2}(e_{t-1} + e_{t+1}) \ge \gamma\right);
                                                                                // remove shallow bumps
1149
               if isLocalMax \land prominent then
1150
                \mid \mathcal{T}_{\star} \leftarrow \mathcal{T}_{\star} \cup \{t\} ;
                                                                                                       // keep peak t
1151
1152
          Non-maximum suppression (window \Delta):;
1153
          Sort \mathcal{T}_{\star} by descending e_t to get (t^{(1)}, t^{(2)}, \dots);
                                                                         // highest entropy first;
1154
           O(|\mathcal{T}_{\star}|\log|\mathcal{T}_{\star}|)
1155
          \mathcal{T} \leftarrow \emptyset;
                                                                                                 // final peak set
1156
          foreach u \in (t^{(1)}, t^{(2)}, \dots) do
1157
              if \forall v \in \mathcal{T}, |u - v| > \Delta then |\mathcal{T} \leftarrow \mathcal{T} \cup \{u\};
1158
                                                                        // suppress neighbors within \Delta
1159
1160
          Segmentation (cut after peaks):;
1161
          Set 1 = \tau_0 < \tau_1 < \dots < \tau_{M-1} < \tau_M = T with \{\tau_1, \dots, \tau_{M-1}\} = \mathcal{T};
1162
            // deterministic cuts
1163
          for k \leftarrow 1 to M do
1164
              S_k \leftarrow \{ t \mid \tau_{k-1} \le t < \tau_k \} ;
                                                                                     // SMILES index segment
1165
              \widehat{\mathcal{S}}_k \leftarrow \{ \pi(t) : t \in \mathcal{S}_k \} ;
                                                                                         // map to graph nodes
1166
          Pooling to dynamic tokens (avg-pool):;
1167
          for k \leftarrow 1 to M do
1168
               z_k \leftarrow \operatorname{AvgPool}ig(\{\,X_i: i \in \widehat{\mathcal{S}}_k\,\}ig) \in \mathbb{R}^d\,; // substructure token; linear in
1169
1170
1171
          return \{z_k\}_{k=1}^{M}, \{S_k\}_{k=1}^{M};
                                                                   // dynamic structure-aware tokens
1172
```

Pre-training objectives. The Dynamic Query Transformer is pre-trained with three complementary objectives on Mol-Llama-Instruct dataset (Kim et al., 2025): (1) cross-modal contrastive alignment to pull together representations of the same molecule across modalities and push apart different molecules; (2) modality-matching to make fixed anchors decode the modality identity, enforcing a shared semantic interface; and (3) masked substructure reconstruction to teach dynamic tokens to preserve fine-grained structural content.

(1) Cross-modal contrastive loss. For a mini-batch \mathcal{B} , and modalities m and m', we use an InfoNCE objective

$$\mathcal{L}_{\text{contrast}} = -\frac{1}{|\mathcal{B}|} \sum_{(i,m)\in\mathcal{B}} \log \frac{\exp(\sin(z^{(i,m)}, z^{(i,m')})/\tau)}{\sum_{j\in\mathcal{B}} \exp(\sin(z^{(i,m)}, z^{(j,m')})/\tau)},\tag{16}$$

where $z^{(i,m)} \in \mathbb{R}^d$ is the pooled query embedding of molecule i in modality m, \sin is cosine similarity, and τ is a temperature.

Table 11: Configurations and pre-training settings of Dynamic Query Transformer. The model configuration is based on SciBERT and trained from scratch.

Item	Value	Item	Value
Base config	SciBERT	Embedding dimension	512
Max dynamic queries	64	Anchor queries	16
Max input tokens	594	Batch size	48
Attention heads	8	Transformer layers	8
Training epochs	10	Total steps	15,599
Precision	bf16-mixed	Scheduler	cosine_lr
Initial learning rate	1.0×10^{-4}	Minimum learning rate	1.0×10^{-5}
Warmup learning rate	1.0×10^{-6}	Warmup steps	500
Weight decay	0.05	Temperature	0.10
Top-k	50	Top-p	1.0
Repetition penalty	1.0	Accumulate grad batches	1

(2) Modality-matching loss. Fixed-length anchors $Q_{\rm fix}$ must predict the modality identity via an M-way classification:

$$\mathcal{L}_{\text{match}} = \frac{1}{|\mathcal{B}|} \sum_{(i,m)\in\mathcal{B}} \text{CE}(\mathbf{W}_m^\top \bar{z}^{(i)}, m), \tag{17}$$

where $\bar{z}^{(i)}$ is the average of the k anchor outputs for molecule i, \mathbf{W}_m are classifier parameters, and CE is cross-entropy.

(3) Masked substructure reconstruction loss. Dynamic queries $Q_{\rm dyn}$ are trained to recover masked fragments within each modality:

$$\mathcal{L}_{\text{recon}} = \frac{1}{|\mathcal{B}| M} \sum_{(i,m) \in \mathcal{B}} \text{CE}(f_{\text{dec}}^{(m)}(z^{(i,m)}), \text{ masked_tokens}^{(i,m)}),$$
(18)

where $f_{\text{dec}}^{(m)}$ is a modality-specific decoder and masked_tokens $^{(i,m)}$ are the targets.

The total objective is a weighted sum:

$$\mathcal{L}_{\text{total}} = \lambda_1 \, \mathcal{L}_{\text{contrast}} + \lambda_2 \, \mathcal{L}_{\text{match}} + \lambda_3 \, \mathcal{L}_{\text{recon}}, \qquad \lambda_1, \lambda_2, \lambda_3 > 0. \tag{19}$$

C.3 EDT-FORMER SETTINGS AND FINETUNING

After pre-training the Dynamic Query Transformer, the next stage is finetuning this connector with frozen LLMs.

Selection of LLMs. 6 After pre-training the Dynamic Query Transformer bridge, we finetune only the connector jointly with graph encoders and the LLM, while keeping both the encoders and the LLM frozen. We screen several 7–8B backbones, including Mistral-8B (Mistral AI, 2024), Qwen3-8B (Yang et al., 2025), Llama2-7B (Touvron et al., 2023b), and Llama3.1-8B (MetaAI, 2024)–covering strong base and reasoning-oriented models. As shown in Tab. 6, Llama3.1-8B yields consistent performance on the MoleculeQA (Lu et al., 2024) benchmark. We therefore adopt Llama3.1-8B as the default backbone for subsequent experiments, as it provides the strongest raw molecular understanding and the best absolute performance once paired with our connector. Notably, EDT-Former delivers consistent gains across all LLM backbones, underscoring the effectiveness of our connector design and its robustness to different architectures. For the frozen backbone, we freeze the main layers and retain the embedding layers trainable to ensure a stable convergence process, as the graph tokens from the Dynamic Query Transformer are mapped into LLM's embedding space.

Finetune configuration and settings. The graph encoders and LLM backbone use the default settings released by the author (Zhou et al., 2023; Liu et al., 2023). The training settings of Dynamic Query Transformer list in the Tab 12.

Table 12: Fine-tuning settings for EDT-Former with frozen Llama-3.1-8B backbone.

Item	Value	Item	Value
Precision	bf16-mixed	Epochs	2
Accumulate_grad_batches	8	Weight decay	0.05
Initial learning rate	1.0×10^{-4}	Minimum learning rate	5.0×10^{-6}
Warmup learning rate	1.0×10^{-6}	Warmup steps	1,000
Scheduler	linear_warmup_cosine_lr	Temperature	0.10
LoRA rank (r)	8	LoRA alpha	32
LoRA dropout	0.10	Trainable parameters	213M
Total parameters	8.3B	_	

Training dataset and prompt examples. As summarized in Tab. 13, our fine-tuning corpus from Mol-Llama-Instruct comprises four instruction types: (i) Detailed Structural Descriptions, that elicit atom-/substructure-level narratives of a molecule; (ii) Structure-to-Chemical Features, linking structure to physicochemical properties; (iii) Structure-to-Biological Features, covering bioactivity/ADMET-style attributes; and (iv) Comprehensive Conversations

Table 13: Fine-tuning dataset. Four types of instruction data from Mol-Llama-Instruct, the data types and amounts are listed.

Category	Amount
Detailed Structural Descriptions	77,239
Structure-to-Chemical Features	73,712
Structure-to-Biological Features	73,645
Comprehensive Conversations	60,147

for multi-turn reasoning and explanation. The dataset contains 284,743 instances in total, providing balanced coverage of descriptive, analytic, and conversational skills. Fig. 8 shows the prompts and examples of the finetuning data.

D EXPERIMENTS DETAILS

Here, we provide additional experimental details to complement the main results. This section first describes the benchmarks and baseline systems used for comparison, followed by the evaluation setup (splits, prompts, metrics, and decoding). We then include representative response examples to illustrate typical model behavior, and then we report our data-contamination checks and screening procedures. Finally, we document reproducibility essentials—configurations and other necessary information to enable faithful re-runs.

D.1 BENCHMARKS AND BASELINES

BENCHMARKS

MoleculeQA. A large-scale multiple-choice QA benchmark for molecular factuality, where each instance pairs a human-written question with one correct option and three distractors from authoritative molecule descriptions (Lu et al., 2024). It covers four aspects: Structure (configuration, functional groups, backbone), Source (isolation, discovery/derivation, metabolites), Property (physicochemical/biological properties, safety, mechanisms), and Application (therapeutic/chemical

Table 14: MoleculeQA split counts using Bemis–Murcko scaffold-based splitting at the molecule level to avoid scaffold overlap across splits.

Category	Train / Valid / Test
Overall	49,993 / 5,795 / 5,786
Structure	32,176 / 3,314 / 3,113
Property	4,838 / 698 / 731
Application	1,917 / 558 / 599
Source	11,062 / 1,225 / 1,343

uses, approvals, research/agricultural agents). Molecules are split by Bemis-Murcko scaffolds into train/dev/test sets of QA pairs as a total of 61,574. Aspect-wise counts are reported in Tab. 14.

Mol-Instructions. Mol-Instructions (Fang et al., 2023) is a large-scale instruction-tuning corpus for the biomolecular domain in instruction–response format for LLM supervision. Each example couples a natural-language prompt with target text for molecule-centric skills. We focus on two tasks: molecular captioning (generating faithful, structure-aware descriptions) and property prediction (in-

System: "You are a helpful assistant specializing in chemistry and biology. The instruction that describes a task is given, paired with molecules. Write a response that appropriately completes the request."

Structural Descriptions

User: "Can you tell me about the basic structure of the given molecule?\nMolecule <mol>."

Assistant: "The given molecule is a complex carbohydrate or polysaccharide derivative that contains several monosaccharide units with multiple functional groups and modifications.\n\n1. **Core Structure**: The backbone is based on cyclic carbohydrates. Specifically, ... "

Structure-to-Chemical

User: "Molecule <mol>.\nOffer a thorough analysis of the chemical characteristics of the compound."

Assistant: "The given molecule has a structure that combines several functional groups, imparting unique chemical properties to it. Centrally, it contains a propenone linkage, which features a carbon-carbon double bond conjugated with a carbonyl group. This conjugated system not only ..."

Structure-to-Biological

User: "Offer an in-depth discussion of the biological traits of the molecule, specifically highlighting the relationship between the core parts of the molecule and its properties.\nMolecule <mol>."

Assistant: "The given molecule is a complex natural product found in Illicium verum, commonly known as star anise. The biological properties of this molecule are influenced by its structural features, which include several characteristic functional groups and substructures.\\\\n## Tetracyclic Core Structure: ..."

Comprehensive Conversations

User: "Can you tell me about the basic structure of the given molecule?\nMolecule <mol>."

Assistant: "Certainly! The given molecule is a lactone, more specifically, a five-membered cyclic ester. It includes a methylene group attached at the 3-position of the ring and two methyl groups at the 5-position.\n"

User: "What are the key functional groups present in the given molecule?"

Assistant: "The key functional group in the given molecule is the lactone, which is a cyclic ester. In addition, there's an exocyclic methylene group at the 3-position.\n"

Figure 8: Prompt setup and instruction data examples of EDT-Former in the finetuning stage with backbone frozen. One instance is shown for each data type; the comprehensive conversation category contains multiple rounds (exceeding ten).

ferring physicochemical or bioactivity attributes from molecular inputs). The corpus spans over two million instructions and includes the captioning and property instruction sets used in our fine-tuning.

BBBP and Pampa.We use two ADME-focused datasets from Therapeutics Data Commons (TDC) (Huang et al., 2021). BBBP (Blood-Brain Barrier Penetration) is a binary classification task predicting a compound's ability to cross the blood-brain barrier. PAMPA (Parallel Artificial Membrane Permeability Assay) predicts passive membrane permeability from an in vitro surrogate of oral absorption. Both datasets are standardized in TDC for model development and fair comparison.

BASELINES

General Language Models. We compare against strong general LLMs spanning proprietary and open releases: GPT-40 (OpenAI, 2023) and GPT-5 (OpenAI, 2025) as OpenAI's flagship multimodal and next-generation reasoning models. From Meta, we include Llama2-7B and Llama3.1-8B, widely used open backbones for downstream fine-tuning and dialogue (Touvron et al., 2023b; MetaAI, 2024). We also consider domain/community models: Galactica for science-centric pretraining (Taylor et al., 2022), BLOOM as a large multilingual open model (Workshop, 2022), and Pythia for scaling/analysis studies (Biderman et al., 2023). Finally, we include popular instruction-tuned chat models derived from LLaMA/Llama-2 (Vicuna-v1.5 (Zheng et al., 2023), Alpaca-7B (Taori et al., 2023), Baize-7B (Xu et al., 2023)) as well as Qwen3-8B as a recent multilingual/reasoning baseline Yang et al. (2025).

Molecular Language Models. For molecular-domain LLMs spanning datasets, graph/3D alignment, and reasoning, Mol-Instructions (Fang et al., 2023) provides a large instruction–response corpus tailored to molecules, while trained Llama-based Mol-Instructions model themselves. Mol-

Table 15: Hyperparameter settings for benchmarks. For all the evaluated models and datasets, identical finetuning or inference parameters are applied for fair comparison.

Parameter	MoleculeQA	Mol-Instructions	PAMPA/BBBP	
Batch size	2	4	32	
Precision	BF16-Mixed	BF16-Mixed	BF16-Mixed	
Max epochs	5	10	_	
Accumulate grad batches	64	32	_	
Weight decay	0.05	0.05	_	
Initial LR	1.00e-04	1.00e-04	_	
Minimum LR	1.00e-05	1.00e-05	_	
Warmup LR	1.00e-06	1.00e-06	_	
Warmup steps	100	100	_	
Scheduler	Cosine LR	Cosine LR	_	
Temperature	0.1	0.1	0.1	
Max input tokens	1024	1024	512	

Llama (Kim et al., 2025) targets general molecular understanding via multi-modal instruction tuning. 3D-MoLM (Li et al., 2024) equips an LM with a 3D molecular encoder through a dedicated projector for 3D-aware captioning and QA. LLaMo (Park et al., 2024) integrates a molecular graph encoder with a multi-level graph projector to bridge graphs and language. BioMedGPT-LM (Luo et al., 2023b) is a biomedical generative LM adapted for molecule- and biology-centric tasks. Mol-Reasoner (Zhao et al., 2025) emphasizes interpretable chemical reasoning with a two-stage procedure (supervised reasoning initialization followed by reinforcement learning).

Benchmarks span factual QA, molecular captioning/property prediction, and ADME tasks, providing complementary coverage of accuracy, reasoning, and domain grounding. Baselines include both strong general LLMs and chemistry-specialized LLMs to isolate EDT-Former 's contribution beyond backbone capacity.

BASELINE RESULTS

The official benchmark protocols and data splits are adopted for all evaluations. For baselines whose reported performance trails ours by more than 10%, we cite the results directly from the benchmark sources (Lu et al., 2024; Fang et al., 2023; Kim et al., 2025). For competitive models (within this margin) and for baselines not covered by the benchmark, we re-run both the baselines and our method under a unified, fair protocol using released code and configurations. Full implementation and evaluation details are provided in App. D.3.

D.2 DATA PROCESS

For structural inputs, we adopt Uni-Mol (Zhou et al., 2023) and MoleculeSTM (Liu et al., 2023). For Mol-LLaMA-Instruct (Kim et al., 2025), we use the provided 3D coordinates for both pretraining and finetuning. For other datasets (e.g., Mol-Instructions (Fang et al., 2023)), we generate eleven conformers with RDKit and perform geometry optimization using MMFF (Merck Molecular Force Field Halgren (1996)), a classical small molecule force field that assigns bonded and nonbonded terms to atoms and bonds and minimizes molecular energy to a locally stable 3D geometry. After optimization, RDKit is used to construct the 2D molecular graph (including atom and bond types, charges, hybridization, aromaticity, and stereochemical flags). Molecules that fail RDKit sanitization are filtered out. All processed data and scripts are released for reproducibility.

D.3 EVALUATION SETTINGS

As summarized in Tab. 15, all models evaluated on MoleculeQA (Lu et al., 2024) and Mol-Instructions (Fang et al., 2023) are finetuned under identical hyperparameter settings to ensure a fair comparison. For PAMPA and BBBP (Huang et al., 2021), all models are evaluated zero-shot with a shared prompt template—BBBP as shown in Fig. 9 and PAMPA following Mol-Llama (Kim et al., 2025). For the 10-shot setting on MoleculeQA, the support set consists of 10 randomly selected examples fixed once and reused for every 10-shot model to control variance.

BBBP: Blood-Brain Barrier Penetration

Direct

System: "You are a drug discovery assistant tasked with predicting whether a molecule can penetrate the blood-brain barrier (BBBP). Output the final decision strictly using the provided answer format string.\nYour final answer should be formatted as either: 'Final answer: Penetrant' or 'Final answer: Non-penetrant' "

User: "Determine the BBBP class of the given molecule.\nMolecule <mol>."

Reasoning

System: "You are a drug discovery assistant tasked with predicting whether a molecule can penetrate the blood-brain barrier (BBBP). Output the final decision strictly using the provided answer format string.\nYour final answer should be formatted as either: 'Final answer: Penetrant' or 'Final answer: Non-penetrant'"

User: "Determine the BBBP class of the given molecule and briefly explain your rationale.'nMolecule <mol>."

Rich Instructions

System: "You are a drug discovery assistant for BBBP prediction. Molecules that are lipophilic, appropriately sized, and not highly polar are more likely to penetrate the BBB. Consider: lipophilicity (logP), polar surface area, hydrogen bonding capacity, ionization state, and presence of efflux transporter motifs. Output the final decision strictly using the provided answer format string.\nYour final answer should be formatted as either: 'Final answer: Penetrant' or 'Final answer: Non-penetrant' "

User: " Determine the BBBP class of the given molecule.\nMolecule <mol>."

Figure 9: Zero-shot prompt templates for BBBP (Blood-Brain Barrier Penetration). Three settings are evaluated: (i) Direct, which asks for a binary decision; (ii) Reasoning, which requests a brief rationale before the decision; and (iii) Rich Instructions, which provides domain context (e.g., lipophilicity, PSA, H-bonding, ionization, efflux motifs) prior to answering. All prompts enforce an identical output format: Final answer: Penetrant or Final answer: Non-penetrant.

Evaluation Metrics. (1) MoleculeQA (Lu et al., 2024): We report task-wise accuracy (fraction of correctly answered questions within each task). Average accuracy is the unweighted mean of the four task accuracies, and total accuracy is the overall correct/total across all questions pooled from the four tasks (count-based). Higher is better. (2) Mol-Instructions (Fang et al., 2023): For property prediction we use Mean Absolute Error (MAE), $\text{MAE} = \frac{1}{n} \sum_{i=1}^{n} |\hat{y}_i - y_i|$ (lower is better). For molecular description generation, we report: BLEU-2/4 (n-gram precision with a brevity penalty; BLEU-2 uses up to bigrams, BLEU-4 up to 4-grams; higher is better), ROUGE-1/2 (recall-oriented overlap of unigrams/bigrams; higher is better), ROUGE-L (recall-oriented longest common subsequence; higher is better), and METEOR (harmonic mean of unigram precision and recall with stemming/synonym matching and fragmentation penalty; higher is better).

D.4 DATA CONTAMINATION ANALYSIS

13-gram overlap (contamination) analysis. Following the GPT-3 practice (Brown et al., 2020), we measure data leakage by character-level n-gram overlap with n=13. For a text string x, let $\mathcal{G}_{13}(x)$ be the multiset of all contiguous 13-grams. For a test item t and a training document d, the overlap score is

$$\operatorname{overlap}(t,d) = \frac{\left| \mathcal{G}_{13}(t) \cap \mathcal{G}_{13}(d) \right|}{\left| \mathcal{G}_{13}(t) \right|}, \tag{20}$$

and the contamination score for t is $\max_{d \in \mathcal{D}_{train}} \operatorname{overlap}(t,d)$. We report the fraction of test items whose contamination score exceeds a small threshold (e.g., any non-zero match or a preset τ). Applying this procedure to our benchmarks, the test-train 13-gram overlap rate is below 5%. Given that we fine-tune only the connector (the LLM remains frozen) and for fewer than two epochs, memorization of instance-specific labels is unlikely; the observed overlap is within common LLM practice and does not affect our conclusions.

Table 16: Extended evaluation results on Pampa dataset. Accuracy and F1 scores are reported for each model across three prompt settings. Top and second-best results are highlighted.

Models	Direct		Reasoning		RichInst.		Average	
	Acc	F1	Acc	F1	Acc	F1	Acc	F1
GPT-4o	48.65	58.78	58.23	70.41	47.17	56.73	51.35	61.97
Mol-InstLlama2-7B	49.63	63.19	31.16	39.84	38.18	46.61	39.66	49.88
Mol-LLaMA2-7B	75.68	85.96	79.61	88.91	67.90	79.41	74.40	84.76
Llama3.1-8B	56.51	68.87	46.19	57.22	63.64	75.88	55.45	67.32
Mol-InstLlama3.1-8B	55.91	69.84	33.50	39.84	70.47	81.46	53.29	63.71
3D-MoLM-8B	46.93	57.81	50.00	62.43	64.86	76.45	53.93	65.56
LLaMo-8B	49.25	61.78	64.37	77.11	48.51	60.66	54.04	66.52
Mol-LLaMA3.1-8.2B	63.55	75.32	64.37	76.72	72.48	83.51	66.80	78.52
EDT-Former-8.3B	81.57	89.92	81.57	89.74	83.78	91.15	82.31	90.27

D.5 REPRODUCIBILITY

Our work is fully open source under the MIT license. Unlike many prior projects, we release end-to-end resources to reproduce every result in the paper: complete training and evaluation code (including all ablations), scripts for baselines and benchmarks, processed pretraining and finetuning datasets, and model weights. The repository is actively maintained and will serve as the foundation for subsequent extensions. All materials are available via the anonymous GitHub link: https://anonymous.4open.science/r/EDT-Former-844D.

D.6 DATASET IMBALANCE ANALYSIS

To account for potential class imbalance in the Pampa dataset, we additionally report F1 scores alongside accuracy (shown in Tab. 16). As expected, EDT-Former achieves the best performance across both metrics under all prompting strategies, confirming that its gains are not driven by skewed label distributions but reflect genuine improvements in molecular property prediction.

E EXTENDED ABLATIONS

Ablation studies on model size, graph encoders, query token length, computational efficiency, and training hyperparameters are provided in this section. These analyses highlight how each factor influences performance and efficiency, underscoring the robustness of our approach.

E.1 IMPACT OF QUERY TOKEN LENGTH

As shown in Tab. 17, increasing the number of anchor tokens improves accuracy, indicating that a larger anchor set provides a higher-bandwidth, more stable interface for aggregating global structural cues before conditioning the LLM. By contrast, expanding the maximum dynamic length from 64 to 128 yields a negligible change (e.g., 66.75 to 66.81 at 8 anchors) and even a slight dip at 16 anchors. This is expected: with a cap of 64, entropy-guided segmentation already covers the substructure count of most molecules, so raising the limit rarely increases the realized number of dynamic tokens and can introduce marginal redundancy. Overall, the best trade-off is achieved with 16 anchors and a max dynamic length of 64.

E.2 IMPACT OF MODEL SIZE

As shown in Tab. 18, Qwen3 (Yang et al., 2025) backbones below 2B parameters perform near random guessing on MoleculeQA, indicating insufficient capacity. Accuracy improves markedly once scale reaches 4B, with further gains at 8B and 14B, albeit with diminishing returns from 8B ot 14B. We use the Qwen3 family because it offers a broad spread of sizes for a controlled comparison, and we did not evaluate larger backbones due to resource constraints.

Table 17: Effect of anchor and dynamic token length. Different settings are evaluated on Mthe olcueleQA dataset, with total accuracy reported for each ablation model. The top and second-best results are highlighted.

st results are nightighted.					
Anchor Length	Max Dyn. Length	Accuracy			
4	64	64.29			
8	64	66.75			
8	128	66.81			
16	64	68 31			

Table 18: Performance across different model sizes of backbone choices from Qwen3. Models are evaluated on the MoleculeQA benchmark with total accuracy reported. The top and second-best are highlighted.

Model Size	Accuracy
0.6B	27.23
1.7B	26.52
4B	36.82
8B	54.79
14B	56.84

Table 19: Ablation on finetuning hyperparameters on MoleculeQA benchmark. Accuracy is reported for different initial learning rates, global batch sizes, and training epochs.

68.03

Init LR	Accuracy.	Global Batchsize	Accuracy	Training Epochs	Accuracy
5.0×10^{-4}	67.92	32	67.97	2	65.23
1.0×10^{-4}	68.34	64	68.08	4	67.92
5.0×10^{-5}	67.82	128	68.34	6	68.34
1.0×10^{-6}	67.09	256	68.11	8	68.13

E.3 ANALYSIS OF HYPERPARAMETER

As summarized in Tab. 19, varying the initial learning rate, global batch size, and training epochs produces only minor fluctuations in accuracy. This stability indicates that performance gains primarily arise from our multimodal fusion design and entropy-guided structural cues rather than aggressive hyperparameter tuning.

E.4 IMPACT OF GRAPH ENCODERS

We evaluated stronger molecular encoders (e.g., Uni-Mol-v2 (Ji et al., 2024)) under matched settings and observed only minor accuracy gains relative to our base encoder. This suggests that, beyond a reasonable encoder quality threshold, the primary bottleneck for molecular understanding lies not in the encoder itself but in the multimodal fusion interface to the LLM. In our framework, the entropyguided substructure tokens and dynamic query—anchor fusion contribute the dominant improvements by exposing chemically salient evidence to the LLM; upgrading the encoder yields diminishing returns compared to better fusion.

F LIMITATIONS AND FUTURE WORK

F.1 LIMITATIONS

While EDT-Former demonstrates broad applicability across molecular tasks, we note several current constraints that highlight opportunities for improvement. As a generalist molecular LLM, it does not yet consistently match highly specialized classifiers that are trained and tuned per dataset—a pattern commonly observed in the Mol-LLM literature. In this study, we also did not scale to larger backbones due to practical compute limits. In addition, we did not incorporate synthetic "reasoning" corpora, as widely used resources are predominantly GPT-generated and lack reliable human-verified rationales, which could confound conclusions. Finally, the present system answers end-to-end without external tools, which can increase the chance of occasional hallucinations or arithmetic/chemistry slips. We emphasize that these constraints are scope choices rather than fundamental barriers and frame clear directions for progress.

F.2 FUTURE WORK

Building on these findings, we plan to close the generalist–specialist gap via lightweight task adapters and, where feasible, moderate scaling of both backbone and bridge with efficient fine-tuning. We will curate and release a human-verified molecular reasoning benchmark to replace purely synthetic rationales, and integrate tool-augmented agents (e.g., RDKit calculators (Landrum, 2013), literature retrieval, unit/constraint checks) with self-verification to reduce hallucinations and enforce domain validity. Finally, we will broaden evaluation to include robustness under distribution shift, calibration, and abstention, and cost/latency reporting, to make EDT-Former both reliable and practical in real-world molecular workflows.

G LLM USAGE

Large language models (LLMs) were used only for light editorial assistance (grammar or wording polish and minor LaTeX phrasing). They did not contribute to research ideation, experimental design, data collection, analysis, or result writing. Separately, LLMs appear in our experiments solely as baselines/backbones within the method; all runs, configurations, and analyses were conducted and verified by the authors. The authors take full responsibility for the paper's contents.