

# EvoMD-LLM: Learning the Language of Species Evolution in Reactive Molecular Dynamics

Anonymous ACL submission

## Abstract

While large language models (LLMs) excel at static scientific reasoning, they struggle to model the temporal structure of dynamic physical processes. We present **EvoMD-LLM** (Evolutionary Molecular Dynamics Large Language Model), a framework that reformulates species-level molecular dynamics as a symbolic temporal language modeling problem. Reactive MD trajectories are discretized into sequences of molecular events, where each token represents a chemical species augmented with its persistence duration, enabling standard autoregressive LLMs to learn compositional evolution over time through efficient fine-tuning. A key component of EvoMD-LLM is temporal scaffolding, which treats event duration as an explicit linguistic token and serves as a structured inductive bias, significantly reducing invalid or hallucinated molecular outputs compared to conventional sequence modeling approaches. We evaluate EvoMD-LLM on multiple temporal prediction tasks, achieving up to 68.66% accuracy and consistently outperforming sequential neural networks and language-based baselines. Beyond quantitative improvements, we qualitatively observe that the model is capable of generating interpretations for its own predictions by incorporating relevant chemical knowledge, even though it was not explicitly supervised with paired trajectory-explanation data. These results demonstrate that symbolic temporal language modeling provides an effective framework for grounding LLMs in dynamic physical simulations.

## 1 Introduction

The convergence of large language models (LLMs) and molecular representations has emerged as a promising direction in AI for Science. Recent paradigms have successfully aligned static molecular encodings, such as SMILES strings (Cavanagh et al., 2024), with natural language, enabling LLMs to support tasks ranging from molecular property

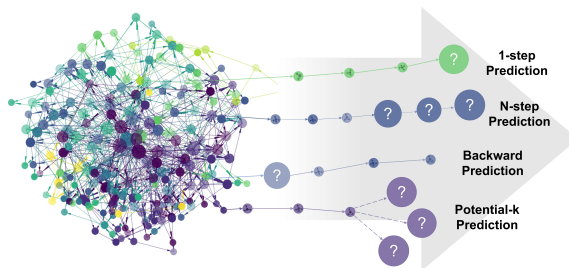


Figure 1: Conceptual overview of EvoMD-LLM. The framework interprets MD trajectories as structured sequences (Nodes: species; Edges: transformations) to reconstruct reaction pathways via four predictive tasks.

prediction (Chithrananda et al., 2020) to retrieval-augmented chemical reasoning (Chen et al., 2025). However, most existing approaches operate on static molecular representations or rely on external tools for reasoning (Boiko et al., 2023). This limits their applicability to physical systems, which evolve over time through sequences of creation, persistence, and transformation events. As a result, enabling LLMs to model temporal physical processes remains a fundamental challenge in AI for Science (Wigh et al., 2022). Molecular dynamics (MD) simulations provide a natural description of temporal physical evolution by recording time-resolved atomic motions (Alder and Wainwright, 1957). Yet, raw MD trajectories consist of high-frequency continuous coordinates that are incompatible with the discrete, symbolic token space of language models. Emerging time-series foundation models (Ansari et al., 2024) remain inapplicable to this challenge, as their numerical quantization schemes destroy the compositional semantics and discrete identity intrinsic to chemical species. Directly aligning MD simulations with LLMs therefore presents a key abstraction challenge: how to represent continuous molecular evolution as symbolic sequences amenable to language modeling. Existing learning-based approaches to MD trajec-

071 tories largely focus on structural dynamics in non-  
072 reactive or weakly reactive systems, such as protein  
073 folding (Tsai et al., 2020; Bera and Mondal, 2025;  
074 Murtada et al., 2024; Hussein Murtada et al., 2025),  
075 and are ill-suited for reactive processes character-  
076 ized by discrete changes in chemical species.

077 To address this gap, we introduce **EvoMD-LLM** (**E**vo-  
078 **M**olecular **D**ynamics **L**arge **L**anguage **M**odel), a framework that reformulates  
079 species-level molecular dynamics as a constrained  
080 generative language task. We propose a modality  
081 alignment scheme that translates continuous trajec-  
082 tories into discrete tokens, where duration serves  
083 as an explicit semantic modifier for each chemi-  
084 cal species. This representation enables standard  
085 autoregressive LLMs to internalize the "grammar"  
086 of chemical evolution directly through fine-tuning,  
087 eliminating the need for external simulators or spe-  
088 cialized architectures.

090 A key component of EvoMD-LLM is temporal  
091 scaffolding, which explicitly encodes event dura-  
092 tion as a linguistic token. By treating persistence  
093 as a semantic attribute, temporal scaffolding pro-  
094 vides a structured inductive bias that encourages  
095 the model to distinguish between long-lived sta-  
096 ble states and short-lived transient intermediates.  
097 Empirical ablation studies show that this design  
098 significantly improves prediction accuracy and re-  
099 duces invalid or hallucinated molecular outputs.

100 We evaluate EvoMD-LLM on a comprehensive  
101 suite of temporal prediction tasks, as illustrated in  
102 Figure 1. Beyond quantitative metrics, we remark-  
103 ably observe that the model exhibits emergent ex-  
104 planatory behaviors: despite lacking explicit super-  
105 vision, it spontaneously produces plausible phys-  
106 ical rationales for kinetic stability. These results  
107 demonstrate that symbolic temporal language mod-  
108 eling serves as an effective framework for learning  
109 species-level dynamics. Our main contributions are  
110 summarized as follows:

- 111 • **EvoMD-LLM Framework:** We propose  
112 a language modeling framework that reform-  
113 ulates species-level molecular dynamics  
114 as symbolic event sequences, enabling stan-  
115 dard autoregressive large language models to  
116 model temporal evolution in reactive systems.
- 117 • **Temporal Scaffolding via Duration Tokens:**  
118 We introduce temporal scaffolding by explic-  
119 itly encoding event persistence as linguistic  
120 tokens. This structured inductive bias signif-  
121 icantly improves prediction accuracy and re-

122 duces invalid molecular outputs, as demon-  
123 strated by extensive ablation studies.

- 124 • **Unified Temporal Prediction Formulation:**  
125 We show that a single instruction-tuned lan-  
126 guage model can flexibly support diverse tem-  
127 poral prediction tasks, including forward fore-  
128 casting and backward inference, without task-  
129 specific architectures.

## 130 2 Methods

131 We propose EvoMD-LLM to treat molecular evo-  
132 lution as a foreign language with its own grammar  
133 of causality and persistence. As illustrated in Fig-  
134 ure 2, our framework operates through a four-stage  
135 pipeline: (1) Dynamic Modality Alignment; (2)  
136 Structured Instruction Formatting; (3) Heteroge-  
137 neous Task Integration; and (4) Model Training  
138 and Inference. In this section, we detail the theoret-  
139 ical formulation and key algorithmic components.

### 140 2.1 Problem Formulation

141 We enable LLMs to learn the dynamics of chemical  
142 reactions by reformulating MD simulations as a  
143 structured symbolic text generation problem.

144 A standard MD simulation produces a raw tra-  
145 jectory  $\mathcal{T}_{\text{raw}}$ , recording atomic positions  $\mathbf{R}$  and mo-  
146 menta  $\mathbf{P}$  at each time step  $\tau$ :

$$147 \mathcal{T}_{\text{raw}} = \{(\mathbf{R}(\tau), \mathbf{P}(\tau)) \mid 0 \leq \tau \leq T_{\text{total}}\}. \quad (1)$$

148 While physically complete, such trajectories are  
149 high-dimensional and dominated by thermal noise,  
150 which obscures long-term reaction patterns. To ob-  
151 tain a representation amenable to language model-  
152 ing, we apply a transformation  $\Phi$  that maps raw tra-  
153 jectories to a discrete sequence of molecular states:

$$154 \mathcal{X} = \Phi(\mathcal{T}_{\text{raw}}) = \{(m_i, \Delta t_i)\}_{i=1}^n, \quad (2)$$

155 where each state consists of a semantic unit  $m_i$   
156 (molecular formula) and its temporal persistence  
157  $\Delta t_i$  (duration). This abstraction suppresses high-  
158 frequency atomic fluctuations while preserving  
159 the causal sequence of chemical transformations.  
160 Unlike standard text generation where tokens are  
161 equidistant, chemical evolution is an irregularly  
162 sampled time series. We treat this sequence directly  
163 as natural language. This allows us to train the  
164 model using standard autoregressive cross-entropy  
165 loss, without requiring specialized regression archi-  
166 tectures.

**Generative Modeling Objective.** We formulate reaction modeling as conditional sequence generation. Given a context sequence  $\mathbf{x}$  and an instruction  $\mathcal{I}$ , the model generates a target sequence  $\mathbf{y}$  according to the factorization:

$$P(\mathbf{y} | \mathbf{x}, \mathcal{I}) = \prod_{j=1}^{|\mathbf{y}|} P(y_j | y_{<j}, \mathbf{x}, \mathcal{I}), \quad (3)$$

where  $\mathbf{y} = ((m'_1, \Delta t'_1), \dots, (m'_L, \Delta t'_L))$ . The instruction  $\mathcal{I}$  specifies the task (e.g., forward or backward prediction), enabling a unified formulation across different reaction reasoning scenarios.

## 2.2 Dynamic Modality Alignment

To bridge the gap between continuous physical simulations and discrete symbolic reasoning, we construct a Dynamic Modality Interface. This process translates raw MD trajectories into a structured "grammar" of reaction events, characterized by semantic identity and temporal persistence.

### From Continuous Trajectories to Discrete

**Events.** Raw MD data consists of high-frequency atomic coordinates dominated by thermal noise. We adopt the ab initio bond-order determination method established by Dang et al. (2025) as the physical ground truth for identifying atomic connectivity. Building upon these snapshots, our framework projects the continuous evolution into a discrete event space by defining molecular formulas as atomic semantic units. Unlike standard NLP approaches that tokenize chemical strings into subword units (e.g., SMILES characters), we treat each distinct molecular formula as an atomic semantic unit. This preserves the integrity of chemical identity, allowing the LLM to reason over species-level transformations rather than character-level statistics.

We define a valid Molecular Event  $\mathcal{E} = (m, \Delta t)$  as a tuple comprising a molecular species  $m$  and its persistence duration  $\Delta t$ . To distill chemically significant states from transient thermal fluctuations, we apply a temporal band-pass filter, retaining only events with durations  $\Delta t \in [10, 500]$  ps. This operation effectively isolates stable reaction intermediates from high-frequency noise. Details about the original dataset scale and filtering statistics are provided in Appendix A.

**Structured Context Construction.** To enable autoregressive forecasting, the discrete event stream is segmented into structured input-output

pairs using a sliding window approach. Each training example consists of a historical context window (3–5 events) and a target future event.

A critical challenge in learning from physical simulations is the long-tail distribution of molecular events. Raw reaction data is often dominated by a few highly stable species (e.g., reactants or final products), while key intermediate transition states appear infrequently. Training directly on such imbalanced data would lead the LLM to degenerate into a trivial identity mapping (predicting the most frequent species). To mitigate this, we employ a stratified sampling strategy that balances the dataset across both molecular types and temporal regimes. This ensures that the model is exposed to a diverse range of reaction pathways and learns to differentiate between rapid intermediates and stable products based on structural context rather than mere occurrence frequency.

Detailed visualizations of the data evolution, species distribution, and the effects of balancing are presented in Appendix A (Figure 5).

## 2.3 Temporal Scaffolding

Standard Transformers, while adept at sequence ordering, remain agnostic to variable time intervals. To bridge this gap, EvoMD-LLM implements Explicit Temporal Tokenization by interleaving duration tokens  $\Delta t_i$  with species tokens. We reinterpret this strategy as a neural implementation of Run-Length Encoding (RLE). Since raw MD trajectories are dominated by high-frequency thermal fluctuations (e.g., non-reactive vibrations), explicitly tokenizing persistence performs semantic compression, effectively decoupling continuous physical time from the logical reaction sequence. This allows the LLM to skip redundant noise and reason directly across chemically significant timescales.

Theoretically, this design aligns with proven paradigms in variable-duration modalities, such as note sustain in Music Transformers (Huang et al., 2018) and phoneme alignment in FastSpeech (Ren et al., 2019). By treating duration as a semantic attribute, this scaffolding introduces a strong inductive bias: it forces the model to differentiate between thermodynamically stable states (long  $\Delta t$ ) and transient intermediates, thereby suppressing "hallucinated" transitions that violate physical kinetic constraints.

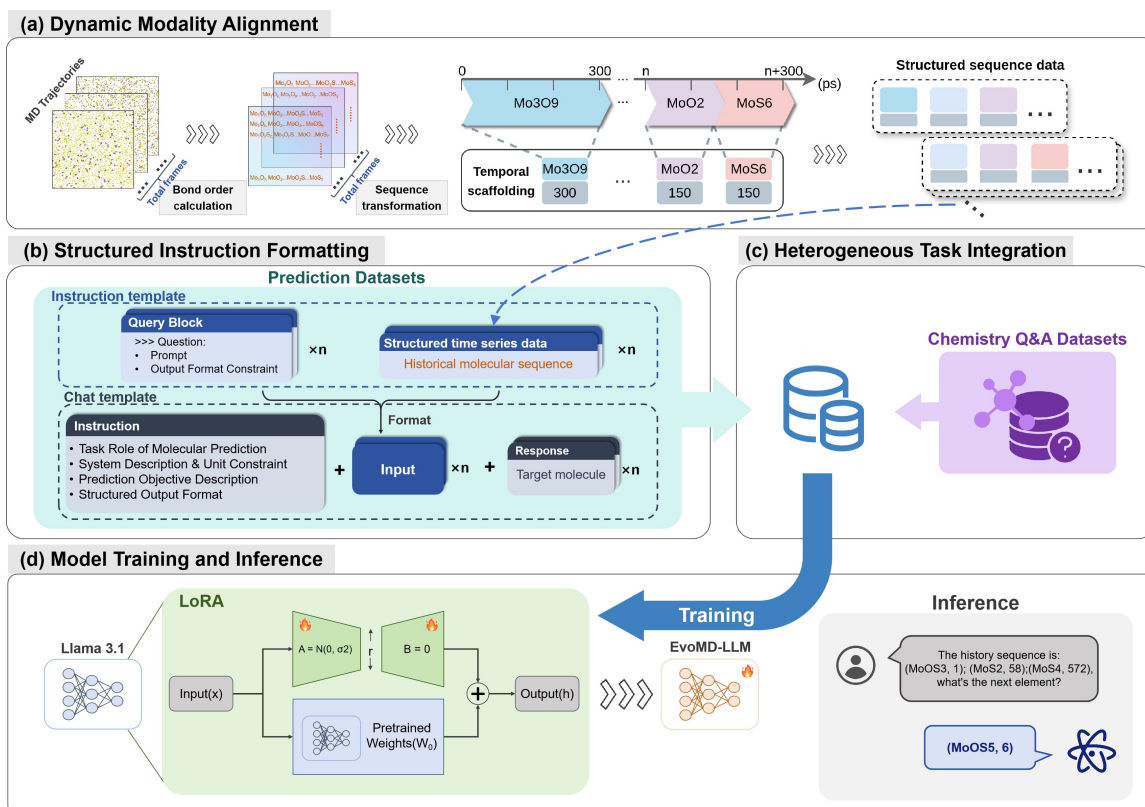


Figure 2: The overall framework of the model. Encompassing dynamic modality alignment, structured instruction formatting, heterogeneous task integration and model training along with inference.

## 2.4 Structured Instruction Formatting

To transform the interleaved event sequences into training samples, we employ a structured instruction-tuning paradigm. As illustrated in Figure 2(b), we design a domain-specific template that enforces strict syntactic constraints on the generative output.

The construction consists of two components:

- System Context (Semantic Definition):** We utilize the system prompt to define the model’s role as a "Scientific Simulator." Crucially, this prompt establishes the semantic mapping for our unified vocabulary, explicitly instructing the model that the output must alternate between molecular formulas (representing state identity) and time tokens (representing kinetic stability).
- Task Instruction (Historical Constraints):** The user prompt encapsulates the historical context window  $x = \mathcal{H}_{<t}$ . Unlike open-ended chat, we inject structural constraints into the instruction, limiting the generation search space to valid physical transitions.

By wrapping the raw sequences in this rigorous

format, we align the stochastic nature of physical dynamics with the deterministic syntax required for language modeling.

## 2.5 Heterogeneous Task Integration

To synergize domain-specific dynamic modeling with general scientific reasoning, we construct a heterogeneous instruction dataset comprising two distinct streams as shown in Figure 2(c).

**Structured Forecasting Stream.** We curate prediction tasks covering 1-step, 2-step, and backward trajectory forecasting. Crucially, we restrict the training horizon to short-term contexts (maximum 2 steps). By mastering local transition rules, the model is forced to acquire temporal inductive reasoning capabilities, enabling it to generalize to long-horizon ( $N$ -step) planning during inference without explicit supervision on long sequences.

**Linguistic Regularization Stream.** While structured forecasting teaches the model the 'syntax' of reaction rules, it risks reducing chemical formulas to arbitrary symbols. To prevent catastrophic forgetting of general capabilities and provide semantic anchoring for chemical tokens, we inter-

309 leave synthetic Chemistry Q&A pairs. This stream  
310 serves as a semantic anchor, forcing the model to  
311 ground the symbolic molecular tokens in its pre-  
312 trained scientific knowledge base. This ensures that  
313 EvoMD-LLM evolves into a dual-capable agent:  
314 structurally grounded in specific physical dynamics  
315 while linguistically aligned with general chemical  
316 principles.

## 317 2.6 Model Training and Inference

318 As illustrated in Figure 2(d), we employ a super-  
319 vised fine-tuning (SFT) framework (Ouyang et al.,  
320 2022; Taori et al., 2023; Daniel Han and team,  
321 2023; Touvron et al., 2023), aligning the model  
322 to predict target molecular events directly from  
323 structured input sequences. SFT enables EvoMD-  
324 LLM to internalize domain-specific transition rules  
325 into its parameters. The training process is explic-  
326 itly designed to balance structural precision with  
327 linguistic generalization.

328 **Input Representation and Architecture** We uti-  
329 lize the Llama 3.1 8B (Meta AI Team, 2024) back-  
330 bone without architectural modifications. Formu-  
331 las are tokenized using the standard Byte-Pair En-  
332 coding (BPE) (Sennrich et al., 2016) vocabulary.  
333 This formatting encourages the model to process  
334 chemical formulas as semantic units, leveraging  
335 pretrained linguistic priors to model statistical reg-  
336 ularities in molecular evolution.

337 **Parameter-Efficient Optimization** To align the  
338 model with reaction dynamics while preserving  
339 general scientific reasoning, we employ Low-Rank  
340 Adaptation (LoRA) (Hu et al., 2022). By inject-  
341 ing trainable low-rank matrices into the attention  
342 and feed-forward layers while freezing pretrained  
343 weights, we achieve two strategic objectives: (1)  
344 prevention of catastrophic forgetting, ensuring the  
345 retention of the base model’s linguistic priors es-  
346 sential for the QA component; and (2) training  
347 efficiency, which allows for rapid convergence on  
348 consumer-grade hardware by focusing capacity ex-  
349 clusively on modeling temporal chemical patterns.

350 **Optimization Strategy** To enable dual compe-  
351 tence in structured forecasting and general reason-  
352 ing, we employ a multi-task sampling strategy:  
353 structured prediction datasets and chemistry Q&A  
354 instructions (Section 2.5) are interleaved during  
355 training. The final loss is computed as the stan-  
356 dard autoregressive cross-entropy over the target  
357 tokens of both tasks, ensuring the model simultane-

ously optimizes for domain-specific dynamics and  
linguistic fluency.

## 360 3 Experiments

361 We evaluate EvoMD-LLM on a range of temporal  
362 prediction tasks derived from molecular dynamics  
363 trajectories to assess its ability to model symbolic  
364 chemical evolution.

### 365 3.1 Tasks and Evaluation Protocol

366 To rigorously assess symbolic dynamic modeling,  
367 we utilize the Mo-S reactive system as a testbed.  
368 This system serves as a challenging benchmark due  
369 to its intrinsic stochasticity and the coexistence of  
370 competing reaction pathways (e.g., simultaneous  
371 growth and etching), demanding reasoning capa-  
372 bilities beyond simple pattern matching. This sec-  
373 tion evaluates whether EvoMD-LLM effectively  
374 addresses the proposed symbolic-temporal abstrac-  
375 tion gap in modeling dynamic chemical systems.  
376 Specifically, we evaluate EvoMD-LLM on four  
377 temporal prediction tasks: 1-step prediction for  
378 short-range consistency, N-step prediction for itera-  
379 tive long-horizon forecasting, backward prediction  
380 for bidirectional reasoning of precursor states, and  
381 Potential- $k$  prediction to capture the stochastic na-  
382 ture of chemical evolution (Coley et al., 2019).

383 Performance is quantified using complementary  
384 metrics. Prediction accuracy and potential- $k$  ac-  
385 curacy track the presence of ground-truth states  
386 within the 1 and  $k$  predictions, respectively. These  
387 serve as proxies for logical consistency, assessing  
388 whether the model captures the valid causal logic of  
389 chemical evolution. Conversely, missing rate calcu-  
390 lates the proportion of generated outputs that fail to  
391 parse as valid molecular formulas, thereby measur-  
392 ing the model’s Syntactic Validity and adherence to  
393 the chemical grammar. These metrics jointly assess  
394 the model’s ability to navigate branching chemical  
395 reaction pathways while maintaining both struc-  
396 tural integrity and instructional adherence.

### 397 3.2 Experimental Setup

398 **Baseline Methods** To evaluate the effectiveness  
399 of EvoMD-LLM, we compare it with four represen-  
400 tative categories of baselines for temporal knowl-  
401 edge integration:

- 402 • **In-context learning (ICL)** (Dong et al., 2024;  
403 Luo et al., 2025): We assess standard prompt-  
404 ing capabilities including **zero-shot (ZS)** ( $k =$   
405 0) and **few-shot (FS)** ( $k = 3$ ) (Brown et al.,

2020). To probe scalability, we also implement **many-shot** (1,000 examples) (Agarwal et al., 2024) and **full-Context** ( $\sim 7,000$  examples) to test the upper bound of long-context reasoning.

- **Retrieval-Augmented generation (RAG)** (Lewis et al., 2020; Zhong et al., 2025): A dynamic memory baseline where a retriever selects the  $k$  most similar historical subsequences from the training set to provide input-dependent context.
- **Sequential baselines (Numerical Modality)**: To assess the necessity of symbolic abstraction over raw numerical fitting, we consider two representative neural baselines that operate on numerical composition vectors (encoding atomic counts) rather than semantic tokens. Baselines include an LSTM (Hochreiter and Schmidhuber, 1997) and a custom encoder-only Transformer (Vaswani et al., 2017), which map trajectories to latent vectors for direct numerical regression.

**Implementation Details** Our final instruction-tuning dataset comprises over 20,000 samples, combining 7,015 stratified trajectory sequences with auxiliary scientific Q&A data. Among these samples, the generated RMD symbolic dataset is split into training and test sets at a 9:1 ratio. EvoMD-LLM is initialized with Llama 3.1 8B. We employ LoRA for parameter-efficient fine-tuning, optimizing approximately 42 M parameters ( $r = 16, \alpha = 16$ ). The model is trained for 2 epochs using the AdamW optimizer (Loshchilov and Hutter, 2019) with a global batch size of 8 and a peak learning rate of  $2e-4$ . We utilize a linear learning rate scheduler and mixed-precision (bfloat16) with a maximum sequence length of 2048. All experiments are conducted on a single consumer-grade GPU (NVIDIA RTX 4090D) to demonstrate accessibility.

### 3.3 Overall Performance Comparison

As shown in Table 1, EvoMD-LLM significantly outperforms all baselines, achieving 68.66% accuracy and a 0% missing rate. In comparison, the RAG baseline achieves only 36.61% accuracy, while zero-shot prompting suffers from a high missing rate (29.49%). Notably, simply extending the context window (Content-1000/All) yields

marginal gains over few-shot prompting, with accuracy plateauing around 18%. This saturation suggests that naive long-context prompting fails to effectively capture complex temporal dependencies. Conversely, our SFT approach enables the model to learn structured temporal correlations directly, leading to superior predictive accuracy and stability.

Figure 3(a) presents the confusion matrix of molecular species predicted by the EvoMD-LLM. The strong diagonal dominance indicates robust discriminative capability, with only minor confusion among chemically similar or temporally adjacent species. This highlights the model’s ability to capture fine-grained symbolic and temporal distinctions within the molecular event space.

### 3.4 Multi-step, Backward and Potential- $k$ Prediction Analysis

We further evaluate EvoMD-LLM on N-step, backward, and potential- $k$  prediction tasks to assess its ability to model long-range temporal dependencies and stochastic reaction dynamics.

As shown in Figure 3(b), accuracy decreases monotonically as the prediction horizon increases. While the average accuracy over the 3-step horizon remains high at 50.09% (Table 2), the point-wise accuracy for the final step drops to 37.74%. This trend reflects error accumulation in autoregressive forecasting, while the remaining performance indicates that fine-tuning improves temporal consistency over multiple steps.

Figure 3(c) shows that EvoMD-LLM substantially outperforms the LLaMA 3.1 base model in both potential- $k$  and backward prediction. Potential-1 accuracy increases from 13.96% to 68.66%, and potential-3 accuracy reaches 77.1%, with a smaller gap between the two metrics, indicating more stable candidate ranking. In addition, EvoMD-LLM achieves 57.12% accuracy in backward prediction, demonstrating its ability to infer plausible precursors from downstream states.

### 3.5 Comparison with Sequential Baseline

Table 2 compares EvoMD-LLM with an LSTM and an encoder-only model. For the 1-step prediction task, EvoMD-LLM achieves the highest accuracy at 68.66%, outperforming both LSTM (40.74%) and the encoder-only baseline (65.95%).

As the prediction horizon increases, all methods exhibit performance degradation. EvoMD-LLM consistently maintains higher accuracy in the 2-step

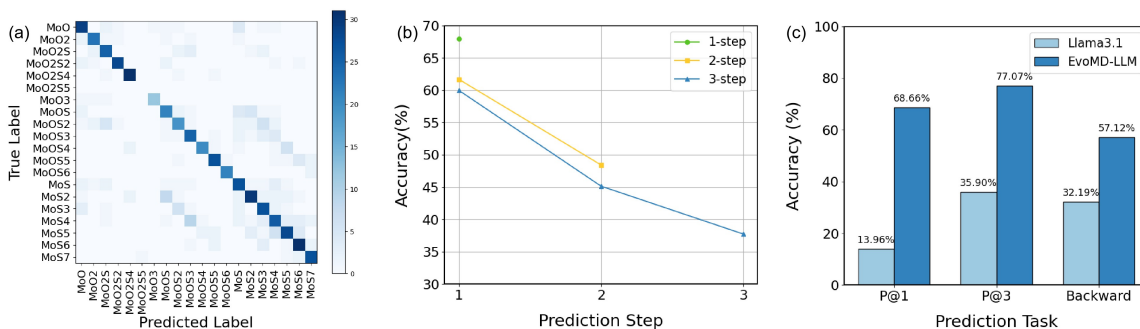


Figure 3: Experimental results. (a) Confusion matrix showing discriminative capability. (b) Accuracy decay over N-step forecasting horizons. (c) Performance comparison against the LLaMA 3.1 base model across three tasks.

Methods	Accuracy $\uparrow$	Missing Rate $\downarrow$
Zero-shot	13.96	29.49
Few-shot	17.02	5.63
Content-1000	17.28	1.05
Content-All	18.39	0.07
RAG ( $k = 5$ )	36.61	1.28
<b>EvoMD-LLM</b>	<b>68.66</b>	<b>0.00</b>

Table 1: Comparison of 1-step accuracy and missing rate between EvoMD-LLM and various LLM-based methods.

Model	1-step	2-step	3-step	Backward
LSTM	40.74	32.77	30.23	30.20
Encoder-only	65.95	51.71	36.84	52.71
<b>EvoMD-LLM</b>	<b>68.66</b>	<b>55.06</b>	<b>50.09</b>	<b>57.12</b>

Table 2: Performance comparison on molecular forecasting tasks. Reported values for N-step tasks represent the **average accuracy** over the entire prediction horizon (1 to N). Best results are **bolded**.

and 3-step settings, with particularly clear gains in the 3-step task. In the backward prediction task, EvoMD-LLM again outperforms both baselines. Overall, these results indicate that EvoMD-LLM provides more robust temporal modeling across different prediction settings.

## 3.6 Ablation Studies and Analysis

### 3.6.1 Impact of Temporal Scaffolding

As detailed in Table 3, excising the duration targets leads to a consistent performance degradation. The 1-step prediction accuracy declines from 68.66% to 55.70%, with a comparable drop in backward reasoning. These results empirically validate that temporal supervision is not merely an auxiliary task but a necessary constraint for correct chemical reasoning.

Model Variant	1-Step	Backward
w/o Temporal	55.70	49.15
w/o Q&A	65.53	55.98
<b>EvoMD-LLM (Full)</b>	<b>68.66</b>	<b>57.12</b>

Table 3: Ablation study. The significant drop in *w/o Temporal* validates the necessity of kinetic scaffolding, while excluding Q&A (*w/o Q&A*) has a relatively minor impact on domain-specific forecasting.

### 3.6.2 Impact of Multi-Task Instruction Tuning

A key design goal of EvoMD-LLM is to improve structured molecular forecasting without degrading general scientific language capabilities. To assess this, we evaluate an ablated variant trained without the chemistry Q&A dataset (*w/o Q&A*).

We adopt both qualitative and quantitative evaluations. For qualitative assessment, we use Qwen3 (Yang et al., 2025) as an automated evaluator to score model responses across multiple scientific capability and quality dimensions. In addition, we report performance on the AI2 Reasoning Challenge (ARC) benchmark (Clark et al., 2018) to measure standardized scientific reasoning, containing a total of 3548 test samples.

As shown in Figure 4(a–b), removing Q&A supervision leads to consistent degradation across all evaluated dimensions, with the largest drops observed in mechanism reasoning and question answering. The ablated model also exhibits reduced coherence and fluency, suggesting that natural language supervision contributes to stable scientific expression.

Figure 4(c) further shows that EvoMD-LLM maintains strong performance on both ARC-e and ARC-c, comparable to the base LLaMA 3.1 model, while the *w/o Q&A* variant performs noticeably worse. These results indicate that the Q&A dataset

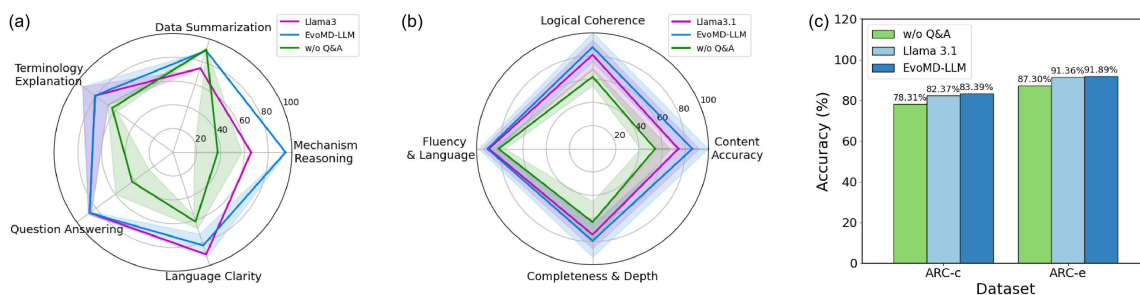


Figure 4: Evaluation of general language understanding. (a–b) present scores from an automated evaluation using Qwen3 as a judge. (a) displays performance across five distinct capability dimensions, while (b) assesses four key criteria for output quality. The shaded areas in the radar charts represent the standard deviation across the evaluation questions. (c) shows the accuracy of the models on the ARC-c (Challenge) and ARC-e (Easy) benchmark datasets.

Input History Sequence	Prediction	Generated Reasoning
<b>Case 1: Trend Following (Oxygen → Sulfur Exchange)</b>		
(MoO <sub>2</sub> , 98); (MoOS <sub>2</sub> , 3); (MoS, 182)	<b>Mol:</b> MoS <sub>3</sub> <b>Time:</b> 106 ps	The transition from MoS to MoS <sub>3</sub> is driven by <b>surface diffusion and adsorption</b> of additional MoS units. The predicted duration suggests the formation of a <b>stable three-dimensional structure</b> favored by thermodynamics.
<b>Case 2: Stability Detection (Distinguishing Stable vs. Transient)</b>		
(MoO <sub>2</sub> , 2); (MoO, 29); (MoOS <sub>2</sub> , 26)	<b>Mol:</b> MoS <b>Time:</b> 277 ps	The transition is driven by the <b>decomposition of MoOS<sub>2</sub></b> due to high temperatures. The predicted duration (277 ps) suggests that MoS is <b>relatively stable</b> , forming a <b>metastable trap</b> at these conditions.
<b>Case 3: Structure Building (Increasing Complexity)</b>		
(MoOS <sub>2</sub> , 43); (MoOS <sub>4</sub> , 6); (MoS <sub>3</sub> , 7)	<b>Mol:</b> MoS <sub>5</sub> <b>Time:</b> 365 ps	The transition involves the <b>addition of two sulfur atoms</b> through a series of elementary reactions. This process reflects the increasing preference for <b>higher-order coordination structures</b> in Mo-S systems.

Table 4: Qualitative examples of reasoning generated by EvoMD-LLM. The selected cases demonstrate the model’s ability to track linear evolutionary pathways without oscillation. Highlights indicate the textual description of **reaction mechanisms (Red)** and **stability/metastability (Blue)** aligned with temporal cues.

549 plays an important role in preserving general scienti-  
550 fic reasoning during domain-specific fine-tuning.

### 551 3.7 Qualitative Analysis

552 Quantitative metrics summarize predictive accuracy but provide limited insight into model behav-  
553 ior. We therefore present a qualitative analysis to examine the how EvoMD-LLM explains its pre-  
554 dictions without explicit supervision on reaction mechanisms. This analysis focuses on the align-  
555 ment of learned sequence patterns with semantic explanations, rather than on validating ab initio  
556 physical correctness.

561 Table 4 demonstrates the model’s context aware-  
562 ness. For instance, it correctly links early-stage  
563 sulfidation to surface diffusion (Case 1) while identifying high-temperature decomposition in inter-  
564 mediate phases (Case 2). Furthermore, it utilizes the duration token as a semantic pivot to differentiate  
565 between kinetically stable products and metastable  
566 traps, as evidenced by the distinct duration predic-  
567 tions. Additional qualitative examples, including  
568  
569

typical failure modes, are provided in Appendix G.

## 570 4 Conclusion

571  
572 In this work, we introduce EvoMD-LLM, a frame-  
573 work that re-frames molecular dynamics as a sym-  
574 bolic language modeling problem, thereby inter-  
575 nalizing the "grammar" of chemical evolution into  
576 LLMs. By aligning continuous physical trajec-  
577 tories with discrete semantic tokens through an Tem-  
578 poral Scaffolding strategy, we enable the model to  
579 treat temporal persistence as a semantic component.  
580 We demonstrate that this design introduces a robust  
581 inductive bias toward temporally consistent genera-  
582 tion, leading to improved forecasting accuracy and  
583 a substantial reduction in invalid molecular states.  
584 More broadly, EvoMD-LLM highlights the poten-  
585 tial of language-based models as general-purpose  
586 sequence learners for scientific simulations, sug-  
587 gesting a promising direction for bridging linguistic  
588 abstraction with time-resolved molecular dynamics  
589 in AI-driven materials discovery.

## 590 Limitations

591 While EvoMD-LLM demonstrates promising capa-  
592 bilities in modeling symbolic chemical evolution,  
593 several limitations remain to be addressed in future  
594 work:

595 **Generalization to Unseen Chemical Spaces.**  
596 Our evaluation focuses on the Mo-S CVD sys-  
597 tem. We selected this system not merely for data  
598 availability, but as a representative "complex proto-  
599 type" of inorganic synthesis: it features high-degree  
600 stochastic branching, reversibility, and multi-phase  
601 transitions (nucleation, etching, growth), which  
602 are often absent in linear organic reaction datasets.  
603 However, extending this framework to diverse  
604 chemical spaces, including heterogeneous biolog-  
605 ical systems, remains an open challenge for future  
606 scaling.

607 **Autoregressive Error Accumulation.** As ob-  
608 served in N-step prediction tasks, the model suf-  
609 fers from error accumulation typical of autoregres-  
610 sive generation, leading to performance degrada-  
611 tion over long horizons. Unlike numerical solvers  
612 that strictly enforce conservation laws, the current  
613 probabilistic generation may occasionally drift into  
614 physically invalid states. Integrating physical con-  
615 straints (e.g., mass conservation or energy consis-  
616 tency) directly into the loss function could mitigate  
617 this issue in future iterations.

618 **Loss of Fine-Grained Geometry.** While our  
619 coarse-grained symbolic representation efficiently  
620 filters thermal noise and captures high-level reac-  
621 tion logic, it inevitably discards fine-grained con-  
622 conformational information (e.g., precise bond lengths  
623 and angles). Consequently, EvoMD-LLM is cur-  
624 rently less suitable for tasks requiring exact geomet-  
625 ric verification. Future work could explore multi-  
626 modal architectures that jointly model symbolic  
627 evolution and geometric deformation to achieve  
628 fully comprehensive dynamic reasoning.

629 **Interpretability and Hallucination Risks.** The  
630 explanatory outputs provided by EvoMD-LLM are  
631 derived from aligning trajectory patterns with scien-  
632 tific knowledge learned during training, rather than  
633 ab initio derivation. Consequently, explanations of-  
634 ten rely on plausible but geometrically ungrounded  
635 terminology, suggesting retrieval-based associa-  
636 tion driven by pre-trained priors. Furthermore,  
637 the model occasionally over-interprets stochastic  
638 cues as deterministic stability guarantees and de-  
639 faults to generic linear narratives for rare intermedi-

ates, reflecting reduced sensitivity in low-frequency  
regimes.

## References

- Rishabh Agarwal, Avi Singh, Lei M. Zhang, Bernd  
Bohnet, Luis Rosias, Stephanie C. Y. Chan, Biao  
Zhang, Ankesh Anand, Zaheer Abbas, Azade Nova,  
John D. Co-Reyes, Eric Chu, Feryal M. P. Behba-  
hani, Aleksandra Faust, and Hugo Larochelle. 2024.  
Many-shot in-context learning. *Advances in Neural  
Information Processing Systems*, 37:76930–76966.
- Berni Julian Alder and Thomas Everett Wainwright.  
1957. Phase transition for a hard sphere system. *The  
Journal of chemical physics*, 27(5):1208.
- Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen,  
Xiyuan Zhang, Pedro Mercado, Huibin Shen, Olek-  
sandr Shchur, Syama Sundar Rangapuram, Sebastian  
Pineda Arango, Shubham Kapoor, Jasper Zschieg-  
ner, Danielle C. Maddix, Hao Wang, Michael W.  
Mahoney, Kari Torkkola, Andrew Gordon Wilson,  
Michael Bohlke-Schneider, and Yuyang Wang. 2024.  
[Chronos: Learning the language of time series.](#)  
*Transactions on Machine Learning Research*.
- Palash Bera and Jagannath Mondal. 2025. Accurate pre-  
diction of the kinetic sequence of physicochemical  
states using generative artificial intelligence. *Chem-  
ical Science*, 16(20):8735–8751.
- Daniil A Boiko, Robert MacKnight, and Gabe Gomes.  
2023. Emergent autonomous scientific research ca-  
pabilities of large language models. *arXiv preprint  
arXiv:2304.05332*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
Subbiah, Jared D. Kaplan, Prafulla Dhariwal, Arvind  
Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
Askell, Sandhini Agarwal, Ariel Herbert-Voss,  
Gretchen Krueger, Tom Henighan, Rewon Child,  
Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,  
Clemens Winter, and 12 others. 2020. Language  
models are few-shot learners. volume 33, pages  
1877–1901.
- J. Cavanagh, K. Sun, A. Gritsevskiy, D. Bagni, T. Head-  
Gordon, and T. D. Bannister. 2024. Smileyllama:  
Modifying large language models for directed chem-  
ical space exploration. In *NeurIPS 2024 Workshop  
on AI for New Drug Modalities*.
- Xiuying Chen, Tairan Wang, Taicheng Guo, Kehan Guo,  
Juexiao Zhou, Haoyang Li, Zirui Song, Xin Gao,  
and Xiangliang Zhang. 2025. Unveiling the power  
of language models in chemical research question  
answering. *Communications Chemistry*, 8(1):4.
- Seyone Chithrananda, Gabriel Grand, and Bharath  
Ramsundar. 2020. [Chemberta: Large-scale self-  
supervised pretraining for molecular property pre-  
diction.](#) *CoRR*, abs/2010.09885.

693	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. <a href="#">Think you have solved question answering? try arc, the ai2 reasoning challenge</a> . <i>CoRR</i> , abs/1803.05457.	
694		
695		
696		
697		
698	Connor W Coley, Wengong Jin, Luke Rogers, Timothy F Jamison, Tommi S Jaakkola, William H Green, Regina Barzilay, and Klavs F Jensen. 2019. A graph-convolutional neural network model for the prediction of chemical reactivity. <i>Chemical science</i> , 10(2):370–377.	
699		
700		
701		
702		
703		
704	OpenCompass Contributors. 2023. Opencompass: A universal evaluation platform for foundation models. <a href="https://github.com/open-compass/opencompass">https://github.com/open-compass/opencompass</a> .	
705		
706		
707		
708	Zhengzheng Dang, Zhichen Tang, Jixin Wu, Yide Chang, and Yanming Wang. 2025. Unraveling the reaction networks and key pathways during the gas phase stage in cvd synthesis of mos2. <i>Chemical Engineering Journal</i> , 503:157957.	
709		
710		
711		
712		
713	Michael Han Daniel Han and Unsloth team. 2023. <a href="#">Unsloth</a> .	
714		
715	Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, Xu Sun, and Zhifang Sui. 2024. A survey on in-context learning. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 1107–1128.	
716		
717		
718		
719		
720		
721	S. Hochreiter and J. Schmidhuber. 1997. Long short-term memory. <i>Neural Computation</i> , 9(8):1735–1780.	
722		
723		
724	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. Lora: Low-rank adaptation of large language models. <i>ICLR</i> , 1(2):3.	
725		
726		
727		
728	Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Ian Simon, Curtis Hawthorne, Noam Shazeer, Andrew M Dai, Matthew D Hoffman, Monica Dinulescu, and Douglas Eck. 2018. <a href="#">Music transformer: Generating music with long-term structure</a> . In <i>International Conference on Learning Representations</i> .	
729		
730		
731		
732		
733		
734	Mhd Hussein Murtada, Z Faidon Brotzakakis, and Michele Vendruscolo. 2025. Md-llm-1: A large language model for molecular dynamics. <i>arXiv e-prints</i> , pages arXiv–2508.	
735		
736		
737		
738	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. <i>Advances in neural information processing systems</i> , 33:9459–9474.	
739		
740		
741		
742		
743		
744		
745	Ilya Loshchilov and Frank Hutter. 2019. <a href="#">Decoupled weight decay regularization</a> . In <i>International Conference on Learning Representations</i> .	
746		
747		
	Feifei Luo, Jinglang Zhang, Qilong Wang, and Chunpeng Yang. 2025. Leveraging prompt engineering in large language models for accelerating chemical research. <i>ACS Central Science</i> , 11(4):511–519.	748
		749
		750
		751
	Meta AI Team. 2024. The llama 3 herd of models: the llama 3.1 family (405b parameters, 128k context window). Technical report. Technical Report.	752
		753
		754
	Mhd Hussein Murtada, Z Faidon Brotzakakis, and Michele Vendruscolo. 2024. Language models for molecular dynamics. <i>bioRxiv</i> , pages 2024–11.	755
		756
		757
	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35:27730–27744.	758
		759
		760
		761
		762
		763
		764
		765
		766
	Yi Ren, Yangjun Ruan, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. 2019. Fastspeech: Fast, robust and controllable text to speech. <i>Advances in neural information processing systems</i> , 32.	767
		768
		769
		770
	Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. <a href="#">Neural machine translation of rare words with subword units</a> . In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.	771
		772
		773
		774
		775
		776
		777
	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. Stanford alpaca: An instruction-following llama model.	778
		779
		780
		781
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 others. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	782
		783
		784
		785
		786
		787
		788
		789
	Sun-Ting Tsai, En-Jui Kuo, and Pratyush Tiwary. 2020. Learning molecular dynamics with simple language model built upon long short-term memory neural network. <i>Nature communications</i> , 11(1):5115.	790
		791
		792
		793
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	794
		795
		796
		797
		798
	Daniel S Wigh, Jonathan M Goodman, and Alexei A Lapkin. 2022. A review of molecular representation in the age of machine learning. <i>Wiley Interdisciplinary Reviews: Computational Molecular Science</i> , 12(5):e1603.	799
		800
		801
		802
		803

804 A. Yang, A. Li, B. Yang, B. Zhang, B. Hui, B. Zheng,  
805 and Z. Qiu. 2025. Qwen3 technical report. *arXiv*  
806 *preprint arXiv:2505.09388*.

807 Xianrui Zhong, Bowen Jin, Siru Ouyang, Yanzhen Shen,  
808 Qiao Jin, Yin Fang, Zhiyong Lu, and Jiawei Han.  
809 2025. [Benchmarking retrieval-augmented generation](#)  
810 [for chemistry](#). In *Second Conference on Language*  
811 *Modeling*.

## 812 A Data Processing and Statistics

813 In this section, we provide a detailed breakdown  
814 of the data processing pipeline, statistical charac-  
815 teristics, and the balancing strategy visualized in  
816 Figure 5.

### 817 A.1 Event Extraction and Filtering Pipeline

818 Our molecular event sequences originate from Re-  
819 active Molecular Dynamics (RMD) simulations of  
820 MoS<sub>2</sub> synthesis, as reported by [Dang et al. \(2025\)](#).  
821 The raw trajectories capture high-frequency atomic  
822 motions that do not directly correspond to sym-  
823 bolic chemical reactions. To align this data with  
824 language modeling, we employed a multi-stage  
825 pipeline:

- 826 • **Raw Extraction (Step 1):** We utilized a  
827 Depth-First Search (DFS) algorithm based  
828 on bond-order cutoffs to identify molecular  
829 species at every frame. This yielded an initial  
830 "Raw MD" dataset comprising 1,648,646  
831 events (Figure 5(a)).
- 832 • **Thermal Noise Reduction (Step 2):** Raw tra-  
833 jectories are dominated by transient thermal  
834 fluctuations where bonds vibrate but do not  
835 break. We filtered out all events with a persis-  
836 tence duration of  $\Delta t < 1$  ps. This resulted in  
837 the "Extracted" dataset of 237,673 events.
- 838 • **Temporal Band-Pass Filtering (Step 3):** To  
839 focus on the primary dynamic scales relevant  
840 for reaction forecasting, we further refined  
841 the dataset by retaining only events with du-  
842 rations  $\Delta t \in [10, 500]$  ps. This step removes  
843 extremely short-lived noise while segmenting  
844 ultra-long stable states, resulting in the "Fil-  
845 tered" dataset of 130,900 events.

### 846 A.2 Handling Data Imbalance

847 A critical challenge in MD data is the long-tail  
848 distribution of species. As shown in the word cloud  
849 in Figure 5(c) and the "Origin Data" histograms in  
850 Figure 5(b), a few dominant species (e.g., reactants

851 like MoS<sub>x</sub> precursors) account for the vast majority  
852 of observations, while critical transition states are  
853 rare.

854 Training a language model directly on this  
855 skewed distribution leads to trivial solutions where  
856 the model simply memorizes the most frequent  
857 tokens. To address this, we applied stratified sam-  
858 pling to balance the dataset across both molecular  
859 identity and duration intervals.

- 860 • **Effect of Balancing:** Figure 5(b) (bottom pan-  
861 els) demonstrates that after sampling, the dis-  
862 tributions of both molecular types and event  
863 durations become significantly more uniform.

- 864 • **Final Dataset:** This process yielded the fi-  
865 nal "Balanced" dataset containing 7,015 high-  
866 quality sequence pairs, which were used for  
867 fine-tuning EvoMD-LLM.

### 868 A.3 Temporal Characteristics

869 Figure 5(d) presents box plots of the existence du-  
870 rations for various molecular species in the final  
871 processed dataset. The distinct temporal distribu-  
872 tions (e.g., some species consistently show shorter  
873 lifetimes than others) confirm that duration is a se-  
874 mantic property intrinsic to each chemical species,  
875 justifying our use of Temporal Scaffolding to cap-  
876 ture these kinetic signatures.

## 877 B Implementation Details

878 To ensure the reproducibility of EvoMD-LLM, we  
879 provide detailed specifications of our software en-  
880 vironment, hardware infrastructure, and training  
881 configurations.

### 882 B.1 Software and Hardware Environment

883 We implemented EvoMD-LLM using the Unsloth  
884 framework, which optimizes memory usage and  
885 training speed for Llama-based models. The core  
886 software dependencies include:

- 887 • **Python:** 3.10
- 888 • **PyTorch:** 2.7.0 (with CUDA 12.6)
- 889 • **Unsloth:** 2025.5.6
- 890 • **Transformers:** 4.51.3
- 891 • **Peft:** 0.15.2

892 All experiments were conducted on a single  
893 consumer-grade NVIDIA RTX 4090D (24GB

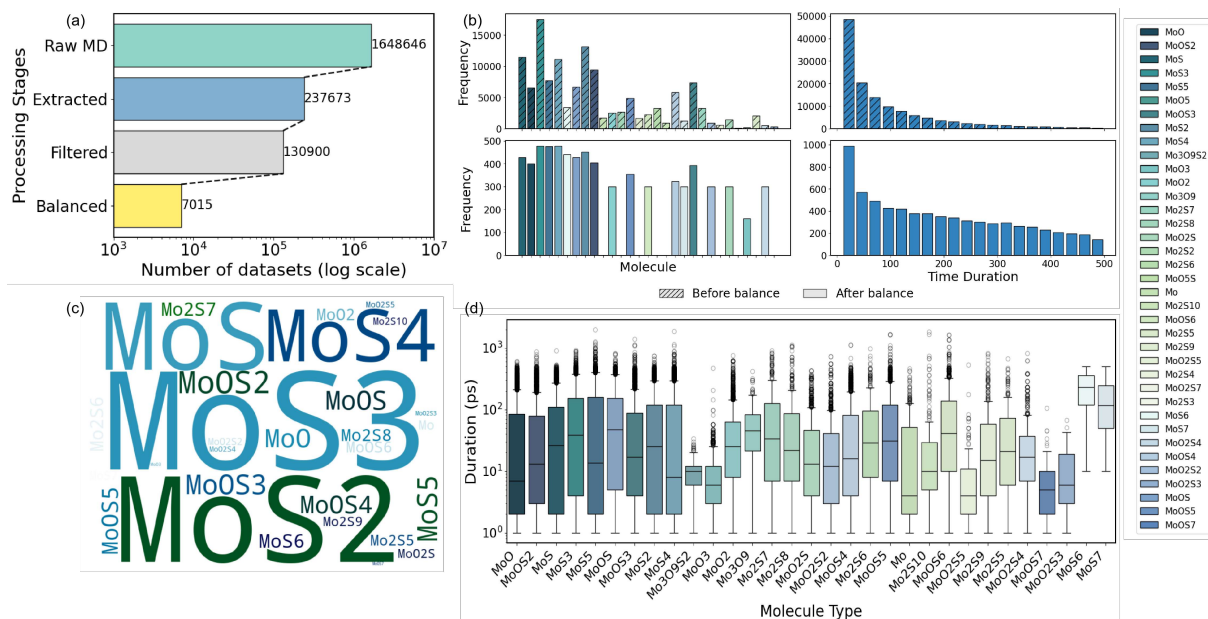


Figure 5: Data processing visualization. (a) Evolution of dataset scale across successive preprocessing stages, showing the reduction from raw events to high-quality balanced sequences. (b) Histograms comparing molecular type and duration distributions before (Origin) and after (Processed) stratified sampling, highlighting the mitigation of the long-tail problem. (c) Word cloud visualizing the dominance of specific species in the raw dataset. (d) Box plots showing the distinct distribution of existence durations (in ps) for different molecular species, reflecting their varying kinetic stabilities.

VRAM) GPU. We utilized mixed-precision training (bfloat16) to maximize computational efficiency without compromising numerical stability.

## B.2 Details on Quantitative Evaluation

All reported quantitative results are computed as the mean over multiple runs with different random seeds to ensure robust and reproducible evaluation. Specifically, for each experiment, we performed 10 runs with distinct random seeds and report the averaged metrics.

For the additional reasoning assessment on the ARC benchmark, we utilized the OpenCompass evaluation platform (Contributors, 2023). The model was evaluated using the HuggingFaceCausalLM.

## B.3 Hyperparameters and Training Costs

The detailed hyperparameters used for fine-tuning are listed in Table 5, optimized by grid search. We employed the LoRA technique, targeting all linear layers in the attention and feed-forward blocks.

With the configuration specified below, the full training process (covering both structured forecasting and Q&A tasks) took approximately 2.5 hours. The peak memory usage was controlled under 16GB thanks to the 4-bit quantization support

from Unsloth during the gradient calculation.

## B.4 Licenses and Terms of Use.

We use the LLaMA 3.1 model released by Meta under the LLaMA community license. The ARC-e and ARC-c benchmarks are publicly available for research purposes. All reactive molecular dynamics simulations and derived symbolic datasets were generated by the authors and do not contain personal or sensitive information. We plan to release the processed datasets and model checkpoints under a permissive research license upon acceptance.

## C Prompt Templates

In this section, we present the exact prompt templates used for training EvoMD-LLM and for eliciting qualitative reasoning. We employed a consistent system prompt to define the model’s role, while task-specific instructions were appended to the user queries to constrain the output format.

### C.1 Training and Prediction Prompts

For the supervised fine-tuning (SFT) stage and standard prediction tasks (1-step, N-step, and Backward), we used the following template structure.

**System Message.** This prompt sets the general behavioral constraints and defines the data format

Table 5: Detailed hyperparameters and training configuration for EvoMD-LLM. The model was fine-tuned using the LoRA method with the Unsloth framework for memory optimization.

Hyperparameter	Value
<i>General Configuration</i>	
Base Model	Llama 3.1 8B Instruct
Framework	Unsloth (TRL)
Precision	bfloat16 (bf16)
Random Seed	3407
Max Sequence Length	2048
<i>Optimization</i>	
Optimizer	AdamW (8-bit)
Learning Rate	$2 \times 10^{-4}$
Weight Decay	0.01
LR Scheduler	Linear
Warmup Steps	400
Num Epochs	2
<i>Batch Size Configuration</i>	
Per-Device Batch Size	2
Gradient Accumulation Steps	4
Effective Batch Size	8
<i>LoRA Configuration</i>	
Rank ( $r$ )	16
Alpha ( $\alpha$ )	16
Dropout	0
Bias	None
Target Modules	q, k, v, o, gate, up, down_proj

(molecule, duration).

You are an AI assistant to help me predict molecular sequence progression based on given molecular compositions and their existence durations and analysis. Each data point consists of a molecule and the duration it persists in the system, the unit of duration is ps. If the question is about predicting molecular sequences, format your answer as (molecule, time). Otherwise, answer normally.

**Task-Specific Instructions.** Different prediction tasks are distinguished by specific suffixes appended to the historical sequence.

### 1. Single-Step Prediction (Forward):

**Input:** The history sequence is {SEQUENCE\_HISTORY}, What is the next element? Output ONLY the next element in the format: (molecule, time). No explanation. No code. No extra words!

### 2. Multi-Step Prediction (N=2):

**Input:** The history sequence is {SEQUENCE\_HISTORY}, What are the next two elements? Output ONLY the next two elements in the format: (molecule, time). No explanation. No code. No extra words!

### 3. Backward Prediction:

**Input:** The history sequence is {SEQUENCE\_HISTORY}, What is the previous element? Output ONLY the previous element in the format: (molecule, time). No explanation. No code. No extra words!

## D Reasoning and Explanation Prompts

To assess the emergent explanatory capabilities of EvoMD-LLM, we utilized a structured prompt designed to constrain the output format.

It is important to note that while the model was fine-tuned with a mixture of symbolic MD sequences and general scientific Q&A pairs (to prevent catastrophic forgetting, see Section 2.5), it was never supervised on paired samples of (trajectory, textual explanation).

The training data for MD trajectories consisted solely of symbolic sequences (e.g., molecule tokens and duration values). Therefore, the detailed reasoning elicited by the prompt below reflects the model’s emergent ability to ground its general chemical knowledge (acquired from pre-training and Q&A regularization) into the specific context of the learned physical dynamics.

The prompt acts as a structural scaffold, directing the model to articulate its learned sequence patterns into explicit linguistic reasoning.

### Expert Simulator System Context.

You are an expert scientific simulator specializing in Reactive Molecular Dynamics (RMD) for Chemical Vapor Deposition (CVD) synthesis.

**System Context:** The reaction system involves the sulfidation of Mo309 precursors by S2 gas. Key dynamics include Oxygen-Sulfur exchange, structural relaxation, and thermal decomposition.

**Task Definition:** Your goal is to forecast the trajectory of chemical evolution. Each data point (Molecule, Duration) represents a distinct chemical state and its kinetic persistence (stability).

- A short duration implies a transient intermediate or transition state.
- A long duration implies a thermodynamically stable product or metastable trap.

**Reasoning Instruction.** After the model generates a prediction, we prompt it to explain the physical rationale using the following template:

**Task:** You are provided with a historical trajectory of molecular species and their durations.

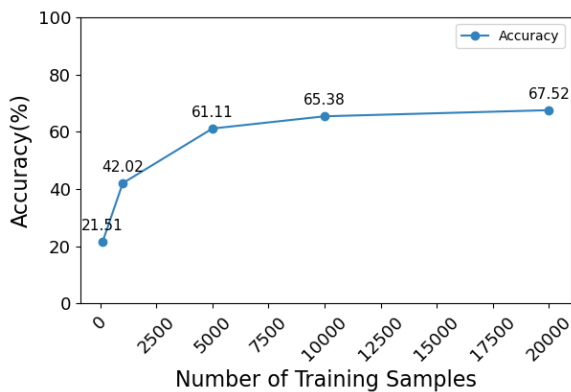
1026 **History Sequence:** {history\_seq} 1059  
1027 **Your Model Prediction:** ({predict\_res}) 1060  
1028 **Instructions:** Provide a scientific 1061  
1029 explanation for this transition. Your 1062  
1030 response must:

- 1031 1. **Mechanism:** Analyze the change in 1063  
1032 stoichiometry from the last history 1064  
1033 step to the predicted step. What 1065  
1034 specific chemical process drives 1066  
1035 this transformation?
- 1036 2. **Stability:** Analyze the predicted 1067  
1037 duration ({duration}). What does 1068  
1038 this specific timescale imply 1069  
1039 about the thermodynamic state or 1070  
1040 kinetic stability of the predicted 1071  
1041 molecule?
- 1042 3. **Format:** Write in strict, concise 1072  
1043 Academic English. 1073

1044 Your answer must be in academic English, 1074  
1045 concise, and only include the reasoning 1075  
1046 (no extra content, no repetition). 1076

## 1047 E Sample Efficiency Analysis 1077

1048 To investigate whether our dataset size is a bot- 1078  
1049 tleneck for performance, we conducted a scaling 1079  
1050 analysis by training EvoMD-LLM on subsets of the 1080  
1051 training data ranging from roughly 100 to 20,000 1081  
1052 samples. Figure 6 illustrates the 1-step prediction 1082  
1053 accuracy as a function of data quantity. 1083



1044 Figure 6: Learning curve of EvoMD-LLM. The plot 1084  
1045 shows 1-step prediction accuracy scaling with training 1085  
1046 data size. The model exhibits strong few-shot general- 1086  
1047 ization, reaching over 60% accuracy with only 5,000 1087  
1048 samples, and shows signs of performance saturation 1088  
1049 beyond 10,000 samples, indicating that the current dataset 1089  
1050 size is sufficient for capturing the core dynamics. 1090

1054 As shown in Figure 6, the model exhibits high 1091  
1055 sample efficiency. 1092

- 1056 • **Rapid Syntax Acquisition (0-5k):** Accuracy 1093  
1057 surges from 21.5% to 61.1% within the first 1094  
1058 5,000 samples. This steep rise suggests that 1095

1059 the LLM, leveraging its pre-trained capabil- 1060  
1061 ities, rapidly aligns with the "grammar" of 1062  
1063 molecular evolution (syntax and basic stoi- 1064  
1065 chiometry) with minimal data. 1066

- 1067 • **Performance Saturation (10k-20k):** As data 1068  
1069 volume doubles from 10,000 to 20,000, accu- 1070  
1071 racy gains moderate (from 65.4% to 67.5%). 1072  
1073 This plateau indicates that the model has effec- 1074  
1075 tively captured the majority of the learnable 1076  
1077 patterns within the current domain. 1078

1079 This analysis confirms that our dataset size 1080  
1081 (~20k total samples) is robust. The constraint on 1082  
1083 further performance improvement is likely not the 1084  
1084 quantity of raw data, but the inherent stochasticity 1085  
1085 of the chemical system itself. 1086

## 1087 F Detailed Error Analysis 1088

1089 To evaluate the reliability of EvoMD-LLM beyond 1090  
1091 standard accuracy metrics, we conducted a fine- 1091  
1092 grained analysis of the prediction errors on the test 1092  
1093 set. 1093

### 1094 F.1 Zero Hallucinations and Chemical 1094 1095 Validity 1095

1096 A critical finding of our analysis is that EvoMD- 1096  
1097 LLM exhibits zero hallucinations regarding chemi- 1097  
1098 cal validity. 100% of the incorrect predictions 1098  
1099 correspond to chemically valid molecular formulas 1099  
1100 that exist within the reaction network (e.g., pre- 1100  
1101 dicting a valid intermediate like  $MoS_3$  but at an 1101  
1102 incorrect time step). 1102

1103 This stands in sharp contrast to generic LLMs, 1103  
1104 which often generate physically impossible stoi- 1104  
1105 chiometry when fine-tuned on scientific data. 1105  
1106 The absence of hallucinations confirms that our 1106  
1107 symbolic tokenization strategy has successfully 1107  
1108 grounded the model in the compositional gram- 1108  
1109 mar of the chemical system, constraining its errors 1109  
1110 to the domain of physical kinetics rather than gen- 1110  
1111 erative syntax. 1111

### 1112 F.2 Physical Symmetry of Kinetic Mismatch 1112

1113 Since all errors are valid "Kinetic Mismatches," we 1113  
1114 further decomposed them to determine if the model 1114  
1115 exhibits systematic bias. Given that the primary re- 1115  
1116 action mechanism is sulfidation (replacing Oxygen 1116  
1117 with Sulfur), we classified the errors based on the 1117  
1118 stoichiometry of the predicted species relative to 1118  
1119 the ground truth: 1119

- 1105 • **Under-Sulfidation (Lagging):** The predicted  
1106 molecule contains fewer Sulfur atoms than the  
1107 ground truth ( $S_{pred} < S_{true}$ ). The model pre-  
1108 dicta a precursor state, effectively "lagging"  
1109 behind the true trajectory.
- 1110 • **Over-Sulfidation (Too Fast):** The pre-  
1111 dicted molecule contains more Sulfur atoms  
1112 ( $S_{pred} > S_{true}$ ). The model anticipates the  
1113 reaction progressing faster than reality.
- 1114 • **Oxygen Deviation:** The Sulfur content is cor-  
1115 rect, but the Oxygen stoichiometry differs, re-  
1116 flecting minor inaccuracies in secondary de-  
1117 oxidation steps.

1118 As illustrated in Figure 7, the error distribution  
1119 is remarkably balanced. The near-perfect symme-  
1120 try between Under-Sulfidation (45.0%) and Over-  
1121 Sulfidation (44.1%) suggests that EvoMD-LLM  
1122 does not suffer from a systematic drift towards  
1123 either accelerating or retarding the reaction kinet-  
1124 ics. Instead, the errors likely reflect the inherent  
1125 stochasticity of Molecular Dynamics simulations,  
1126 where the exact timing of atomic transitions fluctu-  
1127 ates around the mean path. The model effectively  
1128 captures the central tendency of the reaction, with  
1129 errors representing unbiased variance.

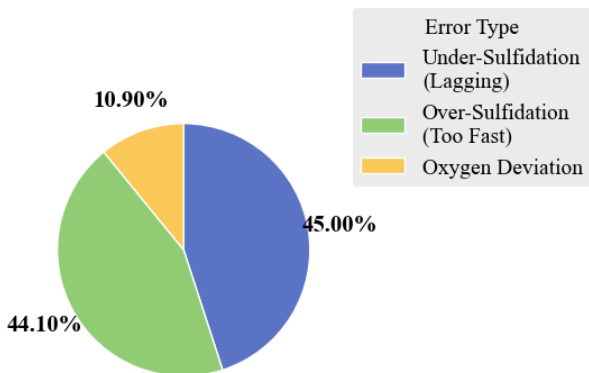


Figure 7: Breakdown of Kinetic Mismatch Errors. The distribution shows a balanced split between Under-Sulfidation (Blue, 45.0%) and Over-Sulfidation (Green, 44.1%). This symmetry indicates that the model's temporal errors are stochastic rather than biased. Oxygen Deviation (Yellow, 10.9%) represents minor stoichiometric noise.

## 1130 G Additional Qualitative Examples

1131 To provide a more comprehensive view of the  
1132 model's qualitative behavior, we present additional

prediction and reasoning examples in Table 6. Unlike the selected cases in the main paper, these examples include both successful predictions and characteristic failure modes, and are intended to illustrate typical patterns rather than exhaustive coverage.

1133  
1134  
1135  
1136  
1137  
1138

Input History Sequence	Prediction	Generated Reasoning
<b>Case 1: Over-sulfidation Drift</b>		
(MoS <sub>4</sub> , 8); (Mo <sub>2</sub> S <sub>7</sub> , 12); (MoS <sub>3</sub> , 41)	<b>Mol:</b> MoS <sub>6</sub> <b>Time:</b> 10 ps	The transition is attributed to <b>continued sulfur adsorption</b> leading to higher sulfur coordination. However, the extremely short predicted duration suggests a <b>transient and unstable configuration</b> .
<b>Case 2: Stability Overestimation</b>		
(Mo <sub>2</sub> S <sub>5</sub> , 7); (MoS <sub>3</sub> , 19); (MoS <sub>4</sub> , 32)	<b>Mol:</b> MoS <sub>4</sub> <b>Time:</b> 125 ps	The model interprets repeated appearance of MoS <sub>4</sub> as evidence of <b>kinetic stability</b> , potentially overestimating persistence due to frequency bias.
<b>Case 3: Reduction Shortcut</b>		
(MoOS, 4); (MoO <sub>2</sub> , 9); (MoO <sub>2</sub> , 12)	<b>Mol:</b> MoO <b>Time:</b> 488 ps	The transition is explained as <b>oxygen loss under reducing conditions</b> . The long duration reflects a <b>stable oxide plateau</b> , though intermediate pathways are omitted.
<b>Case 4: Missed Cluster Growth</b>		
(Mo <sub>2</sub> S <sub>7</sub> , 5); (Mo <sub>2</sub> S <sub>5</sub> , 11); (MoS <sub>4</sub> , 18)	<b>Mol:</b> MoS <sub>4</sub> <b>Time:</b> 495 ps	The model favors <b>structural persistence</b> over further aggregation, overlooking potential <b>cluster growth</b> toward higher sulfur content.
<b>Case 5: Correct Stability Tracking</b>		
(MoS <sub>4</sub> , 2); (MoOS <sub>4</sub> , 9); (MoOS <sub>5</sub> , 14)	<b>Mol:</b> MoOS <sub>5</sub> <b>Time:</b> 62 ps	The prediction reflects <b>incremental sulfur incorporation</b> while correctly identifying a <b>moderately persistent intermediate</b> .
<b>Case 6: Linear Growth Bias</b>		
(MoS <sub>3</sub> , 6); (MoS <sub>4</sub> , 11); (MoS <sub>5</sub> , 9)	<b>Mol:</b> MoS <sub>6</sub> <b>Time:</b> 83 ps	The reasoning assumes monotonic sulfur accumulation, exhibiting a <b>linear growth bias</b> that neglects possible sulfur loss or rearrangement.
<b>Case 7: Transient Misclassification</b>		
(MoO <sub>2</sub> , 3); (MoOS <sub>2</sub> , 7); (MoS <sub>2</sub> , 5)	<b>Mol:</b> MoS <sub>2</sub> <b>Time:</b> 214 ps	Although MoS <sub>2</sub> is predicted as long-lived, the explanation likely overstates <b>kinetic trapping</b> for what may be a transient configuration.
<b>Case 8: Rare Species Generalization</b>		
(Mo <sub>3</sub> S <sub>11</sub> , 4); (Mo <sub>3</sub> S <sub>12</sub> , 6); (Mo <sub>3</sub> S <sub>13</sub> , 3)	<b>Mol:</b> Mo <sub>3</sub> S <sub>13</sub> <b>Time:</b> 57 ps	The model produces a generic explanation invoking <b>coordination saturation</b> , reflecting limited specificity for <b>rare cluster species</b> .
<b>Case 9: Oxygen Retention Bias</b>		
(MoO <sub>3</sub> , 5); (MoO <sub>2</sub> , 18); (MoOS <sub>2</sub> , 7)	<b>Mol:</b> MoO <sub>2</sub> <b>Time:</b> 301 ps	The prediction favors oxygen-rich species, suggesting a bias toward <b>oxide persistence</b> despite emerging sulfidation signals.
<b>Case 10: Competing Pathway Suppression</b>		
(MoS <sub>2</sub> , 14); (MoOS <sub>3</sub> , 9); (MoS <sub>3</sub> , 11)	<b>Mol:</b> MoS <sub>4</sub> <b>Time:</b> 92 ps	The explanation emphasizes sulfur addition while suppressing alternative <b>desulfurization or rearrangement pathways</b> .

Table 6: Additional qualitative examples generated by EvoMD-LLM. These cases complement the main-paper examples by covering a broader range of behaviors, including correct predictions, stability overestimation, linear growth bias, and generic reasoning for rare species. Highlights indicate inferred reaction mechanisms (**Red**) and stability or metastability judgments (**Blue**).