

28 January 2026

# A unified framework for automated transition state generation to accelerate mechanistic exploration in organic synthesis

Li-Cheng Xu<sup>1</sup>, Junyi An<sup>1</sup>, Weiqi Liu<sup>1,2</sup>, Yun-Fei Shi<sup>1</sup>, Chong-Lei Ji<sup>3</sup>, Fenglei Cao<sup>1</sup>, Yuan Qi<sup>1,2</sup>

1. Shanghai Academy of Artificial Intelligence for Science

2. Artificial Intelligence Innovation and Incubation Institute Fudan University

3. School of Physical Science and Technology ShanghaiTech University

## Abstract

Transition states (TS) are pivotal for elucidating reaction mechanisms, but their calculation constitutes a major bottleneck. While generative AI provides a route to automate transition state generation, existing methods remain confined to simplified chemical systems, lacking generalizability to complex reactions like transition metal catalysis. Here, we present a unified framework for general-purpose TS discovery, comprising the UniTS-Lib library of 4,391 high-fidelity structures spanning 42 elements and diverse chemical transformations, coupled with the UniTS-Gen diffusion model that generates 3D TS configurations from 2D reactant graphs using a higher-degree equivariant network. Validation demonstrates UniTS-Gen's superior structural accuracy and robust generalization to unseen chemical systems. We show that UniTS-Gen provides highly reliable initial guesses and actively accelerates discovery by locating kinetically preferred TS conformation. This work provides a scalable and transferable solution for automating mechanistic studies in organic synthesis.

# A unified framework for automated transition state generation to accelerate mechanistic exploration in organic synthesis

Li-Cheng Xu,<sup>1\*</sup> Junyi An,<sup>1</sup> Weiqi Liu,<sup>1,2</sup> Yun-Fei Shi,<sup>1</sup> Chong-Lei Ji,<sup>3</sup> Fenglei Cao,<sup>1</sup> Yuan Qi<sup>1,2</sup>

<sup>1</sup>Shanghai Academy of Artificial Intelligence for Science, Shanghai, 200232, China

<sup>2</sup>Artificial Intelligence Innovation and Incubation Institute, Fudan University, Shanghai, 201203, China

<sup>3</sup>School of Physical Science and Technology, ShanghaiTech University, Shanghai, 200120, China

\*Email: xulicheng@sais.org.cn

**Abstract:** Transition states (TS) are pivotal for elucidating reaction mechanisms, but their calculation constitutes a major bottleneck. While generative AI provides a route to automate transition state generation, existing methods remain confined to simplified chemical systems, lacking generalizability to complex reactions like transition metal catalysis. Here, we present a unified framework for general-purpose TS discovery, comprising the UniTS-Lib library of 4,391 high-fidelity structures spanning 42 elements and diverse chemical transformations, coupled with the UniTS-Gen diffusion model that generates 3D TS configurations from 2D reactant graphs using a higher-degree equivariant network. Validation demonstrates UniTS-Gen's superior structural accuracy and robust generalization to unseen chemical systems. We show that UniTS-Gen provides highly reliable initial guesses and actively accelerates discovery by locating kinetically preferred TS conformation. This work provides a scalable and transferable solution for automating mechanistic studies in organic synthesis.

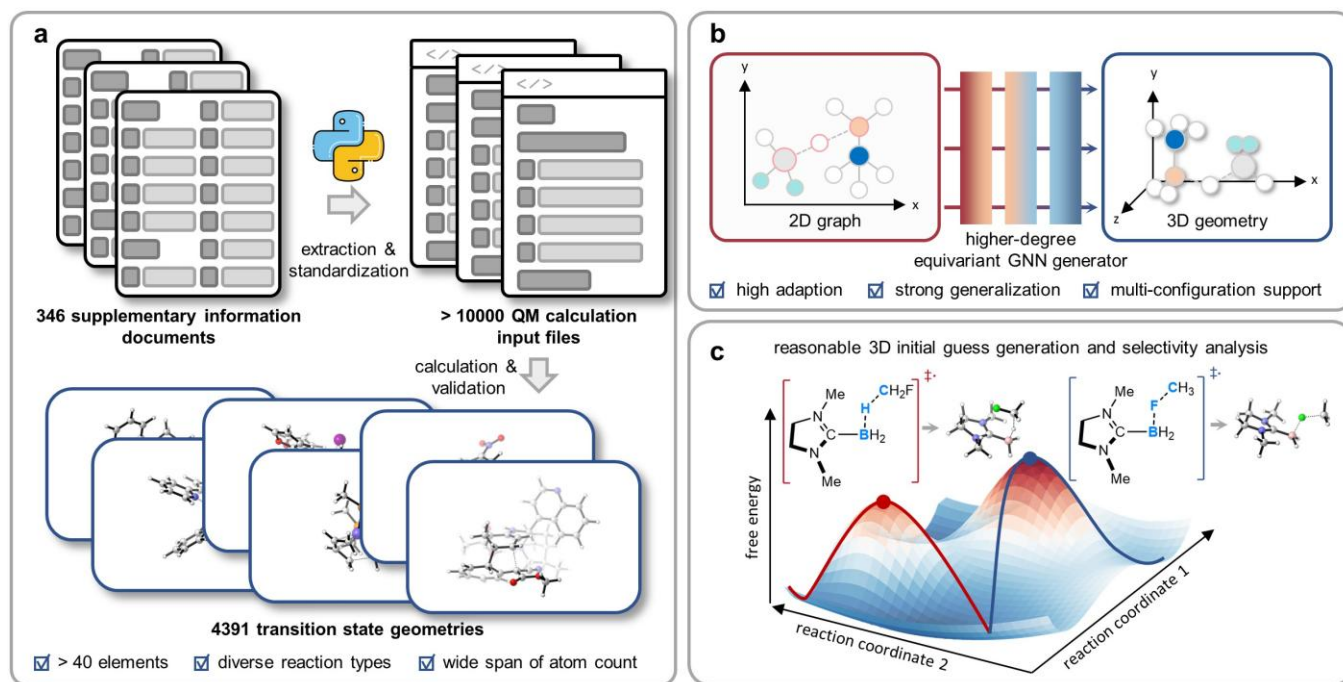
## Introduction

The transition state (TS), representing the saddle point on the reaction coordinate of highest potential energy, dictates the kinetics of chemical transformations and is fundamental to elucidating reaction mechanisms<sup>1-3</sup>. In the rapidly evolving fields of precise synthesis, including asymmetric organometallic catalysis<sup>4</sup>, chiral radical chemistry<sup>5</sup>, and bio-orthogonal chemistry<sup>6</sup>, analyzing TS structures is indispensable for the rational design of effective chemical systems that optimize reactivity and selectivity<sup>7,8</sup>. Despite the critical importance of locating these transient structures, prevailing computational strategies remain heavily reliant on rational initial guesses<sup>9</sup>. Generating plausible TS initial structures typically necessitates labor-intensive trial-and-error and extensive expert supervision. This manual, expertise-dependent paradigm scales poorly and is increasingly unable to keep pace with the rapid growth of novel reaction systems.

Deep learning has precipitated a paradigm shift in structural science, most notably revolutionizing the prediction and design of complex biological macromolecules such as proteins<sup>10,11</sup> and nucleic acids<sup>12</sup>. Drawing inspiration from this transformative success, the chemical community has begun to harness generative AI to automate the identification of transition states<sup>13</sup>, aiming to transcend the constraints of expert intuition. The methodology in this field has evolved rapidly, progressing from early heuristics utilizing geometric distance-based methods<sup>14-17</sup> to sophisticated diffusion-based generative frameworks<sup>18-20</sup>. Representative studies by Kim et al.<sup>19</sup> and Duan et al.<sup>18,20</sup> independently established the feasibility of applying diffusion models to generate plausible TS geometries, using 2D molecular graphs and 3D coordinates as inputs, respectively. Building on this foundation, Liu et al.<sup>21</sup> further expanded the frontier, extending

end-to-end generation from unimolecular rearrangements to a particular class of bimolecular system. Collectively, these contributions have significantly advanced the automation of reaction mechanism investigation<sup>13</sup>.

Despite this progress, current generative frameworks remain confined to simplified chemical transformations, typically characterizing unimolecular systems or singular types of bimolecular reactions. Existing datasets are often limited to fewer than 5 types of heavy elements and restrict molecular size to under 30 atoms. Consequently, developing a general-purpose TS generation model capable of handling the complexity of authentic synthetic scenarios, such as transition metal catalysis<sup>22</sup>, remains a formidable challenge. To deploy generative AI effectively in realistic mechanistic studies, two critical challenges must be addressed: first, the construction of a comprehensive TS dataset covering diverse reaction types, elemental compositions, and broad atom-count distributions; and second, the development of a generative architecture with sufficient expressive power to model the immense structural diversity of transition states from relatively sparse data, while maintaining robust generalization to unseen chemical manifolds.



**Fig. 1: A unified framework for automated transition state discovery.** **a**, The automated workflow for constructing UniTS-Lib. This pipeline processes 346 Supplementary Information documents to standardize over 10,000 initial quantum mechanical (QM) inputs, which are subjected to rigorous DFT optimization and validation to yield 4,391 high-fidelity transition state geometries. **b**, The architecture of the UniTS-Gen model. The model utilizes a higher-degree equivariant GNN as its denoising kernel to generate 3D TS geometries from a 2D molecular graph input, enabling robust generalization and multi-configuration support. **c**, Demonstration of UniTS-Gen in generating reasonable 3D initial guesses and its utility in subsequent selectivity analysis.

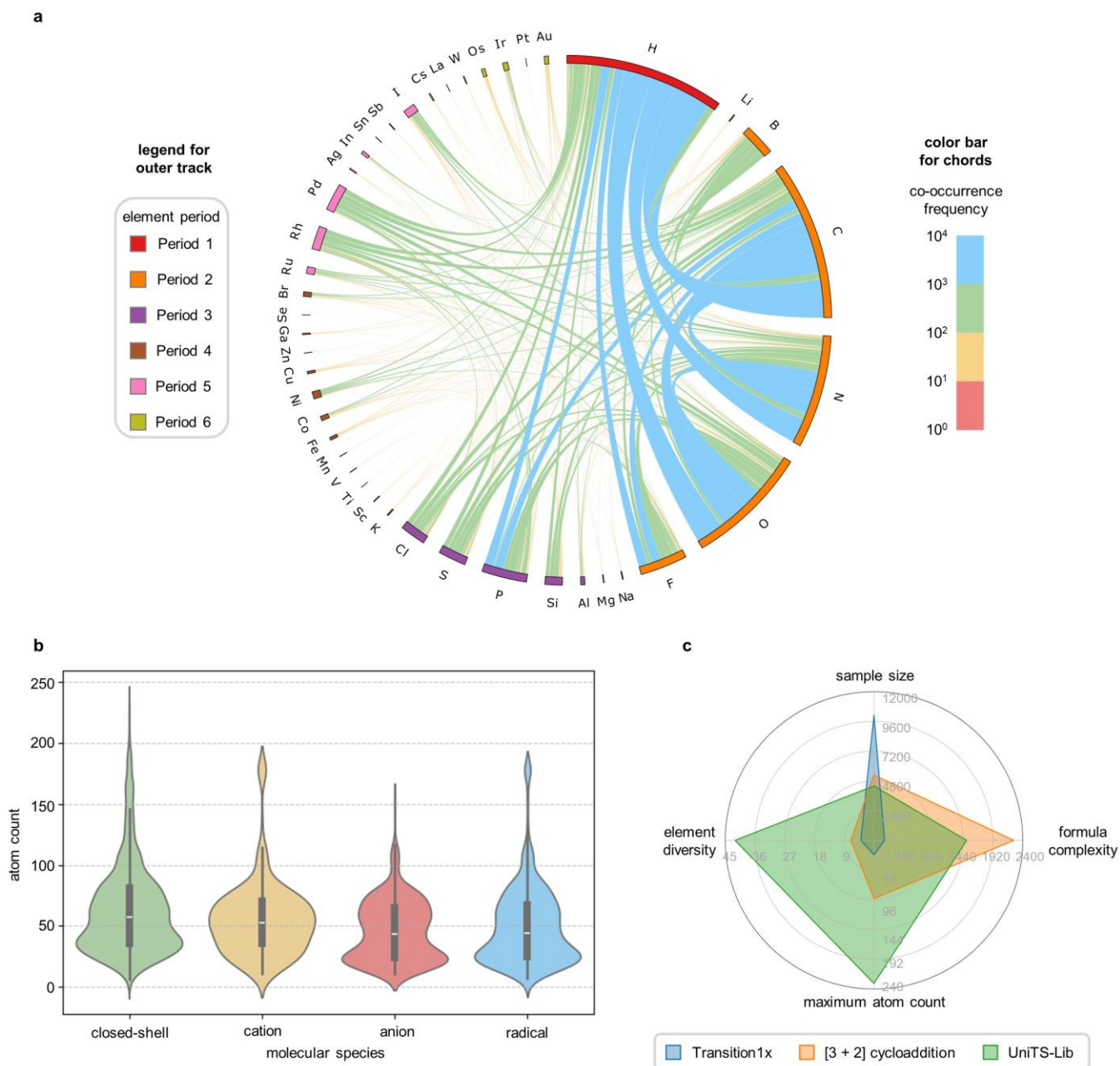
Here, we bridge this gap by introducing a unified framework encompassing the automated construction of a complex TS dataset and the development of a high-accuracy, highly-generalizable end-to-end generation model for predicting high-quality initial guesses of 3D TS structures from 2D molecular graphs. We first utilized the data construction workflow within this framework, a process that extracts coordinates of TS from literature Supplementary Information

documents followed by meticulous density functional theory (DFT) level TS optimization and frequency-based validation, enabling us to establish UniTS-Lib, a library of 4,391 high-fidelity TS structures derived from authentic mechanistic explorations of complex organic synthesis (Fig. 1a). To address the inherent challenge posed by the rich structural diversity and comparative data sparsity of UniTS-Lib, we developed UniTS-Gen, a diffusion model with a novel architecture that centers on a well-designed higher-degree SE(3)-equivariant graph neural network (GNN)<sup>23-25</sup> specifically designed as its denoising kernel (Fig. 1b). By carefully design, UniTS-Gen takes only the 2D molecular graph of the reactants as input and generates the corresponding 3D TS geometries exhibiting rich conformational and configurational diversity. We extensively validated UniTS-Gen across three benchmarks—the classic Transition1x<sup>26</sup>, a [3+2]-cycloaddition dataset<sup>27</sup> known for modeling bio-orthogonal reactions and exhibiting wide elemental compositional variance, and our newly constructed UniTS-Lib. A series of comprehensive experiments robustly demonstrate that UniTS-Gen achieves high structural accuracy, superior generalization capability, and the ability to capture distinct TS structures critical for subsequent kinetic barrier estimation (Fig. 1c). This work thus establishes a robust and unified framework that provides a powerful solution to the multifaceted challenges inherent in automating realistic chemical mechanistic discovery.

## Results

### UniTS-Lib: a versatile and high-fidelity TS library

To establish a robust data foundation for training the general-purpose TS generation model, we developed an automated workflow (Fig. 1a) that efficiently extracts atomic types and coordinates from the Supplementary Information of organic synthesis literature and standardizes them into quantum mechanical (QM) input files. By processing 346 Supplementary Information documents, we successfully extracted and standardized over 10,000 initial QM calculation inputs. Notably, during the file generation process, we explicitly considered the radical, anionic, and cationic states for species that are not stable closed-shell neutrals, thereby ensuring that our data broadly covers charged and radical TS structures commonly encountered in mechanistic studies. We then subjected these initial structures to rigorous DFT-level TS optimization at a consistent accuracy level, followed by frequency analysis combined with graph-matching verification to validate the saddle point identity, yielding a final set of 4,391 high-fidelity TS structures. We deliberately adhered to a standardized protocol rather than employing exhaustive, ad-hoc recovery strategies for non-converging structures. The substantial attrition rate observed from the initial pool to the final dataset highlights the intrinsic difficulty of transition state optimization, a task notoriously sensitive to both the precise initial geometry and the chosen computational methodology. Even when starting from reasonable, literature-derived coordinates, variations in basis sets or functionals can frequently impede convergence. This rigorous filtration process underscores the inherent scarcity of high-fidelity TS data for complex systems and illustrates the formidable challenge of locating correct saddle points.



**Fig. 2: Analysis of the structural and elemental complexity of UniTS-Lib and its comparison to existing TS datasets.** **a**, The element co-occurrence chord diagram illustrates the elemental richness and relationships within UniTS-Lib, which features 42 elements including main-group and transition metals essential for catalysis. The outer track length corresponds to individual element frequency, while the chord thickness and color (cooler tones indicate higher frequency) denote the volume of pairwise co-occurrence. **b**, Violin plots characterize the atom count distribution across four major molecular species (closed-shell neutral, cation, anion, and radical transition states). All species exhibit broad size distributions with median atom counts around 50. **c**, A multi-dimensional radar chart quantitatively compares UniTS-Lib against the classic Transition1x and the representative [3+2] cycloaddition datasets across sample size, element diversity, maximum atom count, and formula complexity.

Based on our unbiased data curation strategy, UniTS-Lib encompasses over 10 distinct transformation types central to mechanistic investigations, including oxidative addition, transmetallation and cycloaddition (see Supplementary Fig. 2 for details). The dataset features an exceptionally broad elemental distribution involving 42 elements, with the maximum atomic number being 79 (gold). As elucidated by the element co-occurrence chord diagram in Fig. 2a, the arc length of the outer track corresponds to the frequency of each element, while the thickness and color intensity (with cooler tones indicating higher frequency) of the internal chords depict the volume of pairwise co-occurrence. This visualization reveals that, beyond main-group elements, transition metals pivotal to catalysis, specifically the iron and platinum groups, are substantially represented.

Beyond elemental diversity, UniTS-Lib exhibits remarkable structural complexity. As shown in Fig. 2b, violin plots characterize the atom count distribution across four primary molecular species: closed-shell neutral, cationic, anionic, and radical transition states. All species display broad distributions with medians centered around 50 atoms. Notably, neutral closed-shell molecules exhibit the widest span, reaching a maximum of 231 atoms; in total, the entire dataset contains 131 structures exceeding 150 atoms. This extensive coverage across diverse electronic and spin states underscores the dataset's capacity to represent complex chemical spaces.

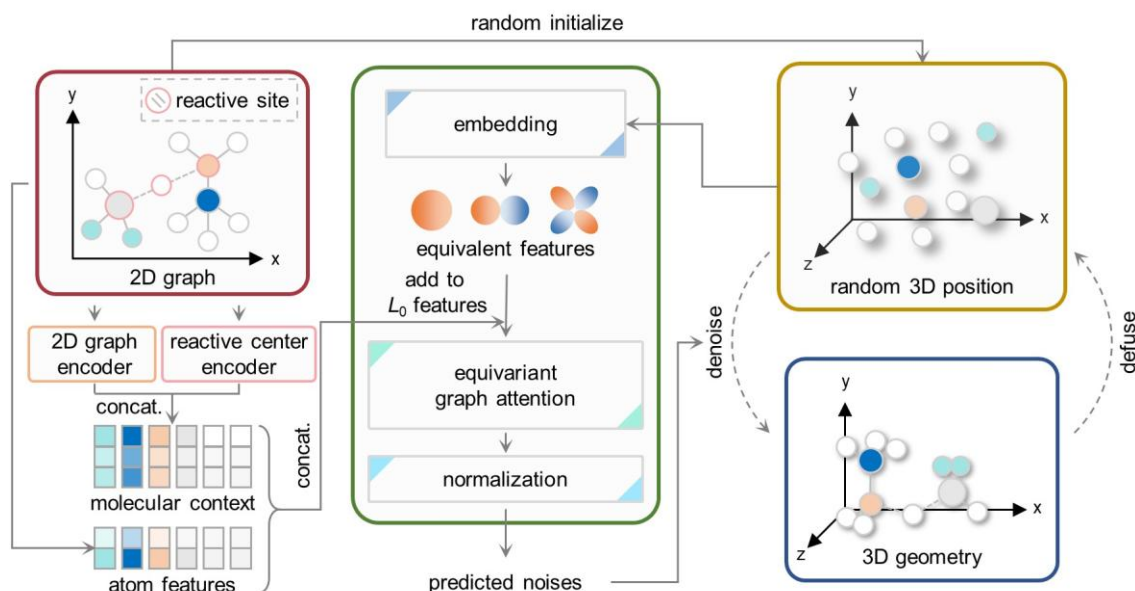
To contextualize the data complexity of UniTS-Lib, we performed a multi-dimensional quantitative benchmark against the classic Transition1x dataset<sup>26</sup> and the representative [3+2] dipolar cycloaddition dataset<sup>27</sup> (Fig. 2c). While Transition1x possesses the largest sample size, its significantly lower scores across the three structural complexity metrics identify it as a baseline for simpler TS generation tasks. The [3+2] dataset achieves the highest chemical formula complexity due to combinatorial dipole-dipolarophile pairs. In contrast, UniTS-Lib leads in both elemental diversity and molecular size breadth while maintaining substantial formula complexity. Although reaction type complexity was excluded from the quantitative plot due to categorization subjectivity, UniTS-Lib inherently encompasses a significantly broader range of reaction types compared to unimolecular or single-reaction benchmarks. Consequently, the integration of these factors positions UniTS-Lib as the most elementally and structurally complex dataset among those compared. Furthermore, the inherent difficulty of acquiring such high-fidelity data results in a naturally smaller sample size, thereby exacerbating the challenge for training general-purpose generative AI models.

### **UniTS-Gen: a higher-degree equivariant generation model**

With UniTS-Lib dataset in hand, we next sought to develop a generative model capable of learning from such structural and elemental diversity. Recent diffusion-based TS generation models have achieved notable accuracy by using 3D reactant and product geometries as inputs<sup>18,20</sup>. While effective for well-defined benchmark systems, this strategy is less practical for complex, mechanistically underexplored reactions, where obtaining reliable 3D product structures is non-trivial and may inadvertently restrict the conformational and configurational exploration of transition states. We therefore designed UniTS-Gen to generate 3D TS geometries directly from 2D reactant graphs, using only atomic indices designating the reaction site. This input scheme preserves flexibility in configurational sampling and lowers the barrier to application in novel reaction systems.

To achieve sufficient structural accuracy and generalization on such complex, diverse, yet relatively sparse TS data, the core architecture must be capable of deeply capturing and learning molecular 3D structural information. Inspired by the exceptional performance of higher-degree equivariant GNNs in molecular property prediction<sup>23-25</sup>, which leverage higher-degree spherical harmonics to precisely capture subtle changes in molecular 3D structure, we posited that this model type serves as an ideal denoising kernel for generating complex 3D TS structures within a diffusion

framework. Accordingly, to address the specific challenge of capturing the intricate 3D geometry of transition states, we designed a novel, higher-degree equivariant denoising kernel, termed HiEGNN, and built the UniTS-Gen framework upon this kernel, as illustrated in Fig. 3.



**Fig. 3: Architecture of the UniTS-Gen conditional diffusion model for 3D TS generation.** The generative framework is based on the higher-degree equivariant GNN as the denoising kernel. The model is engineered as a conditional diffusion process where the inputs are restricted to the 2D molecular graph of the reactants and atomic indices defining the reactive site (red box). These inputs are encoded into a molecular context that guides the reverse diffusion process. The HiEGNN denoising module (green box) iteratively refines the initial random 3D coordinates (yellow box) based on both the noisy structure and the molecular context. The HiEGNN utilizes higher-degree SE(3)-equivariance to accurately model complex 3D structures while inherently respecting rotational and translational equivariances, ultimately predicting the required noise to generate the precise 3D TS geometry (blue box).

UniTS-Gen is engineered as a conditional diffusion model built around the HiEGNN kernel. Its input consists solely of the 2D molecular graph of the reactants (or intermediates), augmented with atomic indices that define the reactive site. These indices serve as an essential conditioning feature for distinguishing among multiple possible active sites. These graph-based inputs are encoded by a 2D molecular encoders<sup>28</sup> to produce a conditioning context that guides the reverse diffusion process. Starting from random coordinates initialized based on the input atom count, the model iteratively denoises the structure toward a chemically meaningful 3D TS geometry.

The core of UniTS-Gen is the HiEGNN denoising module (Fig. 3, green box), which leverages higher-degree SE(3)-equivariance, a vital property that ensures the model inherently respects the rotational and translational equivariances of 3D space, which is important for the accurate and stable prediction of complex coordinates. During the reverse diffusion process, HiEGNN refines the 3D coordinates by integrating the instantaneous noisy structure with the molecular context, resulting in the generation of precise TS geometries that exhibit the rich conformational and configurational diversity required for comprehensive mechanistic analysis (see Supplementary Fig. 4 and Fig. 5 for details of the diffusion and sampling process).

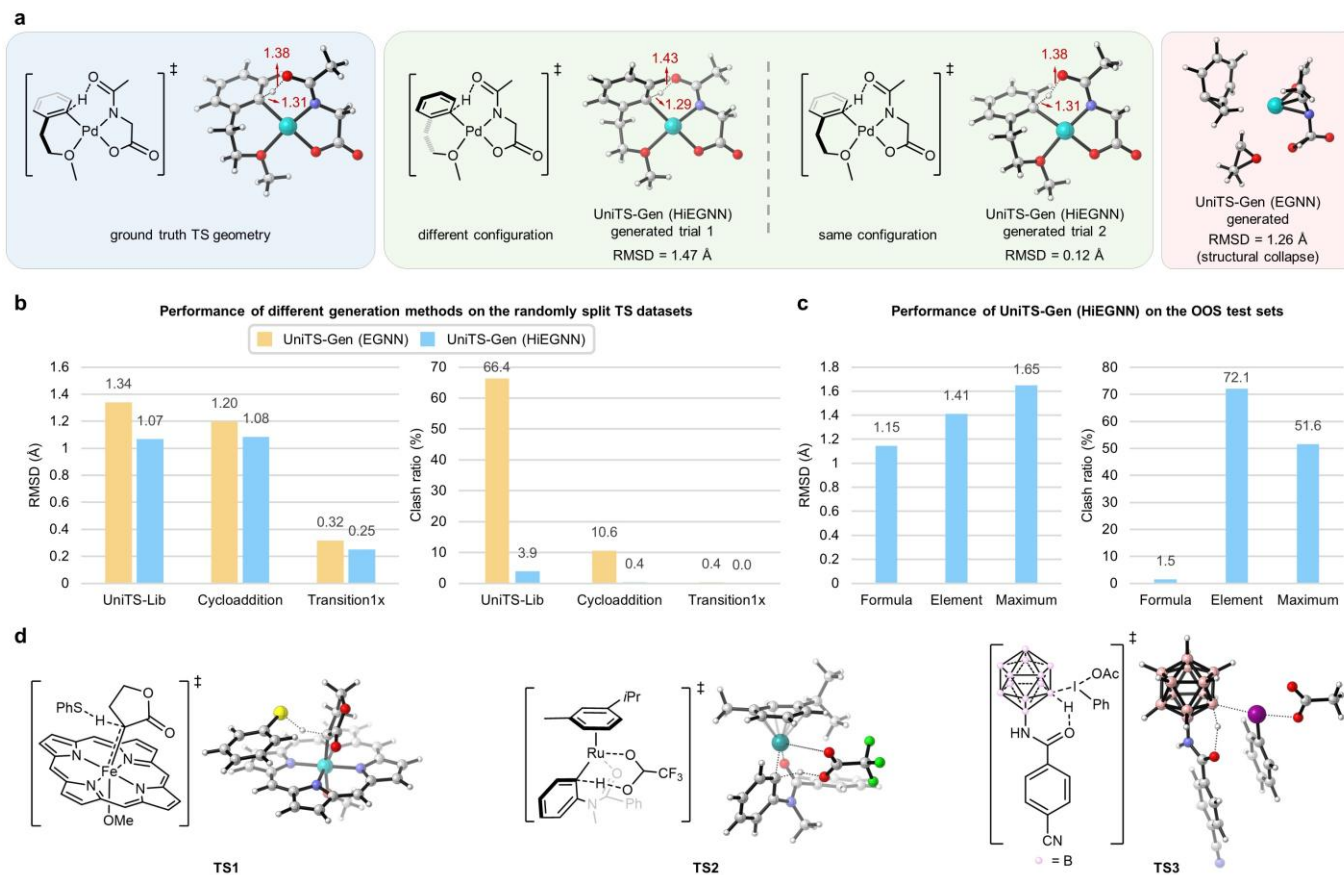
## Structural accuracy and generalization in TS generation

We systematically evaluated the performance of UniTS-Gen in generating 3D TS structures. We first randomly split the UniTS-Lib dataset into training, validation, and test sets (8:1:1) to assess generation accuracy. Fig. 4a visually presents a representative case of a Pd-catalysed C–H activation TS<sup>29</sup> from the test set. The blue panel displays the ground truth structure, while the green panels show two independent structures generated by UniTS-Gen (using HiEGNN as the denoising kernel) from the same input molecular graph and reactive site indices. Crucially, the denoising process in UniTS-Gen transforms random noise into ordered structures guided by molecular graph constraints; while the input graph dictates chemical connectivity, it does not rigidly constrain conformational freedom. Consequently, UniTS-Gen can generate diverse valid conformers and configurations from a single 2D graph. As shown, trial 1 yielded a chemically plausible TS structure but with the phenyl ring oriented differently from the ground truth, resulting in a relatively high root mean square deviation (RMSD) of 1.47 Å. In contrast, trial 2 converged to a configuration nearly identical to the ground truth, achieving a significantly lower RMSD of 0.12 Å. Notably, in both trials, the critical interatomic distances within the reactive center—specifically the C–H and O–H distances—closely matched the ground truth, a factor pivotal for successful downstream DFT optimization toward the correct saddle point.

To validate the necessity of the higher-degree equivariant architecture for the complex UniTS-Lib dataset, we benchmarked our HiEGNN kernel against a standard E(n) Equivariant Graph Neural Network (EGNN)<sup>30</sup> commonly used in 3D molecule generation<sup>31</sup> (Fig. 4a, red panel). Notably, alternative models like React-OT<sup>20</sup> was excluded as it requires 3D structures of reactant and product, which are often unavailable for novel reactions, whereas our method uses only the 2D reactant graph. The structure generated by the EGNN-based model exhibited severe structural collapse, failing to form a chemically valid geometry despite forming a recognizable phenyl substructure. Paradoxically, due to the alignment-based nature of RMSD calculations, this physically invalid, collapsed structure yielded a lower RMSD (1.26 Å) than the chemically valid but configurationally distinct structure from HiEGNN trial 1 (1.47 Å). This observation underscores that RMSD alone is an insufficient metric for evaluating generative models on complex, multimodal TS landscapes, particularly where configurational diversity is high. To address this, we introduced the "clash ratio" metric, which quantifies the percentage of generated structures containing physically implausible interatomic distances (the detailed calculation method of clash ratio can be found in Methods section), providing a robust measure of chemical validity.

To comprehensively assess performance across varying levels of complexity, we evaluated both UniTS-Gen (HiEGNN) and UniTS-Gen (EGNN) on three datasets: the complex UniTS-Lib, the classic Transition1x benchmark, and the [3+2] dipolar cycloaddition dataset known for its diverse substrate combinations. For Transition1x, we adopted the data split reported by Duan et al.<sup>20</sup>, while for the cycloaddition dataset, we used an 8:1:1 split; single-sample inference was performed for all test structures. As shown in Fig. 4b, the performance gap is stark on the complex UniTS-Lib: while the RMSD improvement of HiEGNN over EGNN appears moderate (1.07 Å vs. 1.34 Å), the difference in chemical validity is profound, with HiEGNN achieving a clash ratio of just 3.9% compared to 66.4% for EGNN. This result conclusively demonstrates the superiority of our higher-degree equivariant design for handling complex transition states. For the simpler Transition1x (molecules < 25 atoms, < 5 element types) and cycloaddition datasets, the simpler EGNN suffices, though HiEGNN still maintains a precision advantage. In addition, to account

for the inherent multi-configuration generation capability of UniTS-Gen, we also performed 10-sample evaluations for these datasets, identifying the best-matching conformer; details are provided in Supplementary Table 3–5.



**Fig. 4: Superior structural accuracy and robust generalization of UniTS-Gen using the higher-degree equivariant denoising kernel.** **a**, A representative TS from the test set illustrates UniTS-Gen's ability to generate multiple chemically valid configurations (trial 1 and trial 2) from a single 2D input graph, one of which (trial 2) achieves high accuracy (RMSD 0.12 Å). The visualization contrasts HiEGNN with the standard EGNN kernel, which fails to generate a chemically valid structure (structural collapse). **b**, Benchmark comparison of UniTS-Gen (HiEGNN) and UniTS-Gen (EGNN) across UniTS-Lib, [3+2] cycloaddition, and Transition1x datasets using RMSD (Å) and clash ratio (%). HiEGNN demonstrates profound superiority in chemical validity on the complex UniTS-Lib (clash ratio 3.9% vs. 66.4%). **c**, Performance evaluation of UniTS-Gen (HiEGNN) on three rigorous OOS test sets derived from UniTS-Lib (Formula-OOS, Element-OOS, and Maximum-OOS). The model exhibits remarkable generalization capability on the Formula-OOS set (clash ratio 1.5%), which is most representative of novel mechanistic discovery. **d**, Three generated 3D structures of challenging transformations from the OOS dataset, all successfully optimized and verified as true transition states.

While random splits assess in-distribution performance, the critical test for a general-purpose model is its generalization to unseen chemical spaces. We therefore evaluated the generalizability of UniTS-Gen through three rigorous out-of-sample (OOS) tests on UniTS-Lib. These tests comprised a Formula-OOS set of 403 reactions with unseen elemental compositions, an Element-OOS set of 351 structures containing entirely unseen elements, and a

Maximum-OOS set of 219 structures whose atom counts exceeded the maximum size encountered during training. The results (Fig. 4c) show that the model's performance on the Formula-OOS test (RMSD 1.15 Å, clash ratio 1.5%) closely rivals that of the random split, a finding that underscores its strong extrapolation capability. This unexpected robustness demonstrates UniTS-Gen's ability to generalize effectively to novel organic reactions involving known elements—the most common scenario in mechanistic research<sup>1</sup>. While performance on the highly challenging Element-OOS and Maximum-OOS tasks indicates room for improvement, these limitations will be addressed by the continuous expansion of UniTS-Lib via our automated data workflow.

The practical value of UniTS-Gen lies in its ability to generate high-quality initial guesses that liberate chemists from manual trial-and-error. To verify this utility, we subjected the structures generated from the Formula-OOS test set, which represents a highly challenging and common scenario in practical mechanistic exploration, to DFT-level transition state optimization. Notably, even based on a single random generation without any post-hoc filtering or heuristic constraints, the geometric success rate for converging to a first-order saddle point reached 68.7%. Subsequent manual validation, which filtered out ineffective transition states (e.g., those involving mere group rotation), confirmed that 41.9% of all test cases represented the correct transition state (detailed in Supplementary Information Section 11). The final success rate closely approaches that achieved during the construction of the UniTS-Lib dataset itself.

To further illustrate this practical utility, Fig. 4d presents three representative, structurally complex, and mechanistically important transition states from the Formula-OOS dataset: **TS1**, a hydrogen atom transfer transition state featuring an iron-porphyrin moiety; **TS2**, a Ru-catalyzed C–H activation transition state involving aromatic  $\pi$ -coordination; and **TS3**, a B–H bond activation transition state containing a complex icosahedral borane substructure. UniTS-Gen correctly generated chemically reasonable initial geometries for these challenging systems, accurately positioning the critical bond-forming/breaking distances essential for subsequent TS optimization. The model also successfully reproduced intricate substructures, such as the iron-porphyrin macrocycle and the icosahedral borane cluster. Crucially, all three generated structures were successfully optimized to correct first-order saddle points and verified as valid transition states.

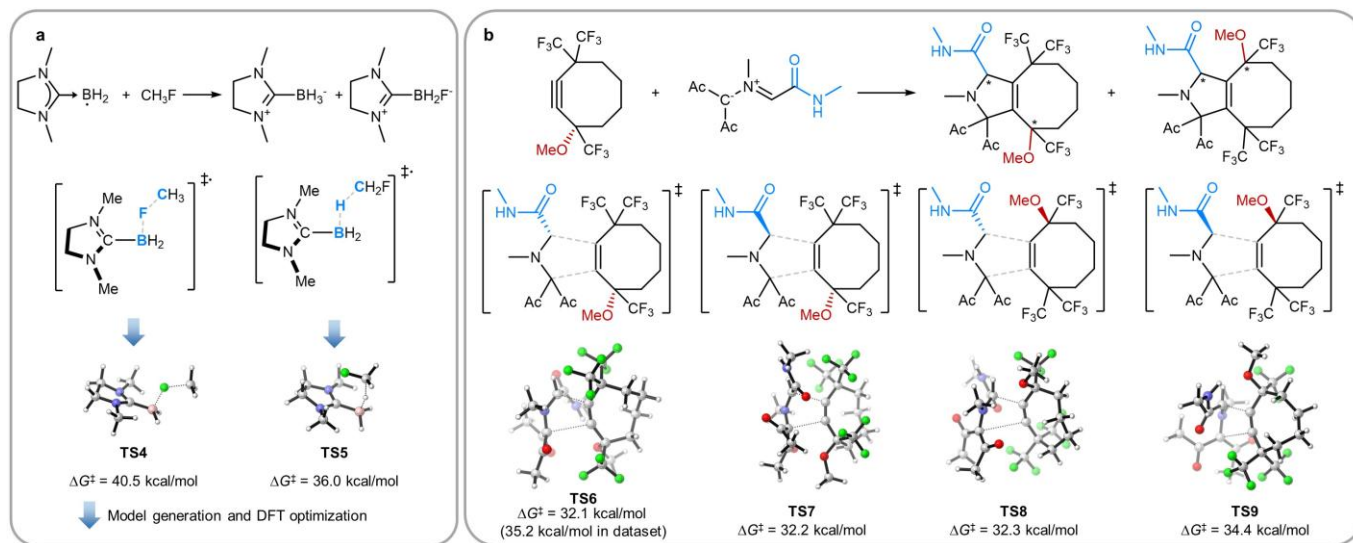
These results on the highly challenging Formula-OOS test underscore the model's potential to significantly accelerate mechanistic exploration by providing reliable starting points for complex reaction systems, thereby mitigating the bottleneck of manual trial-and-error. This result serves as a definitive proof-of-concept, highlighting UniTS-Gen's practical efficacy in accelerating automated mechanistic discovery.

### **Applications in kinetic selectivity analysis and mechanistic discovery**

To further substantiate UniTS-Gen's capability in kinetic and mechanistic studies, we first selected a highly selective bond activation system from the UniTS-Lib Formula-OOS test set, specifically the selective C–F or C–H bond activation of CH<sub>3</sub>F by a Lewis base–boryl radical, as reported by Xue et al.<sup>32</sup> (Fig. 5a). Despite sharing the same reactant graph input, the two pathways yield distinct TS structures (**TS4** for C–F activation and **TS5** for C–H activation) differentiated solely by the reactive site indices provided as conditional inputs. Following structure generation by UniTS-Gen and subsequent DFT optimization using the computational methodology consistent with the original literature, we successfully located the saddle points for both **TS4** and **TS5**. By comparing the free energies of these TS structures with the reactant, we obtained kinetic selectivities consistent with the reported findings, showing that the C–H bond activation pathway ( $\Delta G^\ddagger = 36.0$  kcal/mol) is kinetically favored over the C–F bond

activation pathway ( $\Delta G^\ddagger = 40.5$  kcal/mol) promoted by the boryl radical. This case highlights how UniTS-Gen enables selectivity analysis by enabling the generation of distinct, path-specific transition states controlled solely by the input reaction center.

Furthermore, we demonstrated UniTS-Gen's utility in exploring complex, multi-configurational systems by selecting a representative reaction from the test set of the [3+2] dipolar cycloaddition dataset<sup>27</sup> (Fig. 5b), whose formula appeared only once in the entire dataset, making it a critical Formula-OOS case. Based on the same reactant molecular graph and reactive site indices, UniTS-Gen successfully sampled and generated initial guesses for all four possible diastereomeric TS configurations (**TS6–TS9**). Subsequent DFT optimization successfully located the correct first-order saddle points for all generated structures. Although the original authors performed extensive conformer searches when building the [3+2] dataset and recorded the lowest-barrier configuration, we further found that while our generated **TS6** configuration matched the lowest-barrier stereoisomer, the conformation located by UniTS-Gen yielded a significantly lower energy barrier ( $\Delta G^\ddagger = 32.1$  kcal/mol) compared to the value reported in the dataset (35.2 kcal/mol). This result not only confirms UniTS-Gen's ability to cover the multi-configurational space but also highlights its powerful conformational diversity sampling capability. This high-precision, broad-sampling approach is a substantial advantage over traditional manual or computationally expensive conformer searches, which are often limited by predefined templates, thereby significantly accelerating accurate mechanistic discovery.



**Fig. 5: Demonstration of UniTS-Gen's capability in reaction selectivity analysis and sampling complex multi-configurational space.** **a**, UniTS-Gen accurately predicts the kinetic selectivity in the C–H/C–F bond activation of  $\text{CH}_3\text{F}$  by a Lewis base–boryl radical. By generating initial guesses for both pathways (**TS4** and **TS5**) guided by reactive site indices, subsequent DFT optimization yields free energies ( $\Delta G^\ddagger$ ) consistent with original report, correctly identifying the C–H activation pathway ( $\Delta G^\ddagger = 36.0$  kcal/mol) as kinetically favored over C–F activation ( $\Delta G^\ddagger = 40.5$  kcal/mol). **b**, Utility in exploring complex stereochemical landscapes. UniTS-Gen successfully samples and generates initial guesses for all four possible diastereomeric TS configurations (**TS6–TS9**) in a complex [3+2] dipolar cycloaddition. Furthermore, the model's conformational sampling power located a significantly lower energy barrier for the lowest-barrier stereoisomer (**TS6**,  $\Delta G^\ddagger = 32.1$  kcal/mol) compared to the literature dataset value (35.2 kcal/mol), highlighting the ability to accelerate high-precision mechanistic discovery.

## Discussion

Our results establish a universal strategy for generating initial TS structures in complex reaction systems, which integrates a universally applicable, automated workflow for TS dataset construction and a generative model driven by higher-degree equivariant diffusion. As a proof of concept, we leveraged this data construction pipeline to assemble UniTS-Lib, a library containing over 4,000 TS structures characterized by extensive elemental diversity, broad molecular size distribution, and coverage of diverse elementary steps. Building upon this high-complexity yet inherently sparse dataset, we designed UniTS-Gen, a diffusion-based molecular generation model that employs HiEGNN, a higher-degree equivariant graph neural network, as its denoising kernel. This architecture empowers mechanistic researchers to precisely control the geometry of generated TS structures by simply defining the reactant molecular graph and potential reactive sites.

UniTS-Gen has been rigorously validated across multiple TS benchmarks, exhibiting exceptional structural accuracy, multi-configuration generation capabilities, and robust OOS generalization. On the classic Transition1x benchmark, the highly diverse [3+2] dipolar cycloaddition dataset, and our own highly complex UniTS-Lib, UniTS-Gen consistently delivered high-fidelity initial TS guesses. We further probed the model's capacity to generate "unseen" structures within UniTS-Lib. Remarkably, UniTS-Gen maintained high generation accuracy even for transition states with elemental compositions unseen during training (Formula-OOS), a scenario most representative of novel mechanistic research. While generating structures containing unseen elements or exceeding the maximum atom count of the training set remains challenging, these limitations can be effectively mitigated through the continuous, automated expansion of UniTS-Lib using our established workflow. Notably, although our generated initial structures possess high geometric accuracy, making them suitable for use in certain data-driven modeling of reaction structure-performance relationships, it is crucial to emphasize that they are designed to serve only as reliable starting structures for subsequent DFT TS optimization, and not as the true converged TS. To address this critical utility, we performed detailed DFT-level TS optimization on structures generated from the highly challenging and realistic UniTS-Lib OOS-formula test set. The results showed that 68.7% of the initial guesses successfully converged to a first-order saddle point, with 41.9% of the entire test set being manually validated as the correct, chemically valid TS connecting reactants and products.

Furthermore, we demonstrated the practical utility of UniTS-Gen in kinetic selectivity analysis and mechanistic discovery using two OOS systems: the Lewis base-boryl radical promoted selective C–H/C–F bond activation and a complex [3+2] cycloaddition system. In the radical system, by specifying different reactive site information, UniTS-Gen successfully generated chemically reasonable initial structures for both C–H and C–F activation. Subsequent DFT optimization allowed us to locate the correct first-order saddle points and reproduce the originally reported kinetic selectivity. In the multi-configurational [3+2] cycloaddition system, our model successfully generated the complete set of potential TS configurations. Crucially, among these generated structures, we identified a conformation with a lower energy barrier than the ground truth recorded in the test set.

These findings underscore that our proposed general strategy significantly enhances the automation of organic reaction mechanism research. By providing a scalable and transferable solution, this approach is readily applicable to TS generation across organic synthesis. Furthermore, our work paves the way for next-generation computational tools, as the reliable initial guesses may enable future integration with machine learning potentials<sup>33,34</sup> for end-to-end

optimization and could provide critical 3D TS information for quantitative structure-selectivity modeling<sup>35–37</sup> ultimately accelerating data-driven catalysis research.

## Methods

### Details of the UniTS-Lib dataset construction

#### Data collection

The TS structures comprising UniTS-Lib were harvested from the Supplementary Information of literature reports detailing organic reaction mechanisms. We developed a custom data curation pipeline for the automated extraction and standardization of 3D atomic coordinates from computational output files across 346 distinct literature sources. Detailed statistics regarding the journal sources and the 346 literature documents used to construct UniTS-Lib are provided in Supplementary Table 1.

#### Computational details

All DFT calculation results were obtained using the Gaussian 09 program<sup>38</sup>. Geometry optimizations of all structures in UniTS-Lib were carried using B3LYP functional<sup>39,40</sup>, employing the D3 version of Grimme's dispersion corrections<sup>41</sup> with Becke-Johnson damping<sup>42</sup>, and the Def2SVP basis set<sup>43</sup> for all elements. Frequency analysis was performed at the same level of theory as geometry optimization to confirm whether optimized stationary points were either local minimum or transition states.

### UniTS-Gen architecture

#### 2D molecular graph encoding

The UniTS-Gen framework utilizes 2D molecular graphs as input. The input molecular graph includes nine node features: atom type, explicit atom degree, total charge, spin multiplicity, chiral tag, aromaticity, total valence, hydrogen count, and Cahn–Ingold–Prelog (CIP) code. It also includes five edge features: bond type, bond direction, stereo information, presence in a ring and conjugation status (Supplementary Fig. 4). These features are processed by a 2D graph encoder, adapted from our prior work<sup>28</sup>, which employs embedding layers followed by multiple graph convolutional blocks to generate latent 2D representations. To incorporate reaction-specific information, we extract a subgraph representing the reactive sites, which is encoded and concatenated with the global molecular embedding to form a unified molecular context for conditional denoising.

#### HiEGNN denoising kernel

The core of UniTS-Gen is the HiEGNN denoising kernel that predicts time-step-specific noise for 3D coordinates. The kernel constructs a dynamic 3D molecular graph based on the current atomic positions and types, from which it derives higher-degree SO(3) embeddings through distance and degree embedding. These structural representations are fused with the 2D molecular context and iteratively refined through multiple Transformer blocks comprising SO(2)-equivariant graph attention layers and feedforward networks (Supplementary Fig. 5). Finally, an equivariant prediction head processes these enhanced features to determine the noise, enabling the diffusion model to reconstruct high-fidelity transition state geometries directly from 2D graphs.

## Metric for assessing structural clash (clash ratio)

To evaluate the chemical plausibility of the generated 3D transition state structures and identify physically implausible geometries, we defined a clash ratio metric. A structure is classified as a "clash" if any pair of atoms  $i$  and  $j$  exhibits an interatomic distance  $d_{ij}$  indicating severe steric compression or structural collapse. Specifically, a structural clash is identified based on the following criterion:

$$d_{ij} < 0.75 \times (r_{cov,i} + r_{cov,j})$$

where  $r_{cov,i}$  and  $r_{cov,j}$  represent the respective covalent radii of the atoms. The clash ratio is calculated as the percentage of generated structures within a given test set that contain at least one such physically implausible interatomic distance. This metric serves as a robust measure for assessing the structural integrity of generative models when applied to complex molecular systems.

## Data availability

The UniTS-Lib dataset and preprocessed Transition1x and [3+2]-cycloaddition dataset is available in <https://github.com/licheng-xu-echo/UniTS>.

## Code availability

Codes for data preprocessing, model training, generation and visualization of transition states are available via <https://github.com/licheng-xu-echo/UniTS>.

## Acknowledgements

We thank M.-J. Tang for valuable discussions and for modifying MolOP to handle transition state structures. We are grateful to X.-T. Yu and C.T. Ouyang for their efforts in the manual inspection of ambiguous transition states within UniTS-Lib. Generous support by the National Natural Science Foundation of China (82394432 and 92249302, Y.Q.; 22503056, C.-L.J.), the Shanghai Municipal Science and Technology Major Project (Grant No. 2023SHZDZX02, Y.Q.), and the Science and Technology Commission of Shanghai Municipality (Grant No. 23YF1426700, C.-L.J.). The DFT computations in this research were performed using the CFFF platform at Fudan University and the HPC Platform of ShanghaiTech University. The AI model training was conducted using the CFFF platform at Fudan University and the Inspire platform at SAIS.

## Author contributions

L.-C.X.: conceptualization, methodology, software, validation, investigation, data curation, writing - original draft, writing - review and editing, and visualization. J.A., Q.L., and Y.-F.S.: software. C.-L.J.: resources and funding acquisition. F.C.: supervision and resources. Y.Q.: supervision, resources, and funding acquisition.

## Competing interests

The authors declare no competing interests.

## Reference

1. Cheong, P. H.-Y., Legault, C. Y., Um, J. M., Çelebi-Ölçüm, N. & Houk, K. N. Quantum Mechanical Investigations of Organocatalysis: Mechanisms, Reactivities, and Selectivities. *Chem. Rev.* **111**, 5042–5137 (2011).
2. Zahrt, A. F., Athavale, S. V. & Denmark, S. E. Quantitative Structure–Selectivity Relationships in Enantioselective Catalysis: Past, Present, and Future. *Chem. Rev.* **120**, 1620–1689 (2020).
3. Sunoj, R. B. Transition State Models for Understanding the Origin of Chiral Induction in Asymmetric Catalysis. *Acc. Chem. Res.* **49**, 1019–1028 (2016).
4. Noyori, R. Asymmetric Catalysis: Science and Opportunities (Nobel Lecture). *Angew. Chem. Int. Ed.* **41**, 2008–2022 (2002).
5. Mondal, S. *et al.* Enantioselective Radical Reactions Using Chiral Catalysts. *Chem. Rev.* **122**, 5842–5976 (2022).
6. Scinto, S. L. *et al.* Bioorthogonal chemistry. *Nat. Rev. Methods Primer* **1**, 30 (2021).
7. Bahmanyar, S., Houk, K. N., Martin, H. J. & List, B. Quantum Mechanical Predictions of the Stereoselectivities of Proline-Catalyzed Asymmetric Intermolecular Aldol Reactions. *J. Am. Chem. Soc.* **125**, 2475–2479 (2003).
8. Sperger, T., Sanhueza, I. A. & Schoenebeck, F. Computation and Experiment: A Powerful Combination to Understand and Predict Reactivities. *Acc. Chem. Res.* **49**, 1311–1319 (2016).
9. Cheng, G.-J., Zhang, X., Chung, L. W., Xu, L. & Wu, Y.-D. Computational Organic Chemistry: Bridging Theory and Experiment in Establishing the Mechanisms of Chemical Reactions. *J. Am. Chem. Soc.* **137**, 1706–1725 (2015).
10. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021).
11. Guo, A. B. *et al.* Deep learning–guided design of dynamic proteins. *Science* **388**, eadr7094 (2025).
12. Li, J., Chiu, T.-P. & Rohs, R. Predicting DNA structure using a deep learning method. *Nat. Commun.* **15**, 1243 (2024).

13. Wang, X., Mao, Y. & Wang, Z. Machine learning approaches for transition state prediction. *Chem Catal.* **5**, 101458 (2025).
14. Pattanaik, L., Ingraham, J. B., Grambow, C. A. & Green, W. H. Generating transition states of isomerization reactions with deep learning. *Phys. Chem. Chem. Phys.* **22**, 23618–23626 (2020).
15. Jackson, R., Zhang, W. & Pearson, J. TSNet: predicting transition state structures with tensor field networks and transfer learning. *Chem. Sci.* **12**, 10022–10040 (2021).
16. Lemm, D., Von Rudorff, G. F. & Von Lilienfeld, O. A. Machine learning based energy-free structure predictions of molecules, transition states, and solids. *Nat. Commun.* **12**, 4468 (2021).
17. Choi, S. Prediction of transition state structures of gas-phase chemical reactions via machine learning. *Nat. Commun.* **14**, 1168 (2023).
18. Duan, C., Du, Y., Jia, H. & Kulik, H. J. Accurate transition state generation with an object-aware equivariant elementary reaction diffusion model. *Nat. Comput. Sci.* **3**, 1045–1055 (2023).
19. Kim, S., Woo, J. & Kim, W. Y. Diffusion-based generative AI for exploring transition states from 2D molecular graphs. *Nat. Commun.* **15**, 341 (2024).
20. Duan, C. *et al.* Optimal transport for generating transition states in chemical reactions. *Nat. Mach. Intell.* **7**, 615–626 (2025).
21. Si, Y. *et al.* Transition state structure detection with machine learnings. *Npj Comput. Mater.* **11**, 199 (2025).
22. Zhang, S.-Q. & Hong, X. Mechanism and Selectivity Control in Ni- and Pd-Catalyzed Cross-Couplings Involving Carbon–Oxygen Bond Activation. *Acc. Chem. Res.* **54**, 2158–2171 (2021).
23. Batzner, S. *et al.* E(3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials. *Nat. Commun.* **13**, 2453 (2022).

24. Passaro, S. & Zitnick, C. L. Reducing SO(3) convolutions to SO(2) for efficient equivariant GNNs. in *Proceedings of the 40th International Conference on Machine Learning* (JMLR.org, Honolulu, Hawaii, USA, 2023).
25. Liao, Y.-L., Wood, B. M., Das, A. & Smidt, T. E. EquiformerV2: Improved Equivariant Transformer for Scaling to Higher-Degree Representations. in *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024* (OpenReview.net, 2024).
26. Schreiner, M., Bhowmik, A., Vegge, T., Busk, J. & Winther, O. Transition1x - a dataset for building generalizable reactive machine learning potentials. *Sci. Data* **9**, 779 (2022).
27. Stuyver, T., Jorner, K. & Coley, C. W. Reaction profiles for quantum chemistry-computed [3 + 2] cycloaddition reactions. *Sci. Data* **10**, 66 (2023).
28. Xu, L.-C., Tang, M.-J., An, J., Cao, F. & Qi, Y. A unified pre-trained deep learning framework for cross-task reaction performance prediction and synthesis planning. *Nat. Mach. Intell.* **7**, 1561–1571 (2025).
29. Cheng, G.-J. *et al.* Role of *N*-Acyl Amino Acid Ligands in Pd(II)-Catalyzed Remote C–H Activation of Tethered Arenes. *J. Am. Chem. Soc.* **136**, 894–897 (2014).
30. Satorras, V. G., Hoogeboom, E. & Welling, M. E(n) Equivariant Graph Neural Networks. in *Proceedings of the 38th International Conference on Machine Learning* (eds Meila, M. & Zhang, T.) vol. 139 9323–9332 (PMLR, 2021).
31. Hoogeboom, E., Satorras, V. G., Vignac, C. & Welling, M. Equivariant Diffusion for Molecule Generation in 3D. in *Proceedings of the 39th International Conference on Machine Learning* (eds Chaudhuri, K. *et al.*) vol. 162 8867–8887 (PMLR, 2022).
32. Guo, X., Zhang, Y., Lai, X., Pang, Y. & Xue, X. C(sp<sup>3</sup>)–F Bond Activation by Lewis Base-Boryl Radicals via Concerted Electron-Fluoride Transfer. *Angew. Chem. Int. Ed.* **64**, e202415715 (2025).

33. Li, B., Xiao, J., Gao, Y., Zhang, J. Z. H. & Zhu, T. Transition State Searching Accelerated by Neural Network Potential. *J. Chem. Inf. Model.* **65**, 2297–2303 (2025).
34. Yin, J. *et al.* CaTS: Toward Scalable and Efficient Transition State Screening for Catalyst Discovery. *ACS Catal.* **15**, 15754–15764 (2025).
35. Moskal, M., Beker, W., Szymkuć, S. & Grzybowski, B. A. Scaffold-Directed Face Selectivity Machine-Learned from Vectors of Non-covalent Interactions. *Angew. Chem. Int. Ed.* **60**, 15230–15235 (2021).
36. Xu, L.-C. *et al.* Enantioselectivity prediction of pallada-electrocatalysed C–H activation using transition state knowledge in machine learning. *Nat. Synth.* **2**, 321–330 (2023).
37. Chen, G.-M., Ye, Z.-H., Li, Z.-M. & Zhang, J.-L. Universal descriptors of quasi transition states for small-data-driven asymmetric catalysis prediction in machine learning model. *Cell Rep. Phys. Sci.* **5**, 102043 (2024).
38. Frisch, M. J. *et al.* Gaussian 09 Revision D.01. (2016).
39. Lee, C., Yang, W. & Parr, R. G. Development of the Colle-Salvetti correlation-energy formula into a functional of the electron density. *Phys. Rev. B* **37**, 785–789 (1988).
40. Becke, A. D. Density-functional thermochemistry. III. The role of exact exchange. *J. Chem. Phys.* **98**, 5648–5652 (1993).
41. Grimme, S., Antony, J., Ehrlich, S. & Krieg, H. A consistent and accurate *ab initio* parametrization of density functional dispersion correction (DFT-D) for the 94 elements H-Pu. *J. Chem. Phys.* **132**, 154104 (2010).
42. Grimme, S., Ehrlich, S. & Goerigk, L. Effect of the damping function in dispersion corrected density functional theory. *J. Comput. Chem.* **32**, 1456–1465 (2011).
43. Schäfer, A., Horn, H. & Ahlrichs, R. Fully optimized contracted Gaussian basis sets for atoms Li to Kr. *J. Chem. Phys.* **97**, 2571–2577 (1992).