
SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

Anonymous Authors¹

Abstract

LLMs agents have incredible potential for scientific discovery applications. However, the performance of LLM agents on real-world, small molecule drug design (SMDD) tasks across diverse chemistries and targets is unclear. Current evaluation methods are either ad hoc, too simple for real-world discovery, limited in scale, or restricted to single-turn question answering. In effort to standardize the evaluation of LLM agents on small molecule design, we introduce SMDD-Bench, a challenging, multi-turn, long-horizon agentic benchmark consisting of 502 guaranteed-solvable task instances spanning 5 task types: 2D Pharmacophore Identification, Interaction Point Discovery, Scaffold Hopping, Lead Optimization, and Fragment Assembly. SMDD-Bench tasks span a wide region of chemical space and involve 102 unique protein targets. Completely solving the benchmark would require having strong chemical and biological reasoning and 3D intuition, understanding specialized tool use, and displaying planning expertise over a limited number of oracle calls. We benchmark 7 frontier open and closed source LLMs and find even the most performant LLM, GPT5.4, solves only 40.2% of tasks. We hope SMDD-Bench provides a standardized testbed to invigorate the field towards training and evaluating LLM agents for fully autonomous computational drug design.

1. Introduction

LLMs have seen considerable success in many scientific discovery tasks such as synthesizing scientific literature, orchestrating chemistry experiments, and even performing end to end applied machine learning research (Boiko et al.,

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Submitted to the 2026 Workshop on Generative and Agentic AI for Biology (ICML 2026). Do not distribute.

2023; Lu et al., 2026; Asai et al., 2026). Within the drug discovery and medicinal chemistry space, LLMs have been benchmarked on tasks including single step retrosynthesis prediction, molecular description, de novo molecular generation, and property prediction (Yu et al., 2024; Mirza et al., 2025; Huang et al., 2021; Bartmann et al., 2026). However, these benchmarks are either restricted to single-turn question answering, are rudimentary and unrealistic when compared to the complexity of real-world computational drug discovery, or are limited to only a small non-diverse handful of tasks and protein targets. Although such benchmarks are useful in understanding an LLM’s ability to understand basic chemistry principles, they are not as useful for understanding and improving an LLM **agent’s** ability to handle the geometric and chemical complexity and diversity of drug discovery tasks facing expert human medicinal chemists.

In this work, we present SMDD-Bench, an agentic benchmark designed to evaluate LLM agents, multi-agent systems, and agentic workflows on complex and realistic small molecule drug design (SMDD) tasks facing human medicinal chemists. SMDD-Bench is comprised of 5 task types with a total of 502 task instances. The task types are 2D Pharmacophore Identification, Interaction Point Discovery, Scaffold Hopping, Lead Optimization, and Fragment Assembly. Each task is deeply embedded into real-world drug discovery workflows and requires deep chemical intuition, biological reasoning ability, expertise in selecting and running specialized cheminformatics tools, hypothesis generation and validation, and understanding 3D geometries in 3D space. We also introduce the notion of witness-aware task generation, allowing us to generate a large number of chemically valid and guaranteed-solvable task instances procedurally and without a human in the loop.

Real world small molecule drug discovery requires the LLM to be able to write RDKit code, analyze protein-ligand interactions in 3D space, propose targeted molecular edits, understand how those molecular edits affect predicted properties, binding affinities, and 3D conformations, and run specialized cheminformatics tools. These skills go far beyond the simple single-turn question answering tasks that previous benchmarks study and require the development of

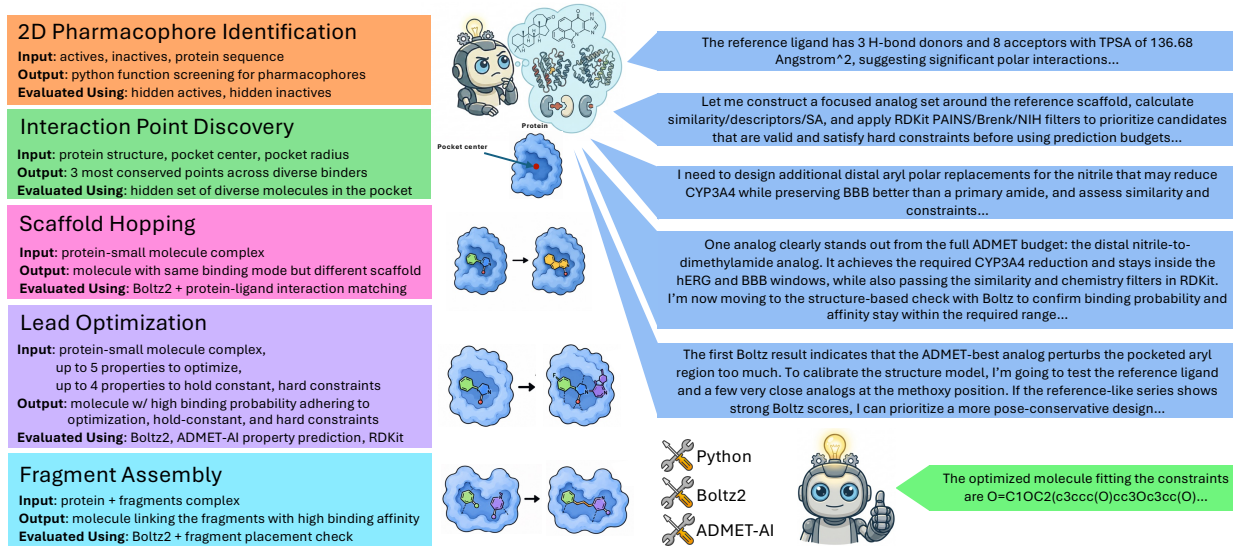


Figure 1. Overview of SMDD-Bench’s task types and example reasoning turns.

a new and challenging agentic benchmark consisting of a diverse spread of chemistries and protein targets.

However, creating real-world drug discovery tasks can be difficult due to the requirement for there to be at least one known solution to the task. For example, in lead optimization, while it is tempting to create a task instance by sampling a protein-ligand pair before sampling a set of properties to optimize as well as a set of properties to hold constant, there is no guarantee there exists a solution that adheres to the set of provided objectives and constraints. As a result, when creating task instances, in SMDD-Bench, we create a *witness-aware* task generation pipeline, allowing us to generate difficult, chemically valid, and guaranteed-solvable task instances while simultaneously generating an example solution molecule – all without requiring a human expert in the loop.

We evaluate 7 frontier open and closed source LLMs using a minimalist agent scaffold on SMDD-Bench and find that, although LLMs achieve up to 50% on simple lead optimization tasks, more challenging lead optimization tasks as well as other task types have near 0% success rates – highlighting significant room for growth for LLM applications as autonomous medicinal chemists. We also introduce **SMDD-Bench Diversity**, a subset of SMDD-Bench along with a principled framework for studying the diversity of SMDD agents’ outputs.

We hope that the release of SMDD-Bench standardizes and invigorates research into training and evaluating LLMs for real-world autonomous drug discovery.

2. SMDD-Bench

SMDD-Bench consists of 502 instances of 5 task types: 2D Pharmacophore Identification (25 tasks), Interaction Point Discovery (25 tasks), Scaffold Hopping (52 tasks), Lead Optimization (340 tasks), and Fragment Assembly (60 tasks). Each task is generated in a **witness-aware** fashion, where a witness molecule that solves the task is simultaneously generated along with the starting molecule – ensuring the task instance is guaranteed to be solvable and has at least one solution. Each task instance is also evaluated purely computationally, without a human in the loop, using tools such as RDKit, PLIP, OpenBabel, Boltz2, and ADMET-AI (Passaro et al., 2025; Swanson et al., 2024; Landrum, 2010; Salentin et al., 2015; O’Boyle et al., 2011). Boltz2 is used for protein-ligand co-folded structure prediction, binding affinity¹, and binding probability prediction while the ADMET-AI model is used for property prediction. An overview of the benchmark can be found in Figure 1.

2.1. Task Types

Further details on the generation and evaluation of each task type can be found in Appendix B. Only the *Scaffold Hopping*, *Lead Optimization*, and *Fragment Assembly* tasks require **witness-aware** generation due to the *2D Pharmacophore Identification* and *Interaction Point Discovery* tasks having ground truth solutions by default.

¹Note that Boltz2 predicts binding affinity as $\log_{10}(\text{IC}_{50} \text{ in } \mu\text{M})$. As a result, a decrease in the Boltz2 predicted affinity value leads to stronger binding (Passaro et al., 2025)

1. 2D Pharmacophore Identification

Task: Given a target protein sequence as well as the SMILES strings associated with 10 experimentally active molecules and 10 experimentally inactive molecules from the ChEMBL database (Mendez et al., 2019), create a 2D pharmacophore in the form of a python function, which returns True if an inputted SMILES string matches the pharmacophore and False otherwise. See Appendix B.1 for more details.

Evaluation: The python function is graded on its recall and specificity performance on a held-out set of hidden, experimentally validated actives and inactives.

Generation: We filtered ChEMBL activity records for high-confidence, direct-binding assays, split the candidates into actives and inactives as determined by potency, and then pruned with MaxMin diversity picking to ensure the visible set spans chemical space and introduces binding mode diversity.

This task assesses the agent’s ability to extract a generalizable, discriminative structural hypothesis from a small set of actives and inactives. In real world applications, the pharmacophore produced by the agent can be used as a first-filter for the screening of additional candidate compounds.

2. Interaction point discovery

Task: Given a PDB file of the protein structure as well as the location of the pocket, specify the 3D coordinates of the 3 interaction points within a protein pocket most likely to be conserved across a diverse set of small molecule binders. The agent must output a csv file of 3D coordinates as well as a label for the type of interaction point. The types are chosen as donor, acceptor, aromatic, hydrophobic, cation or anion interaction points.

Evaluation: The predicted points are assigned to the ground-truth points via the Hungarian algorithm. Each prediction is filtered to a ground truth point whose type is compatible (exact match, or a permitted cross-type match, see Appendix B.2 for more details), then the Hungarian algorithm computes the optimal one-to-one assignment under a distance cutoff. All three predictions must match a distinct conserved ground truth point in order for the task to be considered a success.

Generation: We determine ground truth points by fetching thousands of co-crystal structures per target from the PDB (Berman et al., 2000), aligning binding sites, extracting per-ligand protein interactions, and run density-based clustering with a conservation threshold. Steric clashes with the agent-provided prototype PDB are filtered out and ambiguous type labels are merged. As a result, it is possible for there to be more than three ground truth points of well-conserved interaction points. The agent only needs to output points that match to three of these ground truth points.

This task assesses the agent’s ability to identify valuable hotspots within a protein pocket. In real world applications, these hotspots can then be targeted during fragment-based drug discovery campaigns, serve as guardrails during lead optimization, or used as 3D pharmacophores during pharmacophore-based virtual screening.

3. Scaffold hopping

Task: Given a Boltz2-folded protein and reference ligand complex, propose a new SMILES with a distinct scaffold that reproduces the reference ligand’s binding interactions with the protein pocket.

Evaluation: The submitted SMILES string must have a low Tanimoto similarity with the reference molecule. The Murcko scaffolds of the submitted and reference molecules must also have a maximum common substructure which is only a portion of the overall heavy atom count of the Murcko scaffold of the reference molecule. Both of these checks ensure that a proper scaffold hop has been performed.

In order to check if the submitted molecule adheres to the same binding mode as the reference molecule, Boltz2 is used to co-fold the submitted SMILES string against the target using the same pocket contacts as the reference molecule. Potential stochasticity within the Boltz2 co-folded structure is mitigated by matching binding interactions via an ensemble of 10 Boltz2 co-folded structures.

Witness-Aware Generation: We first mine experimentally known actives for a specific pocket. Afterwards, we consider all potential pairs of actives. Every pair which the two ligands satisfy the Tanimoto distance and scaffold substructure constraint is co-folded via Boltz to the shared target. If both members of the pair share the same interaction fingerprint with the protein pocket and have high predicted binding probability, the pair and target are created into a task instance. One molecule within the pair is labeled as the reference while the other molecule is labeled as the witness. The witness molecule ensures the existence of at least one solution and guarantees the solvability of the task instance.

This task assesses the agent’s ability to understand bioisosterisms as well as internally reason through binding poses before committing to making limited, computationally heavy oracle calls. In real world applications, scaffold hopping is used for rescuing lead molecules with problematic scaffolds, mitigating off-target binding effects, enabling the extension of SAR exploration beyond a depleted search, and avoiding IP infringement.

4. Lead optimization

Task: Given a reference molecule that is Boltz2 co-folded with a target protein, optimize the reference molecule into a new molecule that improves a set of ADMET and binding affinity objectives while holding other properties constant, satisfying drug-likeness constraints, maintaining high bind-

ing probability, and maintaining close Tanimoto similarity to the reference molecule (due to this being a lead optimization task).

Evaluation: The submitted molecule must satisfy the set of specified ADMET and binding affinity objectives. RDKit is used to measure drug-likeness, ADMET properties are calculated using the ADMET-AI library, and Boltz2 is used for binding probability and binding affinity prediction.

Witness-Aware Generation: For each selected protein target, we build a pool of pocket-compatible molecules that pass drug-likeness filters and a Boltz2 binding probability gate. We consider all pairs of molecules from this pool that are similar via Tanimoto similarity. From these pairs, if there exists a molecule whose ADMET/binding affinity profile dominates the other molecule on some subset of the 8 properties in SMDD-Bench (found in Table 7), the task’s optimization objectives and hold-constant constraints are derived from the measured property differences between the two molecules. The dominating molecule becomes the witness while the dominated molecule becomes reference to provide to the agent to start lead optimization. The task instance is only instantiated if the witness’s profile satisfies every objective and constraint. As a result, the witness solves the task instance by construction and the task becomes guaranteed solvable. More details on this can be found in Appendix B.4.

This task assesses the agent’s ability to reason through how certain edits may affect complicated, intertwined relationships between molecular properties. In real world applications, multi-property optimization with optimization objectives and hold-constant constraints define the core loop of any lead optimization campaign. Lead optimization is also where the bulk of time is spent in most drug discovery programs. We envision a setting where hundreds of thousands or millions of high-performing agents are instantiated in parallel with the task of discovering a diverse array of optimized molecules given a single starting reference molecule.

5. Fragment Assembly

Task: Given either one or two 3D molecule fragments placed in a protein pocket, design a single drug-like molecule that contains each fragment as a substructure, covalently links the fragments when more than one fragment is provided, achieves high binding probability and affinity, and preserves the input fragment poses when the submitted molecule is re-docked.

Evaluation: The submitted molecule must incorporate each of the provided fragments as distinct substructures as well as pass drug-likeness checks. The submitted molecule is also co-folded with the protein via Boltz2 and must have a high predicted binding probability and binding affinity

with respect to the witness molecule. Finally, the fragment substructures within the submitted molecule must be similar to the provided 3D coordinates of the input fragments after pocket alignment is performed. This is to ensure a molecule of the same binding mode as the input fragments was designed. Otherwise, it would be possible for a submitted molecule to contain the input fragments but have the fragments be in nonessential locations for binding. More details can be found in Appendix B.5.

Witness-Aware Generation: We co-fold known actives via Boltz2 and keep the co-folded structures which pass a binding probability gate. We perform BRICS decomposition (Degen et al., 2008) on the molecule to separate the molecule into its synthesizable fragments. Each input fragment must be large enough to be a meaningful pharmacophore, have at least one direct interaction with the protein pocket, and be in close proximity to the protein itself. Fragments are then sampled from the same molecule to create the task instance. As a result, the co-folded active molecule in which fragments are sampled from becomes the witness molecule by construction – ensuring the guaranteed solvability of the task instance. Task instances with 2 fragments have the additional constraint for fragments to be sufficiently far and not directly linked together.

This task assesses the agent’s ability to reason through 3D geometry and fragment-pocket complementarity. In real world applications, fragment assembly is the central step of fragment-based drug discovery campaigns, where weakly-binding fragments identified via crystallographic or biophysical screening must be elaborated or linked into a single high-affinity lead molecule.

2.2. Benchmark Setup and Variations

SMDD-Bench Agent Setup. To evaluate the agent’s innate medicinal chemistry reasoning ability, internet access is strictly prohibited, and agents have no access to academic literature, patent filings, textbooks, or any other source of non-parametric knowledge. We benchmark current frontier LLMs using a minimalist ReAct agent (Yao et al., 2022) with a minimal harness, providing access to an isolated filesystem, a python interpreter with relevant biology and chemistry packages installed, and a limited number of Boltz2 and ADMET-AI oracle calls (the same oracles used to evaluate agent submissions). To minimize memorization, all identifying information such as PDB codes, ligand codes, Uniprot IDs, and target names are obfuscated or removed. More details on the minimalist agent can be found in Appendix C.

Benchmark Statistics. Across 502 task instances, SMDD-Bench contains 102 unique protein targets and 855 unique small molecules. Figure 2 shows the distribution of task types and protein targets across the benchmark. Figure 3(a)

and (b) report, respectively, the frequencies with which each of the 9 properties appears in Lead Optimization tasks and the joint distribution of task complexity. Additional details on benchmark diversity are provided in Appendix A.

SMDD-Bench Lite. Due to the heavy computational overhead of running an agent on all 502 tasks, we release SMDD-Bench Lite, a performance-representative subset of 100 task instances designed to encourage broader adoption of the benchmark. SMDD-Bench Lite serves as a testbed for iterating on agent development against a difficulty-diverse and representative sample of the full benchmark. Details on how the SMDD-Bench Lite instances are selected can be found in Appendix E.

3. Results

LLMs do not perform well on real-world small molecule drug discovery tasks. We evaluated 7 frontier LLMs on SMDD-Bench using the minimalist agent harness described in Section 2.2 and present results in Table 1. As described in Appendix C, each agent was provided with 8 Boltz2 oracle calls and 15 ADMET-AI oracle call – each ADMET-AI oracle call returns all 8 predicted ADMET properties. Closed source models such as Claude Sonnet 4.6, GPT-5.4 and Gemini 3.1 pro significantly outperform leading open source models such as Kimi K2.5 Thinking, Deepseek V3.2, and Qwen 3.5-397B-A17B with GPT5.4 setting the state of the art at 40.2% overall success rate.

However, most of the successes across all models are a result of performing well on the Lead Optimization task. Particularly, current LLMs perform poorly on tasks that require a heavy intuition and understanding of 3D geometries and interactions such as the Interaction Point Discovery, Scaffold Hopping and Fragment Assembly tasks. Considering the overwhelming challenges of multi-property, multi-constraint lead optimization, frontier LLMs perform surprisingly well at the Lead Optimization task. More specific results on LLM performance, such as cost breakdown and tool usage can be found in Appendix D.

The chemical novelty of the submissions of each agent, broken down into task type, are found in Table 2. Chemical novelty is defined as the % of total agent submissions which are not found in the ChEMBL36, SureChEMBL, PubChem, of BindingDB databases (Mendez et al., 2019; Papadatos et al., 2016; Kim et al., 2023; Liu et al., 2007; 2025b). These databases are chosen due to their potential inclusion of bioactivity data for small molecules with respect to protein targets. The novelty check we prescribe is a conservative measure of the agent’s ability to design a molecule in the context of a protein target for which there doesn’t exist pre-existing training data on.

LLMs are not diverse in their outputs. In real world ap-

plications, during lead optimization campaigns, researchers may want to instantiate many agents or agentic systems in parallel on the same task instance. The goal would be for the agents to propose a diverse and novel spread of successful solutions rather than converging on a single solution, allowing for a more robust exploration and discovery of lead molecules. In order to test for this behavior, we propose **SMDD-Bench Diversity**, a strict subset of 20 SMDD-Bench task instances for the Lead Optimization task type. The task instances represent the most challenging lead optimization tasks that are still solvable by current frontier LLMs. See Appendix F for more details on selecting SMDD-Bench Diversity.

For each task on SMDD-Bench Diversity, we run each agent 10 times and report the distribution of novel and diverse molecules which pass the evaluation requirements. In order to measure diversity, we measure the average number of unique successful solutions per task instance, the average pairwise tanimoto similarity of successful solutions, the percentage of successful solutions which are novel (using the same criteria as in Table 2), and the average number of successful solutions of each LLM on the subset. Results are found in Table 3.

Although Claude Sonnet 4.6 consistently achieves the highest number of successful submissions, the model converges to the same set of molecules over half the time. Surprisingly, Minimax M2.7 submits the most unique spread of successful molecules on SMDD-Diversity while maintaining strong % novelty across successful submissions. SMDD agents deployed in the wild must consistently achieve a high number of both unique and novel successful solutions when provided the same task instance.

Enumeration vs. Selection. Throughout the solve process, each LLM enumerates many potential candidates, reasons through those candidates, and then runs the most promising of the candidates through an oracle tool call. We investigate whether or not LLMs are capable of enumerating passing molecules but fail at the reasoning and selection step. In order to do so, we use GPT5.4 mini to extract every SMILES string mentioned in each of the agent’s traces for the Scaffold Hopping, Lead Optimization, and Fragment Assembly tasks. We then, using the task instance’s evaluation pipeline, evaluate all of the molecules that the agent did not perform an oracle call on – this set of molecules represents the molecules which the agent hypothesized could be strong candidates, but considered to be inferior. We then calculate the overall recovered success rate on SMDD-Bench had the LLM selected the best candidate from its enumeration and report the values in Table 4. More details on how we performed this experiment can be found in Appendix G.

Scaffold hopping sees a significant increase in recovered rate for Claude Sonnet 4.6, Gemini 3.1 Pro, and GPT5.4 –

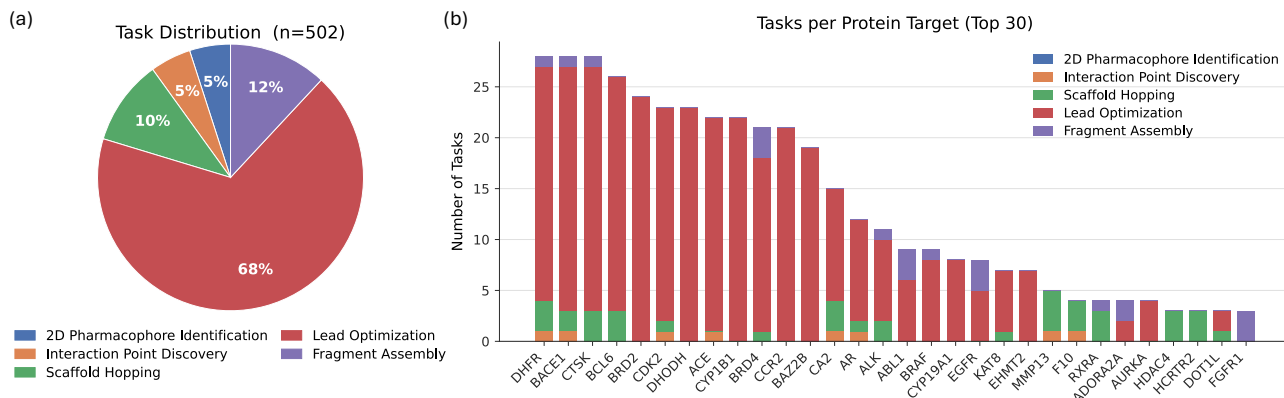


Figure 2. (a) The distribution of SMDD-Bench tasks across the five task types. Each task type is generated to be guaranteed-solvable by construction. (b). The top 30 most frequent protein targets found through all task types, broken down by task type.

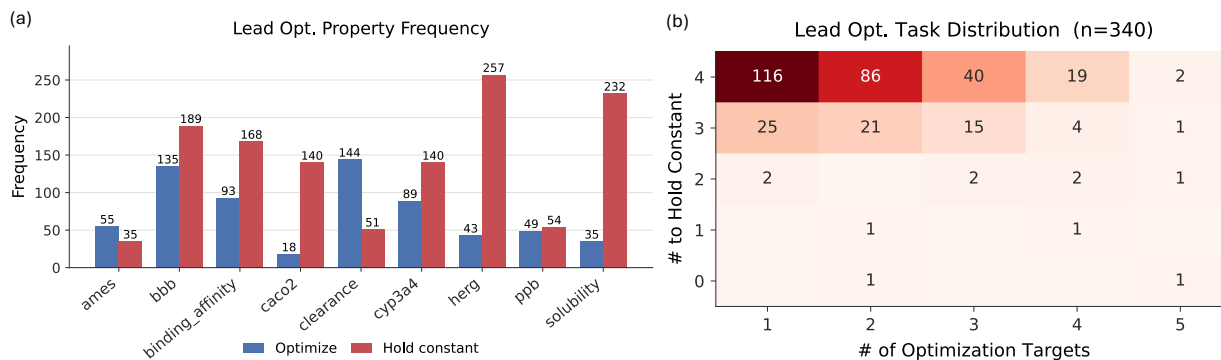


Figure 3. (a) Frequency with which each ADMET and binding affinity property appears as an optimization objective versus a hold-constant constraint across the 340 Lead Optimization instances. (b) Joint distribution of the number of optimization objectives and hold-constant properties across the 340 Lead Optimization task instances. Each cell reports the number of task instances with the given (hold-constant, optimization target) pair.

Table 1. Success rate of 7 frontier closed and open source LLMs on SMDD-Bench, broken down by task type. Each LLM is evaluated using the minimalist agent harness outlined in Section 2.2 and Appendix C with a limited number of oracle calls.

Model	Success Rate (%)					Overall	Cost (\$/instance)
	2D Pharmaco-phore ID	Interaction Point Discovery	Scaffold Hopping	Lead Optimization	Fragment Assembly		
Claude Sonnet 4.6	28.0	0.0	3.8	53.5	0.0	38.0	1.31
GPT-5.4	12.0	0.0	3.8	57.6	1.7	40.2	0.78
Gemini 3.1 Pro	20.0	4.0	0.0	55.6	1.7	39.0	0.63
Kimi K2.5 Thinking	12.0	0.0	1.9	43.5	0.0	30.3	0.40
DeepSeek V3.2	8.0	0.0	3.8	34.7	0.0	24.3	0.43
MiniMax M2.7	16.0	0.0	1.9	27.1	0.0	19.3	0.36
Qwen3.5 397B A17B	4.0	0.0	1.9	40.0	0.0	27.5	0.75

Table 2. Novelty of agent submissions, defined by % of submitted molecules not found in ChEMBL, PubChem, SureChEMBL, or BindingDB.

Model	Scaffold Hopping	Lead Optimization	Fragment Assembly
Claude Sonnet 4.6	88.5	77.9	82.8
GPT-5.4	94.2	75.3	68.3
Gemini 3.1 Pro	71.2	76.5	83.3
Kimi K2.5 Thinking	80.8	77.6	65.0
DeepSeek V3.2	94.2	83.2	69.1
MiniMax M2.7	82.7	86.5	67.2
Qwen3.5 397B A17B	82.7	77.9	61.7

Table 3. Statistics across 10 runs of each of the 20 task instances in the SMDD-Bench Diversity subset. Novelty is defined in Table 2. "Avg. Successful" and "Avg. Unique and Successful" categories are out of 10 due to the 10 runs performed on each task instance.

Model	Avg. Successful	Avg. Unique and Successful	Novel and Successful (%)	Pairwise Tanimoto Similarity
Claude Sonnet 4.6	8.40	3.70	74.0	0.823
GPT-5.4	7.90	2.75	64.6	0.863
Gemini 3.1 Pro	8.00	4.00	67.6	0.809
Kimi K2.5 Thinking	6.00	3.85	65.0	0.786
DeepSeek V3.2	5.35	3.85	68.4	0.763
MiniMax M2.7	6.00	4.05	73.1	0.763
Qwen3.5 397B A17B	7.25	3.55	67.2	0.814

indicating that the agents are not as capable at reasoning through binding interactions of chemically different analogues without first performing oracle calls. On the other hand, fragment assembly, a task requiring a deep understanding of how the binding pocket and surrounding chemistry affects the 3D geometry of the small molecule binding pose, sees very minimal overall improvement. Interestingly, MiniMax M2.7 nearly doubles its success rate for the Lead Optimization task, indicating its ability to enumerate solutions and inability to select from the enumerations.

SMDD-Bench Lite performance. The performance of each agent on the SMDD-Bench Lite subset are reported in Table 5. We also report the performance of Gemini 3.1 pro on high thinking mode.

4. Common Failure Modes

Beyond the issues regarding the novelty rate, diversity rate, and enumeration vs. selection performance gap, we investigated individual traces searching for other common failure modes in the frontier LLMs.

Absence of Cross-Turn Structure–Activity Relationship (SAR) Synthesis. A recurring failure across models is their inability to generalize SAR rules. When a molecule violates predefined constraints, a budget-aware agent should ask whether a specific structural feature caused this failure and apply this exclusion rule to later candidates. Instead, agents repeatedly test variants of already disqualified structures, effectively wasting budget on confirming the same rule multiple times. This reflects a lack of meta-cognitive reasoning: agents fail to convert individual failures into structural generalizations that could narrow down the search space. We provide an annotated reasoning trace in Appendix I.1.

Incoherent planning across multiple turns. We notice that a common failure mode involves candidate-generation turns being treated independently rather than as part of an iterative, sequential process. This results in the agent generating new molecules instead of selecting from candidates that passed earlier-stage criteria. Furthermore, the same molecule may be proposed multiple times across turns even though it had already been tested in a previous turn and did not meet proper thresholds. Refer to a representative trace in Appendix I.2.

Code Execution Errors. Agents also exhibit tool-specific coding mistakes. These include malformed calls to structure-prediction tools (e.g., incorrect argument types or missing required fields for Boltz binding predictions, cheminformatics routines (e.g., MolToSmiles conversion), and interaction-analysis pipelines (e.g., PLIP).

5. Related Work

SMDD-Bench evaluates LLM agents on challenging, real-world small molecule drug design tasks. Existing chemistry and molecular benchmarks – ChemBench (Mirza et al., 2025), MolecularIQ (Bartmann et al., 2026), SmolInstruct

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

Table 4. Original vs recovered success rates (%) and total number of enumerated molecules (#Mols) for all agents on SMDD-Bench. The recovered success rate is the rate the agent would have if it had selected to submit the ground truth best molecule that it enumerated during the reasoning process.

Model	Scaffold Hopping			Lead Optimization			Fragment Assembly			Overall		
	Orig. (%)	Recov. (%)	#Mols	Orig. (%)	Recov. (%)	#Mols	Orig. (%)	Recov. (%)	#Mols	Orig. (%)	Recov. (%)	#Mols
Claude Sonnet 4.6	3.8	32.7	2657	53.5	64.1	5872	0.0	3.3	4982	40.7	52.4	13 511
DeepSeek V3.2	3.8	11.5	1052	34.7	39.4	3413	0.0	0.0	1358	26.5	31.0	5823
Gemini 3.1 Pro	0.0	19.2	735	55.6	62.9	3328	1.7	3.3	1855	42.0	50.0	5918
GPT-5.4	3.8	15.4	850	57.6	61.2	4486	1.7	1.7	1945	44.0	48.0	7281
Kimi K2.5 Thinking	1.9	3.8	986	43.5	53.8	4940	0.0	1.7	2528	33.0	41.2	8454
MiniMax M2.7	1.9	5.8	1203	27.1	47.1	7231	0.0	0.0	2372	20.6	36.1	10 806
Qwen3.5 397B	1.9	3.8	1055	40.0	53.8	4367	0.0	0.0	1915	30.3	40.9	7337

Table 5. Success rate (%) by task type and overall cost per instance of 8 frontier LLMs on the SMDD-Bench Lite subset.

Model	Success Rate (%)						Cost (\$/inst.)
	2D Pharmacophore	Interaction Point	Scaffold Hopping	Lead Opt.	Fragment Assembly	Overall	
Claude Sonnet 4.6	40.0	0.0	0.0	52.9	0.0	38.0	1.29
GPT-5.4 (medium)	0.0	0.0	0.0	54.4	0.0	37.0	0.77
Gemini 3.1 Pro (medium)	20.0	0.0	0.0	54.4	0.0	38.0	0.57
Gemini 3.1 Pro (high)	0.0	0.0	20.0	52.9	8.3	39.0	0.46
Kimi K2.5 Thinking	0.0	0.0	0.0	42.6	0.0	29.0	0.38
DeepSeek V3.2	0.0	0.0	10.0	33.8	0.0	24.0	0.40
MiniMax M2.7	20.0	0.0	0.0	25.0	0.0	18.0	0.34
Qwen3.5 397B A17B	0.0	0.0	0.0	38.2	0.0	26.0	0.79

(Yu et al., 2024), QCBench (Xie et al., 2025), and MuMOInstruct (Dey et al., 2025) – are restricted to single-turn question answering on relatively simple chemistry tasks and do not evaluate complex agentic systems on real-world problems. Prior work on multi-turn iterative molecular optimization (Chennakesavalu et al., 2026; Wang et al., 2024) likewise lacks realistic agentic tool calling, instead querying the LLM repeatedly in a loop over a small number of simple optimization objectives. In contrast, harnessing LLM agents as molecular designers lets them use tools to understand the 3D problem setting, analyze previous failed hypotheses, and run simulations and other cheminformatics tools before committing to oracle calls. Iterative prompting approaches are not faithful representations of the true potential for LLM agents in drug discovery.

Beyond single-turn QA, prior LLM agents for drug design and discovery (Ock et al., 2026; Li et al., 2026a;b; Zhang et al., 2024; Wang et al., 2024) either report results on a few hand-selected successes or use ad hoc, non-standardized evaluation setups, often retrofitting datasets built for de novo molecular generation such as ZINC, Therapeutics Data Commons, GuacaMol, PMO, and MOSES (Tingle et al., 2023; Huang et al., 2021; Brown et al., 2019; Gao et al., 2022; Polykovskiy et al., 2020). Adjacent domains like computational biology have produced agentic benchmarks such as DeepResearch Bench (Du et al., 2025), BioProBench (Liu et al., 2025c), BioLP-Bench (Ivanov, 2024), ABC-Bench (Liu et al., 2025a), Lab-Bench (Laurent et al., 2024), BixBench (Mitchener et al., 2025), and CompBioBench (Nair et al., 2026), but no prior work evaluates LLMs on

agentic small molecule drug design.

The primary contributions of SMDD-Bench compared to previous work are as follows: (1) the first standardized benchmark for evaluating long-horizon, multi-turn, small molecule drug design **agents** on a chemically diverse array of challenging drug discovery tasks; (2) support for five relevant, real-world task types: 2D Pharmacophore Identification, Interaction Point Discovery, Scaffold Hopping, Lead Optimization, and Fragment Assembly; and (3) guaranteed solvability for all tasks, including the Scaffold Hopping, Lead Optimization, and Fragment Assembly tasks via the existence of a held-out witness molecule which passes all evaluation requirements.

6. Discussion and Conclusion

We present SMDD-Bench, the first agentic benchmark evaluating LLM agents on challenging, long-horizon, real-world small molecule drug design. SMDD-Bench tests scientific and chemical reasoning across diverse chemistries, task types, and protein targets, alongside domain-specific tool-calling capabilities, with every task instance generated to be guaranteed-solvable. Our hope is that SMDD-Bench unifies and standardizes the design, training, and evaluation of LLMs for small molecule drug design. The dominant approach to small molecule drug design today is conditional generative diffusion and flow models (Sun et al., 2025; Zeng et al., 2026; Zhang et al., 2025; Morehead et al., 2026; Klein et al., 2023; Hoogeboom et al., 2022), but LLM-based SMDD agents offer distinct advantages in flexibility and generality. An agent can perform lead optimization on properties unseen at training time, understand the failure modes of specific oracles, selecting the right tool for a given task, and adhere to constraints (such as maintaining protein-protein contacts when designing a PROTAC) that were never specified during training. Communicating objectives in natural language further frees the agent from the rigid train-time settings that bind diffusion and flow models.

Reliance on Boltz2 and ADMET-AI. Boltz2 and ADMET-AI are imperfect oracles. However, in SMDD-Bench we

440 treat these oracles as the ground truth within the "world" of
 441 the benchmark, analogous to an agent requesting wet lab
 442 assays. As a result, we deliberately provide a limited number
 443 of oracle calls to reflect the sparsity of real-world wet
 444 lab experiments. To enforce that the Boltz2 and ADMET-
 445 AI predicted distributions are the operative ground truth,
 446 all PDB-sourced complexes are re-folded via Boltz2 and
 447 re-scored by ADMET-AI during task generation, so the
 448 agent only sees structure-activity relationships and property
 449 distributions produced by these oracles. SMDD-Bench is intended
 450 as a testbed for developing training frameworks calibrated
 451 to the Boltz2 and ADMET-AI distributions. The same
 452 frameworks can later be interfaced with an autonomous laboratory,
 453 where the oracle becomes a real wet lab that synthesizes,
 454 characterizes, and tests candidates, **allowing agents to be re-trained under the true distribution of real-world structure-activity relationships.**

458 Impact Statement

460 SMDD-Bench presents a benchmark for the evaluation of
 461 LLM agents on small molecule drug design tasks. The
 462 intended impact is to provide the ML community with a
 463 rigorous testbed for scientific reasoning in a domain where
 464 progress could meaningfully accelerate therapeutic development.
 465 We acknowledge dual-use considerations inherent to drug
 466 discovery tooling: the same capabilities that enable
 467 beneficial therapeutics could in principle inform the design
 468 of harmful molecules. However, we believe marginal uplift
 469 this benchmark provides to a malicious actor is minimal
 470 relative to existing public tooling, while the upside is
 471 substantial: by making agent progress in chemistry legible,
 472 comparable, and grounded in scientifically meaningful
 473 verification, SMDD-Bench can help steer LLM development
 474 toward capabilities that translate into real benefit for patients
 475 and for the broader scientific community.

477 References

479 Asai, A., He, J., Shao, R., Shi, W., Singh, A., Chang, J. C.,
 480 Lo, K., Soldaini, L., Feldman, S., D’Arcy, M., et al.
 481 Synthesizing scientific literature with retrieval-augmented
 482 language models. *Nature*, 2026.

484 Ashton, M., Barnard, J., Casset, F., Charlton, M., Downs, G.,
 485 Gorse, D., Holliday, J., Lahana, R., and Willett, P. Identification
 486 of diverse database subsets using property-based and fragment-based
 487 molecular descriptions. *Quantitative Structure-Activity Relationships*,
 488 21(6):598–604, 2002.

490 Baell, J. B. Feeling nature’s pains: natural products, natural
 491 product drugs, and pan assay interference compounds (pains).
 492 *Journal of natural products*, 79(3):616–628, 2016.

Bartmann, C., Schimunek, J., Ielanskyi, M., Seidl, P., Klam-
 496 bauer, G., and Luukkonen, S. Molculariq: Characterizing
 497 chemical reasoning capabilities through symbolic verification on
 498 molecular graphs. *arXiv preprint arXiv:2601.15279*, 2026.

Bemis, G. W. and Murcko, M. A. The properties of known
 499 drugs. 1. molecular frameworks. *Journal of medicinal chemistry*,
 500 39(15):2887–2893, 1996.

Berman, H. M., Westbrook, J., Feng, Z., Gilliland, G., Bhat,
 501 T. N., Weissig, H., Shindyalov, I. N., and Bourne, P. E. The
 502 protein data bank. *Nucleic Acids Research*, 28(1):235–242,
 503 2000. doi: 10.1093/nar/28.1.235.

Boiko, D. A., MacKnight, R., Kline, B., and Gomes, G. Au-
 504 tonomous chemical research with large language models.
 505 *Nature*, 624(7992):570–578, 2023.

Brenk, R., Schipani, A., James, D., Krasowski, A., Gilbert,
 506 I. H., Frearson, J., and Wyatt, P. G. Lessons learnt from
 507 assembling screening libraries for drug discovery for neglected
 508 diseases. *ChemMedChem: Chemistry Enabling Drug Discovery*,
 509 3(3):435–444, 2008.

Brown, N., Fiscato, M., Segler, M. H., and Vaucher, A. C.
 510 Guacamol: benchmarking models for de novo molecular design.
 511 *Journal of chemical information and modeling*, 59(3):1096–1108,
 512 2019.

Chennakesavalu, S., Shmilovich, K., Weir, H., Grambow,
 513 C., Bradshaw, J., Suriana, P., Cheng, C., and Chuang, K. Evaluating
 514 the progression of large language model capabilities for small-
 515 molecule drug design. *arXiv preprint arXiv:2604.16279*, 2026.

Degen, J., Wegscheid-Gerlach, C., Zaliani, A., and Rarey,
 516 M. On the art of compiling and using ‘drug-like’ chemical
 517 fragment spaces. *ChemMedChem*, 3(10):1503, 2008.

Dey, V., Hu, X., and Ning, X. Gellm3o: Generalizing large
 518 language models for multi-property molecule optimization.
 519 *arXiv preprint arXiv:2502.13398*, 10, 2025.

Du, M., Xu, B., Zhu, C., Wang, X., and Mao, Z. Deep-
 520 research bench: A comprehensive benchmark for deep research
 521 agents. *arXiv preprint arXiv:2506.11763*, 2025.

Ertl, P. and Schuffenhauer, A. Estimation of synthetic access-
 522 ibility score of drug-like molecules based on molecular complexity
 523 and fragment contributions. *Journal of cheminformatics*,
 524 1(1):8, 2009.

Ester, M., Kriegl, H.-P., Sander, J., Xu, X., et al. A density-
 525 based algorithm for discovering clusters in large spatial
 526 databases with noise. In *kdd*, volume 96, pp. 226–231,
 527 1996.

- 495 Gao, W., Fu, T., Sun, J., and Coley, C. Sample efficiency
496 matters: a benchmark for practical molecular optimiza-
497 tion. *Advances in neural information processing systems*,
498 35:21342–21357, 2022.
- 499
500 Gleeson, M. P. Generation of a set of simple, interpretable
501 admet rules of thumb. *Journal of medicinal chemistry*, 51
502 (4):817–834, 2008.
- 503
504 Hoogeboom, E., Satorras, V. G., Vignac, C., and Welling, M.
505 Equivariant diffusion for molecule generation in 3d. In
506 *International conference on machine learning*, pp. 8867–
507 8887. PMLR, 2022.
- 508
509 Huang, K., Fu, T., Gao, W., Zhao, Y., Roohani, Y., Leskovec,
510 J., Coley, C. W., Xiao, C., Sun, J., and Zitnik, M. Ther-
511 apeutics data commons: Machine learning datasets and
512 tasks for drug discovery and development. *arXiv preprint*
513 *arXiv:2102.09548*, 2021.
- 514
515 Ivanov, I. Biolp-bench: Measuring understanding of biolog-
516 ical lab protocols by large language models. *bioRxiv*, pp.
517 2024–08, 2024.
- 518
519 Kim, S., Chen, J., Cheng, T., Gindulyte, A., He, J., He, S.,
520 Li, Q., Shoemaker, B. A., Thiessen, P. A., Yu, B., et al.
521 Pubchem 2023 update. *Nucleic acids research*, 51(D1):
522 D1373–D1380, 2023.
- 523
524 Klein, L., Krämer, A., and Noé, F. Equivariant flow match-
525 ing. *Advances in Neural Information Processing Systems*,
526 36:59886–59910, 2023.
- 527
528 Kuhn, H. W. The hungarian method for the assignment
529 problem. *Naval research logistics quarterly*, 2(1-2):83–
530 97, 1955.
- 531
532 Landrum, G. Rdkit: Open-source cheminformatics. <http://www.rdkit.org>, 2010. Accessed: 2026-04-11.
- 533
534 Laurent, J. M., Janizek, J. D., Ruzo, M., Hinks, M. M.,
535 Hammerling, M. J., Narayanan, S., Ponnampati, M., White,
536 A. D., and Rodrigues, S. G. Lab-bench: Measuring capa-
537 bilities of language models for biology research. *arXiv*
538 *preprint arXiv:2407.10362*, 2024.
- 539
540 Leeson, P. D. and Springthorpe, B. The influence of drug-
541 like concepts on decision-making in medicinal chemistry.
542 *Nature reviews Drug discovery*, 6(11):881–890, 2007.
- 543
544 Li, H., Cao, H., Peng, S., Liu, Z., Feng, B., Wang, Y., Yan,
545 Z., Tian, Y., Li, Y., and Yuan, L. Agentic reinforcement
546 learning empowers next-generation chemical language
547 models for molecular design and synthesis. *arXiv preprint*
548 *arXiv:2601.17687*, 2026a.
- 549
550 Li, T., Hou, K., Vinh, T., Raj, M., Guo, Z., and Yang, C.
551 Reinforcement learning with llm-guided action spaces
552 for synthesizable lead optimization. *arXiv preprint*
553 *arXiv:2604.07669*, 2026b.
- 554
555 Limban, C., Nuță, D. C., Chiriță, C., Negreș, S., Arsene,
556 A. L., Goumenou, M., Karakitsios, S. P., Tsatsakis, A. M.,
557 and Sarigiannis, D. A. The use of structural alerts to avoid
558 the toxicity of pharmaceuticals. *Toxicology reports*, 5:
559 943–953, 2018.
- 560
561 Liu, A. B., Nedungadi, S., Cai, B., Kleinman, A., Bhasin, H.,
562 and Donoughe, S. Abc-bench: An agentic bio-capabilities
563 benchmark for biosecurity. In *NeurIPS 2025 Workshop*
564 *on Biosecurity Safeguards for Generative AI*, 2025a.
- 565
566 Liu, T., Lin, Y., Wen, X., Jorissen, R. N., and Gilson, M. K.
567 Bindingdb: a web-accessible database of experimentally
568 determined protein–ligand binding affinities. *Nucleic*
569 *acids research*, 35(suppl_1):D198–D201, 2007.
- 570
571 Liu, T., Hwang, L., Burley, S. K., Nitsche, C. I., Southan,
572 C., Walters, W. P., and Gilson, M. K. Bindingdb in
573 2024: a fair knowledgebase of protein-small molecule
574 binding data. *Nucleic acids research*, 53(D1):D1633–
575 D1644, 2025b.
- 576
577 Liu, Y., Lv, L., Zhang, X., Yuan, J. W. L., and Tian, Y.
578 Bioprobench: Comprehensive dataset and benchmark in
579 biological protocol understanding and reasoning. *arXiv*
580 *preprint arXiv:2505.07889*, 2025c.
- 581
582 Lu, C., Lu, C., Lange, R. T., Yamada, Y., Hu, S., Foerster,
583 J., Ha, D., and Clune, J. Towards end-to-end automation
584 of ai research. *Nature*, 651(8107):914–919, 2026.
- 585
586 Mendez, D., Gaulton, A., Bento, A. P., Chambers, J.,
587 De Veij, M., Félix, E., Magariños, M. P., Mosquera,
588 J. F., Mutowo, P., Nowotka, M., Gordillo-Marañón, M.,
589 Hunter, F., Junco, L., Mugumbate, G., Rodriguez-Lopez,
590 M., Atkinson, F., Bosc, N., Radoux, C. J., Segura-
591 Cabrera, A., Hersey, A., and Leach, A. R. ChEMBL: to-
592 wards direct deposition of bioassay data. *Nucleic acids*
593 *research*, 47(D1):D930–D940, January 2019. ISSN
594 0305-1048. doi: 10.1093/nar/gky1075. URL <https://europepmc.org/articles/PMC6323927>.
- 595
596 Mirza, A., Alampara, N., Kunchapu, S., Ríos-García,
597 M., Emoekabu, B., Krishnan, A., Gupta, T., Schilling-
598 Wilhelmi, M., Okereke, M., Aneesh, A., et al. A frame-
599 work for evaluating the chemical knowledge and reason-
600 ing abilities of large language models against the exper-
601 tise of chemists. *Nature Chemistry*, 17(7):1027–1034,
602 2025.
- 603
604 Mitchener, L., Laurent, J. M., Andonian, A., Tenmann, B.,
605 Narayanan, S., Wellawatte, G. P., White, A., Sani, L., and
606 Rodrigues, S. G. Bixbench: a comprehensive benchmark
607 for llm-based agents in computational biology. *arXiv*
608 *preprint arXiv:2503.00096*, 2025.

- 550 Morehead, A., Cretu, M., Panescu, A., Anand, R., Weiler,
551 M., Perez, T., Blau, S., Farrell, S., Bhimji, W., Jain,
552 A., et al. Zatom-1: A multimodal flow foundation
553 model for 3d molecules and materials. *arXiv preprint*
554 *arXiv:2602.22251*, 2026.
- 555 Nair, S., Gunsalus, L., Orcutt-Jahns, B., Rossen, J., Lal, A.,
556 Donno, C. D., Celik, M. H., Fletez-Brant, K., Xie, X.,
557 Bravo, H. C., et al. Agentic systems are adept at solv-
558 ing well-scoped, verifiable problems in computational
559 biology. *bioRxiv*, pp. 2026–04, 2026.
- 561 O’Boyle, N. M., Banck, M., James, C. A., Morley, C., Van-
562 dermeersch, T., and Hutchison, G. R. Open babel: An
563 open chemical toolbox. *Journal of Cheminformatics*, 3:
564 33, 2011. doi: 10.1186/1758-2946-3-33.
- 566 Ock, J., Meda, R. S., Badrinarayanan, S., Aluru, N. S., Chan-
567 drasekhar, A., and Barati Farimani, A. Large language
568 model agent for modular task execution in drug discovery.
569 *Journal of Chemical Information and Modeling*, 66(4):
570 2055–2068, 2026.
- 572 Papadatos, G., Davies, M., Dedman, N., Chambers, J.,
573 Gaulton, A., Siddle, J., Koks, R., Irvine, S. A., Petters-
574 son, J., Goncharoff, N., et al. Surechembl: a large-scale,
575 chemically annotated patent document database. *Nucleic*
576 *acids research*, 44(D1):D1220–D1228, 2016.
- 578 Passaro, S., Corso, G., Wohlwend, J., Reveiz, M., Thaler, S.,
579 Somnath, V., Getz, N., Portnoi, T., Roy, J., Stark, H., et al.
580 Boltz-2: Towards accurate and efficient binding affinity
581 prediction. *biorxiv 2025*. *Google Scholar*, 2025.
- 582 Polykovskiy, D., Zhebrak, A., Sanchez-Lengeling, B., Golo-
583 vanov, S., Tatanov, O., Belyaev, S., Kurbanov, R., Arta-
584 monov, A., Aladinskiy, V., Veselov, M., et al. Molecular
585 sets (moses): a benchmarking platform for molecular gen-
586 eration models. *Frontiers in pharmacology*, 11:565644,
587 2020.
- 589 Salentin, S., Schreiber, S., Haupt, V. J., Adasme, M. F., and
590 Schroeder, M. Plip: fully automated protein–ligand inter-
591 action profiler. *Nucleic acids research*, 43(W1):W443–
592 W447, 2015.
- 594 Sun, Y., Chen, L., Jing, Z., Li, Y. Y., Kim, D., Gao, J.-Y.,
595 Noroozi, R., Yi, G. Y., Tetsassi Feugmo, C. G., Klinkova,
596 A., et al. Generative ai for the design of molecules: ad-
597 vances and challenges. *Journal of Chemical Information*
598 *and Modeling*, 65(23):12668–12690, 2025.
- 599 Swanson, K., Walther, P., Leitz, J., Mukherjee, S., Wu,
600 J. C., Shivnaraine, R. V., and Zou, J. Admet-ai:
601 a machine learning admet platform for evaluation of
602 large-scale chemical libraries. *Bioinformatics*, 40(7):
603 btae416, 06 2024. ISSN 1367-4811. doi: 10.1093/
604 bioinformatics/btae416. URL <https://doi.org/10.1093/bioinformatics/btae416>.
- Tingle, B. I., Tang, K. G., Castanon, M., Gutierrez, J. J.,
Khurelbaatar, M., Dandarchuluun, C., Moroz, Y. S., and
Irwin, J. J. Zinc-22– a free multi-billion-scale database
of tangible compounds for ligand discovery. *Journal of*
chemical information and modeling, 63(4):1166–1176,
2023.
- Varma, M. V., Feng, B., Obach, R. S., Troutman, M. D.,
Chupka, J., Miller, H. R., and El-Kattan, A. Physico-
chemical determinants of human renal clearance. *Journal*
of medicinal chemistry, 52(15):4844–4852, 2009.
- Wang, H., Skreta, M., Ser, C.-T., Gao, W., Kong, L., Strieth-
Kalthoff, F., Duan, C., Zhuang, Y., Yu, Y., Zhu, Y., et al.
Efficient evolutionary search over chemical space with
large language models. *arXiv preprint arXiv:2406.16976*,
2024.
- Waring, M. J. Lipophilicity in drug discovery. *Expert*
opinion on drug discovery, 5(3):235–248, 2010.
- Xie, J., Wang, W., Gao, B., Yang, Z., Wan, H., Zhang,
S., Fu, T., and Li, Y. Qcbench: Evaluating large lan-
guage models on domain-specific quantitative chemistry.
Journal of Chemical Information and Modeling, 65(22):
12268–12278, 2025.
- Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan,
K. R., and Cao, Y. React: Synergizing reasoning and
acting in language models. In *The eleventh international*
conference on learning representations, 2022.
- Yu, B., Baker, F. N., Chen, Z., Ning, X., and Sun, H. Llas-
mol: Advancing large language models for chemistry
with a large-scale, comprehensive, high-quality instruc-
tion tuning dataset. *arXiv preprint arXiv:2402.09391*,
2024.
- Zeng, C., Jin, J., Ambrose, C., Karypis, G., Transtrum, M.,
Tadmor, E. B., Hennig, R. G., Roitberg, A., Martiniani, S.,
and Liu, M. Propmolflow: property-guided molecule gen-
eration with geometry-complete flow matching. *Nature*
Computational Science, pp. 1–10, 2026.
- Zhang, O., Jin, J., Wu, Z., Zhang, J., Yuan, P., Yu, Y., Lin,
H., Zhong, H., Zhang, X., Hua, C., et al. Ecloudgen:
leveraging electron clouds as a latent variable to scale up
structure-based molecular design. *Nature Computational*
Science, pp. 1–12, 2025.
- Zhang, Y., Chen, X., Jin, B., Wang, S., Ji, S., Wang, W.,
and Han, J. A comprehensive survey of scientific large
language models and their applications in scientific dis-
covery. In *Proceedings of the 2024 Conference on Em-
pirical Methods in Natural Language Processing*, pp.
8783–8817, 2024.

A. Detailed Benchmark Diversity

The complete list of the 102 unique protein targets broken down into task types can be found in Figure 4. In Figure 5, we plot each of the reference molecule baseline property values for optimization properties and hold-constant properties in the Lead Optimization task. Blood brain barrier penetration can either be optimized to enhance penetration or to decrease penetration depending on the real world application. Both directions of optimization are included in SMDD-Bench.

In Figure 6, we plot the pairwise Tanimoto similarity between all pairs of reference molecules in the Scaffold Hopping and Lead Optimization tasks. The low distribution of Tanimoto similarity indicates that the bulk of the reference molecules aren’t collapsed to similar points in chemical space – highlighting the small molecule diversity of SMDD-Bench.

Figure 7 shows the breakdown of the protein targets within each task instance into their respective families, highlighting the protein diversity of SMDD-Bench. For an agent to be successful, it must be able to generalize and reason across diverse chemistries, pockets, and biological targets.

B. Details on Task Type Generation

In the following section, we provide the specific details for generating the task instances for each task type.

B.1. 2D Pharmacophore Identification

Each instance of this task type asks the agent to identify a 2D pharmacophore that distinguishes a set of active molecules from a set of inactive molecules for a given protein target. Tasks are generated procedurally from ChEMBL.

Activity data extraction We restrict to direct binding assays by requiring `standard_type` \in $\{IC_{50}, K_i\}$, `standard_units` = nM, `standard_relation` = '=', and `confidence_score` \geq 8. Only targets with `target_type` = 'SINGLE PROTEIN' are considered.

Compound normalization and deduplication. SMILES are canonicalized via RDKit. We discard multi-component salts/mixtures, molecules with exact molecular weight greater than 900 Da, and entries that fail RDKit parsing. Multiple measurements of the same canonical compound are collapsed by taking the per-compound median activity in nM.

Active/inactive labeling Compounds with median activity below 1,000 nM ($IC_{50}/K_i < 1 \mu M$) are labeled active; compounds with median activity above 10,000 nM ($> 10 \mu M$) are labeled inactive. Compounds in the 1–10 μM ambiguous zone are discarded, providing a 10 \times potency gap between the two classes that excludes borderline molecules. Targets with fewer than 30 actives or fewer than 30 inactives after filtering are not eligible.

Diverse subset selection For each class, ten visible exemplars are selected via the MaxMin algorithm operating on Morgan fingerprints (radius = 2, 2048 bits, seed = 42) (Ashton et al., 2002). The remaining compounds form a held-out evaluation set that the agent never sees. MaxMin selection forces the visible exemplars to span chemical space, allowing the agent to see the most chemically diverse spread of actives as possible.

Agent inputs and outputs Each generated task directory contains: `actives.smi` and `inactives.smi` (the 10 visible exemplars per class, given to the agent), `protein.fasta` (the target sequence). The agent is required to emit `solution.py` defining `check_pharmacophore(smiles: str) -> bool`, which is run on every molecule of the held-out evaluation set during scoring.

Design rationale Using ChEMBL-measured inactives rather than synthetically constructed decoys (e.g. DUD-E) ensures that negatives are experimentally verified non-binders that may share scaffolds with the actives, ruling out solutions that rely on coarse topological dissimilarity. Providing both classes to the agent and requiring generalization to a held-out set forces extraction of a genuine substructural pattern rather than memorization of the visible exemplars.

B.2. Interaction Point Discovery

Each instance of this task type asks the agent to predict, given a single apo-like protein structure and a pocket center, the 3D coordinates and pharmacophoric type of the three interaction points most likely to be conserved across diverse binders of that pocket. Tasks are generated procedurally from the RCSB PDB (Berman et al., 2000).

Ensemble construction For each specified target from the PDB, we fetch up to 500 unique X-ray crystal structures with

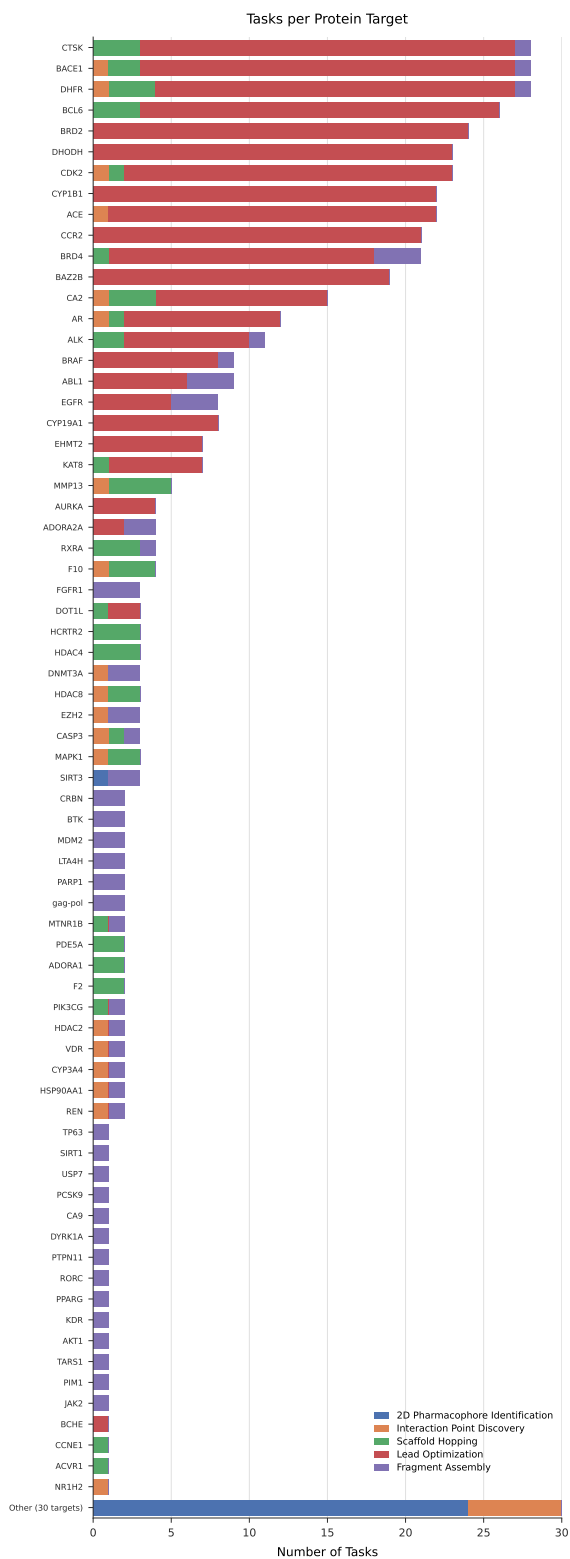


Figure 4. The complete breakdown of task instances into protein targets. There are 102 unique protein targets in SMDD-Bench.

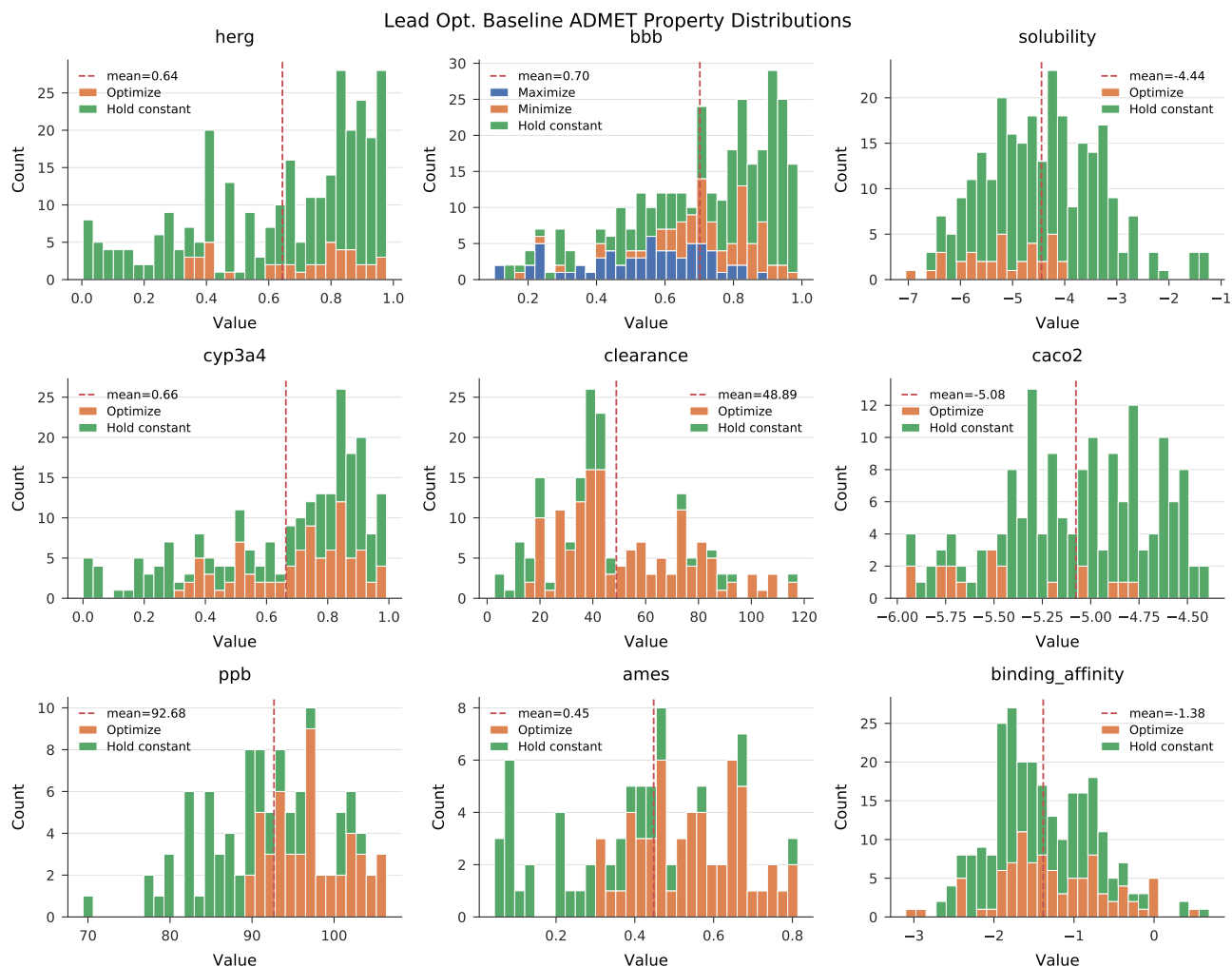


Figure 5. The baseline ADMET property values for all optimization objectives and hold-constant properties for the Lead Optimization task type. These are the values calculated for the reference molecules that the LLM agent must either optimize or hold constant. Note that blood brain barrier penetration (bbb) can either be a property that is minimized or maximized. All other properties are either only minimized or only maximized.

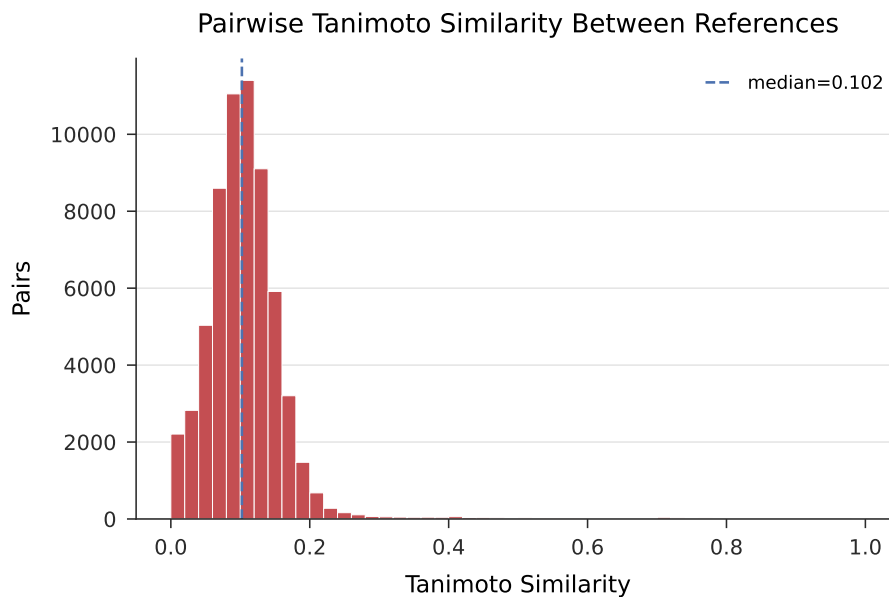


Figure 6. A histogram of the tanimoto similarities between all pairs of reference molecules provided to the LLM agent in the Scaffold Hopping and Lead Optimization tasks. The low median tanimoto similarity between pairs indicates the chemical diversity of SMDD-Bench.

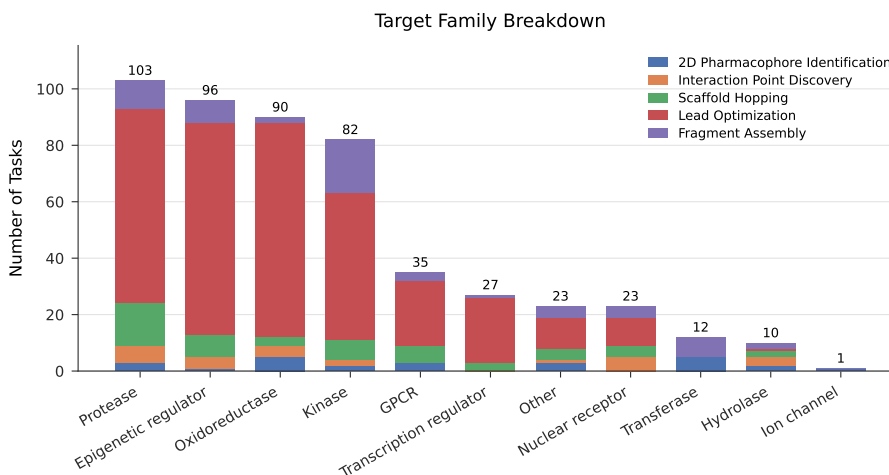


Figure 7. A breakdown of the task types of SMDD-Bench into the families of the protein targets associated with each task instance.

825 resolution better than 2.5 Å. The first structure containing a drug-like ligand is taken as the *reference*; common crystallization
826 additives (glycerol, PEG, DMSO, sulfate, and similar) are explicitly blacklisted. Within that reference structure, the largest
827 non-blacklisted ligand by heavy-atom count is selected as the reference ligand. The reference receptor is stripped of all
828 heteroatoms, waters, and ions to produce `receptor.pdb`, and the centroid of the reference ligand defines the pocket
829 center provided to the agent.

830 **Local structural alignment** All ensemble members are superimposed onto the reference using only the C_{α} atoms of
831 residues within 12 Å of the pocket center. Restricting the alignment to the local pocket suppresses contributions from global
832 protein motion and induced fit – minimizing potential global variance in the 3D reference frame.

833 **Pharmacophore feature extraction** For each aligned structure, we select the ligand whose centroid lies closest to the
834 reference pocket center (within 12 Å); waters, the same blacklist of crystallization additives and common ions used during
835 reference selection, and small fragments (fewer than six heavy atoms) are excluded from candidacy. RDKit is then used to
836 assign each atom or functional group of the chosen ligand to one of six pharmacophoric families: hydrogen-bond donor,
837 hydrogen-bond acceptor, aromatic, hydrophobic, cation, or anion. Features further than 12 Å from the pocket center are
838 discarded to exclude allosteric binders and crystallographic artifacts.

840 **Consensus clustering** Clustering proceeds in two stages. First, DBSCAN (Ester et al., 1996) is run independently per
841 pharmacophore family at $\epsilon = 1.0$ Å on the 3D points of that family alone, so that, e.g., donor and acceptor densities never
842 collapse into a single cluster regardless of spatial proximity. A cluster is retained only if it is supported by at least 15%
843 of the unique structures contributing to the ensemble; the conservation score of a cluster is the fraction of unique PDB
844 entries supporting it. Cross-family relationships between surviving clusters are handled separately during ground-truth point
845 selection.

846 **Ground-truth point selection** Surviving clusters are post-processed to yield a physically valid, non-redundant ground-truth
847 set: (i) clusters within 2 Å of any heavy atom of the reference receptor are discarded as sterically clashing with the reference
848 protein conformation given to the agent; (ii) clusters of different types (e.g. donor and cation arising from a basic amine, or
849 donor and acceptor from a hydroxyl) whose centroids lie within 2 Å are merged into a single point with a set-valued type
850 label, reflecting and allowing for chemical ambiguity in the agent output. The full set of $N \geq 3$ surviving consensus points
851 constitutes the ground truth.

853 **Agent inputs and outputs** The agent is provided `receptor.pdb` (the reference protein with all ligands and solvent
854 removed) the pocket center coordinates and the 12 Å pocket radius where all interaction points are within. The agent
855 must emit `solution.csv` with exactly three rows of the form `x, y, z, type`, where `type` is drawn from {Donor,
856 Acceptor, Aromatic, Hydrophobic, Cation, Anion}.

858 **Evaluation** Predicted points are matched against the N ground-truth points by optimal bipartite assignment (Kuhn, 1955)
859 on Euclidean distance, restricted to type-compatible pairs: a prediction is type-compatible with a ground-truth point if its
860 label is in the ground-truth label set, with an additional aromatic/hydrophobic equivalence at scoring time, reflecting that
861 aromatic ring features and hydrophobic features are partially co-tagged by RDKit and that the ground-truth construction
862 may assign only one of the two labels at a position where both are chemically reasonable. This choice was made out of
863 evaluation leniency. A match is counted if the assigned distance is below 2.5 Å. The task score is (number of matched
864 predictions) / 3, and the task is considered solved only when all three predictions match (score = 1.0).

865 **Design rationale** The local-pocket alignment, conservation threshold, and steric clash filter were designed to ensure that the
866 targets are simultaneously experimentally grounded, chemotype-agnostic, empirically solvable, and medicinally useful in
867 the specific receptor conformation handed to the agent.

869 B.3. Scaffold Hopping

871 Each instance of this task type asks the agent to propose a novel small molecule that binds the same protein pocket as a given
872 reference ligand and reproduces its protein–ligand interaction pattern, while presenting a chemically distinct scaffold. Tasks
873 are generated procedurally from the RCSB PDB and ChEMBL, and every task is constructed around an explicit *witness*
874 molecule that is known and evaluated a priori to satisfy the evaluation criteria – thus guaranteeing the solvability of the task
875 instance.

876 **Target and canonical structure selection** Targets are sourced from a curated pool of ~ 300 (UniProt, PDB, ligand)
877 triplets spanning diverse protein families, grouped by UniProt identifier. For each UniProt, the PDB entry with the best
878

X-ray resolution is chosen as the canonical receptor; its FASTA sequence and a single fixed pocket definition are reused throughout downstream steps.

Pocket definition The binding pocket is defined as the set of residues with any heavy atom within 6 Å of the canonical ligand. This pocket, expressed as (chain, SEQRES position) pairs, is held fixed for all downstream Boltz2 cofolding runs and is also handed to the agent as part of the task specification.

Candidate ligand collection For each target we assemble a candidate ligand pool from two sources: (i) all ligands co-crystallized with the same UniProt across the PDB, with bond orders corrected via Chemical Component Dictionary (CCD) templates and deduplicated by canonical SMILES; and (ii) ChEMBL actives against the target with reported $IC_{50}/K_i/K_d$ below 10 μ M, pre-filtered to remove unparseable structures, molecules above 1,000 Da, and metal-containing species.

Scaffold-diversity pre-filter Bemis–Murcko scaffolds are extracted for every candidate, and for each unordered ligand pair we compute the Maximum Common Substructure (MCS) of the two scaffolds (Bemis & Murcko, 1996). A pair is retained only if the MCS covers less than 65% of the heavy atoms of the larger scaffold; the symmetric denominator ensures the filter does not depend on the eventual reference/witness role assignment.

Crystal-PLIP pre-ranking For pairs in which both members appear as PDB co-crystals, each source structure is superimposed onto the canonical receptor using C_α atoms of pocket residues within 12 Å, and a protein-ligand interaction fingerprint is extracted with PLIP (Salentin et al., 2015). Pairs are then ordered by descending crystal interaction-fingerprint similarity, with PDB-PDB pairs prioritized over pairs involving ChEMBL-only ligands. This ranking concentrates downstream compute on pairs that are most likely to clear the interaction-similarity gate during the co-folding screen.

Co-folding screen For each target, the top-ranked unique candidate ligands are co-folded with the canonical receptor using Boltz2 (Passaro et al., 2025), constrained to the fixed pocket definition with 10 diffusion samples per ligand. From the resulting ensemble of 10 conformations, we extract: (i) the Boltz2-predicted binding probability across samples; (ii) per-sample PLIP interaction fingerprints, where each interaction is coarsened to a (chain, residue, broad type) tuple over six interaction families (HBond, Aromatic, Ionic, Halogen, Metal, Hydrophobic); and (iii) a consensus interaction fingerprint defined as the set of tuples present in at least 60% of the diffusion samples. Any ligand whose predicted binding probability falls below 0.75 is discarded. The consensus interaction fingerprint is required during evaluation to offset variance induced by Boltz2-predicted binding poses.

Witness validation and role assignment A pair is promoted to a task only if both members survive the binding gate and the set Tanimoto between their consensus interaction fingerprints exceeds 0.75. Within a surviving pair, the ligand with the higher predicted binding probability is assigned as the *reference* (the molecule provided to the agent to hop away from) and the other becomes the witness (the hidden molecule whose existence guarantees the solvability of the task instance). We retain at most three tasks per UniProt target and require distinct reference molecules across tasks of the same target, preventing one well-behaved reference from dominating the benchmark.

Agent inputs and outputs The agent is given: `complex.pdb` (the Boltz2-predicted reference complex, with the ligand on chain L), `reference.sdf` (the reference ligand conformer extracted from the Boltz2 complex), `protein.fasta` (the canonical protein sequence), and the target ligand identifier and pocket contact residue list used during the task generation process. The agent is required to emit `solution.smi` containing a single SMILES string. The witness molecule is never exposed to the agent.

Evaluation A submission must clear five gates: (1) the SMILES is parseable by RDKit; (2) the Morgan fingerprint Tanimoto (radius = 2, 2048 bits) between the agent’s molecule and the reference is below 0.5; (3) the Bemis–Murcko scaffold MCS, normalized by the heavy-atom count of the larger of the two scaffolds, is below 0.65; (4) the agent’s molecule, co-folded under the same Boltz protocol used during generation (10 diffusion samples, identical pocket constraint), produces a consensus interaction fingerprint whose set Tanimoto with the reference’s consensus fingerprint exceeds 0.75; and (5) the predicted binding probability of the agent’s molecule across the same diffusion samples exceeds 0.7. Gates 1, 2, and 3 enforce 2D novelty; gates 4 and 5 enforce 3D mimicry, ensuring the scaffold hopped molecule adheres to the same pharmacophore, and binding competence.

Design rationale Every task is anchored to a concrete witness molecule that has already cleared the same evaluation pipeline the agent will be scored against, ensuring the existence of at least one solution for the task instance. The 2D-dissimilarity and 3D-interaction-similarity criteria rule evaluates for a true scaffold hopped molecule and that the submitted molecule matches the binding mode of the original reference molecule. High predicted binding probability of the submitted molecule

Table 6. List of hard-constraints for molecules evaluated in Lead Optimization and Fragment Assembly tasks.

Category	Constraint
SMILES validity	Valid SMILES only
Molecular weight	MW < 600 Da
Lipophilicity	$\log P \in [-1, 5]$
Topological polar surface area	TPSA < 140 Å ²
H-bond donors	≤ 5
H-bond acceptors	≤ 10
Rotatable bonds	≤ 10
Formal charge	∈ [-2, +2]
Synthetic accessibility	SA score < 4.5 (Ertl & Schuffenhauer, 2009)
Structural alerts	No PAINS (Baell, 2016) or Brenk/NIH alerts (Brenk et al., 2008; Limban et al., 2018)

Table 7. List of ADMET properties used as optimization objectives and hold-constant constraints.

Property
hERG inhibition
Blood-brain barrier permeability
Aqueous solubility
CYP3A4 inhibition
Hepatocyte clearance
Caco-2 permeability
Plasma protein binding
Ames mutagenicity

is also enforced. We score interaction similarity against a consensus fingerprint over the 10-sample diffusion ensemble rather than against the single most confident pose to ensure robustness against variance in the Boltz2 poses. Requiring a tuple to appear in at least 60% of samples discards interactions that are artifacts of any one diffusion draw and retains only those the model is confident enough to recover under structural perturbation.

B.4. Lead Optimization

Each instance of this task type asks the agent to modify a known active molecule so that it improves a specified set of pharmacological optimization objectives, holds a separate set of properties within tolerance, and clears a fixed panel of hard drug-likeness constraints. As in the Scaffold Hopping task type, every task instance is constructed around an explicit witness molecule that has already been verified to satisfy every gate the agent will be scored against, guaranteeing the solvability of the task instance.

Target and pocket selection Targets are drawn from a curated set of UniProt entries with high-resolution co-crystal coverage. For each target, a single representative PDB is chosen; its FASTA sequence and the residues within 6 Å of the crystal ligand define a fixed pocket that is reused for all candidate co-folding within the target.

Candidate pool collection For each target we aggregate experimentally reported actives from ChEMBL, BindingDB, and PubChem, deduplicated by canonical SMILES. The pool is then filtered through the same hard-constraint panel that will be applied to the agent’s submission at evaluation time. The hard-constraint panel is found in Table 6.

ADMET annotation Every surviving molecule is scored on eight ADMET properties (hERG inhibition, blood-brain-barrier permeability, aqueous solubility, CYP3A4 inhibition, hepatocyte clearance, Caco-2 permeability, plasma protein binding, Ames mutagenicity) via the ADMET-AI ensemble (Swanson et al., 2024) trained on the corresponding TDC datasets (Huang et al., 2021).

Co-folding screen Each filtered candidate is co-folded against the canonical receptor with Boltz2 (Passaro et al., 2025), constrained to the fixed pocket via fixed pocket contacts with 10 diffusion samples. We retain the predicted binding probability and predicted $\log_{10}(\text{IC}_{50})$ across samples, and discard any molecule whose predicted binding probability falls below 0.75. The pool of molecules surviving every filter constitutes the per-target validated pool from which both reference and witness are drawn.

Pair finding and task generation Within each target’s validated pool, we compute pairwise Morgan fingerprint Tanimoto

(radius = 2, 2048 bits) and consider only pairs with similarity ≥ 0.7 .

For each filtered pair of molecules within the validated pool, we derive the optimization objectives and hold-constant set directly from the measured property differences between the two molecules. A property is eligible to be an optimization objective if the witness improves over the reference by at least a property-specific threshold. A property is eligible to be a hold-constant constraint only if the witness’s value already lies within the property’s tolerance band relative to the reference.

Up to five objectives and four hold-constant properties are selected greedily. The direction of BBB optimization is inferred from the witness properties rather than imposed externally. Furthermore, for the task instances without binding affinity as optimization objectives or hold-constant constraints (only a handful in SMDD-Bench), Boltz2 calculated binding probability is enforced as a hard constraint for the agent’s submitted molecule to prevent the agent from throwing away binding in the pursuit of optimizing ADMET properties.

Agent inputs and outputs The agent is given: `complex.pdb` (the Boltz2-predicted reference complex, with the ligand on chain L), `reference.sdf` (the reference ligand conformer extracted from the complex), `protein.fasta` (the canonical sequence), and the pocket residue list, the target ligand identifier, the list of optimization objectives with directions and thresholds, the list of hold-constant properties with tolerance bands, the full hard-constraint panel, the BBB direction, and the baseline property values measured on the reference. The agent is required to emit `solution.smi` containing a single SMILES string. The witness molecule is never exposed to the agent.

Evaluation A submission must clear five gates: (1) the SMILES is valid and sanitizable; (2) every hard constraint in the panel is satisfied, including a Tanimoto similarity to the reference ≥ 0.7 that anchors the submission to the same lead series; (3) every hold-constant property lies within its tolerance band relative to the reference baseline; (4) every optimization objective improves over the reference by at least its required threshold; and (5) Boltz2 co-folding of the agent’s molecule under the same protocol used during generation produces a predicted binding probability above 0.7 together with the affinity checks implied by the binding-affinity objective or hold-constant constraint. ADMET properties are recomputed from the agent’s SMILES with the same ADMET-AI ensemble used in generation.

Design rationale Deriving objectives and hold-constant sets from the actual reference–witness property delta, rather than imposing a fixed objective list per target, produces tasks that provably solvable. The high-Tanimoto pairing constraint enforces that the task captures local lead-series chemistry rather than de novo design, which is the regime in which the multi-objective trade-offs are most relevant. Because the entire hard-constraint panel and the ADMET model are applied identically at generation and at evaluation, the agent is scored under the exact filter that defined task feasibility, and any submission that clears the gates is at least as good as a known real molecule on every property the task names.

We investigate the indicators of the difficulty of the Lead Optimization tasks in Section D.1.

B.5. Fragment Assembly

Each instance of this task type asks the agent to design a single drug-like molecule that incorporates one or two given 3D molecular fragments at their specified poses inside a protein binding pocket, covalently links them when more than one is provided, and binds the target well. As in the previous two task types, every task instance is constructed around an explicit witness molecule – here, the known active from which the fragments were carved – that has already been verified to satisfy every gate the agent will be scored against, guaranteeing the solvability of the task instance.

Target and active selection Targets are sourced from UniProt entries that have at least two PDB entries containing distinct, non-buffer ligands of at least 10 heavy atoms, after filtering crystallization additives. For each (UniProt, PDB, ligand) triplet, the ligand SMILES is reconstructed via Chemical Component Dictionary templates and the binding pocket is defined as the protein residues with any heavy atom within 6 Å of the ligand.

Witness co-folding Each candidate active is co-folded against its source protein with Boltz2 (Passaro et al., 2025) under pocket constraints derived from the crystal pose, with the affinity head enabled and 10 diffusion samples. Only chains that participate in pocket contacts are retained for the co-fold; co-ligands and unrelated chains are excluded. Actives are kept only if their predicted binding probability exceeds 0.75. The resulting predicted $\log_{10}(\text{IC}_{50})$ is recorded as the witness affinity baseline against which the agent’s submission will be scored.

Fragment generation Surviving actives are decomposed via BRICS (Degen et al., 2008) into chemically reasonable fragments at retrosynthetically meaningful bonds. Each fragment is retained only if it (i) contains at least 6 heavy atoms after dummy-atom removal, (ii) makes at least one PLIP (Salentin et al., 2015) protein-ligand interaction with the receptor

when scored on the witness’s docked pose, and (iii) lies within 4 Å of the protein. The first criterion ensures the fragment is large enough to be a meaningful pharmacophore rather than a methyl-scale stub; the second and third together guarantee that the fragment is genuinely engaged with the pocket rather than solvent-exposed.

Pair sampling for two-fragment tasks When a task instance is to contain two fragments, the pair is drawn from a single source active and is required to have a closest-heavy-atom distance of at least 6 Å, this prevents the fragments from being spatially close enough to admit a short or simple linker.

Agent inputs and outputs The agent is given: `pocket.pdb` (the witness’s Boltz2-predicted protein structure with the witness ligand stripped), one `fragment_i.sdf` per fragment (the 3D pose of each fragment as it sat in the witness’s Boltz2 docked complex), the binding affinity baseline target for the submitted molecule, and the full hard-constraint panel also used in the Lead Optimization task (see Table 6). The agent is required to emit `solution.smi` containing a single SMILES string. The witness is never exposed to the agent.

Evaluation A submission must clear five gates: (1) the SMILES is valid; (2) every fragment SMARTS pattern is a substructure of the agent’s molecule; (3) the fragments admit a disjoint substructure assignment (no shared atoms across fragments) and are pairwise connected by a covalent path in the agent’s molecule; (4) the agent’s molecule clears the same hard-constraint panel specified in Table 6; and (5) Boltz2 co-folding of the agent’s molecule against the source receptor under the same pocket constraint and 10 diffusion samples produces (a) a predicted binding probability above 0.7, (b) a predicted $\log_{10}(\text{IC}_{50})$ no more than 0.1 log units worse than the witness baseline, and (c) a per-fragment 3D pose RMSD below 2.0 Å. RMSD is computed by SMARTS-matching each fragment against the docked agent ligand, extracting the 3D coordinates of the matched atoms, aligning the docked complex to `pocket.pdb` by Kabsch superposition on the C_{α} atoms of the pocket residues, and selecting the best-aligned diffusion sample for the comparison.

Design rationale The single-active witness construction guarantees that for every task there exists at least one molecule – the original active itself – that contains all fragments at exactly the required poses, links them through known chemistry, and already satisfies the binding and drug-likeness panels under the same Boltz2 protocol used at evaluation time. The 3D pose RMSD threshold enforces the docked geometry to recapitulate the fragment poses given as input, which forces the linker chemistry to be physically compatible with the relative orientation of the fragments rather than just being topologically valid or bind well.

C. Minimalist Agent Setup

To isolate model capability from harness engineering, all baseline runs use a deliberately minimalist ReAct loop. The agent sees only the obfuscated task description and the tool schemas. All tools run in an isolate docker container with an isolated filesystem. The standard agent configuration can be found in Table 8.

Tools. The agent is provided four tools: (i) `python(code)`: executes Python in the sandbox with RDKit, OpenBabel, PLIP, BioPython, NumPy, SciPy, and scikit-learn preinstalled; (ii) `predict_boltz(smiles, protein_sequences, ...)`: runs Boltz-2 co-folding with affinity prediction; (iii) `predict_admet(smiles)`: runs ADMET-AI property predictors and returns all the values of all properties of the molecule in a single call; The Boltz/ADMET oracles live in an isolated venv that the agent’s `python` tool cannot import, forcing all oracle calls through the budgeted interface.

Sandbox. Each task runs in a fresh `smdd-agent` container on an internal Docker network with no external routing; the only reachable endpoint is `api.colabfold.com:443` via an `nginx` TCP-passthrough sidecar (required by Boltz’s MSA server). Task files (e.g. `complex.pdb`, `reference.sdf`) are obfuscated to remove target names, PDB codes, and ligand three letter codes to suppress data-leakage shortcuts via memorized identifiers.

C.1. Selecting Oracle Calls

By default, across all evaluations, we set the total number of available Boltz2 calls to 8 and the number of ADMET-AI calls to 15. Note that a single ADMET-AI call returns all 8 ADMET properties for the inputted SMILES string. The complete reasoning behind only providing the LLM with a limited number of oracle calls is outlined in Section 6. In short, our goal is to treat the oracle calls as though they were ground truth, real-world wet lab experiments, hence the limited number of oracle calls. Ultimately, SMDD-Bench can be considered a testbed for developing the right training strategies and agent harnesses for SMDD agents – even though agents trained on SMDD-Bench adhere to the Boltz2/ADMET-AI distribution of binding affinities, structures, and properties. Afterwards, the agents can be trained using those same strategies with a

Table 8. Baseline agent harness configuration. All harness values are constants across tasks and models. Boltz diffusion samples was set to 10 to mitigate stochasticity in Boltz outputs. We empirically find settings above 3 diffusion samples to sufficiently reduce stochasticity across multiple runs.

Parameter	Value	Notes
<i>Sampling</i>		
Temperature	1.0	—
Sampling seed	42	—
Max ReAct turns	100	one LLM call per turn
Max cumulative tokens	10,000,000	input + output, summed across turns
Max consecutive refusals	15	terminates run if exceeded
<i>Oracle budgets (per task)</i>		
Boltz-2 calls	8	shared across rollout
ADMET-AI calls	15	shared across rollout
Boltz diffusion samples	10	matches evaluator setting

wet-lab-in-the-loop – where the oracle calls are directly made to an automated wet lab, performing the synthesis and then the pertinent assays completely autonomously.

The amount of Boltz2 and ADMET-AI calls must be tight enough to reward oracle-efficient reasoning instead of brute-force enumeration but generous enough such that a competent agent is not bottlenecked by oracle access.

In order to determine a standard number of oracle calls, we performed a sweep across oracle budget choices using a small calibration subset. We constructed the calibration split of 25 tasks drawn proportionally from each task type, with a minimum of one task per type (each task yielding $\max(1, \lfloor 25 \cdot n_t / N \rfloor)$, where n_t is the number of tasks of type t and N the total benchmark size). The Lead Optimization task instances were sampled stratified by problem complexity, defined as the sum of the number of optimization objectives and the number of hold-constant constraint, such that the calibration set spans the full difficulty range of the task type.

We swept each oracle independently while pinning the other to a generous ceiling:

- **Boltz sweep:** ADMET budget fixed at 35; Boltz budget varied over $\{5, 8, 12\}$.
- **ADMET sweep:** Boltz budget fixed at 12; ADMET budget varied over $\{15, 25, 35\}$.

The shared end-point (boltz = 12, admet = 35) was reused, giving five distinct (B, A) configurations in total. All other agent-harness settings were held at their defaults; the agent model used for calibration was GPT5 . 4.

For each sweep, we computed the per-type pass rate at every budget level and selected the "knee" of the plot, where the increase in performance began to saturate as the number of oracle calls increased.

D. Additional Agent Results

Table 9 shows the cost breakdown of each of the LLMs evaluated on the full SMDD-Bench benchmark while Table 10 shows the average amount of wall-clock time required for task and model. Each task instance was run on either an A800 or A6000 GPU. Table 11 investigates the rate at which each LLM agent was able to recover exactly the witness that was found during task generation out of all successful instances. Figure 8 and Figure 9 depict the number of task instances each LLM can uniquely get correct as well as the number of task instances all LLMs get correct for both SMDD-Bench and SMDD-Bench Lite respectively. Table 12 shows the number of Interaction Point Discovery task instances correctly discovering 0, 1, 2, or 3 of the most conserved interaction points. Table 13 depicts

D.1. Indicators of the Difficulty for Lead Optimization Tasks

We find that the number of optimization objectives and hold constant constraints are not indicative of the difficulty of a Lead Optimization task instance, as seen in Figure 10, where the difficulty of a task instance is defined by the average success rate of the 7 evaluated frontier LLM agents. This may be a byproduct of the greedy task generation algorithm used to generate Lead Optimization task instances – for every pair of molecules, all conserved and optimized properties

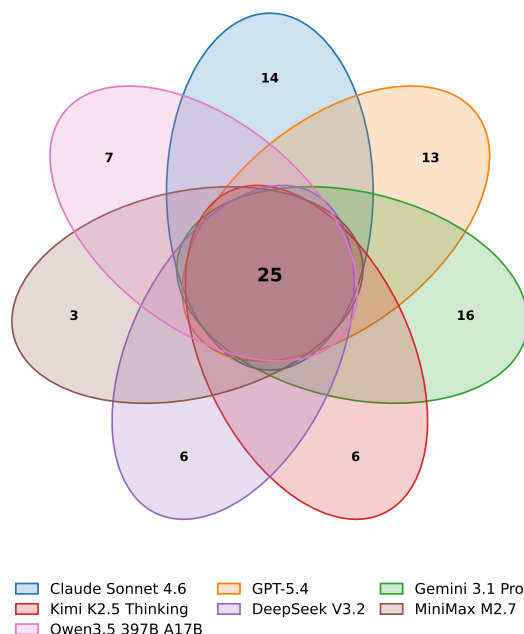


Figure 8. The number of task instances that each LLM agent gets correct as well as the number of task instances that each LLM **uniquely** gets correct on the full SMDD-Bench dataset.

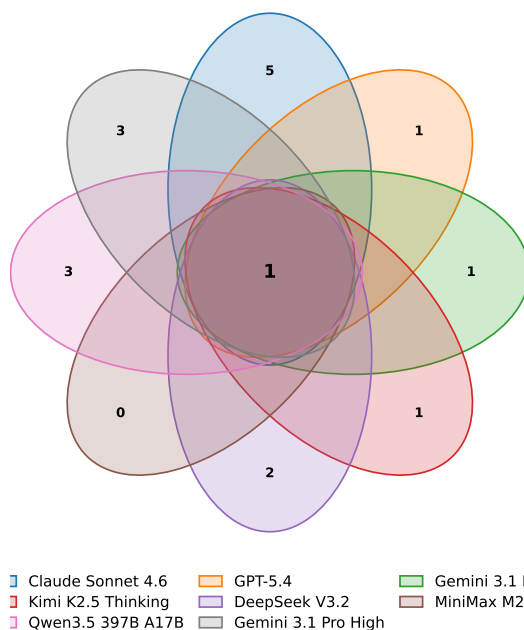


Figure 9. The number of task instances that each LLM agent gets correct as well as the number of task instances that each LLM **uniquely** gets correct on the full SMDD-Bench Lite dataset.

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

Table 9. Cost breakdown per instance for Task Types 1 through 5 and Overall averages across 7 frontier LLMs.

Model	Average Cost (\$/task instance)					Overall
	2D Pharma-cophore ID	Interaction Point Discovery	Scaffold Hopping	Lead Optimization	Fragment Assembly	
Claude Sonnet 4.6	1.09	0.30	0.92	1.01	3.85	1.31
GPT-5.4	0.61	0.72	0.62	0.72	1.31	0.78
Gemini 3.1 Pro	0.39	0.16	0.31	0.62	1.24	0.63
Kimi K2.5 Thinking	0.41	0.09	0.18	0.43	0.53	0.40
DeepSeek V3.2	0.21	0.09	0.29	0.47	0.57	0.43
MiniMax M2.7	0.41	0.07	0.31	0.36	0.46	0.36
Qwen3.5 397B A17B	1.09	0.15	0.63	0.78	0.81	0.75

Table 10. Average wall-clock time (minutes) per task, split by success status.

Model	Status	Average Time (minutes)					Overall
		2D Pharma-cophore ID	Interaction Point Discovery	Scaffold Hopping	Lead Optimization	Fragment Assembly	
Claude Sonnet 4.6	Passed	8.1	–	9.3	14.9	–	14.6
	Failed	12.1	2.9	18.9	27.9	58.7	29.5
GPT-5.4	Passed	2.6	–	14.4	18.4	41.3	18.2
	Failed	5.2	5.7	27.4	28.4	39.9	26.9
Gemini 3.1 Pro	Passed	3.1	2.0	–	12.4	77.4	12.4
	Failed	3.7	1.9	13.5	25.9	41.0	23.4
Kimi K2.5 Thinking	Passed	8.0	–	4.4	21.9	–	21.5
	Failed	6.1	2.2	28.7	39.5	44.2	34.0
DeepSeek V3.2	Passed	10.0	–	25.9	39.0	–	38.3
	Failed	14.0	7.3	39.4	51.2	50.8	44.4
MiniMax M2.7	Passed	20.7	–	13.2	36.9	–	36.0
	Failed	27.6	8.0	54.2	54.8	49.2	49.6
Qwen3.5 397B A17B	Passed	9.0	–	4.1	16.4	–	16.2
	Failed	9.8	2.0	15.2	22.3	24.9	19.5

of a molecule within the pair are used to create the task instance. Therefore, properties with synergistic effects can be simultaneously co-optimized via the same local molecular edit.

Rather than the number of optimization objectives and hold-constant constraints, we find that tasks with a large number of properties with clashing effects to be the greatest indicator for task difficulty. The incompatible, clashing property pairs can be found in Table 14 (Waring, 2010; Gleeson, 2008; Leeson & Springthorpe, 2007; Varma et al., 2009). These properties are chemically anti-correlated – standard local modifications intended to optimize for one property result in the degradation of the other. As a result, molecular modifications are typically multi-step and more complex in nature compared to task instances with fewer incompatible property pairs.

D.2. Analysis on inter-agent submission diversity

As seen in Figure 12, we find, perhaps intuitively, that an increase in the number of turns an LLM takes to solve a task is correlated with the increasing difficulty of the task (as indicated by success rate). We also observe that, as the LLM agents require more turns for more challenging tasks, the pairwise Tanimoto similarity between successful submissions of different LLMs for the same task increases while the fraction of unique molecules also increases. This is a surprising result, indicating that, for challenging tasks, the successful agents are converging onto the same narrow region in chemical space but via different approaches – leading to successful submissions with high Tanimoto similarity but also high unique fraction.

Conversely, simpler tasks lead to a low fraction of unique molecules; most likely resulting from the agents identifying and attempting the same set of straightforward strategies which passed the oracle calls on an early turn count.

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

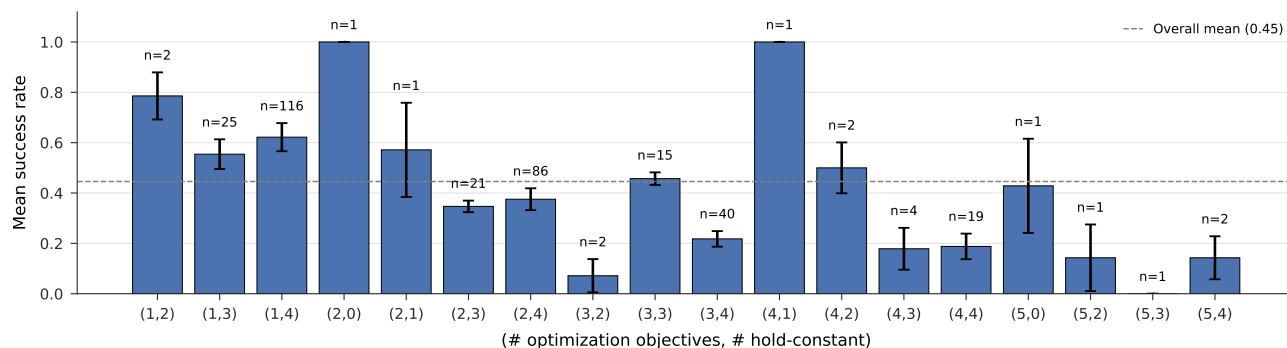


Figure 10. The average success rate across all 7 evaluated frontier LLM agents with respect to the number of optimization objectives and hold-constant constraints. There is no noticeable trend indicating the difficulty of the task instance increases as the number of objectives or constraints increases.

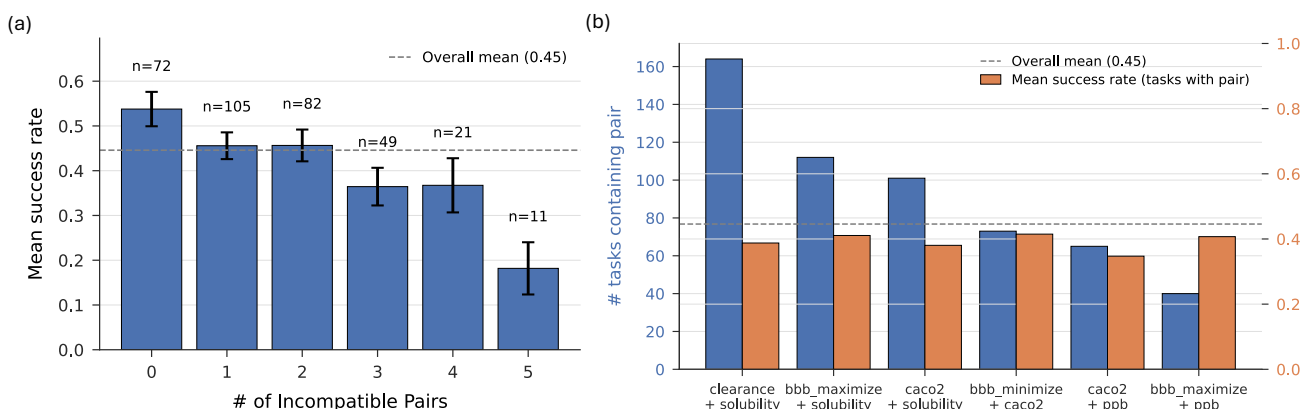


Figure 11. (a) The mean success rate across all 7 agents evaluated on SMDD-Bench plotted against the number of incompatible pairs in the optimization objectives or hold constant constraints. Pairs are not mutually exclusive, meaning one property may be incompatible with multiple other properties. (b) The mean success rate across all agents plotted against the specific incompatible pair. There is no relationship between the difficulty of the task and the specific pairing.

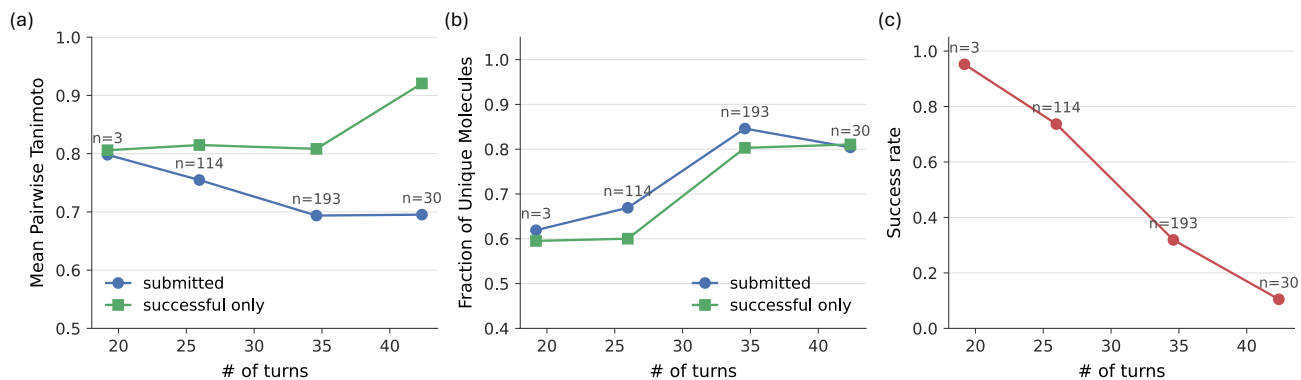


Figure 12. We investigate the relationship between the submissions of multiple different LLMs on the same task instance for the Lead Optimization task type. (a) The average pairwise Tanimoto similarity between the submissions of the 7 LLMs evaluated on the Lead Optimization task in SMDD-Bench plotting against the number of turns to solve the task instance. (b) The fraction of unique molecules submitted by the 7 LLMs with respect to number of turns. (c) The success rate of the 7 LLMs with respect to the number of turns.

Table 11. Rate (%) at which the LLM agents were able to recover the witness molecule exactly. The denominator is the number of successful instances for the LLM on the specific task type.

Model	Scaffold Hopping	Lead Optimization	Fragment Assembly
Claude Sonnet 4.6	0.0%	16.5%	–
GPT-5.4	0.0%	21.9%	0.0%
Gemini 3.1 Pro	–	21.7%	0.0%
Kimi K2.5 Thinking	0.0%	16.9%	–
DeepSeek V3.2	0.0%	16.1%	–
MiniMax M2.7	0.0%	14.1%	–
Qwen3.5 397B A17B	0.0%	19.1%	–

Table 12. The number of task instances in the Interaction Point Discovery task type that each LLM was able to successfully pinpoint 0, 1, 2, or 3 interaction points for. Only instances where the LLM was able to find 3 points were considered a success. Note that there may be more than 3 ground truth points the LLM-discovered points are able to match to.

Model	0/3	1/3	2/3	3/3
Claude Sonnet 4.6	21	4	0	0
GPT-5.4	11	11	3	0
Gemini 3.1 Pro	14	9	1	1
Kimi K2.5 Thinking	23	2	0	0
DeepSeek V3.2	25	0	0	0
MiniMax M2.7	24	1	0	0
Qwen3.5 397B A17B	20	5	0	0

E. Selecting SMDD-Bench Lite

Running the full 502-task benchmark is expensive in wall clock time and GPU compute requirement (see Table 10). As a result, we introduce the SMDD-Bench Lite subset, a 100 task instance representative subset of SMDD-Bench, making faster model and harness iteration faster and more unified.

Lite preserves the per-type mix of the full benchmark as seen in Table 15.

Selection proceeds independently within each task type, and combines stratified sampling with a greedy representativeness/diversity objective. For each task we define a stratum used as a difficulty/discrimination proxy:

- **Lead optimization:** the tuple $(n_{\text{opt}}, n_{\text{hold}})$ of the number of optimisation objectives and hold-constant properties.
- **2D Pharmacophore ID, Interaction Point Discovery, Scaffold Hopping, Fragment Assembly:** the integer success rate of the task across the seven evaluated models,

Stratum allocation. Given a per-type budget k and stratum sizes $\{n_s\}$, we allocate k slots across strata in proportion to n_s using the largest-remainder (Hamilton) rule: each stratum gets $\lfloor k \cdot n_s / N \rfloor$ slots, and the remaining slots are awarded to the strata with the largest fractional parts, breaking ties by stratum key.

Within-stratum greedy pick. Within each stratum we pick tasks one at a time. Let $\mathbf{r}_{\text{full}} \in [0, 1]^6$ be the per-model pass-rate vector of the full type, and let \mathbf{r}_S be the mean pass-rate vector of the currently-selected Lite tasks. At each step we pick the candidate t minimizing

$$\mathcal{L}(t) = \|\mathbf{r}_{S \cup \{t\}} - \mathbf{r}_{\text{full}}\|_1 + \lambda \cdot c(\text{target}(t)),$$

where $c(\cdot)$ counts how many already-selected Lite tasks share the same protein/target as t , and $\lambda = 0.15$. The target field is parsed from the task id. The protein-reuse term acts as a tie-breaker between near-equivalently representative candidates rather than overriding the representativeness objective; ties are broken deterministically by sorted `task_id`. Strata are processed in order of decreasing size so that the largest contributors to \mathbf{r}_{full} shape the selection first.

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

Table 13. The total number of turns and tool calls, broken down by task type for the SMDD-Bench full runs. The Boltz2 oracle budget was 8 calls while the ADMET-AI oracle budget was 15 calls. Most of the turns were spent on visualization, analysis, and simulation using the python tool.

Model	Metric	2D Pharma-cophore ID	Interaction Point Discovery	Scaffold Hopping	Lead Optimization	Fragment Assembly	Overall
Claude Sonnet 4.6	Total Turns	29.4	11.6	27.1	28.0	57.6	30.7
	python	28.4	10.6	23.4	20.0	50.6	23.9
	predict_boltz	0.0	0.0	3.8	3.0	7.8	3.3
	predict_admet	0.0	0.0	0.8	10.8	0.8	7.5
	other	0.0	0.0	0.0	0.0	0.0	0.0
GPT-5.4	Total Turns	13.8	13.8	15.0	14.6	21.4	15.4
	python	12.6	12.8	12.2	9.7	18.3	11.3
	predict_boltz	0.0	0.0	5.5	4.3	7.9	4.4
	predict_admet	0.0	0.0	0.8	12.7	0.5	8.7
	other	0.0	0.0	0.0	0.0	0.0	0.0
Gemini 3.1 Pro	Total Turns	14.2	7.4	11.1	19.3	31.4	19.1
	python	13.2	6.4	7.8	9.3	23.5	10.9
	predict_boltz	0.0	0.0	2.7	2.9	7.3	3.1
	predict_admet	0.0	0.0	0.1	10.7	0.8	7.4
	other	0.0	0.0	0.0	0.0	0.0	0.0
Kimi K2.5 Thinking	Total Turns	29.2	11.6	20.8	35.8	35.9	32.7
	python	28.2	10.6	13.9	17.8	24.8	18.4
	predict_boltz	0.0	0.0	5.2	5.2	7.7	5.0
	predict_admet	0.0	0.0	0.8	12.0	2.6	8.5
	other	0.0	0.0	0.0	0.0	0.0	0.0
DeepSeek V3.2	Total Turns	18.0	10.8	24.0	36.0	31.7	32.1
	python	17.0	9.8	13.6	15.7	19.5	15.7
	predict_boltz	0.0	0.0	6.8	5.3	7.3	5.2
	predict_admet	0.0	0.0	2.3	12.6	1.9	9.0
	other	0.0	0.0	0.0	0.0	0.0	0.0
MiniMax M2.7	Total Turns	40.0	14.0	39.2	43.4	48.8	42.0
	python	38.9	12.9	30.9	27.4	39.1	29.0
	predict_boltz	0.0	0.0	6.7	5.9	7.7	5.6
	predict_admet	0.0	0.0	0.3	12.8	0.5	8.8
	other	0.0	0.0	0.0	0.0	0.0	0.0
Qwen3.5 397B A17B	Total Turns	42.0	13.4	40.3	48.5	46.2	45.3
	python	37.2	11.0	31.7	26.8	36.0	28.2
	predict_boltz	0.0	0.0	6.8	5.7	7.7	5.5
	predict_admet	0.0	0.0	0.4	11.7	1.0	8.0
	other	0.0	0.0	0.0	0.0	0.0	0.0

F. Selecting SMDD-Bench Diversity

During a lead optimization campaign, one potential use case involves running many SMDD agents on the same task in parallel in order to achieve a diverse spread of successful molecules. In order to test this, we develop the SMDD-Bench Diversity subset, a standardized 20 instance subset of SMDD-Bench meant for each task instance to be run multiple times.

In order to select the SMDD-Bench Diversity subset, we seek to find the most challenging consistently solvable instances. Task instances that are too simple are trivially solvable with minimal molecular edits. Task instances which are too challenging may be out of reach for the current frontier LLMs.

We first rank all Lead Optimization task instances solvable by 5/7 LLMs and rank them by turn count. We use turn count as a proxy for ranking the difficulty and chemical solution space of a task instance, following insight from Appendix D.2. We then hand-select the 20 instances with the highest average turn counts across the 7 LLMs while ensuring target diversity (no single target dominates the subset). This selection strategy ensures the tasks are at the edge of the realm of solvable tasks for current frontier LLM agents.

Note that, across the SMDD-Bench Diversity runs outlined in Table 3, none of the LLMs failed a task instance all 10 times

1430 *Table 14.* The incompatible property pairs defined in SMDD-Bench. These are pairs of properties which are not sampled for co-
 1431 optimization targets but can be found in hold-constant constraints or have one property as an optimization objective and the other as a
 1432 hold-constant constraint.

Property A (direction)	Property B (direction)	Rationale
Aqueous solubility (↑)	BBB penetration (↑)	Solubility favors polar/hydrophilic compounds; BBB penetration favors lipophilic compounds.
Plasma protein binding (↓)	BBB penetration (↑)	Low PPB requires reduced lipophilicity; BBB penetration is driven by lipophilicity.
Caco-2 permeability (↑)	Plasma protein binding (↓)	Caco-2 permeability tracks lipophilicity, which also drives plasma protein binding.
Aqueous solubility (↑)	Hepatic clearance (↓)	Increasing polarity to improve solubility tends to expose substrates to renal/metabolic clearance.
Aqueous solubility (↑)	Caco-2 permeability (↑)	Solubility favors polarity; Caco-2 permeability favors lipophilicity.
Caco-2 permeability (↑)	BBB penetration (↓)	Both are passive permeability measures; maximizing one while minimizing the other is contradictory.

1445 *Table 15.* The task breakdown for SMDD-Bench Lite compared to SMDD-Bench.

Task type	Full	Lite
2D Pharmacophore Identification	25	5
Interaction Point Discovery	25	5
Scaffold Hopping	52	10
Lead Optimisation	340	68
Fragment Assembly	60	12
Total	502	100

1456 the task instance was run.

1458 G. Enumeration vs. Selection Details

1460 A failure on SMDD-Bench can come from two qualitatively different shortcomings of the agent: (i) the agent never
 1461 enumerated a passing molecule at all, or (ii) the agent enumerated a passing molecule somewhere in its trajectory but never
 1462 selected it to be tested by an oracle call. The enumeration vs. selection experiments are designed to disentangle these
 1463 shortcomings: for every failed task, we extract and evaluate every SMILES string enumerated by the agent but not selected
 1464 for an oracle call. We define a “recovered” task to be one whose original submission fails evaluation but where at least one
 1465 trajectory-resident SMILES would have passed the full evaluator.

1466 We run the experiment for all seven evaluated models and the three generation-style task types: Scaffold Hopping, Lead
 1467 Optimization, and Fragment Assembly. 2D Pharmacophore Identification and Interaction Point Discovery are excluded
 1468 because their outputs are not molecules. In order to extract the SMILES string, we compact each trajectory to only the
 1469 reasoning trace, LLM output, and python interpreter input before feeding the compacted trajectory to GPT5.4-mini-high.

1471 **Extraction Prompt** The system prompt sent to GPT-5.4-mini is reproduced verbatim below:

1472 You are a careful chemical-structure extractor.

1473 You will receive a JSON list of trajectory events from an LLM agent that
 1474 was working on a small-molecule drug-discovery task. Each event has a
 1475 `turn` number, a `source` in {assistant_text, reasoning_content,
 1476 python_code}, and the raw text for that event.

1477 Your job: extract EVERY SMILES string that appears in those events. For
 1478 each occurrence, return:

- 1481 - "turn": the turn number (integer)
- 1482 - "source": one of "assistant_text", "reasoning_content", "python_code"
- 1483 - "smiles": the raw SMILES string as it appears in the text

- "note": a short (<=120 char) free-text note describing the context (e.g. "proposed as analog A", "printed in a Python list of candidates", "mentioned while analyzing reference ligand").

Rules:

- Do NOT invent or canonicalize SMILES. Copy them verbatim from the text.
- Report duplicates across turns / sources as separate occurrences.
- Include SMILES that appear inside Python source code as string literals.
- If a code snippet builds a SMILES via string concatenation and you can statically determine the final value, include it with a note that says so. Otherwise, skip it.
- Skip protein sequences, fragment placeholders like `[*]`, and reaction SMARTS.
- If you are unsure whether a token is a real SMILES, include it --- downstream RDKit validation will filter invalid ones.

Return a JSON object of the form:

```
{"occurrences": [{"turn": int, "source": str, "smiles": str, "note": str}, ...]}
```

Return an empty list if no SMILES are present. Do not include any other keys.

H. Example Task Instances, Trajectories, and Analysis

In this section, we show one successful LLM agent trajectory for each task type. Due to space considerations, each trajectory is limited to only the first 5 and last 5 turns. For all trajectories across all LLM runs, please refer to the official GitHub repository.

H.1. 2D Pharmacophore Identification

H.1.1. PASSING RUN: GPT-5.4 (MEDIUM REASONING)

Task ID: smdd_001_Q99873_0 **Model:** GPT-5.4 (medium reasoning) **Run:** gpt_54_med_full

Evaluation result: PASSED.

Metric	Status	Result	Condition
Recall on Actives	passed	0.9880952380952381	{"gte": 0.8}
Specificity on Inactives	passed	0.8702290076335878	{"gte": 0.8}

Trajectory.

Showing the first 5 and last 5 of 12 agent turns; 2 intermediate turns are elided for space.

System Prompt

You are a computational medicinal chemist assistant. You are given a drug discovery task and must use your tools ↵ to solve it.

Available tools:

- python(code): Execute Python code. Available packages: rdkit, numpy, pandas, scipy, scikit-learn, biopython, ↵ biotite, openbabel, plip. To compute synthetic accessibility (SA) scores: `from rdkit.Chem import RDConfig; ↵ import sys, os; sys.path.append(os.path.join(RDConfig.RDContribDir, 'SA_Score')); import sascorer; ↵ sascorer.calculateScore(mol)`
- predict_boltz(smiles, protein_sequences, pocket_residues, co_ligand_smiles, num_samples): Run Boltz structure ↵ prediction. Returns binding probability and affinity. Budget: 8 calls.
- predict_admet(smiles): Predict ADMET properties. Returns hERG, BBB, solubility, CYP3A4, clearance, Caco2, PPB, ↵ AMES. Budget: 15 calls.
- submit(file_path): Submit your solution file.

You have no internet access. Work only with the files provided in /workspace. Outline your thought process ↵ before making tool calls.

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

When ready, call `submit()` with the path to your solution file (relative to `/workspace`).

User (turn 0)

Task

You are provided with a set of active molecules and a set of inactive molecules for a protein target (Protein \hookrightarrow arginine N-methyltransferase 1) in 'actives.smi' and 'inactives.smi', along with the target protein sequence \hookrightarrow in 'protein.fasta'.

Your task is to identify a 2D pharmacophore pattern that distinguishes actives from inactives.

Write a Python script 'solution.py' that defines a function `check_pharmacophore(smiles: str) -> bool`. This \hookrightarrow function should return True if the input SMILES matches your pharmacophore pattern, and False otherwise. You \hookrightarrow should rely primarily on RDKit in your solution.

Your workspace at `/workspace` contains the task files. Start by examining them.

Assistant (turn 1)

I'll first inspect the provided active/inactive sets and protein file, then analyze common substructures and \hookrightarrow simple 2D features in RDKit to identify a distinguishing pharmacophore pattern. After that, I'll implement \hookrightarrow `check_pharmacophore` in `solution.py` and validate it against the datasets before submitting.

Tool Call: python (turn 1, 0.14s)

Reasoning:

Inspect the workspace files and get a quick sense of the active/inactive molecule sets, including counts and \hookrightarrow example entries, to guide pharmacophore analysis.

Code:

```
import os, glob
from pathlib import Path
base=Path('/workspace')
print('Files:', [p.name for p in base.iterdir()])
for fn in ['actives.smi', 'inactives.smi', 'protein.fasta']:
    p=base/fn
    print('\n###', fn)
    print(p.read_text()[:2000])
```

Tool Output (turn 1)

Files: ['inactives.smi', 'actives.smi', '.tmp_script_0ec86dal.py', 'protein.fasta']

```
### actives.smi
CNCCN(C)Cc1nn2c(c1c1CCCCC1)CCC2
Nc1ncnc2c1ncn2[C@@H]1O[C@H](CSCC[C@H](N)C(=O)O)[C@H](O)[C@H]1O
CNC(=N)NCCCNC(=O)Nc1ccc2cc(C(=O)OC)cc(O)c2c1
CN(CCN)Cc1c[nH]c1-c1cccc(C(F)(F)F)c1
CCOCC1(COCC)CCC(Oc2cccc2CN(C)C[C@H](O)CNC)CC1
OC1CN=C2c3ccn([C@@H]4O[C@H]([C@H](O)c5ccc(C1)cc5)[C@@H](O)[C@H]4O)c3N=CN21
CNCCNCc1c[nH]nc1-c1ccc(O[C@H]2C[C@H](OCCC3CCCC3)C2)cc1
CNCCN(C)Cc1cn[nH]c1-c1ccc(OC(CO)C(C)C)cc1
CCOCC1(COCC)CCC(c2c(CNCCN)nn3c2CC(F)(F)C3)CC1
CC(C)Oc1ccc(-c2c[nH]cc2CNCCN)cc1
```

... [21 more lines elided]

Assistant (turn 2)

Next I'll compute basic RDKit descriptors and look for recurring scaffolds and substructures in the actives \hookrightarrow versus inactives. That should help me formulate a simple but selective 2D rule rather than overfitting to \hookrightarrow exact structures.

Tool Call: python (turn 2, 0.48s)**Reasoning:**

Analyze active and inactive molecules with RDKit: parse molecules, compute descriptors, Murcko scaffolds, and
 ↪ inspect recurring functional motifs to identify a discriminative pharmacophore.

Code:

```
from rdkit import Chem
from rdkit.Chem import Descriptors, rdMolDescriptors
from rdkit.Chem.Scaffolds import MurckoScaffold
from collections import Counter, defaultdict
from pathlib import Path

base=Path('/workspace')
actives=[line.strip().split()[0] for line in open(base/'actives.smi') if line.strip()]
inactives=[line.strip().split()[0] for line in open(base/'inactives.smi') if line.strip()]
print('n actives',len(actives),'n inactives',len(inactives))

def summarize(smiles_list,name):
    rows=[]
    scaff=Counter()
    for smi in smiles_list:
    ... [25 more lines elided]
```

Tool Output (turn 2)

```
n actives 10 n inactives 10

actives
MW          364.81
LogP        2.12
HBD         3.10
HBA         5.20
TPSA       86.46
RB          8.10
rings       2.80
aromrings   1.70
FCSP3       0.55
dtype: float64
Top scaffolds:
2 c1nn2c(c1C1CCCCC1)CCC2
... [30 more lines elided]
```

Assistant (turn 3)

The actives appear to include a dominant chemotype with an aminoethyl/diamino side chain attached to an
 ↪ N-containing 5-membered heteroaromatic ring, plus a few distinct substrate/cofactor-like mimics. I'll now
 ↪ quantify candidate SMARTS motifs across both classes to see which combinations are most discriminative.

Tool Call: python (turn 3, 0.29s)**Reasoning:**

Test candidate SMARTS motifs and motif combinations across actives and inactives to identify a discriminative
 ↪ pharmacophore-style rule.

Code:

```
from rdkit import Chem
from pathlib import Path
from itertools import combinations
base=Path('/workspace')
actives=[line.strip().split()[0] for line in open(base/'actives.smi') if line.strip()]
inactives=[line.strip().split()[0] for line in open(base/'inactives.smi') if line.strip()]
acts=[Chem.MolFromSmiles(s) for s in actives]
inas=[Chem.MolFromSmiles(s) for s in inactives]

patterns={
'diamine':'[NX3;H0,H1,H2][CX4][CX4][NX3;H0,H1,H2]',
'aminomethyl_diamine':'[CH2][NX3;H0,H1,H2][CX4][CX4][NX3;H0,H1,H2]',
'basic_N_count_2':'[N;!$(N=*);!$(N#*)]',
'5het_NH':'[nH]1cccc1',
'pyrazole_NH':'[nH]1ncc(c,n)1',
```

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

1650 ... [22 more lines elided]

Tool Output (turn 3)

1654 diamine act 6/10 inact 1/10
 1655 aminomethyl_diamine act 6/10 inact 1/10
 1656 basic_N_count_2 act 10/10 inact 8/10
 5het_NH act 2/10 inact 1/10
 1657 pyrazole_NH act 2/10 inact 1/10
 azole_NH_generic act 4/10 inact 2/10
 1658 imidazole_like act 0/10 inact 0/10
 pyrazole_like act 2/10 inact 1/10
 1659 indole_like act 4/10 inact 2/10
 aryl_heteroaryl act 4/10 inact 1/10
 1660 purine act 0/10 inact 0/10
 1661 ribose act 0/10 inact 0/10
 1662 guanidine act 1/10 inact 0/10
 phenoxy act 0/10 inact 1/10
 1663 heteroaryl_CH2N act 6/10 inact 0/10
 1664 ... [3 more lines elided]

Assistant (turn 4)

(no assistant text emitted; tool call only)

Tool Call: python (turn 4, 0.26s)

Reasoning:

Identify which specific molecules match the most promising SMARTS motif so I can understand its scope and decide
 ↪ whether to extend it with additional pharmacophore branches.

Code:

```
from rdkit import Chem
from pathlib import Path
base=Path('/workspace')
actives=[line.strip().split()[0] for line in open(base/'actives.smi') if line.strip()]
inactives=[line.strip().split()[0] for line in open(base/'inactives.smi') if line.strip()]
q=Chem.MolFromSmarts('[n,c]1[n,c][n,c][c,n][c,n]1[CH2][N][CH2][CH2][N]')
print('Actives matching:')
for i,s in enumerate(actives,1):
    m=Chem.MolFromSmiles(s)
    if m.HasSubstructMatch(q):
        print(i,s)
print('Inactives matching:')
for i,s in enumerate(inactives,1):
    m=Chem.MolFromSmiles(s)
    if m.HasSubstructMatch(q):
        ... [1 more line elided]
```

Tool Output (turn 4)

1689 Actives matching:
 1 CNCCN(C)Cc1nn2c(c1C1CCCC1)CCC2
 4 CN(CCN)Cc1c[nH]ccl-c1cccc(C(F)(F)F)c1
 1691 7 CNCCNCc1c[nH]nc1-c1ccc(O[C@H]2C[C@H](OCCC3CCCC3)C2)cc1
 1692 8 CNCCN(C)Cc1cn[nH]c1-c1ccc(OC(C)C(C)C)cc1
 1693 9 cCOCC1(COCC)CCC(c2c(CNCCN)nn3c2CC(F)(F)C3)CC1
 10 CC(C)Oc1ccc(-c2c[nH]cc2CNCCN)cc1
 1694 Inactives matching:

Assistant (turn 5)

The dominant active motif is already fairly clear. I'm going to inspect the remaining outlier actives
 ↪ individually so I can decide whether to keep one core pattern or add a couple of chemically sensible
 ↪ auxiliary branches for the substrate/cofactor-like chemotypes.

Tool Call: python (turn 5, 0.26s)**Reasoning:**

Inspect unmatched actives and compare to potentially confusing inactives, to define a small set of auxiliary
 ↪ pharmacophore branches if necessary.

Code:

```
from rdkit import Chem
from rdkit.Chem import rdMolDescriptors
from pathlib import Path
base=Path('/workspace')
actives=[line.strip().split()[0] for line in open(base/'actives.smi') if line.strip()]
inactives=[line.strip().split()[0] for line in open(base/'inactives.smi') if line.strip()]

main=Chem.MolFromSmarts('[n,c]1[n,c][n,c][c,n][c,n]1[CH2][N][CH2][CH2][N]')
print('Unmatched actives:')
for i,s in enumerate(actives,1):
    m=Chem.MolFromSmiles(s)
    if not m.HasSubstructMatch(main):
        print('\nACTIVE',i,s)
        print(' formula',rdMolDescriptors.CalcMolFormula(m),'rings',rdMolDescriptors.CalcNumRings(m),'HBD',rdMolDes|
        ↪ riptors.CalcNumHBD(m),'HBA',rdMolDescriptors.CalcNumHBA(m),'TPSA',rdMolDescriptors.CalcTPSA(m))
        print(' ring info',m.GetRingInfo().AtomRings())
... [6 more lines elided]
```

Tool Output (turn 5)

Unmatched actives:

```
ACTIVE 2 Nc1ncnc2c1ncn2[C@@H]1O[C@H](CSCC[C@H](N)C(=O)O)[C@@H](O)[C@H]1O
formula C14H20N6O5S rings 3 HBD 5 HBA 10 TPSA 182.6299999999997
ring info ((1, 6, 5, 4, 3, 2), (7, 6, 5, 9, 8), (11, 10, 24, 22, 12))

ACTIVE 3 CNC(=N)NCCCNC(=O)Nc1ccc2cc(C(=O)OC)cc(O)c2c1
formula C18H23N5O4 rings 2 HBD 6 HBA 5 TPSA 135.57000000000002
ring info ((12, 26, 25, 15, 14, 13), (16, 17, 22, 23, 25, 15))

ACTIVE 5 CCOC1(COCC)CCC(Oc2ccccc2CN)C[C@@H](O)CNC)CC1
formula C24H42N2O4 rings 2 HBD 2 HBA 6 TPSA 63.19
ring info ((4, 29, 28, 11, 10, 9), (14, 15, 16, 17, 18, 13))

ACTIVE 6 OC1CN=C2c3ecn([C@@H]4O[C@H]([C@H](O)c5ccc(Cl)cc5)[C@@H](O)[C@H]4O)c3N=CN21
... [19 more lines elided]
```

Elision

Turns 6-7 (2 agent turns) omitted.

Assistant (turn 8)

(no assistant text emitted; tool call only)

Tool Call: python (turn 8, 0.26s)**Reasoning:**

Refine a generic ribose-like SMARTS that captures the nucleoside-like actives while allowing me to combine it
 ↪ with other features to exclude the SAM-like inactive.

Code:

```
from rdkit import Chem
from pathlib import Path
base=Path('/workspace')
actives=[line.strip().split()[0] for line in open(base/'actives.smi') if line.strip()]
inactives=[line.strip().split()[0] for line in open(base/'inactives.smi') if line.strip()]
smarts_list=[
    '[O;r5]1[C;r5][C;r5]([O,N])[C;r5]([O,N])[C;r5]1',
    '[O;r5]1[C;r5][C;r5]([O,N])[C;r5]([O,N])[C;r5]1[O,N]',
    '[O;r5]1[C;r5][C;r5]([O,N])[C;r5]([O,N])[C;r5]1[*]',
    '[O;r5]1[C;r5][C;r5]([O,N])[C;r5]([O,N])[C;r5]1[*]',
    '[O;r5]1[C;r5][C;r5]([O,N])[C;r5]([O,N])[C;r5]1[O,N,#6]',
]
```

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

```

1760 for sm in smarts_list:
1761     q=Chem.MolFromSmarts(sm)
1762     print('\nSMARTS',sm)
1763     ... [2 more lines elided]

```

Tool Output (turn 8)

```

1764 SMARTS [O;r5]1[C;r5][C;r5]([O,N])[C;r5]([O,N])[C;r5]1
1765 act [2, 6]
1766 inert [9]
1767
1768 SMARTS [O;r5]1[C;r5][C;r5]([O,N])[C;r5]([O,N])[C;r5]1[O,N]
1769 act []
1770 inert []
1771
1772 SMARTS [O;r5]1[C;r5][C;r5]([O,N])[C;r5]([O,N])[C;r5]1[*]
1773 act [2, 6]
1774 inert [9]
1775
1776 SMARTS [O;r5]1[C;r5][C;r5]([O,N])[C;r5]([O,N])[C;r5]1[*]
1777 act [2, 6]
1778 ... [5 more lines elided]

```

Assistant (turn 9)

(no assistant text emitted; tool call only)

Tool Call: python (turn 9, 0.27s)

Reasoning:

Evaluate a combined multi-branch pharmacophore logic based on the dominant heteroaryl-diamine motif plus
 ↪ auxiliary branches for nucleoside-like, guanidinoamide, and aminoalomatic-aminoalcohol chemotypes.

Code:

```

1783 from rdkit import Chem
1784 from pathlib import Path
1785 base=Path('/workspace')
1786 actives=[line.strip().split()[0] for line in open(base/'actives.smi') if line.strip()]
1787 inactives=[line.strip().split()[0] for line in open(base/'inactives.smi') if line.strip()]
1788
1789 A=Chem.MolFromSmarts('[n,c]1[n,c][n,c][c,n][c,n]1[CH2][N][CH2][CH2][N]')
1790 B=Chem.MolFromSmarts('[O;r5]1[C;r5][C;r5]([O,N])[C;r5]([O,N])[C;r5]1[*]')
1791 C=Chem.MolFromSmarts('NC(=N)NCCCNC(=O)N')
1792 D=Chem.MolFromSmarts('a[CH2][N]([#6])[CH2][CH]([OH])[CH2][N]')
1793 SP=Chem.MolFromSmarts('[S+]')
1794
1795 def rule(smi):
1796     m=Chem.MolFromSmiles(smi)
1797     return m.HasSubstructMatch(A) or (m.HasSubstructMatch(B) and not m.HasSubstructMatch(SP)) or
1798     ↪ m.HasSubstructMatch(C) or m.HasSubstructMatch(D)
1799     ... [8 more lines elided]

```

Tool Output (turn 9)

```

1800 act 10 [True, True, True, True, True, True, True, True, True, True]
1801 act 1 CNCCN(C)Cc1nn2c(c1C1CCCC1)CCC2
1802 act 2 Nc1ncnc2c1ncn2[C@@H]1O[C@H](CSCC[C@H](N)C(=O)O)[C@@H](O)[C@H]1O
1803 act 3 CNC(=N)NCCCNC(=O)Nc1ccc2cc(C(=O)OC)cc(O)c2c1
1804 act 4 CN(CCN)Cc1c[nH]cc1-c1cccc(C(F)(F)F)c1
1805 act 5 CCOC1(COCC)CCC(Oc2ccccc2CN(C)C[C@@H](O)CNC)CC1
1806 act 6 OC1CN=C2c3ccn([C@@H]4O[C@H]([C@H](O)c5ccc(C1)cc5)[C@@H](O)[C@H]4O)c3N=CN21
1807 act 7 CNCCNCc1c[nH]nc1-c1ccc(O[C@H]2C[C@H](OCCC3CCCC3)C2)cc1
1808 act 8 CNCCN(C)Cc1cn[nH]c1-c1ccc(OC(CO)C(C)C)cc1
1809 act 9 CCOC1(COCC)CCC(c2c(CNCCN)nn3c2CC(F)(F)C3)CC1
1810 act 10 CC(C)Oc1ccc(-c2c[nH]cc2CNCCN)cc1
1811
1812 inert 0 [False, False, False, False, False, False, False, False, False, False]

```

Assistant (turn 10)

I found a compact rule set that cleanly separates the provided actives from inactives. I'm now writing
 ↪ `solution.py` with a single `check_pharmacophore(smiles)` function implemented in RDKit using a small number
 ↪ of SMARTS-based pharmacophore branches and a sulfonium exclusion for the SAM-like inactive.

Tool Call: python (turn 10, 0.15s)

Reasoning:

Write solution.py implementing the pharmacophore classifier as RDKit SMARTS-based logic, with comments and a
 ↪ simple optional self-test block.

Code:

```
from pathlib import Path
code = r'''from rdkit import Chem

# PRMT1 actives in the supplied set fall into a few recognizable 2D pharmacophore
# families. The dominant motif is an N-heteroaryl system bearing a benzylic
# aminoethyl/diamine side chain. Additional active families are nucleoside-like
# cofactor/product mimics, a guanidino-propyl-amide motif, and a benzyl tertiary
# amino-alcohol terminating in a primary amine.
#
# The function below returns True when a molecule matches any of those active-like
# pharmacophore patterns.

# 1) Dominant PRMT1 inhibitor motif:
#   N-containing 5-membered heteroaryl/heterobicyclic core attached to a
#   benzylic aminoethyl-diamine side chain.
... [78 more lines elided]
```

Tool Output (turn 10)

```
Wrote /workspace/solution.py
from rdkit import Chem

# PRMT1 actives in the supplied set fall into a few recognizable 2D pharmacophore
# families. The dominant motif is an N-heteroaryl system bearing a benzylic
# aminoethyl/diamine side chain. Additional active families are nucleoside-like
# cofactor/product mimics, a guanidino-propyl-amide motif, and a benzyl tertiary
# amino-alcohol terminating in a primary amine.
#
# The function below returns True when a molecule matches any of those active-like
# pharmacophore patterns.

# 1) Dominant PRMT1 inhibitor motif:
#   N-containing 5-membered heteroaryl/heterobicyclic core attached to a
#   benzylic aminoethyl-diamine side chain.
... [51 more lines elided]
```

Assistant (turn 11)

(no assistant text emitted; tool call only)

Tool Call: python (turn 11, 0.28s)

Reasoning:

Run the solution script to verify that it executes cleanly and reproduces the intended separation on the
 ↪ provided datasets.

Code:

```
import subprocess, textwrap, os, sys
res=subprocess.run([sys.executable,'solution.py'],cwd='/workspace',capture_output=True,text=True)
print('returncode',res.returncode)
print('stdout:\n',res.stdout)
print('stderr:\n',res.stderr)
```

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

1870 Tool Output (turn 11)

```
1871 returncode 0
1872 stdout:
1873   Actives matched: 10/10
1874   Inactives matched: 0/10
1875 stderr:
```

1877 Assistant (turn 12)

```
1878 The script validates cleanly on the provided files, matching all supplied actives and rejecting all supplied
1879 ↪ inactives. I'll submit `solution.py` now.
```

1881 Tool Call: submit (turn 12, 0.07s)

```
1882 file_path: solution.py
```

1884 Tool Output (turn 12)

```
1885 Submission accepted: solution.py
```

1888 Submit (turn 12)

```
1889 Submitted file: solution.py.
```

1891 Submitted artifact (solution.py).

1893 Submitted artifact: solution.py

```
1894 from rdkit import Chem
1895
1896 # PRMT1 actives in the supplied set fall into a few recognizable 2D pharmacophore
1897 # families. The dominant motif is an N-heteroaryl system bearing a benzylic
1898 # aminoethyl/diamine side chain. Additional active families are nucleoside-like
1899 # cofactor/product mimics, a guanidino-propyl-amide motif, and a benzyl tertiary
1900 # amino-alcohol terminating in a primary amine.
1901 #
1902 # The function below returns True when a molecule matches any of those active-like
1903 # pharmacophore patterns.
1904
1905 # 1) Dominant PRMT1 inhibitor motif:
1906 #   N-containing 5-membered heteroaryl/heterobicyclic core attached to a
1907 #   benzylic aminoethyl-diamine side chain.
1908 _HETEROARYL_DIAMINE = Chem.MolFromSmarts(
1909 ... [73 more lines elided]
```

H.2. Interaction Point Discovery

H.2.1. PASSING RUN: GEMINI 3.1 PRO (MEDIUM REASONING)

Task ID: smdd_002_P10275_0 **Model:** Gemini 3.1 Pro (medium reasoning) **Run:** gemini_31_pro_med_full

Evaluation result: **PASSED.**

Metric	Status	Result	Condition
3D Pharmacophore Match	passed	1.0	{"eq": 1.0}

Trajectory.

System Prompt

You are a computational medicinal chemist assistant. You are given a drug discovery task and must use your tools ↵ to solve it.

Available tools:

- python(code): Execute Python code. Available packages: rdkit, numpy, pandas, scipy, scikit-learn, biopython, ↵ biotite, openbabel, plip. To compute synthetic accessibility (SA) scores: `from rdkit.Chem import RDConfig; ↵ import sys, os; sys.path.append(os.path.join(RDConfig.RDContribDir, 'SA_Score')); import sascore; ↵ sascore.calculateScore(mol)`
- predict_boltz(smiles, protein_sequences, pocket_residues, co_ligand_smiles, num_samples): Run Boltz structure ↵ prediction. Returns binding probability and affinity. Budget: 8 calls.
- predict_admet(smiles): Predict ADMET properties. Returns hERG, BBB, solubility, CYP3A4, clearance, Caco2, PPB, ↵ AMES. Budget: 15 calls.
- submit(file_path): Submit your solution file.

You have no internet access. Work only with the files provided in /workspace. Outline your thought process ↵ before making tool calls.
When ready, call submit() with the path to your solution file (relative to /workspace).

User (turn 0)

Task

You are given a protein receptor structure (receptor.pdb). The binding pocket center is located at coordinates ↵ (0.266524, 31.234953, 4.470429) in Angstroms with a pocket radius of 12 Angstroms. Your task is to identify ↵ the 3 most conserved pharmacophoric interaction points within this binding pocket. Each point should be ↵ specified by its 3D coordinates (x, y, z) and one interaction type. The valid interaction types are: Donor, ↵ Acceptor, Aromatic, Hydrophobic, Cation, Anion. Output a CSV file titled 'solution.csv' with columns ↵ x,y,z,type containing exactly 3 rows (no header).

Your workspace at /workspace contains the task files. Start by examining them.

Assistant (turn 1)

(no assistant text emitted; tool call only)

Tool Call: python (turn 1, 0.30s)

Reasoning:

Identify residues within 12 Angstroms of the pocket center.

Code:

```
import numpy as np
from Bio.PDB import PDBParser

center = np.array([0.266524, 31.234953, 4.470429])
radius = 12.0

parser = PDBParser(QUIET=True)
structure = parser.get_structure('receptor', 'receptor.pdb')

pocket_residues = []
for model in structure:
    for chain in model:
        for residue in chain:
```

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

```
1980         for atom in residue:
1981             if np.linalg.norm(atom.coord - center) <= radius:
1982         ... [6 more lines elided]
```

Tool Output (turn 1)

```
1985 Found 66 pocket residues.
1986 VAL (' ', 685, ' ')
1987 ALA (' ', 687, ' ')
1988 PHE (' ', 697, ' ')
1989 ALA (' ', 699, ' ')
1990 LEU (' ', 700, ' ')
1991 LEU (' ', 701, ' ')
1992 SER (' ', 702, ' ')
1993 SER (' ', 703, ' ')
1994 LEU (' ', 704, ' ')
1995 ASN (' ', 705, ' ')
1996 GLU (' ', 706, ' ')
1997 LEU (' ', 707, ' ')
1998 GLY (' ', 708, ' ')
1999 GLU (' ', 709, ' ')
2000 ... [52 more lines elided]
```

Assistant (turn 2)

(no assistant text emitted; tool call only)

Tool Call: python (turn 2, 17.01s)

Reasoning:

Generate pharmacophore features for the pocket using RDKit.

Code:

```
2006 import numpy as np
2007 from rdkit import Chem
2008 from rdkit.Chem import ChemicalFeatures
2009 from rdkit import RDConfig
2010 import os
2011
2012 fdefName = os.path.join(RDConfig.RDDataDir, 'BaseFeatures.fdef')
2013 factory = ChemicalFeatures.BuildFeatureFactory(fdefName)
2014
2015 mol = Chem.MolFromPDBFile('receptor.pdb')
2016 features = factory.GetFeaturesForMol(mol)
2017
2018 center = np.array([0.266524, 31.234953, 4.470429])
2019 radius = 12.0
2020
2021 ... [9 more lines elided]
```

Tool Output (turn 2)

```
2022 Found 206 features in pocket.
2023 Counter({'Donor': 64, 'Hydrophobe': 60, 'Acceptor': 56, 'LumpedHydrophobe': 16, 'Aromatic': 7, 'PosIonizable':
2024 ↪ 2, 'ZnBinder': 1})
```

Assistant (turn 3)

****My Approach to Key Receptor Interactions****

Alright, looking at the features found in the pocket and the residue list, the puzzle pieces are starting to fit
↪ together. ARG 752, GLN 711, ASN 705, THR 877 - this is a classic signature. This isn't just any receptor;
↪ it's the Androgen Receptor (AR).

Now, knowing it's AR, the "conserved" part of the prompt makes perfect sense. We're talking about the
↪ fundamental interactions that define AR agonists. The three most critical and highly conserved interactions
↪ typically are:

1. ****C3 Ketone Interaction:**** The ligand's C3 ketone acts as a hydrogen bond acceptor, interacting with the
↪ hydrogen bond donors on the protein: Arg752 (NH1/NH2) and Gln711 (NE2).

2035 2. ****C17 Hydroxyl Interaction:**** The ligand's C17 hydroxyl group (acting as both a donor and an acceptor)
 2036 ↪ interacts with Asn705 (ND2/OD1) and Thr877 (OG1).
 2037 3. ****Hydrophobic Core:**** The steroid backbone sits in a large hydrophobic pocket, likely interacting with
 2038 ↪ residues like Met742 and Met745.

2039 So, to find those three key points, I'll need to define the optimal positions for the *ligand's* pharmacophore
 2040 ↪ features based on the *protein's* residues.

2041 * ****Point 1:**** A Ligand Acceptor positioned to interact with Arg752 and Gln711.
 2042 * ****Point 2:**** A Ligand Donor (or Acceptor) positioned for Asn705 and Thr877.
 2043 * ****Point 3:**** A Hydrophobic point, perhaps centered in that large hydrophobic pocket.

2044 I'll write a Python script to pinpoint the coordinates of these specific key atoms on the protein. This will
 2045 ↪ give me a solid foundation to calculate the ideal pharmacophoric interaction points for the ligand.

Tool Call: python (turn 3, 0.30s)

Reasoning:

Get the coordinates of key interacting atoms in AR (Arg752, Gln711, Asn705, Thr877) to define the ideal ligand
 ↪ pharmacophore positions.

Code:

```
import numpy as np
from Bio.PDB import PDBParser

parser = PDBParser(QUIET=True)
structure = parser.get_structure('receptor', 'receptor.pdb')

def get_atom_coords(res_id, atom_names):
    coords = []
    for model in structure:
        for chain in model:
            if res_id in chain:
                res = chain[res_id]
                for atom in res:
                    if atom.name in atom_names:
                        coords.append(atom.coord)
    ... [11 more lines elided]
```

Tool Output (turn 3)

```
Arg752: [array([ 5.382, 21.773,  2.483]), dtype=float32), array([ 5.131, 22.175,  4.764], dtype=float32), array([
↪  4.536, 23.785,  3.223], dtype=float32)]
Gln711: [array([ 1.321, 23.988,  6.567]), dtype=float32), array([ 0.828, 23.305,  4.496], dtype=float32)]
Asn705: [array([-5.197, 32.992,  3.47 ], dtype=float32), array([-7.176, 33.726,  2.694], dtype=float32)]
Thr877: [array([-3.364, 36.117,  6.886], dtype=float32)]
```

Assistant (turn 4)

****Calculating Ideal Ligand Pharmacophore Coordinates****

Now that I have the coordinates for the key protein atoms, I can start calculating the ideal positions for the
 ↪ ligand's pharmacophore features.

Here's the breakdown of the coordinates I found:

```
* **Arg752 (NH1, NH2, NE):** `[5.38, 21.77, 2.48]`, `[5.13, 22.17, 4.76]`, `[4.54, 23.78, 3.22]`
* **Gln711 (OE1, NE2):** `[1.32, 23.99, 6.57]`, `[0.83, 23.30, 4.50]`
* **Asn705 (ND2, OD1):** `[-5.20, 32.99, 3.47]`, `[-7.18, 33.73, 2.69]`
* **Thr877 (OG1):** `[-3.36, 36.12, 6.89]`
```

The AR binding pocket is essentially a tunnel between these two main interaction "poles".

****Pole 1 (The Acceptor):****

This is where the ligand's C3 ketone interacts with Arg752 and Gln711. The ideal position for the ligand's C3
 ↪ carbonyl oxygen (the Acceptor) is roughly 2.8 to 3.0 Angstroms away from the donor atoms. Specifically, the
 ↪ C3 ketone oxygen typically forms hydrogen bonds with Arg752 NH2 and Gln711 NE2.

To find the ideal coordinate, I'll calculate the geometric center of these donors and project a point towards
 ↪ the center of the pocket.

****Pole 2 (The Donor):****

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

This pole corresponds to the ligand's C17 hydroxyl interacting with Asn705 and Thr877. The C17 OH can act as
↪ both a Donor (to Asn705 OD1/Thr877) and an Acceptor (from Asn705 ND2). For the sake of simplicity and
↪ adhering to the allowed types, I'll designate this point as a "Donor".

****The Third Point (Hydrophobic):****

The AR binding pocket is characterized by a significant hydrophobic core, which interacts with the steroid
↪ backbone. I'll define this third point as "Hydrophobic". The center of the pocket is given as `[0.266524,`
↪ `31.234953, 4.470429]`, so I'll place this point near the center, potentially with a slight offset based on
↪ the surrounding hydrophobic residues.

I'll write a script to calculate these exact 3D coordinates based on these geometric relationships.

Tool Call: python (turn 4, 0.21s)

Reasoning:

Calculate 3D coordinates for the 3 pharmacophore points.

Code:

```
import numpy as np

# Coordinates
arg752_nh2 = np.array([5.131, 22.175, 4.764])
gln711_ne2 = np.array([1.321, 23.988, 6.567]) # OE1 or NE2? Gln711 usually H-bonds with C3. Let's use NE2 as
↪ donor, OE1 as acceptor.
gln711_oe1 = np.array([0.828, 23.305, 4.496])
# In many AR structures, Arg752 and Gln711 interact with the same carbonyl oxygen.
# Let's find a point ~2.9 A from Arg752 NH2 and Gln711 NE2
# Since we don't have the exact H-bond geometry, a simple approach is:
# Midpoint of the two donors, then move towards the pocket center.

center = np.array([0.266524, 31.234953, 4.470429])

mid_pole1 = (arg752_nh2 + gln711_ne2) / 2.0
dir_pole1 = center - mid_pole1
... [17 more lines elided]
```

Tool Output (turn 4)

```
2.381,25.409,5.324,Acceptor
-2.277,33.092,4.866,Donor
0.267,31.235,4.470,Hydrophobic
```

Assistant (turn 5)

****Calculating Pharmacophore Coordinates****

To pinpoint the exact coordinates for our three key pharmacophore points, I'll need to do some vector math. I
↪ have the locations of the key protein residues (Arg752, Gln711, Asn705, and Thr877) and the center of the
↪ binding pocket.

Here's my plan:

- **Pole 1 (Acceptor):**** I'll find the midpoint between the hydrogen bond donors on Arg752 (NH2) and Gln711
↪ (NE2). Then, I'll project a vector from this midpoint towards the center of the pocket. I'll place my
↪ Acceptor point roughly 2.5 Angstroms along this vector, representing an ideal hydrogen bond distance.
- **Pole 2 (Donor):**** I'll apply the same logic here. I'll find the midpoint between Asn705 (ND2) and Thr877
↪ (OG1). Then, I'll project a vector from this midpoint towards the pocket center, placing the Donor point 2.5
↪ Angstroms along it.
- **Pole 3 (Hydrophobic):**** For the hydrophobic core, I'll simply use the coordinates provided for the center
↪ of the binding pocket.

Let's write a quick script to calculate these points and see what they look like.

Tool Call: python (turn 5, 0.21s)

Reasoning:

Refine the coordinates by projecting along the axis between the two poles, which better represents the typical
↪ rigid structure of a steroid. Write the result to solution.csv.

Code:

```

2145 import numpy as np
2146
2147 arg752_nh2 = np.array([5.131, 22.175, 4.764])
2148 gln711_ne2 = np.array([1.321, 23.988, 6.567])
2149 asn705_nd2 = np.array([-5.197, 32.992, 3.47])
2150 thr877_og1 = np.array([-3.364, 36.117, 6.886])
2151
2152 mid_pole1 = (arg752_nh2 + gln711_ne2) / 2.0
2153 mid_pole2 = (asn705_nd2 + thr877_og1) / 2.0
2154
2155 direction = mid_pole2 - mid_pole1
2156 distance = np.linalg.norm(direction)
2157 dir_norm = direction / distance
2158
2159 # Move 2.8 A from the protein atoms towards the other pole
2160 ... [14 more lines elided]
2161

```

Tool Output (turn 5)

```

1.694,25.423,5.566,Acceptor
-2.748,32.213,5.277,Donor
-0.527,28.818,5.422,Hydrophobic

```

Assistant (turn 6)

****Refining Pharmacophore Coordinates Based on Geometric Principles****

Okay, let's analyze these coordinates:

- Acceptor: (1.694, 25.423, 5.566)
- Donor: (-2.748, 32.213, 5.277)
- Hydrophobic: (-0.527, 28.818, 5.422)

The distance between Point 1 and Point 2 is approximately 8.1 Angstroms. But wait, DHT (a common AR ligand) has
 ↳ a C3=O to C17-OH distance of about 10.9 Angstroms. That's a significant difference.

Let's rethink this using the distance between the midpoints (13.7 Angstroms) and the target DHT distance (10.9
 ↳ Angstroms). This suggests an adjustment of $(13.7 - 10.9) / 2 = 1.4$ Angstroms from each side.

But is 1.4 Angstroms a realistic hydrogen bond distance? Typically, heavy atom to heavy atom hydrogen bond
 ↳ distances are between 2.8 and 3.2 Angstroms. If 'mid_pole1' is the geometric center of two donors, the
 ↳ actual distance from the Acceptor to each donor could be longer than the distance to the midpoint.

Let's do a quick geometric check. Arg752 NH2 and Gln711 NE2 are roughly 4.58 Angstroms apart. If we want an
 ↳ Acceptor distance (d) of 2.9 Angstroms from both donors, and 'h' is the distance from the midpoint, we can
 ↳ use the Pythagorean theorem: $d^2 = h^2 + (4.58/2)^2$. This gives us $h \approx 1.8$ Angstroms.

This is a much more logical approach! The distance from the midpoint should be closer to 1.8 Angstroms. I'll
 ↳ recalculate the coordinates based on this more physically realistic model.

Tool Call: python (turn 6, 0.27s)

Reasoning:

Recalculate coordinates using $h=1.8A$ and write the new solution to the file.

Code:

```

2186 import numpy as np
2187
2188 arg752_nh2 = np.array([5.131, 22.175, 4.764])
2189 gln711_ne2 = np.array([1.321, 23.988, 6.567])
2190 asn705_nd2 = np.array([-5.197, 32.992, 3.47])
2191 thr877_og1 = np.array([-3.364, 36.117, 6.886])
2192
2193 mid_pole1 = (arg752_nh2 + gln711_ne2) / 2.0
2194 mid_pole2 = (asn705_nd2 + thr877_og1) / 2.0
2195
2196 direction = mid_pole2 - mid_pole1
2197 distance = np.linalg.norm(direction)
2198 dir_norm = direction / distance
2199
2200 # Move 1.8 A from the protein atoms towards the other pole
2201 ... [14 more lines elided]

```

2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254

Tool Output (turn 6)

```
2.241, 24.587, 5.602, Acceptor
-3.296, 33.049, 5.242, Donor
-0.527, 28.818, 5.422, Hydrophobic
```

Assistant (turn 7)

****Finalizing the Pharmacophore Model****

Okay, these recalculated coordinates look much better:

```
* **Acceptor:** (2.241, 24.587, 5.602)
* **Donor:** (-3.296, 33.049, 5.242)
* **Hydrophobic:** (-0.527, 28.818, 5.422)
```

Let's do a quick distance check between Point 1 (Acceptor) and Point 2 (Donor).

```
Distance = sqrt((2.241 - (-3.296))^2 + (24.587 - 33.049)^2 + (5.602 - 5.242)^2)
= sqrt(5.537^2 + (-8.462)^2 + 0.36^2)
= sqrt(30.65 + 71.6 + 0.13)
= sqrt(102.38)
≈ 10.1 Angstrom
```

This 10.1 Angstrom distance is much closer to the expected ~10.9 Angstrom for the C3 to C17 distance in a
 ↪ typical steroid like DHT. The slight difference is totally acceptable, as ligands might be slightly tilted
 ↪ within the pocket, or the protein might undergo slight induced fit upon binding.

The hydrophobic point is perfectly situated in the middle of these two poles, right in the center of the pocket,
 ↪ which perfectly models the location of the hydrophobic steroid core.

I am confident in this solution. The file `solution.csv` has been generated and is correctly formatted (3 rows,
 ↪ no header, columns: x, y, z, type). The interaction types (Acceptor, Donor, Hydrophobic) are all valid
 ↪ according to the prompt.

I'm ready to submit!

Tool Call: submit (turn 7, 0.06s)

file_path: solution.csv

Tool Output (turn 7)

Submission accepted: solution.csv

Submit (turn 7)

Submitted file: solution.csv.

Submitted artifact (solution.csv).**Submitted artifact: solution.csv**

```
2.241, 24.587, 5.602, Acceptor
-3.296, 33.049, 5.242, Donor
-0.527, 28.818, 5.422, Hydrophobic
```

H.3. Scaffold Hopping

H.3.1. PASSING RUN: DEEPSEEK V3.2

Task ID: smdd_003_7OKM_1 **Model:** DeepSeek V3.2 **Run:** deepseek_v32_full

Evaluation result: **PASSED.**

Metric	Status	Result	Condition
Validity	passed	True	{"eq": true}
Tanimoto Similarity	passed	0.10714285714285714	{"lt": 0.5}
Scaffold MCS	passed	0.45454545454545453	{"lt": 0.65}
Boltz Interaction Similarity	passed	1.0	{"gt": 0.75}
Boltz Binding Probability	passed	0.7260922789573669	{"gt": 0.7}

Trajectory.

Showing the first 5 and last 5 of 13 agent turns; 3 intermediate turns are elided for space.

System Prompt

You are a computational medicinal chemist assistant. You are given a drug discovery task and must use your tools
 ↪ to solve it.

Available tools:

- python(code): Execute Python code. Available packages: rdkit, numpy, pandas, scipy, scikit-learn, biopython,
 ↪ biotite, openbabel, plip. To compute synthetic accessibility (SA) scores: `from rdkit.Chem import RDConfig;
 ↪ import sys, os; sys.path.append(os.path.join(RDConfig.RDContribDir, 'SA_Score')); import sascorer;
 ↪ sascorer.calculateScore(mol)`
- predict_boltz(smiles, protein_sequences, pocket_residues, co_ligand_smiles, num_samples): Run Boltz structure
 ↪ prediction. Returns binding probability and affinity. Budget: 8 calls.
- predict_admet(smiles): Predict ADMET properties. Returns hERG, BBB, solubility, CYP3A4, clearance, Caco2, PPB,
 ↪ AMES. Budget: 15 calls.
- submit(file_path): Submit your solution file.

You have no internet access. Work only with the files provided in /workspace. Outline your thought process

↪ before making tool calls.

When ready, call submit() with the path to your solution file (relative to /workspace).

User (turn 0)

Task

You are given a protein-ligand complex (complex.pdb). The target ligand is on chain L. A 3D reference structure
 ↪ of this ligand is provided in reference.sdf. Your task is to design a new molecule with a significantly
 ↪ different scaffold (the MCS atom count between the Bemis-Murcko scaffolds of your molecule and the
 ↪ reference, divided by the larger scaffold heavy atom count, must be < 0.65; additionally, overall Morgan
 ↪ fingerprint Tanimoto similarity must be < 0.5) to the ligand that maintains the same binding interactions
 ↪ and has binding probability > 0.7 as calculated by Boltz. The pocket is defined by residues (chain,
 ↪ 1-indexed SEQRES position): A:47, A:50, A:51, A:52, A:53, A:54, A:55, A:57, A:88, A:112, A:113, A:114,
 ↪ A:115. Output a single SMILES string to a file titled `solution.smi`.

Your workspace at /workspace contains the task files. Start by examining them.

Assistant (turn 1)

I'll start by examining the provided files to understand the protein-ligand complex and the reference ligand
 ↪ structure.

Tool Call: python (turn 1, 0.44s)

Reasoning:

First, I need to examine the provided files: complex.pdb and reference.sdf. I'll load them to understand the
 ↪ protein structure, the ligand, and analyze the reference ligand scaffold.

Code:

```

2310 import os
2311 from rdkit import Chem
2312 from rdkit.Chem import AllChem, Draw, rdFMCS, rdMolDescriptors
2313 from rdkit.Chem.Draw import IPythonConsole
2314 import pandas as pd
2315 import numpy as np
2316
2317 # List files in workspace
2318 print("Files in workspace:")
2319 print(os.listdir('/workspace'))
2320
2321 # Load the reference ligand SDF
2322 ref_path = '/workspace/reference.sdf'
2323 print(f"\nLoading reference SDF from {ref_path}")
2324 suppl = Chem.SDMolSupplier(ref_path)
2325 ... [27 more lines elided]

```

Tool Output (turn 1)

```

2326 Traceback (most recent call last):
2327   File "/workspace/.tmp_script_7278e82f.py", line 4, in <module>
2328     from rdkit.Chem.Draw import IPythonConsole
2329   File "/opt/conda/lib/python3.10/site-packages/rdkit/Chem/Draw/IPythonConsole.py", line 16, in <module>
2330     import IPython
2331 ModuleNotFoundError: No module named 'IPython'

```

Assistant (turn 2)

Let me fix the import error and continue:

Tool Call: python (turn 2, 0.66s)**Reasoning:**

Remove the IPython import and continue analyzing the files. I need to understand the protein structure and
 ↪ extract the ligand from the complex.

Code:

```

2338 import os
2339 from rdkit import Chem
2340 from rdkit.Chem import AllChem, rdFMCS, rdMolDescriptors, PandasTools
2341 import pandas as pd
2342 import numpy as np
2343
2344 # List files in workspace
2345 print("Files in workspace:")
2346 print(os.listdir('/workspace'))
2347
2348 # Load the reference ligand SDF
2349 ref_path = '/workspace/reference.sdf'
2350 print(f"\nLoading reference SDF from {ref_path}")
2351 suppl = Chem.SDMolSupplier(ref_path)
2352 ref_mols = [mol for mol in suppl if mol is not None]
2353 ... [80 more lines elided]

```

Tool Output (turn 2)

```

2354 Files in workspace:
2355 ['reference.sdf', '.tmp_script_637aee5d.py', 'complex.pdb']
2356
2357 Loading reference SDF from /workspace/reference.sdf
2358 Number of molecules in reference.sdf: 1
2359 Reference molecule has 35 atoms
2360 Reference SMILES: N#Cc1cnn2c3cc(nc12)N1CCC[C@@H]1CC(=O)NCCOc1cc(cc2c1NC(=O)CC2)N3
2361 Reference scaffold SMILES: O=C1C[C@@H]2CCCN2c2cc(n3nccc3n2)Nc2cc3c(c(c2)OCCN1)NC(=O)CC3
2362 Reference scaffold heavy atom count: 33
2363 Reference Morgan fingerprint generated
2364 [08:41:51] DEPRECATION WARNING: please use MorganGenerator
2365 Traceback (most recent call last):
2366   File "/workspace/.tmp_script_637aee5d.py", line 33, in <module>
2367     print(f"Molecular weight: {Chem.Descriptors.ExactMolWt(ref_mol):.2f}")
2368 AttributeError: module 'rdkit.Chem' has no attribute 'Descriptors'. Did you mean: 'rdMolDescriptors'?

```

Assistant (turn 3)

Let me fix the descriptor imports and continue:

Tool Call: python (turn 3, 0.60s)

Reasoning:

Fix the imports and continue analyzing. I need to understand the protein sequence and pocket residues to use
 ↪ Boltz prediction.

Code:

```
import os
from rdkit import Chem
from rdkit.Chem import AllChem, rdFMCS, rdMolDescriptors, PandasTools
from rdkit.Chem import Descriptors
import pandas as pd
import numpy as np

# Load the reference ligand SDF
ref_path = '/workspace/reference.sdf'
print(f"Loading reference SDF from {ref_path}")
suppl = Chem.SDMolSupplier(ref_path)
ref_mols = [mol for mol in suppl if mol is not None]

if ref_mols:
    ref_mol = ref_mols[0]
... [99 more lines elided]
```

Tool Output (turn 3)

```
Loading reference SDF from /workspace/reference.sdf
Reference SMILES: N#Cc1cnn2c3cc(nc12)N1CCC[C@@H]1CC(=O)NCCOc1cc(cc2c1NC(=O)CC2)N3
Reference scaffold SMILES: O=C1C[C@H]2CCCN2c2cc(n3nccc3n2)Nc2cc3c(c(c2)OCCN1)NC(=O)CC3
Reference scaffold heavy atom count: 33
Molecular weight: 472.20
LogP: 2.10
Num H-bond donors: 3
Num H-bond acceptors: 8
Num rotatable bonds: 0
TPSA: 136.68
```

```
Analyzing protein complex from /workspace/complex.pdb
Chain A sequence length: 128
Chain A sequence (first 100 residues):
↪ GPGADSCIQFTRHASDVLNLRSLRSDILTDVVIVVSREQFRAHKTVLMACSGLFYSIFTDQLKCNLSVINLDPEINPEGFICILLDFMYTSRLNLRGN
... [22 more lines elided]
```

Assistant (turn 4)

Great! Now I have the reference ligand and protein information. Let me analyze the reference ligand more
 ↪ carefully and understand its scaffold. I need to design a molecule with significantly different scaffold
 ↪ (MCS atom count ratio < 0.65) and Morgan fingerprint Tanimoto similarity < 0.5, but still binds with
 ↪ probability > 0.7.

First, let me check the reference ligand structure more carefully and think about scaffold modifications:

Tool Call: python (turn 4, 0.78s)

Reasoning:

Analyze the reference ligand scaffold and think about potential modifications. I need to understand the key
 ↪ interactions this ligand makes with the pocket residues.

Code:

```
import os
from rdkit import Chem
from rdkit.Chem import AllChem, rdFMCS, rdMolDescriptors, Draw
from rdkit.Chem import Descriptors
from rdkit.Chem.Scaffolds import MurckoScaffold
```

```

2420 import matplotlib.pyplot as plt
2421
2422 # Load reference molecule
2423 ref_path = '/workspace/reference.sdf'
2424 suppl = Chem.SDMolSupplier(ref_path)
2425 ref_mol = [mol for mol in suppl if mol is not None][0]
2426
2427 print("=== REFERENCE LIGAND ANALYSIS ===")
2428 print(f"SMILES: {Chem.MolToSmiles(ref_mol)}")
2429 print(f"Heavy atoms: {ref_mol.GetNumHeavyAtoms()}")
2430 ... [61 more lines elided]

```

Tool Output (turn 4)

```

2430 === REFERENCE LIGAND ANALYSIS ===
2431 SMILES: N#Cc1cnn2c3cc(nc12)N1CCC[C@@H]1CC(=O)NCCOc1cc(cc2c1NC(=O)CC2)N3
2432 Heavy atoms: 35
2433
2434 Reference Bemis-Murcko scaffold:
2435 SMILES: O=C1C[C@H]2CCCN2c2cc(n3nccc3n2)Nc2cc3c(c(c2)OCCN1)NC(=O)CC3
2436 Heavy atoms: 33
2437
2438 Reference ligand saved to /workspace/reference.png
2439 Reference scaffold saved to /workspace/reference_scaffold.png
2440
2441 === FUNCTIONAL GROUP ANALYSIS ===
2442 Contains nitrile (C#N): False
2443 Contains amide (C(=O)N): True
2444 Contains aromatic N: 3
2445 ... [21 more lines elided]

```

Assistant (turn 5)

```

2442 Good! The simple test molecule already meets the scaffold difference criteria (ratio 0.182 < 0.65, similarity
2443 ↪ 0.056 < 0.5). But I need to design something that maintains binding interactions. Let me think about the
2444 ↪ pocket residues and what interactions the reference ligand likely makes.
2445
2446 Looking at the pocket residues: A:47-THR, A:50-MET, A:51-ALA, A:52-CYS, A:53-SER, A:54-GLY, A:55-LEU, A:57-TYR,
2447 ↪ A:88-PHE, A:112-GLN, A:113-MET, A:114-GLU, A:115-HIS.
2448
2449 These suggest potential interactions:
2450 - Aromatic residues (TYR57, PHE88, HIS115) could have pi-pi stacking or hydrophobic interactions
2451 - Polar residues (THR47, SER53, GLN112, GLU114, HIS115) could form hydrogen bonds
2452 - CYS52 could form disulfide or hydrophobic interactions
2453
2454 The reference ligand has 3 H-bond donors and 8 acceptors with TPSA of 136.68 Angstrom^2, suggesting significant
2455 ↪ polar interactions.
2456
2457 Let me design a new scaffold that maintains similar pharmacophores but with different core structure. I'll start
2458 ↪ by designing a molecule, then test it with Boltz.

```

Tool Call: python (turn 5, 0.49s)**Reasoning:**

```

2459 Design a new molecule with different scaffold but similar pharmacophore properties. I'll create a molecule with
2460 ↪ aromatic rings, amide groups, and similar H-bond donors/acceptors, but with a completely different ring
2461 ↪ system.

```

Code:

```

2462
2463 from rdkit import Chem
2464 from rdkit.Chem import AllChem, rdFMCS, rdMolDescriptors
2465 from rdkit.Chem import Descriptors
2466 from rdkit.Chem.Scaffolds import MurckoScaffold
2467 import numpy as np
2468
2469 # Load reference
2470 ref_path = '/workspace/reference.sdf'
2471 suppl = Chem.SDMolSupplier(ref_path)
2472 ref_mol = [mol for mol in suppl if mol is not None][0]
2473 ref_scaffold = MurckoScaffold.GetScaffoldForMol(ref_mol)
2474 ref_fp = AllChem.GetMorganFingerprint(ref_mol, 2)
2475
2476 # Design 1: Biphenyl-based scaffold with amide linkages
2477 design1_smiles = "CC(=O)Nc1ccc(cc1)C1=NC(=O)NC(=O)C1c1ccc(NC(=O)C)cc1"
2478

```

2475 ... [68 more lines elided]
 2476

2477 **Tool Output (turn 5)**

2478 === MOLECULE DESIGNS ===
 2479
 2480 Design1:
 2481 SMILES: CC(=O)Nc1ccc(cc1)C1=NC(=O)NC(=O)C1c1ccc(NC(=O)C)cc1
 2482 MW: 378.1, LogP: 2.43, TPSA: 116.7
 2483 HBD: 3, HBA: 4
 2484 Scaffold atoms: 20
 2485 MCS ratio: 0.333 (GOOD)
 2486 FP similarity: 0.118 (GOOD)
 2487 [ok] Meets scaffold difference criteria!
 2488
 2489 Design2:
 2490 SMILES: CC(=O)Nc1ccc(cc1)C1=NC2=CC=CC=C2C(=O)N1c1ccc(NC(=O)C)cc1
 2491 MW: 412.2, LogP: 3.97, TPSA: 93.1
 2492 HBD: 2, HBA: 4
 2493 ... [31 more lines elided]

2494 **Elision**

2495 *Turns 6-8 (3 agent turns) omitted.*

2496 **Assistant (turn 9)**

2497 Better at 0.555, but still below 0.7. Let me analyze the reference ligand more carefully and try to design
 2498 ↪ something closer to its properties. The reference has TPSA 136.7, HBD=3, HBA=8. Let me check Design4 which
 2499 ↪ has TPSA 122.3, HBD=4, HBA=6.

2500 **Tool Call: predict_boltz (turn 9, 135.56s)**

2501 **Arguments:**

2502 {
 2503 "reasoning": "Test Design4 which has properties closer to the reference ligand (higher TPSA, more H-bond
 2504 ↪ features). This molecule has a piperazine ring that might provide good geometry for binding.",
 2505 "smiles": "O=C(Nc1ccc(cc1)c1nc2ccc(cc2o1)Nc1ccc(cc1)C(=O)N)C1CCNCC1",
 2506 "protein_sequences": [
 2507 {
 2508 "id": "A",
 2509 "sequence": "GPGADSCIQFTRHASDVLLNLRSLRDLTDVVIVVSREQFRAHKTVLMACSGLFYSIFTDQLKCNLSVINLDPEINPEGFCILLDFMYTS]
 2510 ↪ RLNLREGNIMAVMATMYLQMEHVVDTCRKFIKASE"
 2511 }
 2512],
 2513 "pocket_residues": [
 2514 ["A",
 2515 47
 2516],
 2517]
 2518 ... [49 more lines elided]

2519 **Tool Output (turn 9)**

2520 {
 2521 "call_number": 3,
 2522 "budget_remaining": 5,
 2523 "cif_dir": "boltz_results/call_3/boltz_results_boltz_call_3/predictions/boltz_call_3",
 2524 "affinity_data": {
 2525 "affinity_pred_value": -0.5056207180023193,
 2526 "affinity_probability_binary": 0.7028120756149292,
 2527 "affinity_pred_value1": -0.7198585867881775,
 2528 "affinity_probability_binary1": 0.8170371055603027,
 2529 "affinity_pred_value2": -0.2913827896118164,
 2530 "affinity_probability_binary2": 0.5885871052742004
 2531 },
 2532 "binding_probability": 0.7028120756149292,
 2533 "affinity_pred_value": -0.5056207180023193
 2534 }

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

2530 **Assistant (turn 10)**

2531 Excellent! Design4 has binding probability of 0.703, which meets the >0.7 requirement! Now I need to verify it
2532 ↪ still meets the scaffold difference criteria. Let me double-check and also test its ADMET properties.
2533

2534 **Tool Call: python (turn 10, 0.39s)**

Reasoning:

2537 Double-check that Design4 meets all criteria: scaffold difference (MCS ratio < 0.65, fingerprint similarity <
2538 ↪ 0.5) and has binding probability > 0.7. Also check ADMET properties.
2539

Code:

```
2541 from rdkit import Chem
2542 from rdkit.Chem import AllChem, rdFMCS, rdMolDescriptors
2543 from rdkit.Chem import Descriptors
2544 from rdkit.Chem.Scaffolds import MurckoScaffold
2545
2546 # Load reference
2547 ref_path = '/workspace/reference.sdf'
2548 suppl = Chem.SDMolSupplier(ref_path)
2549 ref_mol = [mol for mol in suppl if mol is not None][0]
2550 ref_scaffold = MurckoScaffold.GetScaffoldForMol(ref_mol)
2551 ref_fp = AllChem.GetMorganFingerprint(ref_mol, 2)
2552
2553 # Design4 molecule
2554 design4_smiles = "O=C(Nc1ccc(cc1)clnc2ccc(cc2o1)Nc1ccc(cc1)C(=O)N)C1CCNCC1"
2555 design4_mol = Chem.MolFromSmiles(design4_smiles)
2556 ... [34 more lines elided]
```

2553 **Tool Output (turn 10)**

```
2554 === FINAL VERIFICATION FOR Design4 ===
2555 SMILES: O=C(Nc1ccc(cc1)clnc2ccc(cc2o1)Nc1ccc(cc1)C(=O)N)C1CCNCC1
2556
2557 Scaffold difference criteria:
2558   Reference scaffold atoms: 33
2559   Design4 scaffold atoms: 31
2560   MCS atoms: 15
2561   MCS ratio: 0.455 [ok]
2562   FP similarity: 0.149 [ok]
2563
2564 Binding probability from Boltz: 0.703 [ok]
2565
2566 Molecular properties:
2567   MW: 455.2
2568   LogP: 4.28
2569   ... [9 more lines elided]
```

2566 **Assistant (turn 11)**

2567 Perfect! Design4 meets all criteria. Let me also check its ADMET properties to ensure it's a reasonable
2568 ↪ drug-like molecule, though this isn't explicitly required by the task.
2569

2571 **Tool Call: predict_admet (turn 11, 7.90s)**

Arguments:

```
2572 {
2573   "reasoning": "Check ADMET properties for Design4 to ensure it's a reasonable drug-like molecule, even though
2574   ↪ not explicitly required. This will help confirm it's a viable candidate.",
2575   "smiles": "O=C(Nc1ccc(cc1)clnc2ccc(cc2o1)Nc1ccc(cc1)C(=O)N)C1CCNCC1"
2576 }
2577
```

2578 **Tool Output (turn 11)**

```
2580 {
2581   "herg": 0.7713615775108338,
2582   "bbb": 0.5673878788948059,
2583   "solubility": -3.329554869252309,
2584   "cyp3a4": 0.13602183759212494,
```

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

```
2585 "clearance": -19.668116752268368,  
2586 "caco2": -5.866392531878797,  
2587 "ppb": 86.2324968967286,  
2588 "ames": 0.31036943197250366  
2589 }
```

Assistant (turn 12)

```
2591 The ADMET properties look reasonable:  
2592 - hERG inhibition risk: 0.77 (moderate risk, common for many drugs)  
2593 - BBB penetration: 0.57 (moderate, could cross BBB)  
2594 - Solubility: -3.33 logS (moderately soluble)  
2595 - CYP3A4 inhibition: 0.14 (low risk)  
2596 - Clearance: -19.7 (moderate clearance)  
2597 - Caco2 permeability: -5.87 (moderate)  
2598 - Plasma protein binding: 86% (high, but acceptable)  
2599 - AMES mutagenicity: 0.31 (low risk)
```

The molecule appears drug-like. Now let me create the solution file with the SMILES string.

Tool Call: python (turn 12, 0.28s)

Reasoning:

Create the solution file with the SMILES string of Design4 molecule that meets all criteria.

Code:

```
2606 # Write the solution SMILES to file  
2607 solution_smiles = "O=C(Nc1ccc(cc1)c1nc2ccc(cc2o1)Nc1ccc(cc1)C(=O)N)C1CCNCC1"  
2608  
2609 with open('/workspace/solution.smi', 'w') as f:  
2610     f.write(solution_smiles)  
2611  
2612 print(f"Solution SMILES written to /workspace/solution.smi")  
2613 print(f"SMILES: {solution_smiles}")  
2614  
2615 # Verify file was created  
2616 import os  
2617 print(f"\nFile exists: {os.path.exists('/workspace/solution.smi')}")  
2618 print(f"File size: {os.path.getsize('/workspace/solution.smi')} bytes")  
2619  
2620 # Read back to verify  
2621 ... [3 more lines elided]
```

Tool Output (turn 12)

```
2620 Solution SMILES written to /workspace/solution.smi  
2621 SMILES: O=C(Nc1ccc(cc1)c1nc2ccc(cc2o1)Nc1ccc(cc1)C(=O)N)C1CCNCC1  
2622  
2623 File exists: True  
2624 File size: 56 bytes  
2625 File content: 'O=C(Nc1ccc(cc1)c1nc2ccc(cc2o1)Nc1ccc(cc1)C(=O)N)C1CCNCC1'
```

Assistant (turn 13)

Perfect! Now let me submit the solution.

Tool Call: submit (turn 13, 0.09s)

file_path: solution.smi

Tool Output (turn 13)

Submission accepted: solution.smi

Submit (turn 13)

Submitted file: solution.smi.

2640 Submitted artifact (solution.smi).

2641
2642
2643
2644
2645
2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694

Submitted artifact: solution.smi

```
O=C(Nc1ccc(cc1)c1nc2ccc(cc2o1)Nc1ccc(cc1)C(=O)N)C1CCNCC1
```

H.4. Lead Optimization

H.4.1. PASSING RUN: MINIMAX M2.7

Task ID: smdd_004_CYP19A1_8 **Model:** MiniMax M2.7 **Run:** minimax_m27

Evaluation result: PASSED.

Metric	Status	Result	Condition
Validity	passed	True	{"eq": true}
Hard Constraints	passed	True	{"eq": true}
Hold-Constant Properties	passed	True	{"eq": true}
Optimization Objectives	passed	True	{"eq": true}
Boltz Binding & Affinity	passed	True	{"eq": true}

Trajectory.

Showing the first 5 and last 5 of 19 agent turns; 9 intermediate turns are elided for space.

System Prompt

You are a computational medicinal chemist assistant. You are given a drug discovery task and must use your tools ↪ to solve it.

Available tools:

- python(code): Execute Python code. Available packages: rdkit, numpy, pandas, scipy, scikit-learn, biopython, ↪ biotite, openbabel, plip. To compute synthetic accessibility (SA) scores: `from rdkit.Chem import RDConfig; ↪ import sys, os; sys.path.append(os.path.join(RDConfig.RDContribDir, 'SA_Score')); import sascorer; ↪ sascorer.calculateScore(mol)`
- predict_boltz(smiles, protein_sequences, pocket_residues, co_ligand_smiles, num_samples): Run Boltz structure ↪ prediction. Returns binding probability and affinity. Budget: 8 calls.
- predict_admet(smiles): Predict ADMET properties. Returns hERG, BBB, solubility, CYP3A4, clearance, Caco2, PPB, ↪ AMES. Budget: 15 calls.
- submit(file_path): Submit your solution file.

You have no internet access. Work only with the files provided in /workspace. Outline your thought process ↪ before making tool calls. When ready, call submit() with the path to your solution file (relative to /workspace).

User (turn 0)

Task

You are given a co-folded protein-ligand complex (complex.pdb) and a reference ligand (reference.sdf). Modify ↪ the reference ligand to optimize the following properties:

- cyp3a4: decrease by at least 0.1 from baseline 0.4468
 - bbb: decrease by at least 0.1 from baseline 0.9828
- While holding these properties approximately constant:
- binding_affinity: keep within 0.3 of baseline (-1.3447) (increase tolerance)
 - hERG: keep within 0.1 of baseline (0.2067) (increase tolerance)
 - solubility: keep within 0.5 of baseline (-4.7943) (decrease tolerance)

The output molecule must satisfy ALL of the following hard constraints:

- SMILES validity: must be parseable and sanitizable by RDKit (no radicals, no disconnected fragments, no ↪ unreasonable valences)
 - Molecular weight: < 600 Da
 - LogP: between -1 and 5
 - TPSA: < 140 Å²
 - H-bond donors: <= 5
 - H-bond acceptors: <= 10
 - Rotatable bonds: <= 10
 - Net formal charge: between -2 and +2
 - Synthetic accessibility: SA score < 4.5
 - PAINS filter: must pass (no pan-assay interference substructures)
 - Structural alerts: must pass Brenk/NIH reactive group filters
 - Tanimoto similarity to reference: >= 0.7 (Morgan fingerprint, radius=2, 2048 bits)
 - Boltz binding probability: affinity_probability_binary > 0.7
- The pocket is defined by residues (chain, 1-indexed SEQRES position): A:81, A:99, A:100, A:187, A:190, A:268, ↪ A:271, A:272, A:273, A:275, A:276, A:336, A:338, A:339, A:340, A:341, A:443, A:444
- Output a single SMILES string to solution.smi.

Your workspace at /workspace contains the task files. Start by examining them.

2750
2751
2752
2753
2754
2755
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804

Assistant (turn 1)

(no assistant text emitted; tool call only)

Tool Call: python (turn 1, 0.45s)

Reasoning:

First, I need to examine the workspace files to understand the protein-ligand complex and reference ligand I'm
↪ working with.

Code:

```
import os

# List workspace files
print("Workspace files:")
for f in os.listdir('/workspace'):
    print(f" {f}")

# Read the reference ligand
from rdkit import Chem

print("\n--- Reference ligand (SDF) ---")
suppl = Chem.SDMolSupplier('/workspace/reference.sdf')
ref_mols = list(suppl)
print(f"Number of molecules in SDF: {len(ref_mols)}")
for i, mol in enumerate(ref_mols):
    ... [13 more lines elided]
```

Tool Output (turn 1)

```
Workspace files:
.tmp_script_11a51954.py
reference.sdf
complex.pdb

--- Reference ligand (SDF) ---
Number of molecules in SDF: 1
Molecule 0: reference_ligand
SMILES: CC[C@@H]1C[C@@H]2[C@H](CC[C@]3(C)C(=O)CC[C@@H]23)[C@]2(C)C=CC(=O)C=C12
Traceback (most recent call last):
  File "/workspace/.tmp_script_11a51954.py", line 19, in <module>
    print(f" MW: {AllChem.CalcExactMolWt(mol):.4f}")
NameError: name 'AllChem' is not defined
```

Assistant (turn 2)

(no assistant text emitted; tool call only)

Tool Call: python (turn 2, 0.75s)

Reasoning:

I need to fix the import for AllChem and properly examine the reference ligand.

Code:

```
from rdkit import Chem
from rdkit.Chem import AllChem, Descriptors, Lipinski, rdMolDescriptors
from rdkit import DataStructs

# Read the reference ligand
print("--- Reference ligand (SDF) ---")
suppl = Chem.SDMolSupplier('/workspace/reference.sdf')
ref_mol = list(suppl)[0]
ref_smiles = Chem.MolToSmiles(ref_mol)
print(f"SMILES: {ref_smiles}")

# Calculate properties
print(f"\nMolecular properties:")
print(f" MW: {AllChem.CalcExactMolWt(ref_mol):.4f}")
```

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

```
2805 print(f" LogP (MLP): {Descriptors.MolLogP(ref_mol):.4f}")
2806 ... [55 more lines elided]
2807
```

Tool Output (turn 2)

```
2808 --- Reference ligand (SDF) ---
2809 SMILES: CC[C@@H]1C[C@@H]2[C@H](CC[C@]3(C)C(=O)CC[C@@H]23)[C@@]2(C)C=CC(=O)C=C12
2810
2811 Molecular properties:
2812 MW: 312.2089
2813 LogP (MLP): 4.4995
2814 TPSA: 34.1400
2815 HBD: 0
2816 HBA: 2
2817 Rotatable bonds: 1
2818 Net charge: 0
2819 Formal charge: 0
2820 SA score: 4.3399
2821 Traceback (most recent call last):
2822   File "/workspace/.tmp_script_cl688f62.py", line 33, in <module>
2823 ... [2 more lines elided]
2824
```

Assistant (turn 3)

(no assistant text emitted; tool call only)

Tool Call: python (turn 3, 0.44s)

Reasoning:

I have the reference SMILES. Now let me properly check PAINS filter and get the ADMET properties, and extract → the protein sequence.

Code:

```
2825 from rdkit import Chem
2826 from rdkit.Chem import AllChem, Descriptors, Lipinski
2827 from rdkit.Chem.FilterCatalog import FilterCatalogParams, FilterCatalog
2828
2829 # Read reference ligand
2830 suppl = Chem.SDMolSupplier('/workspace/reference.sdf')
2831 ref_mol = list(suppl)[0]
2832 ref_smiles = Chem.MolToSmiles(ref_mol)
2833 print(f"Reference SMILES: {ref_smiles}")
2834
2835 # Check PAINS - try different approach
2836 try:
2837     params = FilterCatalogParams()
2838     params.AddCatalog(FilterCatalogParams.PAINS)
2839     catalog = FilterCatalog(params)
2840     ... [35 more lines elided]
2841
2842
```

Tool Output (turn 3)

```
2843 Reference SMILES: CC[C@@H]1C[C@@H]2[C@H](CC[C@]3(C)C(=O)CC[C@@H]23)[C@@]2(C)C=CC(=O)C=C12
2844 PAINS check error: type object 'FilterCatalogParams' has no attribute 'PAINS'
2845
2846 SEQRES chain lengths:
2847
2848
```

Assistant (turn 4)

(no assistant text emitted; tool call only)

Tool Call: python (turn 4, 0.43s)

Reasoning:

Let me look more closely at the PDB file structure and extract the protein sequence properly.

Code:

```
2849 from rdkit import Chem
2850
```

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

```

2860 from rdkit.Chem import AllChem, Descriptors, Lipinski
2861
2862 # Read and parse complex PDB for protein sequence
2863 with open('/workspace/complex.pdb', 'r') as f:
2864     pdb_content = f.read()
2865
2866 print("First 200 lines of PDB:")
2867 for i, line in enumerate(pdb_content.split('\n')[:200]):
2868     print(f"{i}: {line.rstrip()}")

```

Tool Output (turn 4)

```

2869 First 200 lines of PDB:
2870 0: ATOM      1  N  MET  A  1      10.329  51.314  11.107  1.00  0.00      N
2871 1: ATOM      2  CA  MET  A  1      10.138  49.975  10.539  1.00  0.00      C
2872 2: ATOM      3  C   MET  A  1      10.169  50.023   9.013  1.00  0.00      C
2873 3: ATOM      4  O   MET  A  1       9.598  50.924   8.406  1.00  0.00      O
2874 4: ATOM      5  CB  MET  A  1       8.816  49.353  11.010  1.00  0.00      C
2875 5: ATOM      6  CG  MET  A  1       8.685  47.857  10.741  1.00  0.00      C
2876 6: ATOM      7  SD  MET  A  1       7.043  47.336  10.234  1.00  0.00      S
2877 7: ATOM      8  CE  MET  A  1       6.102  47.552  11.740  1.00  0.00      C
2878 8: ATOM      9  N   ALA  A  2      10.848  49.055   8.417  1.00  0.00      N
2879 9: ATOM     10  CA  ALA  A  2      10.945  48.961   6.965  1.00  0.00      C
2880 10: ATOM     11  C   ALA  A  2       9.642  48.435   6.361  1.00  0.00      C
2881 11: ATOM     12  O   ALA  A  2       8.875  47.736   7.023  1.00  0.00
2882 ... [middle truncated] ...
2883 ... [13 more lines elided]

```

Assistant (turn 5)

(no assistant text emitted; tool call only)

Tool Call: python (turn 5, 0.39s)

Reasoning:

Let me extract the ligand from the complex PDB file and get the full protein sequence.

Code:

```

2890 from rdkit import Chem
2891 from rdkit.Chem import AllChem, Descriptors, Lipinski
2892
2893 # Read the complex PDB
2894 with open('/workspace/complex.pdb', 'r') as f:
2895     pdb_lines = f.readlines()
2896
2897 # Find all chains and their sequences from ATOM lines
2898 chain_res = {}
2899 for line in pdb_lines:
2900     if line.startswith('ATOM') or line.startswith('HETATM'):
2901         chain = line[21]
2902         resnum = int(line[22:26].strip())
2903         resname = line[17:20].strip()
2904         if chain not in chain_res:
2905             chain_res[chain] = []
2906         chain_res[chain].append(resname)
2907 ... [26 more lines elided]

```

Tool Output (turn 5)

```

2908 Chains found: ['A', 'L']
2909 Chain A: 473 residues
2910 Sequence: MAKKTSSKGRSSIPGPGYCMGIGPLISHGRFLWMGIGSACNYNRYVGEFMRVWISGEETLIISKSSSMFHIMKHNHYSSRFGSKLGLQCIGMHEKGI I
2911 ↳ FNNNPPELWKTTRPFFMKALSGPGLVRMVTVCAESLKTHLDRLEEVTVNESGYVDVLTLLRVRMLDTSNTLFLRIPLDSEAIUVVKIQGYFDQWQALLIKPDIFFKIS I
2912 ↳ WLYKYEKSVKDLKDAIEVLIAEKRRRISTEEKLEECMDFATELILAEKRGDLTRENVNQCILEMLIAAPDTMSVSLFFMLFLIAKHPNVVEEAIKEIQTVIGER I
2913 ↳ DIKIDDIQKLKVMENFIYESMRYQPVDLVMRKALEDDVIDGYPVKKGTNIILNIGRMHRELFPPKPNFLENFAKNVPYRYFQPFQPGFGRGCAGKYIAMVMMK I
2914 ↳ AILVTLRFRFHVKTLLQGCQVESIQIKIHDLSLHPDETNRNMLEMIFTPRNSDRCLHEHHHH
2915 First 100:
2916 ↳ MAKKTSSKGRSSIPGPGYCMGIGPLISHGRFLWMGIGSACNYNRYVGEFMRVWISGEETLIISKSSSMFHIMKHNHYSSRFGSKLGLQCIGMHEKGIIF
2917 Chain L: 1 residues
2918 Sequence: X
2919 First 100: X
2920 Full chain A length: 473

```

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

Chain A residue numbers: [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, ↪ 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50]...
 Last residues: [424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442, ↪ 443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, ↪ 464, 465, 466, 467, 468, 469, 470, 471, 472, 473]

Elision

Turns 6-14 (9 agent turns) omitted.

Assistant (turn 15)

(no assistant text emitted; tool call only)

Tool Call: predict_boltz (turn 15, 266.93s)

Arguments:

```
{
  "reasoning": "Testing the binding probability of candidate Mod 3 (shortened side chain). Need to provide the
  ↪ protein sequence and pocket residues.",
  "smiles": "C[C@@H]1C[C@@H]2[C@H](CC[C@]3(C)C(=O)CC[C@@H]23)[C@@]2(C)C=CC(=O)C=C12",
  "protein_sequences": [
    {
      "id": "A",
      "sequence": "MAKKTSSKGRSSIPGPYCMGIGPLISHGRFLWMGIGSACNYNRYVGEFMRVWISGEETLIISKSSSMFHMKHNHYSSRFGSKLGLQCIG ]
  ↪ MHEKGIIFNPNPELWKTTRPFFMKALSGPGLVRMVTCAESLKTHLDRLEEVNTEGSDVVDVTLRRVMDLTSNTLFLRIPLESAIVVKIQGYFDAQAL ]
  ↪ LIKPDIFFKISWLYKYEKSVKDLKDAIEVLIAEKRRRISTEEKLEECMDFATELILAEKRGDLTRENVNQCILEMLIAAPDTMSVSLFFMLFLIAKHPNV ]
  ↪ EEATIKIQTIVIGERDIKIDDIQKLVEMENFIYESMRYQPVDLVMRKALEDDVIDGYPVKKGTNIIILNIGRMHRLEFFPKPNEFTLENFAKNVPHYRFPQ ]
  ↪ FGFGRGCAGKYIAMVMMKAILVTLRRFHVKTLQGCQVESIQKIHDLSLHPDETKNMLEMIFTPRNSDRCLEHHHHH"
    }
  ],
  "pocket_residues": [
    [
      "A",
      81
    ],
    [
      ... [69 more lines elided]
    ]
  ]
}
```

Tool Output (turn 15)

```
{
  "call_number": 1,
  "budget_remaining": 7,
  "cif_dir": "boltz_results/call_1/boltz_results_boltz_call_1/predictions/boltz_call_1",
  "affinity_data": {
    "affinity_pred_value": -1.272486686706543,
    "affinity_probability_binary": 0.889390766620636,
    "affinity_pred_value1": -1.1789921522140503,
    "affinity_probability_binary1": 0.9195749163627625,
    "affinity_pred_value2": -1.365981101989746,
    "affinity_probability_binary2": 0.8592066168785095
  },
  "binding_probability": 0.889390766620636,
  "affinity_pred_value": -1.272486686706543
}
```

Assistant (turn 16)

(no assistant text emitted; tool call only)

Tool Call: python (turn 16, 0.41s)

Reasoning:

The binding probability is 0.889 > 0.7, which meets the constraint. Now let me check PAINS filter and structural
 ↪ alerts before writing the solution.

Code:

```
from rdkit import Chem
from rdkit.Chem import AllChem, Descriptors, Lipinski
```

```

2970 from rdkit.Chem import FilterCatalog
2971
2972 smi = "C[C@@H]1C[C@@H]2[C@H](CC[C@]3(C)C(=O)CC[C@@H]23)[C@@]2(C)C=CC(=O)C=C12"
2973 mol = Chem.MolFromSmiles(smi)
2974
2975 # Check PAINS filter
2976 try:
2977     # Try to use FilterCatalog
2978     from rdkit.Chem import rdfilterlibs
2979
2980     # PAINS filter using predefined patterns
2981     pains = rdfilterlibs.PAINS()
2982     matches = pains.RemoveFilter(mol)
2983     ... [49 more lines elided]

```

Tool Output (turn 16)

```

2984 PAINS check: cannot import name 'rdfilterlibs' from 'rdkit.Chem'
2985 ↪ (/opt/conda/lib/python3.10/site-packages/rdkit/Chem/__init__.py)
2986 Brenk check: cannot import name 'rdfilterlibs' from 'rdkit.Chem'
2987 ↪ (/opt/conda/lib/python3.10/site-packages/rdkit/Chem/__init__.py)
2988
2989 Running basic structural alerts check...
2990
2991 === FINAL SUMMARY ===
2992 SMILES: C[C@@H]1C[C@@H]2[C@H](CC[C@]3(C)C(=O)CC[C@@H]23)[C@@]2(C)C=CC(=O)C=C12
2993 All ADMET targets met:
2994   cyp3a4: 0.225 < 0.3468 [ok]
2995   bbb: 0.827 < 0.8828 [ok]
2996   herg: 0.116 (within 0.1 of baseline) [ok]
2997   solubility: -4.437 (within 0.5 of baseline) [ok]
2998 All hard constraints met:
2999   MW < 600: 298.2 [ok]
3000   LogP between -1 and 5: 4.11 [ok]
3001   ... [8 more lines elided]

```

Assistant (turn 17)

(no assistant text emitted; tool call only)

Tool Call: python (turn 17, 0.45s)**Reasoning:**

Let me try checking PAINS filter using an alternative approach and make sure the solution is valid.

Code:

```

3003 from rdkit import Chem
3004 from rdkit.Chem import AllChem, Descriptors, Lipinski
3005 from rdkit.Chem import FilterCatalog
3006
3007 smi = "C[C@@H]1C[C@@H]2[C@H](CC[C@]3(C)C(=O)CC[C@@H]23)[C@@]2(C)C=CC(=O)C=C12"
3008 mol = Chem.MolFromSmiles(smi)
3009
3010 # Try different approach for PAINS
3011 try:
3012     from rdkit.Chem import FilterCatalog
3013     params = FilterCatalog.FilterCatalogParams()
3014     # Try adding catalogs
3015     params.AddCatalog(FilterCatalog.FilterCatalogs.PAINS)
3016     print("PAINS catalog available")
3017 except Exception as e:
3018     ... [45 more lines elided]

```

Tool Output (turn 17)

```

3019 PAINS catalog error: module 'rdkit.Chem.FilterCatalog' has no attribute 'FilterCatalogs'
3020
3021 Checking for structural alerts...
3022
3023 No problematic structural alerts found.
3024 Writing solution to file...
3025 Solution written: C[C@@H]1C[C@@H]2[C@H](CC[C@]3(C)C(=O)CC[C@@H]23)[C@@]2(C)C=CC(=O)C=C12
3026 File content: C[C@@H]1C[C@@H]2[C@H](CC[C@]3(C)C(=O)CC[C@@H]23)[C@@]2(C)C=CC(=O)C=C12

```

3025
3026
3027
3028
3029
3030
3031
3032
3033
3034
3035
3036
3037
3038
3039
3040
3041
3042
3043
3044
3045
3046
3047
3048
3049
3050
3051
3052
3053
3054
3055
3056
3057
3058
3059
3060
3061
3062
3063
3064
3065
3066
3067
3068
3069
3070
3071
3072
3073
3074
3075
3076
3077
3078
3079

Assistant (turn 18)

(no assistant text emitted; tool call only)

Tool Call: python (turn 18, 0.78s)

Reasoning:

Let me do a final comprehensive check before submitting.

Code:

```
from rdkit import Chem
from rdkit.Chem import AllChem, Descriptors, Lipinski, rdMolDescriptors
from rdkit import DataStructs
import sys, os
sys.path.append(os.path.join(Chem.RDConfig.RDContribDir, 'SA_Score'))
import sascorer

# Read and verify solution
with open('/workspace/solution.smi', 'r') as f:
    solution_smiles = f.read().strip()

print("=== FINAL VERIFICATION ===")
print(f"SMILES: {solution_smiles}")

# Parse and verify
... [46 more lines elided]
```

Tool Output (turn 18)

```
=== FINAL VERIFICATION ===
SMILES: C[C@@H]1C[C@@H]2[C@H](CC[C@]3(C)C(=O)CC[C@@H]23)[C@@]2(C)C=CC(=O)C=C12
Valid: [ok]

--- Property Summary ---
MW: 298.19 (< 600) [ok]
LogP: 4.11 (-1 to 5) [ok]
TPSA: 34.14 (< 140) [ok]
HBD: 0 (<= 5) [ok]
HBA: 2 (<= 10) [ok]
RotB: 0 (<= 10) [ok]
Charge: 0 (-2 to 2) [ok]
SA: 4.22 (< 4.5) [ok]
Tanimoto: 0.729 (>= 0.7) [ok]

... [13 more lines elided]
```

Assistant (turn 19)

(no assistant text emitted; tool call only)

Tool Call: submit (turn 19, 0.09s)

file_path: solution.smi

Tool Output (turn 19)

Submission accepted: solution.smi

Submit (turn 19)

Submitted file: solution.smi.

Submitted artifact (solution.smi).

Submitted artifact: solution.smi

```
C[C@@H]1C[C@@H]2[C@H](CC[C@]3(C)C(=O)CC[C@@H]23)[C@@]2(C)C=CC(=O)C=C12
```

H.5. Fragment Assembly

H.5.1. PASSING RUN: GEMINI 3.1 PRO (HIGH REASONING)

Task ID: smdd_005_P04585_0 **Model:** Gemini 3.1 Pro (high reasoning) **Run:** gemini_31_pro_high_lite

Evaluation result: PASSED.

Metric	Status	Result	Condition
Validity	passed	True	{"eq": true}
Fragment Incorporation	passed	True	{"eq": true}
Fragment Linking	passed	True	{"eq": true}
Hard Constraints	passed	True	{"eq": true}
Boltz Binding & Pose	passed	True	{"eq": true}

Trajectory.

Showing the first 5 and last 5 of 35 agent turns; 25 intermediate turns are elided for space.

System Prompt

You are a computational medicinal chemist assistant. You are given a drug discovery task and must use your tools
↪ to solve it.

Available tools:

- python(code): Execute Python code. Available packages: rdkit, numpy, pandas, scipy, scikit-learn, biopython,
↪ biotite, openbabel, plip. To compute synthetic accessibility (SA) scores: `from rdkit.Chem import RDConfig;
↪ import sys, os; sys.path.append(os.path.join(RDConfig.RDContribDir, 'SA_Score')); import sascorer;
↪ sascorer.calculateScore(mol)`
- predict_boltz(smiles, protein_sequences, pocket_residues, co_ligand_smiles, num_samples): Run Boltz structure
↪ prediction. Returns binding probability and affinity. Budget: 8 calls.
- predict_admet(smiles): Predict ADMET properties. Returns hERG, BBB, solubility, CYP3A4, clearance, Caco2, PPB,
↪ AMES. Budget: 15 calls.
- submit(file_path): Submit your solution file.

You have no internet access. Work only with the files provided in /workspace. Outline your thought process
↪ before making tool calls.
When ready, call submit() with the path to your solution file (relative to /workspace).

User (turn 0)

Task

You are given 1 molecular fragments with 3D poses inside a protein binding pocket. Design a single drug-like
↪ molecule that covalently links all fragments and binds to the target. The incorporated fragments must
↪ maintain their 3D poses in the binding pocket when docked.

Fragments:

- Fragment 0: see fragment_0.sdf

Protein pocket: see pocket.pdb

The pocket is defined by residues (chain, 1-indexed SEQRES position): A:8, A:23, A:25, A:26, A:27, A:28, A:29,
↪ A:30, A:31, A:32, A:45, A:47, A:48, A:49, A:50, A:52, A:76, A:80, A:81, A:82, A:84, B:23, B:25, B:26, B:27,
↪ B:28, B:29, B:30, B:31, B:32, B:45, B:46, B:47, B:48, B:49, B:50, B:51, B:52, B:53, B:76, B:80, B:81, B:82,
↪ B:84

The output molecule must satisfy ALL of the following hard constraints:

- SMILES validity: must be parseable and sanitizable by RDKit (no radicals, no disconnected fragments, no
↪ unreasonable valences)
- Molecular weight: < 600 Da
- LogP: between -1 and 5
- TPSA: < 140 Å²
- H-bond donors: <= 5
- H-bond acceptors: <= 10
- Rotatable bonds: <= 10
- Net formal charge: between -2 and +2
- Synthetic accessibility: SA score < 4.5
- PAINS filter: must pass (no pan-assay interference substructures)
- Structural alerts: must pass Brenk/NIH reactive group filters
- Fragment incorporation: all fragment SMARTS must be present as substructures
- Fragment linking: all fragments must be covalently connected
- Boltz binding probability: affinity_probability_binary > 0.7
- Boltz binding affinity: affinity_pred_value < -2.2303 (lower is better)
- Fragment 3D pose RMSD: < 2.0 Å for each fragment (docked pose vs. input fragment pose)

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

3135 Output a single SMILES string to solution.smi.
 3136 Your workspace at /workspace contains the task files. Start by examining them.
 3137

3138 **Assistant (turn 1)**

3139 *(no assistant text emitted; tool call only)*

3140 **Tool Call: python (turn 1, 0.13s)**

3141 **Reasoning:**

3142 Check available files in the workspace

3143 **Code:**

```
3144 import os
3145 print("Files in /workspace:", os.listdir())
```

3146 **Tool Output (turn 1)**

3147 Files in /workspace: ['fragment_0.sdf', 'pocket.pdb', '.tmp_script_8735ed5a.py']

3148 **Assistant (turn 2)**

3149 *(no assistant text emitted; tool call only)*

3150 **Tool Call: python (turn 2, 0.30s)**

3151 **Reasoning:**

3152 Examine the fragment SMILES and the protein pocket structure.

3153 **Code:**

```
3154 from rdkit import Chem
3155 frag0 = Chem.SDMolSupplier('fragment_0.sdf')[0]
3156 smiles = Chem.MolToSmiles(frag0)
3157 print("Fragment 0 SMILES:", smiles)
3158
3159 with open('pocket.pdb') as f:
3160     pdb_lines = f.readlines()
3161     print("Pocket PDB head:", "".join(pdb_lines[:20]))
3162     print("Pocket PDB tail:", "".join(pdb_lines[-20:]))
```

3163 **Tool Output (turn 2)**

```
3164 Fragment 0 SMILES: c1c0cc0c1
3165 Pocket PDB head: ATOM      1  N      PRO A      1      8.976 -13.075  5.421  1.00  0.00      N
3166 ATOM      2  CA      PRO A      1      8.509 -13.590  4.130  1.00  0.00      C
3167 ATOM      3  C      PRO A      1      8.707 -12.580  3.004  1.00  0.00      C
3168 ATOM      4  O      PRO A      1      8.956 -11.397  3.253  1.00  0.00      O
3169 ATOM      5  CB      PRO A      1      7.009 -13.842  4.347  1.00  0.00      C
3170 ATOM      6  CG      PRO A      1      6.634 -12.863  5.409  1.00  0.00      C
3171 ATOM      7  CD      PRO A      1      7.844 -12.790  6.324  1.00  0.00      C
3172 ATOM      8  N      GLN A      2      8.583 -13.064  1.779  1.00  0.00      N
3173 ATOM      9  CA      GLN A      2      8.529 -12.201  0.615  1.00  0.00      C
3174 ATOM     10  C      GLN A      2      7.105 -12.276  0.079  1.00  0.00      C
3175 ATOM     11  O      GLN A      2      6.634 -13.354 -0.276  1.00  0.00      O
3176 ATOM     12  CB      GLN A      2      9.527 -12.633 -0.461  1.00  0.
3177
3178 ... [middle truncated] ...
3179 ... [15 more lines elided]
```

3180 **Assistant (turn 3)**

3181 *(no assistant text emitted; tool call only)*

Tool Call: python (turn 3, 0.46s)**Reasoning:**

Extract the full protein sequence from pocket.pdb to be used in Boltz predictions and examine fragment_0 3D
 ↪ structure.

Code:

```
import urllib.request
import json

# Can we get the sequence from PDB directly using biopython?
from Bio import PDB
parser = PDB.PDBParser()
structure = parser.get_structure('pocket', 'pocket.pdb')
from Bio.PDB.Polypeptide import PPBuilder
ppb = PPBuilder()
seqs = {}
for pp in ppb.build_peptides(structure):
    chain_id = pp[0].get_parent().id
    seq = pp.get_sequence()
    if chain_id not in seqs:
        seqs[chain_id] = str(seq)
... [9 more lines elided]
```

Tool Output (turn 3)

```
Sequences: {'A':
↪ 'PQITLWKRPLVTIKIGGQLKEALLDTGADDTVIEEMSLPGRWPKMIGGIGGFVKVRQYDQIIIEIAGHKAIGTVLVGPTPVNIIGRNLLTQIGATLNF', 'B':
↪ 'PQITLWKRPLVTIKIGGQLKEALLDTGADDTVIEEMSLPGRWPKMIGGIGGFVKVRQYDQIIIEIAGHKAIGTVLVGPTPVNIIGRNLLTQIGATLNF'}
Fragment 0 SDF:
reference_ligand      3D
RDKit

11 11  0  0  0  0  0  0  0  0  0999 V2000
-7.2030  1.7580  1.2190 C  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
-7.5100  0.5480  0.5420 C  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
-8.8450  0.1370  0.4340 C  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
-9.8700  0.9200  0.9900 C  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
-9.5560  2.1060  1.6750 C  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
-8.2160  2.5290  1.7850 C  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
-6.7055 -0.0621  0.1056 H  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
-9.0908 -0.8001 -0.0869 H  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
-10.9190  0.6046  0.8891 H  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
... [8 more lines elided]
```

Assistant (turn 4)

(no assistant text emitted; tool call only)

Tool Call: python (turn 4, 0.29s)**Reasoning:**

Identify the attachment point on fragment_0

Code:

```
from rdkit import Chem
mol = Chem.SDMolSupplier('fragment_0.sdf', removeHs=False)[0]
for atom in mol.GetAtoms():
    if atom.GetSymbol() == 'C' and atom.GetDegree() == 2:
        print(f"Attachment point at C index {atom.GetIdx()} (no H attached)")
```

Tool Output (turn 4)

Attachment point at C index 0 (no H attached)

Assistant (turn 5)

(no assistant text emitted; tool call only)

Tool Call: python (turn 5, 0.52s)**Reasoning:**

Check standard HIV protease inhibitors to see if any meet all constraints (especially Rotatable bonds ≤ 10) and \leftrightarrow contain the phenyl fragment.

Code:

```
from rdkit import Chem
from rdkit.Chem import Descriptors, rdMolDescriptors
import sys, os
from rdkit.Chem import RDConfig
sys.path.append(os.path.join(RDConfig.RDContribDir, 'SA_Score'))
import sasorerer

candidates = {
    "Darunavir": "CC(C)CN(CC(O)C(Cc1ccccc1)NC(=O)OC2COC2O)S(=O)(=O)c3ccc(N)cc3",
    "Amprenavir": "CC(C)CN(CC(O)C(Cc1ccccc1)NC(=O)OC2COC2)S(=O)(=O)c3ccc(N)cc3",
    "Fosamprenavir_active": "CC(C)CN(CC(O)C(Cc1ccccc1)NC(=O)OC2COC2)S(=O)(=O)c3ccc(N)cc3", # same as amprenavir
    "Indinavir": "CC(C)(C)NC(=O)C1CN(CC2CC(O)C(Cc3ccccc3)NC2=O)Cc4cncc(c4)C1",
    "Atazanavir": "COC(=O)NC(C(=O)NN(Cc1ccc(cc1)c2ccccc2)CC(O)C(Cc3ccccc3)NC(=O)C(NC(=O)OC)C(C)C(C)C(C)C"
}

... [12 more lines elided]
```

Tool Output (turn 5)

```
Darunavir: MW=521.6, LogP=1.4, TPSA=151.4, HBD=4, HBA=8, RotB=11, SA=3.8, HasPh=True
Amprenavir: MW=505.6, LogP=2.4, TPSA=131.2, HBD=3, HBA=7, RotB=11, SA=3.5, HasPh=True
Fosamprenavir_active: MW=505.6, LogP=2.4, TPSA=131.2, HBD=3, HBA=7, RotB=11, SA=3.5, HasPh=True
Indinavir: MW=464.6, LogP=2.1, TPSA=94.6, HBD=3, HBA=5, RotB=5, SA=5.0, HasPh=True
Atazanavir: MW=704.9, LogP=4.2, TPSA=171.2, HBD=5, HBA=9, RotB=14, SA=4.1, HasPh=True
```

Elision

Turns 6-30 (25 agent turns) omitted.

Assistant (turn 31)

(no assistant text emitted; tool call only)

Tool Call: python (turn 31, 0.54s)**Reasoning:**

Check if benzothiazole analog violates constraints.

Code:

```
from rdkit import Chem
from rdkit.Chem import Descriptors, rdMolDescriptors
import sys, os
from rdkit.Chem import RDConfig
sys.path.append(os.path.join(RDConfig.RDContribDir, 'SA_Score'))
import sasorerer

smi = "C1CCCC1N(C[C@@H](O)[C@H](Cc2ccccc2)NC(=O)O[C@H]3CO[C@@H]4OCC[C@H]43)S(=O)(=O)c5nc6ccccc6s5"
mol = Chem.MolFromSmiles(smi)
mw = Descriptors.MolWt(mol)
logp = Descriptors.MolLogP(mol)
tpsa = rdMolDescriptors.CalcTPSA(mol)
hbd = rdMolDescriptors.CalcNumHBD(mol)
hba = rdMolDescriptors.CalcNumHBA(mol)
rotb = rdMolDescriptors.CalcNumRotatableBonds(mol)
... [2 more lines elided]
```

Tool Output (turn 31)

```
Benzothiazole: MW=601.7, LogP=3.7, TPSA=127.3, HBD=2, HBA=9, RotB=10, SA=4.2
```

3300 Assistant (turn 32)

3301 (no assistant text emitted; tool call only)

3302
3303 Tool Call: python (turn 32, 0.49s)

3304 Reasoning:

3305 Check p-chloro for final submission.

3306 Code:

```

3307 from rdkit import Chem
3308 from rdkit.Chem import Descriptors, rdMolDescriptors
3309 import sys, os
3310 from rdkit.Chem import RDConfig
3311 sys.path.append(os.path.join(RDConfig.RDContribDir, 'SA_Score'))
3312 import sascore

3313 smi = "ClCCCC1N(C[C@@H](O)[C@H](Cc2ccccc2)NC(=O)O[C@H]3CO[C@@H]4OCC[C@H]43)S(=O)(=O)c5ccc(Cl)cc5"
3314 mol = Chem.MolFromSmiles(smi)
3315 mw = Descriptors.MolWt(mol)
3316 logp = Descriptors.MolLogP(mol)
3317 tpsa = rdMolDescriptors.CalcTPSA(mol)
3318 hbd = rdMolDescriptors.CalcNumHBD(mol)
3319 hba = rdMolDescriptors.CalcNumHBA(mol)
3320 rotb = rdMolDescriptors.CalcNumRotatableBonds(mol)
3321 ... [2 more lines elided]

```

3322 Tool Output (turn 32)

3323 p-Chloro: MW=579.1, LogP=3.7, TPSA=114.4, HBD=2, HBA=7, RotB=10, SA=4.0

3324 Assistant (turn 33)

3325 (no assistant text emitted; tool call only)

3326
3327 Tool Call: python (turn 33, 0.48s)

3328 Reasoning:

3329 Check p-OH analog properties.

3330 Code:

```

3331 from rdkit import Chem
3332 from rdkit.Chem import Descriptors, rdMolDescriptors
3333 import sys, os
3334 from rdkit.Chem import RDConfig
3335 sys.path.append(os.path.join(RDConfig.RDContribDir, 'SA_Score'))
3336 import sascore

3337 smi = "ClCCCC1N(C[C@@H](O)[C@H](Cc2ccccc2)NC(=O)O[C@H]3CO[C@@H]4OCC[C@H]43)S(=O)(=O)c5ccc(O)cc5"
3338 mol = Chem.MolFromSmiles(smi)
3339 mw = Descriptors.MolWt(mol)
3340 logp = Descriptors.MolLogP(mol)
3341 tpsa = rdMolDescriptors.CalcTPSA(mol)
3342 hbd = rdMolDescriptors.CalcNumHBD(mol)
3343 hba = rdMolDescriptors.CalcNumHBA(mol)
3344 rotb = rdMolDescriptors.CalcNumRotatableBonds(mol)
3345 ... [2 more lines elided]

```

3346 Tool Output (turn 33)

3347 p-OH: MW=560.7, LogP=2.8, TPSA=134.6, HBD=3, HBA=8, RotB=10, SA=4.0

3348 Assistant (turn 34)

3349 (no assistant text emitted; tool call only)

3350
3351
3352
3353
3354

3355 **Tool Call: python** (turn 34, 0.30s)

Reasoning:

Check PAINS/BRENK and save the best candidate to solution.smi

Code:

```
from rdkit import Chem
from rdkit.Chem import FilterCatalog

smi = "C1CCCC1N(C[C@@H](O)[C@H](Cc2ccccc2)NC(=O)O[C@H]3CO[C@@H]4OCC[C@H]43)S(=O)(=O)c5ccc(O)cc5"
mol = Chem.MolFromSmiles(smi)

# PAINS
params = FilterCatalog.FilterCatalogParams()
params.AddCatalog(FilterCatalog.FilterCatalogParams.FilterCatalogs.PAINS)
catalog = FilterCatalog.FilterCatalog(params)
if catalog.HasMatch(mol):
    print("PAINS matched!")

# BRENK
params_brenk = FilterCatalog.FilterCatalogParams()
... [11 more lines elided]
```

Tool Output (turn 34)

Passes PAINS and BRENK
Saved to solution.smi

Assistant (turn 35)

(no assistant text emitted; tool call only)

Tool Call: submit (turn 35, 0.06s)

file_path: solution.smi

Tool Output (turn 35)

Submission accepted: solution.smi

Submit (turn 35)

Submitted file: solution.smi.

Submitted artifact (solution.smi).**Submitted artifact: solution.smi**

```
C1CCCC1N(C[C@@H](O)[C@H](Cc2ccccc2)NC(=O)O[C@H]3CO[C@@H]4OCC[C@H]43)S(=O)(=O)c5ccc(O)cc5
```

3356
3357
3358
3359
3360
3361
3362
3363
3364
3365
3366
3367
3368
3369
3370
3371
3372
3373
3374
3375
3376
3377
3378
3379
3380
3381
3382
3383
3384
3385
3386
3387
3388
3389
3390
3391
3392
3393
3394
3395
3396
3397
3398
3399
3400
3401
3402
3403
3404
3405
3406
3407
3408
3409

I. Failure Mode Analysis

I.1. Absence of Cross-Turn SAR Synthesis: Repeated Testing of a Disqualified Structure

Task ID: smdd_004_ACE_8 **Model:** Claude Sonnet 4.6 **Run:** claude_sonnet_46_med_eval

Evaluation result: **FAILED.**

Metric	Status	Result	Condition
Validity	passed	True	{"eq": true}
Hard Constraints	passed	True	{"eq": true}
Hold-Constant Properties	failed	False	{"eq": true}
Optimization Objectives	passed	True	{"eq": true}
Boltz Binding & Affinity	failed	False	{"eq": true}

Summary of error. The agent records individual molecule failures but never extracts the structural generalization that a free para-OH on the phenyl ring is incompatible with the clearance and Caco-2 constraints of this task. After Turn 17 it knows that para-OH on phenyl drives clearance too high (e.g., candidate: 87.95, target < 72.12) and Caco-2 too low (e.g., candidate: -5.74, target > -5.58). After Turn 18 it knows that methyl ester + para-OH reproduces the same failure, but it still proposes a molecule made up of methyl ester + para-OH in Turn 30 that fails the same constraint requirements.

A budget-aware agent would have recognized after two data points that any free phenol on the aryl ring is a metabolic soft-spot and stopped proposing para-OH variants. Three of fifteen ADMET calls are spent confirming the same structural liability.

It is also worth noting that the witness molecule's SMILES string is O=C(O)CN1C(=O)[C@@H](N[C@@H](CSc2ccccc2)C(=O)O)CCc2ccccc2 which does not contain any para-OH functional groups to start with.

Trajectory.

Showing 5 of 51 agent turns; turns 1–16 and 19–29 are elided for space.

System Prompt

You are a computational medicinal chemist assistant. You are given a drug discovery task and must use your tools to solve it.

Available tools:

- python(code): Execute Python code. [...]
- predict_boltz(...): Run Boltz structure prediction. Budget: 8 calls.
- predict_admet(smiles): Predict ADMET properties. Budget: 15 calls.
- submit(file_path): Submit your solution file.

User (turn 0)

Task

You are given a co-folded protein-ligand complex (complex.pdb) and a reference ligand (reference.sdf).

Modify the reference ligand to optimize the following properties:

- clearance: decrease by at least 5.0 from baseline 77.1211
- binding_affinity: decrease by at least 0.3 from baseline -0.0448
- bbb: decrease by at least 0.1 from baseline 0.7589

While holding these properties approximately constant:

- herg: keep within 0.1 of baseline (0.1145) (increase tolerance)
- solubility: keep within 0.5 of baseline (-3.7941) (decrease tolerance)
- caco2: keep within 0.3 of baseline (-5.2772) (decrease tolerance)
- ppb: keep within 5.0 of baseline (84.0794) (increase tolerance)

The output molecule must satisfy ALL of the following hard constraints:

- SMILES validity: must be parseable and sanitizable by RDKit (no radicals, no disconnected fragments, no ↔ unreasonable valences)
- Molecular weight: < 600 Da
- LogP: between -1 and 5
- TPSA: < 140 Å²
- H-bond donors: <= 5
- H-bond acceptors: <= 10

3465 - Rotatable bonds: <= 10
 3466 - Net formal charge: between -2 and +2
 3467 - Synthetic accessibility: SA score < 4.5
 3468 - PAINS filter: must pass (no pan-assay interference substructures)
 3469 - Structural alerts: must pass Brenk/NIH reactive group filters
 3470 - Tanimoto similarity to reference: >= 0.7 (Morgan fingerprint, radius=2, 2048 bits)
 3471 - Boltz binding probability: affinity_probability_binary > 0.7

3472 The pocket is defined by residues (chain, 1-indexed SEQRES position): A:30, A:31, A:34, A:245, A:246, A:315,
 3473 ↪ A:317, A:318, A:319, A:320, A:321, A:322, A:332, A:340, A:343, A:344, A:345, A:347, A:348, A:351, A:355,
 3474 ↪ A:375, A:379, A:382, A:417, A:418, A:421, A:424, A:475, A:476, A:477, A:480, A:482, A:484, A:486, A:487,
 3475 ↪ A:490, A:491, A:494

3476 Output a single SMILES string to solution.smi.

3477 Your workspace at /workspace contains the task files. Start by examining them.

Elision

3478 *Turns 1–14 and 16 (workspace inspection, reference ligand analysis, baseline ADMET properties prediction, baseline binding affinity*
 3479 *computation using Boltz calls, protein sequence reconstruction, initial candidate generation) omitted.*

Assistant (turn 15)

3481 Now I understand the reference molecule. It's a complex structure with:
 3482 - Ethyl ester (CCOC(=O))
 3483 - Phenyl-ethyl side chain
 3484 - Amide linkage
 3485 - Carboxylic acid
 3486 - Fused bicyclic system (pyrrolidine fused with cyclohexane)

3487 To optimize the properties:
 3488 1. **Decrease clearance**: Typically, reducing lipophilicity, adding polar groups, or reducing metabolic
 3489 ↪ liability can help
 3490 2. **Decrease binding affinity**: This is tricky - we need to make the molecule bind less well (more negative
 3491 ↪ affinity means weaker binding)
 3492 3. **Decrease BBB**: Reduce lipophilicity, add polar groups, reduce LogP

3493 Let me think about modifications:
 3494 - Replace ethyl ester with something more polar (like methyl ester or remove it)
 3495 - Add polar groups to reduce LogP and BBB
 3496 - Modify the phenyl ring to reduce lipophilicity

3497 Let me try some modifications and test them:

Tool Call: predict_admet (turn 15)

Reasoning:

3499 Let me propose some modified structures and test their ADMET properties. I'll start by trying to reduce
 3500 ↪ lipophilicity and clearance by modifying the ester group and the phenyl ring.

Input:

```
3501 {
3502 [...]
3503
3504 # Reference SMILES
3505 ref_smiles = "CCOC(=O)[C@H](CC1CCCC1)O[C@H](C)C(=O)N1[C@H](C(=O)O)C[C@H]2CCCC[C@H]21"
3506
3507 # Proposed modifications:
3508 # 1. Replace ethyl ester with methyl ester (reduce MW, slightly reduce LogP)
3509 # 2. Add hydroxyl to phenyl ring (reduce LogP, BBB, clearance)
3510 # 3. Remove phenyl ring entirely or replace with smaller group
3511
3512 # Test modification 1: methyl ester instead of ethyl
3513 mod1 = "COC(=O)[C@H](CC1CCCC1)O[C@H](C)C(=O)N1[C@H](C(=O)O)C[C@H]2CCCC[C@H]21"
3514
3515 # Test modification 2: add OH to phenyl (para position)
3516 mod2 = "CCOC(=O)[C@H](CC1ccc(O)cc1)O[C@H](C)C(=O)N1[C@H](C(=O)O)C[C@H]2CCCC[C@H]21"
3517
3518 # Test modification 3: both methyl ester and OH on phenyl
3519 mod3 = "COC(=O)[C@H](CC1ccc(O)cc1)O[C@H](C)C(=O)N1[C@H](C(=O)O)C[C@H]2CCCC[C@H]21"
3520
3521 # Test modification 4: replace phenyl with smaller group (methyl)
3522 mod4 = "CCOC(=O)[C@H](CC)O[C@H](C)C(=O)N1[C@H](C(=O)O)C[C@H]2CCCC[C@H]21"
3523
3524 # Test modification 5: remove ester entirely (just H)
3525 mod5 = "O[C@H](CC1CCCC1)O[C@H](C)C(=O)N1[C@H](C(=O)O)C[C@H]2CCCC[C@H]21"
```

```

3520 test_smiles = [mod1, mod2, mod3, mod4, mod5]
3521
3522 for i, smiles in enumerate(test_smiles, 1):
3523     mol = Chem.MolFromSmiles(smiles)
3524     if mol is None:
3525         print(f"Mod {i}: Invalid SMILES")
3526         continue
3527
3528     # Check validity
3529     try:
3530         Chem.SanitizeMol(mol)
3531         print(f"\nMod {i}: {smiles}")
3532         print(f"  MW: {Descriptors.MolWt(mol)}")
3533         print(f"  LogP: {Descriptors.MolLogP(mol)}")
3534         print(f"  TPSA: {rdMolDescriptors.CalcTPSA(mol)}")
3535         print(f"  HBD: {rdMolDescriptors.CalcNumHBD(mol)}")
3536         print(f"  HBA: {rdMolDescriptors.CalcNumHBA(mol)}")
3537         print(f"  RotBonds: {rdMolDescriptors.CalcNumRotatableBonds(mol)}")
3538         print(f"  Charge: {Chem.GetFormalCharge(mol)}")
3539         print(f"  SA: {sascorer.calculateScore(mol)}")
3540     except Exception as e:
3541         print(f"Mod {i}: Error - {e}")
3542 }

```

Tool Output (turn 15)

```

3538 {
3539 Mod 1: COC(=O)[C@H](CCc1cccc1)O[C@@H](C)C(=O)N1[C@H](C(=O)O)C[C@H]2CCCC[C@@H]21
3540 MW: 417.50200000000003
3541 LogP: 2.8102000000000001
3542 TPSA: 93.14
3543 HBD: 1
3544 HBA: 5
3545 RotBonds: 8
3546 Charge: 0
3547 SA: 3.768463126675983
3548
3549 Mod 2: CCOC(=O)[C@H](CCc1ccc(O)cc1)O[C@@H](C)C(=O)N1[C@H](C(=O)O)C[C@H]2CCCC[C@@H]21
3550 MW: 447.528000000000036
3551 LogP: 2.9059000000000001
3552 TPSA: 113.370000000000002
3553 HBD: 2
3554 HBA: 6
3555 RotBonds: 9
3556 Charge: 0
3557 SA: 3.882935729530418
3558
3559 Mod 3: COC(=O)[C@H](CCc1ccc(O)cc1)O[C@@H](C)C(=O)N1[C@H](C(=O)O)C[C@H]2CCCC[C@@H]21
3560 MW: 433.50100000000003
3561 LogP: 2.5158000000000005
3562 TPSA: 113.370000000000002
3563 HBD: 2
3564 HBA: 6
3565 RotBonds: 8
3566 Charge: 0
3567 SA: 3.8811270806988905
3568
3569 Mod 4: CCOC(=O)[C@H](CC)O[C@@H](C)C(=O)N1[C@H](C(=O)O)C[C@H]2CCCC[C@@H]21
3570 MW: 355.43100000000002
3571 LogP: 1.9774999999999994
3572 TPSA: 93.14
3573 HBD: 1
3574 HBA: 5
3575 RotBonds: 7
3576 Charge: 0
3577 SA: 3.8841009312779766
3578
3579 Mod 5: O[C@H](CCc1cccc1)O[C@@H](C)C(=O)N1[C@H](C(=O)O)C[C@H]2CCCC[C@@H]21
3580 MW: 375.46500000000002
3581 LogP: 2.5870000000000006
3582 TPSA: 87.07000000000001
3583 HBD: 2
3584 HBA: 4
3585 RotBonds: 7
3586 Charge: 0
3587 SA: 3.75890866479837
3588 }

```

Assistant (turn 17)

Mod 1 results:

- clearance: 74.87 (need < 72.12) - NOT MET
- bbb: 0.787 (need < 0.66) - NOT MET
- herg: 0.190 (need 0.0145-0.2145) - OK
- solubility: -3.46 (need -4.29 to -3.29) - OK
- caco2: -5.24 (need -5.58 to -4.98) - OK
- ppb: 78.51 (need 79.08-89.08) - NOT MET (too low)

Let me test mod 2 (ethyl ester + OH on phenyl):

Tool Call: predict_admet (turn 17)

Reasoning:

Test ADMET properties for modification 2 (ethyl ester with OH on phenyl) to see if it meets the optimization goals.

Input:

```
{
  "smiles":
    "CCOC(=O)[C@H](CCc1ccc(O)cc1)O[C@@H](C)C(=O)N1[C@H](C(=O)O)C[C@@H]2CCCC[C@@H]21"
}
```

Tool Output (turn 17)

```
{
  "clearance": 87.95, <- FAIL (need < 72.12)
  "caco2": -5.741, <- FAIL (need >= -5.58)
  "bbb": 0.615, <- OK
  "herg": 0.180, <- OK
  "ppb": 84.57, <- OK
  ...
}
```

[ACTION]: Para-OH on phenyl (ethyl ester) is evaluated for the first time. It **fails** on two properties: clearance (87.95 > 72.12) and Caco-2 (-5.74 < -5.58). The agent records the failure at the molecule level but draws no structural inference — it does not conclude that the free para-OH group is likely the metabolic soft-spot driving both failures (glucuronidation/sulfation via Phase II

Assistant (turn 18)

Mod 2 results:

- clearance: 87.95 (need < 72.12) - NOT MET (too high)
- bbb: 0.615 (need < 0.66) - OK!
- herg: 0.180 (need 0.0145-0.2145) - OK
- solubility: -3.91 (need -4.29 to -3.29) - OK
- caco2: -5.74 (need -5.58 to -4.98) - NOT MET (too low)
- ppb: 84.57 (need 79.08-89.08) - OK

Let me test mod 3 (methyl ester + OH on phenyl):

[ERROR — Pattern Not Recognised]: The agent changes only the ester chain length, leaving para-OH intact. A budget-aware agent would have noted that the ester variant is orthogonal to the structural driver (the free phenol) and stopped proposing para-OH.

Tool Call: predict_admet (turn 18)

Reasoning:

Test ADMET properties for modification 3 (methyl ester with OH on phenyl).

Input:

```
{
  "smiles":
    "COC(=O)[C@H](CCc1ccc(O)cc1)O[C@@H](C)C(=O)N1[C@H](C(=O)O)C[C@@H]2CCCC[C@@H]21"
}
```

Tool Output (turn 18)

```
{
  "clearance": 81.65, <- FAIL (need < 72.12) [2nd confirmation]
  "caco2": -5.704, <- FAIL (need >= -5.58) [2nd confirmation]
  "bbb": 0.634, <- OK
  "herg": 0.176, <- OK
  "ppb": 79.93, <- OK
  ...
}
```

[ACTION]: Methyl ester + para-OH reproduces the same failure pattern as Turn 17: **clearance too high, Caco-2 too low**. The ester has changed from ethyl to methyl yet the failure persists. No pattern-level rule is formed or applied; the agent continues to other candidates without writing off the phenol series.

Elision

Turns 19–29 (testing of fluorine, nitrile, dimethylamide, shorter-chain, methoxy, and other variants; none involve para-OH) omitted.

Assistant (turn 30)

[...] mod16 is methyl ester with OH on phenyl --- this should have lower BBB than the previous methoxy variant
 ↳ because the OH is a stronger H-bond donor [...]

[ERROR — Third Wasted Call]: Mod16 is methyl ester + para-OH on phenyl — structurally a minor variant of the molecule tested in Turn 18 (methyl ester + para-OH). The agent has now spent three of its fifteen ADMET calls on para-OH variants, each time receiving clearance ≈ 82 –88 and Caco-2 ≈ -5.7 . No structural exclusion rule has been formed or applied.

Tool Call: predict_admet (turn 30)

Reasoning:

Test ADMET for mod16 (methyl ester with OH on phenyl) --- this should have lower BBB than methoxy.

Input:

```
{
  "smiles":
    "COC(=O)[C@H](CCc1ccc(O)cc1)O[C@@H](C)C(=O)N1[C@H](C(=O)O)C[C@H]2CCCC[C@H]21"
}
```

Tool Output (turn 30)

```
{
  "clearance": 82.54, <- FAIL (need < 72.12) [3rd confirmation]
  "caco2": -5.687, <- FAIL (need >= -5.58) [3rd confirmation]
  "bbb": 0.649, <- OK
  "herg": 0.174, <- OK
  "ppb": 80.76, <- OK
  ...
}
```

[ACTION]: Mod16 is methyl ester + para-OH on phenyl — structurally identical to the compound tested in Turn 18. The agent has now spent three of its fifteen ADMET budget calls on para-OH variants and received the same failure pattern each time: clearance ≈ 82 –88, Caco-2 ≈ -5.7 . The agent subsequently exhausts its budget at Turn 35 and submits CCOC(=O)[C@H](CCc1ccc(OC)cc1)... (the methoxy variant), which passes hard constraints but fails the hERG hold-constant criterion ($0.238 > 0.2145$) and Boltz binding probability ($0.675 < 0.7$).

I.2. Repeated Proposal of a Failed Candidate Across Tool Invocations

Task ID: smdd_003_3WV3_0 **Model:** Claude Sonnet 4.6 **Run:** claude_sonnet_46_med

Evaluation result: **FAILED.**

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

Metric	Status	Result	Condition
Validity	passed	True	{"eq": true}
Tanimoto Similarity	passed	0.455	{"lt": 0.5}
Scaffold MCS	passed	0.467	{"lt": 0.65}
Boltz Interaction Similarity	failed	0.667	{"gt": 0.75}
Boltz Binding Probability	passed	0.733	{"gt": 0.7}

Summary of error. The agent introduces a tetrahydroquinoline-hydroxamic acid scaffold (THQ_hydroxamic) in Turn 8, tests it with Boltz in Turn 16, and receives a binding probability of 0.614 — below the required threshold of 0.70. No explicit record of this failure is retained. In Turn 20, the identical SMILES string is re-proposed under a new alias (Indol2_para), with only a parenthetical code comment acknowledging the duplication. The molecule is again included in the candidate list for further Boltz evaluation.

Trajectory.

Showing 3 of the agent's turns; turns 1–7, 9–15, and 17–19 are elided for space.

System Prompt

You are a computational medicinal chemist assistant. You are given a drug discovery task and must use your tools to solve it.

Available tools:

- python(code): Execute Python code. [...]
- predict_boltz(...): Run Boltz structure prediction. Budget: 8 calls.
- predict_admet(smiles): Predict ADMET properties. Budget: 15 calls.
- submit(file_path): Submit your solution file.

Elision

Turns 1–7 (workspace inspection, reference ligand analysis, baseline ADMET prediction, initial candidate generation, Tanimoto and scaffold MCS filtering) omitted.

Assistant (turn 8)

[...] I'll now check scaffold similarity and Morgan Tanimoto of the candidate molecules to see which ones are sufficiently divergent from the reference scaffold while still drug-like [...]

Tool Call: python (turn 8)

Reasoning:

Check scaffold similarity and Morgan Tanimoto of candidate molecules.

Input:

```

candidates_to_test = [
    ("THQ_hydroxamic", "O=C(NO)C1CCN(Cc2ccccc2)c2ccc(OCc3ccccc3)cc21"),
    ...
]
    
```

Output:

```

Final candidate verification:
THQ_hydroxamic  2.63  388.5  0.467  0.412  True
SMILES: O=C(NO)C1CCN(Cc2ccccc2)c2ccc(OCc3ccccc3)cc21
    
```

[ACTION]: The molecule is introduced for the first time as THQ_hydroxamic, a tetrahydroquinoline scaffold hop with a hydroxamic acid zinc-binding group. It passes both the MCS/Larger ratio threshold (0.467 < 0.65) and the Tanimoto threshold (0.412 < 0.50), and is flagged as a valid candidate for Boltz evaluation.

Elision

Turns 9–15 (additional candidate filtering, ADMET predictions, and Boltz call setup) omitted.

3740 **Assistant (turn 16)**

3741 [...] Let me test THQ_hydroxamic with the correctly formatted pocket residues to get the binding probability.

3742 ↪ [...]

3743 -----

3744

3745 **Tool Call: predict_boltz (turn 16)**

3746 **Reasoning:**

3747 Testing THQ_hydroxamic with correctly formatted pocket residues.

3748

3749 **Input:**

3750

3751 { "smiles": "O=C(NO)C1CCN(Cc2ccccc2) c2ccc (OCc3ccccc3) cc21",

3752 ...

3753 }

3754 **Result:**

3755

3756 { "binding_probability": 0.6138086318969727,

3757 ...

3758 }

3759 **[ACTION]:** THQ_hydroxamic is submitted to the Boltz structure predictor. The returned binding probability is 0.614, which is

3760 **below the required threshold of 0.70.** The candidate is implicitly discarded, but no explicit record of this failure is stored or referenced

3761 in subsequent reasoning.

3762

3763

3764

3765 **Elision**

3766 *Turns 17–19 (analysis of Boltz result, generation of further candidate variants) omitted.*

3767

3768 **Assistant (turn 20)**

3769 [...] Let me run Boltz on the most promising candidates from the current pool. [...]

3770 -----

3771 **[ERROR — Candidate Recycled Under New Alias]:** The agent has no memory that THQ_hydroxamic already received

3772 a binding probability of 0.614 in Turn 16 and was therefore disqualified. It re-proposes the identical SMILES under the alias

3773 Indol2_para.

3774

3775 **Tool Call: python (turn 20)**

3776 **Reasoning:**

3777 Let me run Boltz on the most promising candidates.

3778

3779 **Input:**

3780 candidates = [

3781 ...

3782 ("Indol2_para", "O=C(NO)C1CCN(Cc2ccccc2) c2ccc (OCc3ccccc3) cc21")

3783 ...

3784]

3785 **[ACTION]:** **The identical SMILES string is re-proposed under the new alias Indol2_para.** The agent possesses no persistent

3786 memory of tested molecules and their outcomes across Python tool invocations, so a molecule that already **failed** the binding-probability

3787 threshold ($0.614 < 0.70$) in Turn 16 can be silently recycled, wasting one of the limited Boltz budget calls.

3788

3789

3790 **I.3. Recall of Known Molecules as a Substitute for De Novo Design**

3791 **Task ID:** smdd_003_3WV3_0 **Model:** Gemini 3.1 Pro **Run:** gemini_31_pro_med_full

3792

3793 **Evaluation result: FAILED.**

3794

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

Metric	Status	Result	Condition
Validity	passed	True	{"eq": true}
Tanimoto Similarity	passed	0.2923076923076923	{"lt": 0.5}
Scaffold MCS	passed	0.3	{"lt": 0.65}
Boltz Interaction Similarity	failed	0.45454545454545453	{"gt": 0.75}
Boltz Binding Probability	passed	0.9476999938488007	{"gt": 0.7}

Summary of error. The agent bypassed de novo molecular design entirely. At Turn 5 it surfaced a known PubChem compound (O=C(NO)C(Cc1ccccc1)NS(=O)(=O)c2ccc(c3ccccc3)cc2) from memory, framed it as a design idea, and used one Boltz call to confirm it exceeded all thresholds. No alternative scaffolds were generated or evaluated. Of 8 Boltz calls and 15 ADMET calls budgeted, only 1 was consumed.

Trajectory.

Showing 4 of 8 agent turns; the first 4 turns are elided for space.

System Prompt

You are a computational medicinal chemist assistant. You are given a drug discovery task and must use your tools ↪ to solve it.

Available tools:

- python(code): Execute Python code. [...]
- predict_boltz(smiles, protein_sequences, pocket_residues, ...): Run Boltz structure prediction. Returns ↪ binding probability and affinity. Budget: 8 calls.
- predict_admet(smiles): Predict ADMET properties. Budget: 15 calls.
- submit(file_path): Submit your solution file.

You have no internet access. Work only with the files provided in /workspace.

User (turn 0)

Task

You are given a protein-ligand complex (complex.pdb). The target ligand is on chain L. A 3D reference structure ↪ of this ligand is provided in reference.sdf. Your task is to design a new molecule with a significantly ↪ different scaffold (the MCS atom count between the Bemis-Murcko scaffolds of your molecule and the ↪ reference, divided by the larger scaffold heavy atom count, must be < 0.65; additionally, overall Morgan ↪ fingerprint Tanimoto similarity must be < 0.5) to the ligand that maintains the same binding interactions ↪ and has binding probability > 0.7 as calculated by Boltz. The pocket is defined by residues (chain, ↪ 1-indexed SEQRES position): A:82, A:83, A:115, A:116, A:117, A:119, A:120, A:129, A:134, A:135, A:136, ↪ A:137, A:138, A:139, A:140, A:141, A:142, A:143, A:146, A:149, A:150, A:151, A:152, A:155. Output a single ↪ SMILES string to a file titled `solution.smi`.

elision

Turns 1-4 (workspace inspection, extractions of reference ligand SMILES and protein sequence, scaffold analysis) omitted.

Assistant (turn 5)

[...] The existing scaffold looks like it's built around a piperidine core, a benzyl group, and a ↪ 4-benzyloxy-phenylsulfonyl group. To stay active against the MMP, we need to preserve that zinc-binding ↪ interaction, likely keeping the hydroxamate group (^C(=O)NO^).

[...] What if we shifted to a simpler, more flexible acyclic backbone, like an alpha-amino acid derivative? We ↪ could retain the key pharmacophores--the sulfonamide, the deep pocket binder, and the hydroxamate--but ↪ change the central connective tissue entirely.

Let's test an idea: an amino acid-based sulfonamide. Something like ↪ O=C(NO)C(Cc1ccccc1)NS(=O)(=O)c2ccc(c3ccccc3)cc2. The scaffold for this would be ↪ ^O=S(=O)(c1ccc(c2ccccc2)cc1)NC(Cc1ccccc1)C^. Let's calculate the scaffold and Morgan similarity of this ↪ proposed structure against our reference to see if we're on the right track.

[ERROR — Direct Recall]: The agent comes up with O=C(NO)C(Cc1ccccc1)NS(=O)(=O)c2ccc(c3ccccc3)cc2 without any iterative design process. This molecule is a known PubChem compound retrieved directly from memory.

3850 **Tool Call: python (turn 5)****Reasoning:**

Test a proposed simple amino acid derived scaffold against the constraints.

Code:

```
[...]

ref_smiles = 'O=C(NO)C1(S(=O)(=O)c2ccc(OCC3CCCC3)cc2)CCN(Cc2CCCC2)CC1'
ref_mol = Chem.MolFromSmiles(ref_smiles)
ref_scaffold = MurckoScaffold.GetScaffoldForMol(ref_mol)

test_smiles = 'O=C(NO)C(Cc1cccc1)NS(=O)(=O)c2ccc(c3cccc3)cc2'
test_mol = Chem.MolFromSmiles(test_smiles)
test_scaffold = MurckoScaffold.GetScaffoldForMol(test_mol)

# Calculate scaffold similarity
res = rdFMCS.FindMCS([ref_scaffold, test_scaffold], completeRingsOnly=True)
mcs_atoms = res.numAtoms
max_atoms = max(ref_scaffold.GetNumHeavyAtoms(), test_scaffold.GetNumHeavyAtoms())
scaffold_sim = mcs_atoms / max_atoms

# Calculate Morgan similarity
fp1 = AllChem.GetMorganFingerprintAsBitVect(ref_mol, 2, nBits=2048)
fp2 = AllChem.GetMorganFingerprintAsBitVect(test_mol, 2, nBits=2048)
morgan_sim = DataStructs.TanimotoSimilarity(fp1, fp2)

[...]
```

3870 **Tool Output (turn 5)**

```
Scaffold max atoms: 30
MCS atoms: 9
Scaffold similarity: 0.3 <-- passes < 0.65 threshold
Morgan similarity: 0.2923... <-- passes < 0.5 threshold
```

3876 **Assistant (turn 6)**

The initial tests look promising. The scaffold similarity is well below the 0.65 threshold at 0.3, and the Morgan similarity is comfortably under 0.5 at roughly 0.29. [...] Now we need to
 ↪ determine if this new structure will actually bind to the MMP pocket. [...] I'll pass the SMILES for this
 ↪ test compound (O=C(NO)C(Cc1cccc1)NS(=O)(=O)c2ccc(c3cccc3)cc2) to the Boltz prediction tool [...]

[ERROR — Premature Commitment]: Having confirmed that the molecule clears threshold tests, the agent sends it to binding prediction without additional reasoning or considering any alternative candidates. No de novo design occurs.

3884 **Tool Call: predict_boltz (turn 6)****Reasoning:**

Predict binding of the proposed amino acid derived MMP inhibitor.

Input (abbreviated):

```
{
  "smiles": "O=C(NO)C(Cc1cccc1)NS(=O)(=O)c2ccc(c3cccc3)cc2",
  "protein_sequences": [{"sequence": "YNVFPTLKWSK[...]DEDPN", "id": "A"}],
  ...
}
```

3894 **Tool Output (turn 6)**

```
{
  "binding_probability": 0.9568, <-- well above > 0.7 threshold
  "affinity_pred_value": -2.007,
  ...
}
```

3900 **Assistant (turn 7)**

Okay, wow, a binding probability of roughly 0.95! That is a very strong result.
 [...] Everything is perfectly aligned.

SMDD-Bench: Can LLMs Solve Real-World Small Molecule Drug Design Tasks?

- 3905 1. Binding Probability: ~0.95 (Requirement: > 0.7). Check.
3906 2. Scaffold Similarity: [...] 9/30 = 0.3. (Requirement: < 0.65). Check.
3907 3. Morgan Similarity: 0.29. (Requirement: < 0.5). Check.

3908 [...] I'm going to proceed with the SMILES string I've generated.
3909

3910 **[ERROR — No Exploration]:** The agent treats the first passing molecule as the final solution. The entire “design” process
3911 consumed a single Boltz call on a molecule recalled from memory.
3912

3913 **Tool Call: python** (turn 7)

Reasoning:

3914 Write the successful SMILES to solution.smi.
3915

Code:

3916
3917
3918 with open('solution.smi', 'w') as f:
3919 f.write('O=C(NO)C(Cc1ccccc1)NS(=O)(=O)c2ccc(c3ccccc3)cc2\n')

3920 **Tool Call: submit** (turn 8)

3921 **file_path:** solution.smi
3922

3923 **Tool Output** (turn 8)

3924 Submission accepted: solution.smi
3925

3926 **Submit** (turn 8)

3927 **Submitted file:** solution.smi (O=C(NO)C(Cc1ccccc1)NS(=O)(=O)c2ccc(c3ccccc3)cc2).
3928
3929
3930
3931
3932
3933
3934
3935
3936
3937
3938
3939
3940
3941
3942
3943
3944
3945
3946
3947
3948
3949
3950
3951
3952
3953
3954
3955
3956
3957
3958
3959