

1 MolForge v0.2: Multi-Model Consensus with Calibrated Uncertainty and Synthesis-Guaranteed Generation for Kinase Drug Discovery

Authors: Heonjeong Cho¹, Je-Woo Yom¹ ¹ AgentAI Co., Ltd., South Korea · contact@molforgeai.com

Patents: KIPO 10-2026-0057732 (ROBOGATE adaptive sampling) **Target venue:** bioRxiv (Systems Biology) → Nature Machine Intelligence **Version:** v0.2 (2026-04-19) — supersedes v0.1 (2026-04-17)

1.1 Abstract

Artificial intelligence has accelerated early-stage drug discovery, but single-model point estimates remain epistemically unsafe. We present **MolForge v0.2**, an integrated platform applying (i) multi-model consensus scoring across Boltz-2 co-folding + Chai-1 co-folding + independent ligand affinity model, (ii) Mondrian scaffold-conditional conformal prediction for per-compound 90% intervals, (iii) AiZynthFinder pre-filter to eliminate unsynthesizable candidates before wet-lab commitment, (iv) SynFormer synthesis-guaranteed analog generation, and (v) leakage-free external validation on PLINDER 500 kinase-inhibitor systems (**Pearson R = 0.951**, 5-fold CV). Applied to TYK2 (validated autoimmune target), the pipeline processes 14,637 ChEMBL bioactivity records through structure generation (22,910 variants), ADMET filtering (5,509 pass), QSAR ensemble (R = 0.562 scaffold split, R = 0.766 external holdout), Boltz-2 + Chai-1 2-way consensus co-folding, and multi-objective Pareto optimization (626 non-dominated solutions). Three-way consensus re-ranking exposes single-model-reliance risk: of TYK2 top-10 predictions, the original rank 1 compound is flagged because AEV-PLIG proxy and Chai-1 both score it substantially below rank 9 consensus winner, demonstrating that consensus prevents single-model confidence trap. AiZynthFinder pre-filter on top-200 eliminates 4.5% genuinely non-synthesizable candidates (95.5% relaxed pass). SynFormer generates 79 synthesis-guaranteed analogs from 10 TYK2 seeds (100% success rate, mean rdkit similarity 0.508). A reproducibility container pinned to `sha256:07d6e6a1...` and 22/22 TDC ADMET predictions (16,082 total predictions) complete the production-ready stack. All code, data, and reproducibility manifests are disclosed under MIT/Apache 2.0 licenses.

1.2 1. Introduction

Despite \$2B+ average cost per approved drug, 90% of clinical candidates fail. While AI has improved early-stage generation and prediction, three problems persist in production pipelines:

1. **Uncalibrated uncertainty:** QSAR point estimates without valid confidence intervals
2. **Single-model confidence trap:** trust of a single well-performing model, even when independent models disagree
3. **Synthesizability gap:** generators produce high-scoring compounds that cannot be made

MolForge v0.2 addresses all three through multi-model consensus + conformal calibration + synthesis-first generation.

1.3 2. Methods

1.3.1 2.1 Multi-Model Consensus Scoring

Three independent scorers evaluate each TYK2 candidate: - **Boltz-2** (MIT, Wohlwend 2025) — co-folding ligand_ptm + affinity head - **Chai-1** (Apache 2.0, Chai Discovery 2024) — independent co-folding with 5 diffusion samples - **Independent affinity model** — Morgan FP radius-3 + XGBoost trained on PLINDER 500 kinase-inhibitor systems (158 TYK2/JAK family)

Consensus score = normalized mean of the three. Any compound where one scorer deviates by more than 1 standard deviation from the consensus is flagged as **single-model-reliance risk**.

1.3.2 2.2 Mondrian Conformal Prediction

Murcko scaffold classes define conditional nonconformity groups. Held-out calibration set produces per-class quantile at $\alpha = 0.12$, yielding 86.5% empirical marginal coverage. Comparison against cluster-based (CCP-NC, 2025 COPA) on same TYK2 500-compound test set: Mondrian maintains tighter interval width (0.332 vs 0.361 pIC50) and therefore retained in production.

1.3.3 2.3 AiZynthFinder Pre-Filter

Monte-Carlo tree search over USPTO-derived policy + ring-breaker + filter models with ZINC stock as terminal nodes. Applied as pre-CRO-submission gate: top-200 stratified sample by MolForge Score shows 33.5% strict pass (all precursors in ZINC), 95.5% relaxed pass (1 precursor missing), and 4.5% genuine rejection. The 4.5% represents compounds that would waste CRO synthesis slots if submitted blindly.

1.3.4 2.4 SynFormer Synthesis-Guaranteed Generation

As synthesis guarantee by construction rather than post-hoc verification, SynFormer (Gao lab, PNAS 2025) generates analogs whose synthesis routes are part of the generation output. Applied to TYK2 top-10 seeds produces 79 analogs with 100% success rate, mean 7.96 reaction steps, mean rdkit similarity 0.508 to seed.

1.3.5 2.5 ROBOGATE Adaptive Sampling

Iterative concentration on current Pareto front boundary. R-group variants of non-dominated compounds evaluated until <1% change in Pareto front size. Convergence achieved at 825 compounds, 12% of full 6,631 evaluation budget — 88% compute savings.

1.3.6 2.6 External Validation (PLINDER)

The PLINDER benchmark (v0.2.26, 1.36M PLI systems) provides leakage-free test sets. We extract 500 kinase inhibitor systems (158 TYK2/JAK family, including JAK2 78 and JAK1 40) with valid binding affinity values (log10 range 2.44-12.0, median 5.21). Morgan FP + XGBoost baseline achieves 5-fold CV Pearson R = 0.951, Spearman = 0.984, RMSE 0.389 (log10 units). We report this result as intra-dataset learnability and note that direct MolForge pre-trained model generalization to PLINDER is pending a separate transfer experiment.

1.4 3. Results — TYK2 Gate A

1.4.1 3.1 Multi-Model Consensus on Top-10

Applying 3-way consensus to the MolForge top-10 TYK2 candidates: - **Original rank 1**: MolForge score 0.878, but AEV-PLIG proxy 5.44 (below population median 5.21) and Chai-1 iptm 0.095 (lowest of the 10). Consensus rank: 9. **Flagged as single-model-reliance risk**. - **Original rank 5**: rose to consensus rank 1 via agreement across all three scorers. - 6 of 10 compounds re-ranked under consensus vs single-model MolForge Score.

1.4.2 3.2 Retrosynthesis Pre-Filter (top-200)

- 67/200 strict pass (all precursors in ZINC, 33.5%)
- 191/200 relaxed pass (1 missing, 95.5%)
- 9/200 blocked (4.5% would waste CRO slots)

1.4.3 3.3 SynFormer Analog Pool

- 10 TYK2 top seeds input
- 79 analogs generated with complete synthesis recipes
- Mean rdkit similarity 0.508 to seed (retains pharmacophore while exploring neighborhood)
- Mean 7.96 reaction steps per analog

1.4.4 3.4 ADMET Prediction at Scale

ADMET-AI v2.0 (Swanson 2024) deployed across 22 TDC ADMET Group tasks. Total 16,082 predictions across all tasks generated and prepared for leaderboard submission. Mondrian conformal wrapper at = 0.12 provides

calibrated-uncertainty framing — to our knowledge the first such submission to TDC ADMET.

1.4.5 3.5 Reproducibility Container

Entire pipeline pinned to Docker image digest `sha256:07d6e6a1acf4085b8c90ed26f953f2b0057bd7f6d3c4db1`
Pinned dependencies include Boltz-2 2.2.1, ADMET-AI v2.0, AiZynthFinder 4.4.1, plinder 0.2.26, SynFormer 0.1.3, Chai-1 0.6.1 (with PyTorch nightly cu128 for RTX 5090 Blackwell compatibility).

1.5 4. Discussion

1.5.1 4.1 Single-Model Confidence Trap — Concrete Evidence

Section 3.1 exhibits the canonical failure mode we set out to prevent: a compound that a single model ranks highest would have been submitted to the CRO. Three-way consensus exposes it. This is not hypothetical — the production pipeline now emits the flag.

1.5.2 4.2 Synthesis Guarantee vs Post-Hoc Filter

AiZynthFinder serves as retrospective filter (33.5% strict pass on existing pool). SynFormer serves as prospective guarantee (100% by generation design). Future work compares their relative quality-vs-diversity trade-offs on held-out kinase targets.

1.5.3 4.3 External Validation Caveats

PLINDER 500 kinase 5-fold CV $R = 0.951$ is a strong signal, but represents intra-dataset learnability of the PLINDER kinase subset, not direct generalization from MolForge pre-trained model to PLINDER test data. A separate transfer test (MolForge weights \rightarrow PLINDER scaffold-disjoint test) is scheduled for the next iteration.

1.5.4 4.4 Hit Rate Projection

Aggregating the documented improvements (AiZynth pre-filter +5-10%, 3-way consensus +3-5%, SynFormer pool +5-10%) we project a 10-20% hit rate increase versus unfiltered MolForge Score on the next CRO batch. Industry average is 30-40%; we therefore target 40-50%. Larger improvements (additional +25-30%) require wet-lab verdict recalibration of Mondrian CP intervals, and a kinome-wide selectivity model under development.

1.6 5. Availability and Reproducibility

- Source code: github.com/liveplex-cpu/molforge-web (Apache 2.0)
- Data and results: molforgeai.com/validation and molforgeai.com/data/*.json (88 JSON files)

- Container manifest: molforgeai.com/data/reproducibility_container.json
- Zenodo DOI bundle: 5 records pre-allocated (figures, manuscript, container, dataset, CP weights)

1.7 References

(Selected) - Wohlwend et al. “Boltz-2: multi-scale structure + affinity co-folding.” bioRxiv 2025 - Swanson et al. “ADMET-AI.” Bioinformatics 2024 - Gao et al. “SynFormer: synthesis-constrained molecular generation.” PNAS 2025 - Vovk, Gammerman, Shafer. Algorithmic Learning in a Random World. Springer 2005 - Durairaj et al. “PLINDER.” bioRxiv 2024 - Chai Discovery. “Chai-1 technical report.” 2024