

Literature Survey: Conformal Prediction and ADMET in Toxicity Modeling

Rahul Goel (2022388)
Swapnil Verma (2022523)
Rishit Gupta (2022406)
Pooja (2022356)
Sanyam Garg (2022448)
Dharmender (2022158)

Computer Science and Bioscience Department,
Indraprastha Institute of Information Technology, Delhi

1 Zhang et al. (2021) – *Deep Learning–Based Conformal Prediction of Toxicity*

About the Paper

Zhang et al. present one of the first deep-learning conformal frameworks for molecular toxicity, combining per-endpoint feedforward neural networks (FFNNs) with inductive conformal calibration. They split the Tox21 dataset into a fitting set (to train each FFNN), a calibration set (to compute non-conformity scores), and a held-out test set. Nonconformity is defined as one minus the softmax probability assigned to the true label, and for each endpoint and error rate α , they select the quantile threshold on these calibration scores that guarantees coverage $\geq 1 - \alpha$.

Beyond the core methodology, they carefully analyze the trade-off between average prediction-set size and coverage. By demonstrating $\sim 90\%$ coverage at $\alpha = 0.10$ with average set sizes of ~ 1.8 labels per compound, they establish that conformal sets can remain small yet valid. They also explore how conformal guarantees hold under various model architectures and data splits, providing a blueprint for rigorous uncertainty quantification in toxicity prediction.

Key Results & Application

Their main quantitative finding is stable endpoint coverage between 0.89 and 0.92 at $\alpha = 0.10$ across all 12 Tox21 assays, showing that roughly 9 out

of 10 true labels lie within the model’s conformal sets. Average set sizes near 1.8 demonstrate that the model rarely resorts to trivial two-label sets, preserving informativeness.

At the compound level, they apply an OR-rule—flagging a compound toxic if any endpoint’s set contains the “active” label—and report $\sim 60\%$ accuracy. Although this amplifies both true positives and false positives, it provides an overall metric of pipeline performance when decisions must consider all toxicity assays simultaneously.

Usage in Code

In our Phase I implementation we followed Zhang et al.’s splitting strategy exactly: we randomly allocate 70% of the training data to model fitting, 15% to calibration, and retain 15% as a test set. After training each per-endpoint FFNN, we compute nonconformity scores on the calibration split, sort them, and select the $\lceil (1-\alpha)(n_{\text{cal}}+1) \rceil$ -th largest score as the threshold.

Our code modularizes their routine into a `CalibrationRunner` class, which accepts any PyTorch model and returns α -indexed threshold maps. We also vectorized the nonconformity computation to process all calibration examples in one pass, speeding up threshold estimation by $\sim 5\times$ compared to a Python loop.

How the Paper Helped Us

Zhang et al.’s clear separation of fitting, calibration, and testing guided our entire project structure. Their formal argument for coverage guarantees underpins every subsequent calibration phase—ensuring that when we introduce dropout, ensembles, or quantile heads, we maintain the same rigorous framework.

Moreover, their reported metrics serve as a constant benchmark: every time our Phase II–V results deviate, we investigate whether preprocessing, feature selection, or model tweaks broke the conformal assumptions. Their sensitivity analyses on assay imbalance motivated us to include mask-aware loss functions and to monitor per-endpoint coverage, not just the average.

2 Swanson et al. (2024) – *ADMET-AI: A Machine Learning ADMET Platform*

About the Paper

Swanson et al. introduce **ADMET-AI**, a multi-task deep network that simultaneously predicts absorption, distribution, metabolism, excretion, and toxicity endpoints using a shared FFNN backbone and task-specific heads. They train on diverse public ADMET datasets (e.g., Caco-2 permeability,

microsomal clearance) alongside toxicity assays, demonstrating that shared molecular embeddings capture chemical features beneficial across related tasks. Their network includes descriptor inputs (logP, TPSA, MW) concatenated with molecular fingerprints, and they apply weighted cross-entropy losses to balance tasks.

A key innovation is dynamic task weighting: during training, each task’s loss weight is adjusted based on its validation-set gradient magnitude, preventing large tasks from dominating the shared representation. This leads to robust performance even on rare or small datasets (< 500 compounds), where single-task models typically underperform.

Key Results & Application

On five pharmacokinetic tasks, ADMET-AI achieves AUROC gains of 3–5% over single-task baselines, with the largest improvements on low-data tasks (up to 10% gain). For toxicity endpoints, they observe 2–3% AUROC increases and slightly tighter calibration errors (ECE reductions of 0.01–0.02). These results underscore the synergy between pharmacokinetic and toxicity predictions: descriptors that correlate with solubility or clearance also inform toxicity risk.

They further show that learned embeddings cluster chemical series with similar ADMET profiles, visually demonstrating that multi-task training yields chemically meaningful latent spaces. This clustering supports downstream tasks such as lead-optimization, where chemists can navigate embedding space to find molecules with balanced ADMET and safety profiles.

Usage in Code

In Phase II we replicated their multi-task architecture by concatenating the same 50 physicochemical descriptors with 2048-bit Morgan fingerprints, passing them through a shared three-layer FFNN (512→256→128 neurons) and then branching into 12 classification heads. We implemented dynamic task weighting as in ADMET-AI: after each epoch, we compute per-head gradient norms on a validation subset and adjust the loss weights for the next epoch.

This dynamic weighting stabilized training: without it, our multi-task model sometimes overfit large endpoints and underfit rare active endpoints. With weighting, we saw a 4% lift in OR-rule compound accuracy (0.902 vs. 0.866) and a 0.5% improvement in average endpoint coverage stability across α levels.

How the Paper Helped Us

Swanson et al. convinced us that parameter sharing is not only computationally efficient but also statistically beneficial, especially when coupling

toxicity with pharmacokinetic descriptors. Their dynamic task-weighting scheme inspired our own adaptation in Phase II, which directly translated to higher compound-level accuracy under the OR rule.

Beyond code, their embedding-space visualizations shaped our interpretation modules: we borrowed their t-SNE/UMAP pipelines to inspect how our conformal intervals vary across chemical clusters, identifying regions where multi-task learning yields the tightest intervals.

3 Yan et al. (2024) – *MolToxPred: Small Molecule Toxicity Prediction by Ensemble ML*

About the Paper

Yan et al. propose **MolToxPred**, an ensemble of heterogeneous learners—gradient-boosted decision trees (GBDT), support vector machines (SVM), and shallow neural networks—for small-molecule toxicity classification. Each model type offers complementary strengths: GBDTs excel on sparse fingerprint features, SVMs handle small-sample regimes well, and NNs capture non-linear interactions. They fuse predictions via a trainable stacking layer, calibrate the ensemble with isotonic regression, and demonstrate robust performance across internal and external test sets.

Crucially, they show that ensemble diversity—quantified by pairwise disagreement metrics—correlates strongly with improved calibration (lower ECE) and higher AUROC. This systematic analysis of model heterogeneity provides a roadmap for assembling maximally effective ensembles.

Key Results & Application

MolToxPred achieves an average AUROC of 0.92 on the 12 Tox21 assays, significantly outperforming single-model baselines ($p < 0.01$). Calibration errors (ECE) drop below 0.05 for all endpoints, ensuring that predicted probabilities align closely with observed frequencies—a prerequisite for reliable conformal calibration. They also demonstrate that ensemble stacking improves out-of-distribution robustness, reducing error increases from +7% to +2% on an external drug-like compound set.

Their ablation studies reveal that removing any one model type degrades performance, highlighting the synergy of heterogeneous learners. They further benchmark different stacking strategies (e.g. logistic regression vs. neural meta-learners), finding that simple logistic stacking often suffices, reducing computational overhead.

Usage in Code

For Phase IV, we adapted the MolToxPred stacking concept by training *five independent* FFNNs (rather than heterogeneous models) with different seeds and hyperparameters. We average their temperature-scaled softmax outputs and then apply conformal thresholds to the ensemble probabilities. Temperature scaling is optimized on a held-out validation split via minimizing negative log-likelihood, following their isotonic-calibration rationale.

Although we did not implement full stacking, we measured ensemble diversity (via disagreement score) and confirmed that our seeded-FFNN ensemble achieved diversity metrics comparable to MolToxPred’s GBDT+SVM+NN ensemble. This guided our choice of five models as the sweet spot between calibration gains and compute cost.

How the Paper Helped Us

Yan et al.’s systematic exploration of ensemble heterogeneity underscored the importance of diversity for calibration. While we lacked the compute to train GBDTs and SVMs in parallel, their insights motivated our seeded-FFNN approach and our experiments with varied architectures.

Moreover, their use of stacking layers inspired our temperature-scaling step before calibration: by aligning ensemble confidences with real-world frequencies, we could derive tighter conformal thresholds (3% narrower intervals) without sacrificing coverage.

4 Guo et al. (2023) – *Review of ML and DL Models for Toxicity Prediction*

About the Paper

Guo et al. survey a decade of toxicity-prediction methods, from classical QSAR and random forests to modern graph neural networks (GNNs) and deep FFNNs. They discuss feature engineering (fingerprints, descriptors), handling of missing data, class-imbalance remedies, and evaluation metrics (AUROC, ECE, AUPRC). A key takeaway is that while GNNs often outperform FFNNs on very large datasets, FFNNs coupled with rigorous preprocessing and data augmentation remain highly competitive, especially in low-resource settings.

They also dedicate a section to *uncertainty estimation*, covering Bayesian NN approximations (e.g. MC-Dropout), deep ensembles, and conformal methods—highlighting best practices and pitfalls in calibration, such as over-reliance on softmax confidences.

Key Results & Application

Although primarily a review, Guo et al. compile quantitative benchmarks showing that descriptor scaling and SMILES augmentation can yield 2–3% absolute accuracy gains and reduce calibration error by 0.01–0.02. They also note that variance-threshold feature selection removes noisy descriptors ($\leq 1\%$ variance) with no performance loss, simplifying models and reducing overfitting risk.

Their meta-analysis reveals that MC-Dropout and deep ensembles produce comparable uncertainty estimates, but ensembles tend to be more stable at a higher computational cost. This informs trade-off decisions for pipelines that must run on limited hardware.

Usage in Code

We incorporated nearly all of Guo et al.’s preprocessing recommendations:

- **Descriptor Scaling & Normalization:** Mean-centering and unit-variance scaling of all 50 physicochemical descriptors, implemented via scikit-learn’s `StandardScaler`.
- **Variance Thresholding:** Automatic removal of descriptors with variance ≤ 0.01 , reducing the descriptor set from 50→40 without re-training.
- **SMILES Augmentation:** On-the-fly generation of randomized SMILES strings during training epochs using RDKit’s `MolToRandomSmiles`.

How the Paper Helped Us

Guo et al. served as our **methodological checklist**, validating each preprocessing, augmentation, and calibration choice before we wrote a single line of code. Their review convinced us to prune low-variance features, to track ECE during training, and to implement both MC-Dropout and ensembles so we could compare cost–benefit trade-offs empirically.

Their critique of GNN vs. FFNN motivated our decision to prioritize FFNNs—enabling faster calibration phases—while still experimenting with a slimmed-down message-passing network that we ultimately shelved due to marginal gains.

5 Sharma et al. (2023) – *Clinical Toxicity Prediction Using Multi-Task Nets and Contrastive Molecular Explanations*

About the Paper

Sharma et al. combine a multi-task DNN for clinical toxicity with a **contrastive explanation module** that generates counterfactual molecules. They train a masked-graph autoencoder (GAE) alongside the classifier, learning latent representations that can be perturbed to flip toxicity predictions. Their framework not only predicts toxicity but also recommends minimal structural modifications to reduce risk, bridging prediction and interpretability in medicinal chemistry.

They report a two-step pipeline: first, compute conformal or softmax-based uncertainty sets; second, for ambiguous cases, mask specific atoms/bonds and decode new SMILES via the GAE’s generative head, yielding actionable molecular suggestions.

Key Results & Application

On five clinical toxicity assays, their multi-task model achieves 85% OR-rule accuracy—substantially higher than single-task FFNNs. Contrastive explanations generate chemically valid counterfactuals 70% of the time, with an average of 2–3 atom changes needed to flip the prediction, demonstrating both efficacy and interpretability.

They further analyze the fidelity and sparsity of explanations, showing that contrastive masks highlight key functional groups (e.g. removing a nitro group flips toxicity), offering direct guidance for lead-optimization.

Usage in Code

We ported their masked-graph autoencoder into our `explanations/` module, training it on the same Tox21 graphs. During Phase V, for any compound whose conformal set contains both labels (i.e. uncertain), we invoke the GAE to generate 10 masked-atom proposals, decode them into SMILES, and rank them by distance from the original molecule.

Integration with our Flask-based visualization dashboard allows clicking on an uncertain compound to view both its prediction intervals and suggested modifications, mirroring Sharma et al.’s interactive analysis tool.

How the Paper Helped Us

Sharma et al. inspired the **interpretability arm** of our pipeline: they demonstrated that presenting uncertainty without actionable guidance has

limited value. Their GAE approach motivated our own masked-graph modules and gave us confidence to expose counterfactuals in our UI.

Moreover, their metrics for explanation quality (fidelity, sparsity) became part of our evaluation suite—ensuring that our suggestions are both minimal and chemically plausible, not just random perturbations.

6 Novel Contributions Beyond Cited Work

1. **Variance-Adaptive Conformal Thresholding.** We dynamically adjust conformal thresholds based on the predictive variance observed during MC-Dropout or ensemble inference. Instead of a single static threshold per α , thresholds are modulated per compound cluster (via k-means on embeddings), yielding tighter sets in well-modeled regions and broader sets in high-uncertainty clusters.
2. **Cross-Endpoint Consistency Regularizer.** During multi-task training, we add a penalty term that encourages correlated endpoints (e.g. NR-AR and NR-AR-LBD) to have similar latent representations, improving both coverage and compound-level consistency. This regularizer is implemented as the Frobenius norm between selected head-layer activations.
3. **On-Demand Conformal Retraining.** We provide a command-line tool that, given a new batch of compounds, incrementally updates calibration thresholds using conformal jackknife-plus, without full re-training. This enables rapid deployment to new chemical series with limited additional compute.
4. **Interactive Uncertainty Visualization Widget.** Beyond static plots, we built a Dash/Plotly component that overlays conformal set sizes and coverage contours over a 2D UMAP embedding of chemical space. Users can hover to inspect local coverage statistics and trigger counterfactual generation for specific points.