

Improved Detection of Drug-Induced Liver Injury by Integrating Predicted *In Vivo* and *In Vitro* Data

Srijit Seal, Dominic Williams, Layla Hosseini-Gerami, Manas Mahale, Anne E. Carpenter, Ola Spjuth,* and Andreas Bender*



Cite This: *Chem. Res. Toxicol.* 2024, 37, 1290–1305



Read Online

ACCESS |



Metrics & More



Article Recommendations



Supporting Information

ABSTRACT: Drug-induced liver injury (DILI) has been a significant challenge in drug discovery, often leading to clinical trial failures and necessitating drug withdrawals. Over the last decade, the existing suite of *in vitro* proxy-DILI assays has generally improved at identifying compounds with hepatotoxicity. However, there is considerable interest in enhancing the *in silico* prediction of DILI because it allows for evaluating large sets of compounds more quickly and cost-effectively, particularly in the early stages of projects. In this study, we aim to study ML models for DILI prediction that first predict nine proxy-DILI labels and then use them as features in addition to chemical structural features to predict DILI. The features include *in vitro* (e.g., mitochondrial toxicity, bile salt export pump inhibition) data, *in vivo* (e.g., preclinical rat hepatotoxicity studies) data, pharmacokinetic parameters of maximum concentration, structural fingerprints, and physicochemical parameters. We trained DILI-prediction models on 888 compounds from the DILI data set (composed of DIL1st and DIL1rank) and tested them on a held-out external test set of 223 compounds from the DILI data set. The best model, DILIPredictor, attained an AUC-PR of 0.79. This model enabled the detection of the top 25 toxic compounds (2.68 LR+, positive likelihood ratio) compared to models using only structural features (1.65 LR+ score). Using feature interpretation from DILIPredictor, we identified the chemical substructures causing DILI and differentiated cases of DILI caused by compounds in animals but not in humans. For example, DILIPredictor correctly recognized 2-butoxyethanol as nontoxic in humans despite its hepatotoxicity in mice models. Overall, the DILIPredictor model improves the detection of compounds causing DILI with an improved differentiation between animal and human sensitivity and the potential for mechanism evaluation. DILIPredictor required only chemical structures as input for prediction and is publicly available at <https://broad.io/DILIPredictor> for use via web interface and with all code available for download.



DILI PR predictor

INTRODUCTION

The liver is a major organ of drug metabolism in the human body and thus it is vulnerable to not just drugs but also their reactive metabolites.^{1,2} Drug-induced liver injury (DILI) has been a leading cause of acute liver failure³ (causing over 50% of such cases⁴), accounting for a significant proportion of drug-related adverse events. DILI is detectable generally in phase III clinical trials and a leading cause of postmarket drug withdrawals.⁵ Two common types of DILI are intrinsic and idiosyncratic.⁶ While intrinsic DILI is generally dose-dependent and predictable, idiosyncratic DILI is rare, generally not dependent on dosage, typically unpredictable and undetectable in early drug development using standard preclinical models, with a variable onset time and several phenotypes.

The mechanisms underlying DILI are multifactorial⁷ and not completely understood. These include cellular toxicities such as mitochondrial impairment,⁸ inhibition of biliary efflux,⁹ oxidative stress,¹⁰ and more. Additionally, DILI can be influenced by dose variations, pharmacokinetics (PK), and biological variations, such as variations in cytochrome P450 (CYP) expression.¹¹ Today, there are several established *in vitro* assays and *in vivo* experiments (proxy-DILI data) that are relatively good at detecting DILI risk.¹² The current battery of *in vitro* proxy-DILI assays is generally effective at detecting

hepatotoxic compounds. However, most human-relevant *in vitro* models for DILI risk assessments can still fail to accurately predict patient safety at therapeutic doses due to differing toxicity mechanisms at varying concentrations while animal *in vivo* models are not completely human-relevant.¹³ At the time of writing this research article (February–June 2024), at least three reports show clinical trials in phase 1/2 suffered setbacks or ended due to unexpected liver damage, liver function abnormalities that were not predicted by early *in vitro* models.^{14–16} This underscores the critical need for the development of more accurate and reliable *in vitro* systems that can better simulate human liver responses to drug exposure, thereby improving the safety assessment and reducing the risk of late-stage clinical failures. Thus, there is significant interest in improving *in silico* DILI prediction due to its ability to assess large numbers of compounds more quickly

Received: January 10, 2024

Revised: June 27, 2024

Accepted: July 1, 2024

Published: July 9, 2024



Table 1. Sources of Liver-Safety and Toxicity Data Used in This Study, the DILIst and DILIrak Datasets are Used Together as the Gost Standard DILI Dataset Used in This Study

| data source | assay type | used in this study | total number of compounds | number of compounds in active class | description | reference (data retrieved from) |
|----------------------------|------------------------------|--------------------|---------------------------|-------------------------------------|---|---------------------------------|
| human hepatotoxicity | human hepatotoxicity | training data | 1163 | 588 | human hepatotoxicity | Mulliner et al. |
| animal hepatotoxicity A | animal hepatotoxicity | training data | 542 | 184 | chronic oral administration, hepatic histopathologic effects, ToxRefDB | Liu et al. |
| animal hepatotoxicity B | animal hepatotoxicity | training data | 671 | 369 | hepatocellular hypertrophy, rats, ORAD, HESS, | Ambe et al. |
| preclinical hepatotoxicity | animal hepatotoxicity | training data | 2204 | 1642 | preclinical hepatotoxicity | Mulliner et al. |
| diverse DILI A | heterogenous data | training data | 1106 | 382 | large-scale and diverse DILI data set, | He et al. |
| diverse DILI C | heterogenous data | training data | 445 | 208 | transient liver function abnormalities, adverse hepatic effects, U.S. FDA Orange Book, Micromedex | Mora et al. |
| BESP | mechanisms of liver toxicity | training data | 446 | 240 | bile salt export pump inhibition | McLoughlin et al. |
| Mitotox | mechanisms of liver toxicity | training data | 5239 | 689 | mitochondrial toxicity | Hemmerich et al. |
| reactive metabolite | mechanisms of liver toxicity | training data | 317 | 81 | reactive metabolite | Mazzolari et al. |
| C_{max} (total) | pharmacokinetic properties | predicted property | 718 | N/A | maximum total concentration in plasma | Smith et al. |
| C_{max} (unbound) | pharmacokinetic properties | predicted property | 515 | N/A | maximum unbound concentration in plasma | Smith et al. |
| DILIst | DILI | test data (DILI) | 990 | 619 | DILIst classification | Tong et al. |
| DILIrak | DILI | test data (DILI) | 121 | 97 | DILIrak data set | Chen et al., Chavan et al. |

and cost-efficiently, especially in the early stages of drug discovery.¹⁷

In the drug discovery pipeline, hepatotoxicity assessment encompasses a variety of *in vitro* and *in vivo* experimental models as well as *in silico* models. Several *in vitro* models for liver toxicity testing employ proxy end points (hepatotoxicity assays) with liver slices and cell lines such as primary animal and human hepatocytes¹⁸ or even three-dimensional systems with the dynamic flow for the primary cell and/or stem cell cultures.¹⁹ However, the ideal hepatocyte-like cell model system depends on the evaluation of particular cellular functions given that there are substantial differences among various human liver-derived single-cell culture models as previously explored in the context of drug disposition, bioactivation, and detoxification.²⁰ The agreement between *in vitro* data and human *in vivo* data is also low.²¹ For example, methapyrilene is known to cause changes to the level of iron metabolism in the human hepatic HepaRG cell line²² and oxidative stress and mitochondrial dysfunction in rats²³ but has not been reported to cause hepatotoxicity in humans.^{24,25} On the other hand, *in vivo* animal models also have low concordance as shown by recent studies using the eTOX database where organ toxicities were rarely concordant between species.²⁶ The concordance between animal and human data for liver toxicity, specifically, is often low (with some studies indicating rates as low as 40%²⁷ and others in the range of 39–44%²⁸), which makes extrapolating safety assessments from animals to humans a challenging endeavor.^{29,30} For example, 2-butoxyethanol causes hepatic toxicity in mice via an oxidative stress mechanism but not in humans given humans have higher levels of liver vitamin E (and a high resistance to iron accumulation) compared to mice.³¹ Overall, this leads to a greater need for improved DILI prediction,

especially when translating knowledge from preclinical stage and animal studies to human clinical studies.¹⁷

DILIst³² and DILIrak³³ are lists of compounds that have been classified as inducing DILI or not and were developed from FDA-approved drug labels. Binary classification from labeling documents is challenging, and this is evident in the fact that many DILIrak compounds are labeled ambiguous although the DILI for some of these compounds has been reported in the literature. For *in silico* models, these ambiguous compounds are generally removed. Machine learning models are being increasingly used to model biological systems and identify complex patterns in data sets.³⁴ Generally, *in silico* models rely on identifying chemical structural alerts³⁵ or use a range of chemical or physicochemical features. Ye et al. employed Random Forest algorithms and Morgan fingerprints for DILI prediction, achieving an AUC of 0.75 with random splitting (70% training, 30% testing).³⁶ Liu et al. utilized Support Vector Machines and obtained a 76% balanced accuracy on an external test set using Morgan fingerprints; however, their predicted protein target descriptors provided less accurate predictions (balanced accuracy of 59%) but offered better interpretability.³⁷ Mora et al. employed QuBiLS-MAS 0–2.5D molecular descriptors to predict DILI (labels from various sources) on an external test set comprising 554 compounds, achieving a 77% balanced accuracy.³⁸ Predicting organ-level toxicity solely based on chemical structure is challenging and the use of biological data helps improve toxicity prediction.^{39,40} More recently, predicted off-target effects and experimental P450-inhibitory activity have also been considered to improve DILI prediction.^{41,42} Moving away from binary predictions, Aleo et al. developed the hepatic risk matrix (HRM) to assess the potential for human drug-induced liver injury (DILI) among lead clinical and back-up drug

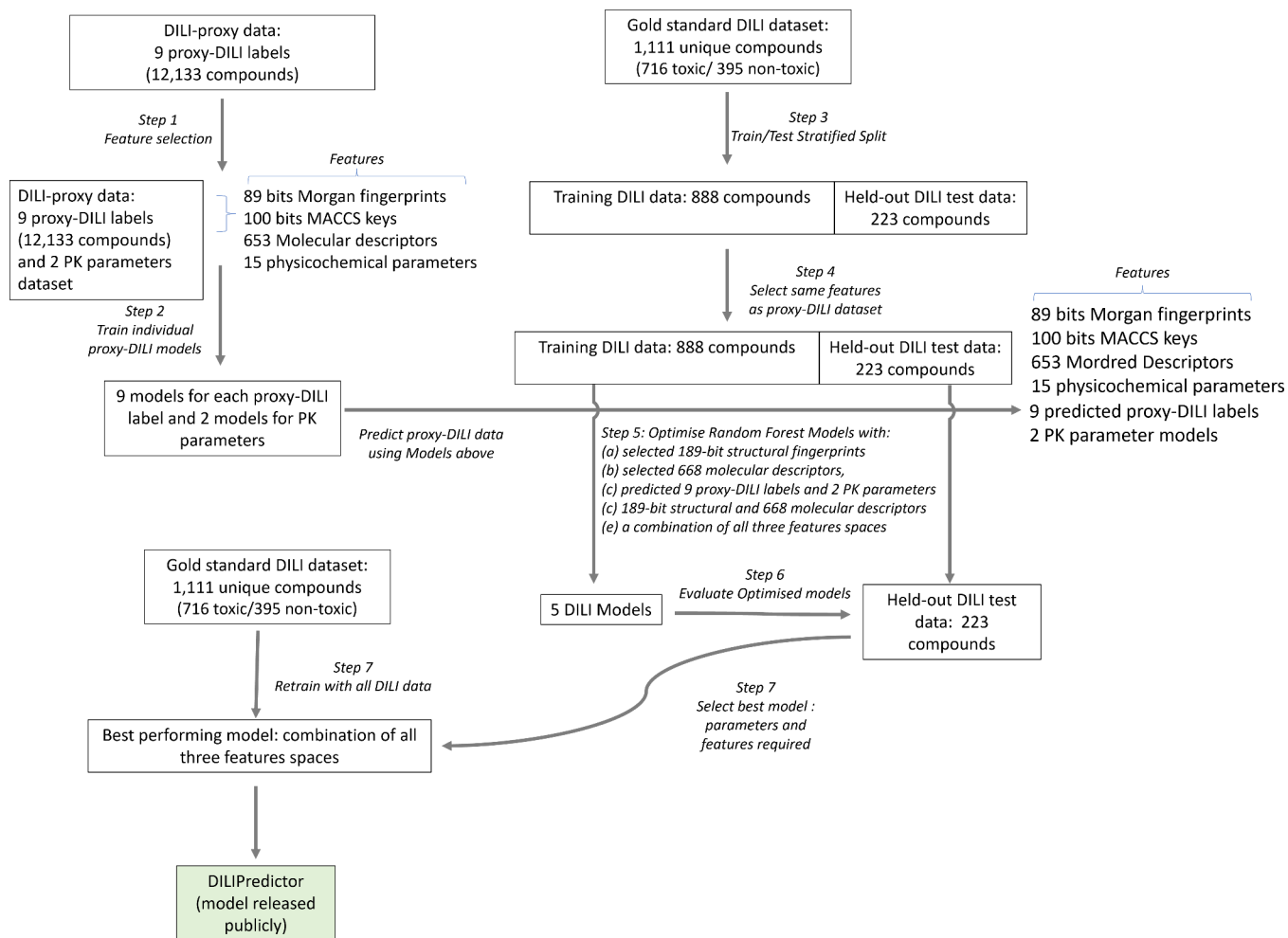


Figure 1. Workflow of the current study. Individual models for 9 *in vivo* and *in vitro* assays in the proxy-DILI data set and 2 PK parameters were used to predict these end points for compounds in the gold standard DILI data set. A combination of these predictions along with chemical structure and molecular descriptors were then used to train and evaluate the models on DILI compounds.

candidates. By integrating physicochemical properties and common toxicity mechanisms, the HRM stratifies drug candidates based on safety margins relative to clinical $C_{\max, \text{total}}$. This study identified 70–80% most-DILI-concern drugs and effectively differentiated successful from unsuccessful drug candidates for liver safety.⁴³ Chavan et al. integrated high-content imaging features with chemical features for DILI label prediction, resulting in a 0.74 AUC.⁴⁴ Previously, the authors of this work explored this in the case of mitochondrial toxicity⁴⁵ (which at high doses is one of the mechanisms known to cause DILI), cytotoxicity,⁴⁶ and also cardiotoxicity.⁴⁷

In this study, we significantly extended the use of different data sources to several *in vivo* and *in vitro* data types in developing the DILIPredictor model presented here. We identified liver injury end points such as human hepatotoxicity,⁴⁸ preclinical hepatotoxicity and animal hepatotoxicity^{48–50} and DILI data sets compiled by various studies^{38,51} (Table 1) These data sets provide the *in vivo* labels for DILI for different species at various stages of the drug discovery pipeline, from preclinical to postmarket withdrawals. We identified three *in vitro* assays that could be indicative of liver toxicity and with public data:⁷ mitochondrial toxicity,⁵² bile salt export pump inhibition (BSEP),⁵³ and the formation of reactive metabolites.⁵⁴ Mitochondria account for 13–20% of the liver, and mitochondrial dysfunction can impact ATP

synthesis, increase ROS generation, and trigger liver injury.⁵⁵ The majority of the mitochondrial toxicity data in Hemmerich et al. originates from a Tox21 assay assessing mitochondrial membrane depolarization in HepG2 cells (which provides a distinct perspective compared to *in vitro* data derived from primary hepatocytes), thereby introducing additional biological information. When BSEP function is inhibited, bile salts accumulate within liver cells, causing hepatocyte injury and a risk of liver failure.⁵⁶ Metabolic processes can form reactive metabolites that bind covalently to hepatic proteins, altering their function and leading to damage in liver tissues.⁵⁷ We also included PK parameters, which have been predicted before using machine learning models based on chemical structures.⁵⁸

Overall, in this study, we hypothesized that these proxy-DILI labels along with chemical structure and physicochemical parameters would lead to improved predictivity in identifying potential liver injury end points while differentiating between sensitivities observed in human and animal proxy-DILI labels, allowing for interpretations of hepatotoxicity data across species. An objective of our study is not just to achieve high overall predictive performance but to understand how individual *in vitro* and *in vivo* proxy end points provide predictive value for DILI outcomes in humans. The prediction of individual proxy end points is valuable, as it aligns with the specific experiments, thereby forming testable hypotheses. It

should be noted that *in vivo* DILI is not easily testable without significant effort, and in the case of humans, clinical trials. Hence, in this work, a FeatureNet approach was adopted, that is, models are trained on predictions from other individual models. This allows us to interpret the importance of the individual predictions to the final prediction with previous studies showing comparable performance to multitask learning.⁵⁹ Furthermore, the nature of our data, characterized by small sizes for *in vivo* relevant data and sparse matrices, presents additional challenges for implementing multitask learning effectively. Finally, by including *in vitro* proxy-DILI labels, the models developed in this study have the potential for mechanistic evaluation and facilitating a comprehensive understanding of the underlying biochemical and cellular processes associated with drug-induced liver injuries.

MATERIALS AND METHODS

The workflow followed in this study is shown in Figure 1 and described in more detail in the following.

Drug-Induced Liver Toxicity Data Sets: DIList and DILrank. The human *in vivo* data set for liver toxicity was collected by combining DIList³² (714 toxic and 440 nontoxic compounds) and DILrank³³ (268 toxic and 76 nontoxic compounds from Chavan et al.³⁹) data sets. The DIList data set classifies compounds into two classes based on their potential for causing DILI. The DILrank data set was released by the FDA prior to DIList. This data set analyzed the hepatotoxic descriptions from FDA-approved drugs and assessed causality evidence from the literature and classified compounds into four groups: ^vMost-, ^vLess-, ^vNo-DILI concern, and Ambiguous-DILI-concern drugs. For the DILrank data set, we retrieved data from Chavan et al.³⁹ We treated ^vMost- and ^vLess as DILI-positive and those labeled with ^vNo-DILI-concern as DILI-negative. Ambiguous-DILI-concern drugs were removed. Together, these data sets form the largest drug list with DILI classification to date.

Proxy-DILI Data Sets: *In Vivo* and *In Vitro* Assays. The data sets we considered include one proxy-DILI label from studies on human hepatotoxicity,⁴⁸ and two proxy-DILI labels from animal hepatotoxicity studies (animal hepatotoxicity A and B, and preclinical hepatotoxicity as detailed in Table 1).^{48–50} Animal hepatotoxicity data sets mentioned above consisted of data compiled by the authors from ToxRefDB,⁶⁰ ORAD,⁵⁰ and HESS⁶¹ as well as hepatic histopathologic effects. Two diverse DILI data sets contain heterogeneous data collected by other studies^{38,51} (Diverse DILI A and C as detailed in Table 1). These data sets consisted of data from the drugs known to cause transient liver function abnormalities and adverse hepatic effects as well as compounds from the U.S. FDA Orange Book and Micromedex. We included three *in vitro* assays related to proposed or known mechanisms of liver injury, namely, mitochondrial toxicity,⁵² bile salt export pump inhibition (BSEP),⁵³ and the formation of reactive metabolites⁵⁴ (as detailed in Table 1). The labels for these *in vitro* data sets were the assay hit calls defined by the original studies. Supplementary Table S1 lists the SMILES, chemical names, and CASRN for all molecules, where available from the original sources; note that not all original sources provide complete information on chemical names or CASRN. Previous studies indicated that mitochondrial toxicity and BSEP are reasonable predictors for cholestatic and mitochondrial toxins; however, they fail when applied to a wider chemical space for drugs with different mechanisms.⁶² Many assay hits screened from chemical libraries often have unfavorable drug metabolism and pharmacokinetics presenting development challenges.⁶³ Thus, we considered pharmacokinetics as two of the proxy-DILI labels and compiled pharmacokinetic parameters of maximum concentration (C_{max}) from Smith et al.⁶⁴ This data set contains maximum unbound concentration in plasma for 515 unique compounds and maximum total concentration in plasma for 718 unique compounds (Supplementary Table S2). Together, as shown in Table 1, we obtained nine *in vivo* and *in vitro* assays (proxy-

DILI labels) related to liver injury and two pharmacokinetic parameters.

Data Set Preprocessing and Compound Standardization. So that the focus of this work remains only on traditional drug like molecules, we dropped compounds if there were no carbon atoms present, if they were disconnected compounds, if they contained only metals, or if the molecular weight was above 1500 Da. Then, compound SMILES were standardized using RDKit,⁶⁵ which included disconnection of metal atoms, normalization, and reionization, isolating the principal molecular fragment, uncharging and charge normalization, standardization to most common isotopic form, removing stereochemical configuration and tautomer standardization. The entire process was iteratively applied up to five times or until the SMILES representation stabilized. If the standardization process did not converge to a single representation, the most frequently occurring SMILES string across the iterations was selected as the standardized form. Finally, the standardized SMILES were protonated to reflect their form at physiological pH of the liver (pH = 7.0) implemented using Dimorphite-DL.⁶⁶ To ensure our data set did not contain duplicate chemistry, we used the first 14 characters of InChIKeys (known as the “hash layer”), which represents the chemical structure excluding stereochemical and isotopic variations. By comparing these truncated InChIKeys, we identified duplicates. In cases where DILrank and DIList contained the same compound with conflicting toxicity labels, we retained the toxic annotation (given there was some evidence of toxicity in either DIList or DILrank). Finally, we obtained a data set of 1,111 unique compounds and associated DILI labels (716 toxic and 395 nontoxic compounds). This data set is henceforth referred to as the gold standard DILI data set (Supplementary Table S3).

For the proxy-DILI data set, in the case of any compounds with conflicting toxicity labels within a particular data set after SMILES standardization, we retained the compound as toxic/active (hence preferring the evidence of toxicity/activity which is a usual practice in drug discovery) resulting in a data set of 13,703 compounds. For each of the nine labels besides PK parameters (as detailed in Table 1), if a compound was already present in the gold standard DILI data set above (compared using the InChIKey hash layer), we removed the compound from the proxy-DILI data set. This was done to avoid any information leaks in the models developed in this study. Finally, we obtained a data set of 12,133 compounds in total for nine proxy-DILI labels, which are henceforth called the proxy-DILI data set in this study (Supplementary Table S4).

Assay Concordance with Experimental Values. To evaluate the concordance of the nine proxy-DILI labels and the gold standard DILI data set with each other, we used all 13,703 compounds in the proxy-DILI data set and compared them to the 1,111 compounds in the gold standard DILI data set. To evaluate concordance, we used Cohen’s kappa (as defined in scikit-learn v1.1.1⁶⁷) to measure the level of agreement between activity values for each pair of labels which were present in the data set.

Exploring the Physicochemical Space. Physicochemical space was explored using six characteristic physicochemical descriptors of molecular weight, TPSA, number of rotatable bonds, number of H donors, number of H acceptors and log P, (as implemented in RDKit⁶⁵ v.2022.09.5). We used a *t*-distributed stochastic neighbor embedding (t-SNE from scikit-learn v1.1.1⁶⁷) to obtain a map of the physicochemical space for all compounds in the gold standard DILI data set and proxy-DILI data set with a high explained variance (PCA: 85.15% using two components).

Structural Fingerprints, Mordred, and Physicochemical Descriptors. We used Morgan fingerprints⁶⁸ of radius 2 and 2048 bits and 166-bit MACCS Keys,⁶⁹ as implemented in RDKit⁶⁵ (v2022.09.5), as structural features for all compounds in the DILI data set and proxy-DILI data set. This resulted in 2,214-bit vector structural fingerprints.

We used molecular descriptors (as implemented in the Mordred⁷⁰ python package) and physicochemical properties (such as topological polar surface area TPSA, partition coefficient log P, etc. as implemented in RDKit⁶⁵ v2022.09.5) for all compounds in the gold

standard DILI data set and proxy-DILI data set. We dropped descriptors with missing values, which resulted in 1,016 molecular descriptors for each compound.

Feature Selection. We first used feature selection on the compounds in the proxy-DILI data set using a variance threshold (as implemented in scikit-learn v1.1.1⁶⁷) to filter features (Figure 1 Step 1). We used a low variance threshold of 0.05 for Morgan fingerprints resulting in 89 selected bits, a threshold of 0.10 for MACCS keys resulting in 100 selected keys, and a threshold of 0.10 for Mordred descriptors resulting in 653 selected descriptors. Lower thresholds for variance ensured strict selection criteria, leading to fewer selected features to strike a balance between the length of all fingerprints and physicochemical parameters. An additional 15 calculated physicochemical parameters (as implemented in RDKit⁶⁵ v2022.09.5: topological polar surface area, hydrogen bond acceptors and donors, fraction of sp³ carbons, log P, and the number of rotatable bonds, rings, assembled rings, aromatic rings, hetero atoms, stereocenters, positive and negatively charged atoms, and the counts of NHOH and NO) were also added. This resulted in 189 bit-vector structural fingerprints and 668 molecular descriptors for each compound in the proxy-DILI data set. The same selected features were used for the gold standard DILI data set (Figure 1 Step 4) to avoid any information leaks.

Evaluation of Predictions from Individual Proxy-DILI Models. First, we trained individual models for each of the nine proxy-DILI end points for all of the other proxy-DILI end points. For each proxy-DILI end point, we trained individual Random Forest models (Figure 1 Step 2) with a 5-fold stratified cross-validation and random halving search hyperparameter optimization (as implemented in scikit-learn v1.1.1⁶⁷ with hyperparameter space given in Supplementary Table S5). We used this hyperparameter-optimized model to obtain predicted probabilities for all compounds for the other proxy-DILI end points for every 9 × 9 combination. For each model built on a proxy-DILI end point, we chose an optimal decision threshold based on the J-statistic value (see released code for implementation) by comparing the predicted probabilities to the true values. We obtained final binary predictions using this threshold, thereby choosing the best-case scenario where the balanced accuracy is optimized from the AUC-ROC curve. Next, we compared how well each proxy-DILI model was at predicting other proxy-DILI labels by comparing the F1 score and likelihood ratios.

Evaluating Predictivity of Individual Proxy-DILI models for the Gold Standard DILI Data Set. To train and evaluate models for DILI, we first split our gold-standard DILI data set (containing 1,111 unique compounds) using ButinaSplitter based on the Butina clustering of a bulk Tanimoto fingerprint matrix split with a cutoff threshold of 0.70 (as implemented in DeepChem,^{71,72} Figure 1 Step 3). This led to a training DILI data of 888 unique compounds (560 toxic and 328 nontoxic compounds) and a held-out DILI test set of 223 unique compounds (156 toxic and 67 nontoxic). This ensures that the DILI test set included a wide variety of compounds that are structurally less similar to those used in training. Therefore, the held-out DILI test set is a more challenging representation (compared to random splits) as encountered in real-world drug discovery: predicting DILI outcomes for new compounds away from the chemical space of the known compounds. We evaluated the performance of individual models built on each of the nine proxy-DILI end points on the held-out DILI test set (223 compounds). First, for each of the nine individual models, we obtained out-of-fold predicted probabilities on the DILI training data (888 compounds) using cross-validation with a 5-fold stratified split. We used these out-of-fold predicted probabilities and true values to obtain an optimal decision threshold based on the J-statistic value. Finally, we used each of the individual models and the corresponding optimal decision threshold to obtain predictions of the held-out DILI test set. We used the Jaccard similarity coefficient score (as implemented in scikit-learn v1.1.1⁶⁷) to compare the similarity of predictions, that is, the predicted DILI vectors from each model. The Jaccard similarity coefficient measures the similarity between two sets of data counting mutual presence (positives/toxic) as matches but not the absences.

Models for Prediction of C_{max}. Next, we trained two Random Forest regressor models to predict the median pMolar unbound plasma concentration and median pMolar total plasma concentration for 515 and 718 compounds, respectively (Figure 1 Step 2). We used the selected 189 bit-vector structural fingerprints and 668 molecular descriptors as features to train the models with a 5-fold stratified cross-validation and random halving search hyperparameter optimization as described above. The best estimator was refit on the entire data set, the final model was used to generate predictions for compounds, and these predicted features were used for training DILI models.

Models for Prediction of DILI. In this study, we built models (Figure 1 Step 5) using (a) selected 189-bit structural fingerprints, (b) selected 688 molecular descriptors, (c) a combination of selected 189-bit structural fingerprints and selected 668 molecular descriptors, (d) predicted nine proxy-DILI labels and two predicted pharmacokinetic parameters which refers to a FeatureNet approach, and (e) a combination of all three features spaces.

For each feature space, we used repeated nested cross-validation. First, the DILI training data was split into 5-folds. One of these folds was used as a validation set while the data from the remaining 4 folds were used to train and hyperparameter optimize a Random Forest Classifier (as implemented in scikit-learn v1.1.1⁶⁷). We optimized the classifier model using a random halving search (as implemented in scikit-learn v1.1.1⁶⁷) and 4-fold cross-validation (see Supplementary Table S5 for hyperparameter space). Once hyperparameters were optimized, we then used the fitted model to generate 4-fold cross-validated estimates for each compound in the fitted data. These predicted probabilities along with the real data were used to generate an optimal threshold using the J statistic value (see released code for implementation). Finally, we predicted the DILI end point for the validation set and used the optimal threshold to determine the DILI toxicity. The process was repeated 5 times in total until all 888 compounds in the DILI training data were used as a validation set. This entire nested-cross validation setup was repeated ten times with different splits. The model with the highest AUC was fit on the entire DILI training data and we obtained the optimal threshold using the J statistic value on the 4-fold cross-validated estimates for each of these compounds. Finally, this threshold was used to evaluate our models (Figure 1 Step 6) on the held-out DILI test set (223 unique compounds). Thus, for each model using a feature space (or the combination), we obtained evaluation metrics on (a) the nested cross-validation (on training data) and (b) the held-out test set. The best-performing model (Figure 1 Step 7), as shown in the Results and Discussion section, was the combination of all three feature spaces. This model was retrained (Figure 1 Step 8) on the complete gold-standard DILI data set consisting of 1,111 distinct compounds. This model, DILIPredictor, can be accessed through a web application <https://broad.io/DILIPredictor> and have all code available for local use on GitHub at <https://github.com/srijitseal/DILI>.

To calculate the structural similarity of the held-out test to training data, we first calculated pairwise Tanimoto similarity (using 2048-bit Morgan fingerprint, see released code for implementation) for each test compound to each training compound. Finally, we calculated the mean of the three highest Tanimoto similarities (that is the three nearest neighbors), which was used to define the structural similarity of the particular test compound.

Evaluation Metrics. All predictions (nested-cross validation and held-out test set) were evaluated using sensitivity, specificity, balanced accuracy (BA), Mathew's correlation constant (MCC), F1 scores, positive predictive value (PPV), likelihood ratio (LR+),⁷³ average precision score (AP), and area under curve-receiver operating characteristic (AUC-ROC) as implemented in scikit-learn v1.1.1.⁶⁷

Feature Importance Measures to Understand the Chemistry and Biological Mechanisms for Common DILI Compounds. For the final model released publicly that used a combination of all feature spaces, we used SHAP values (as implemented in the shap python package⁷⁴) to obtain feature importance for each input compound. This included proxy-DILI data, pharmacokinetic parameters, physicochemical features, and also

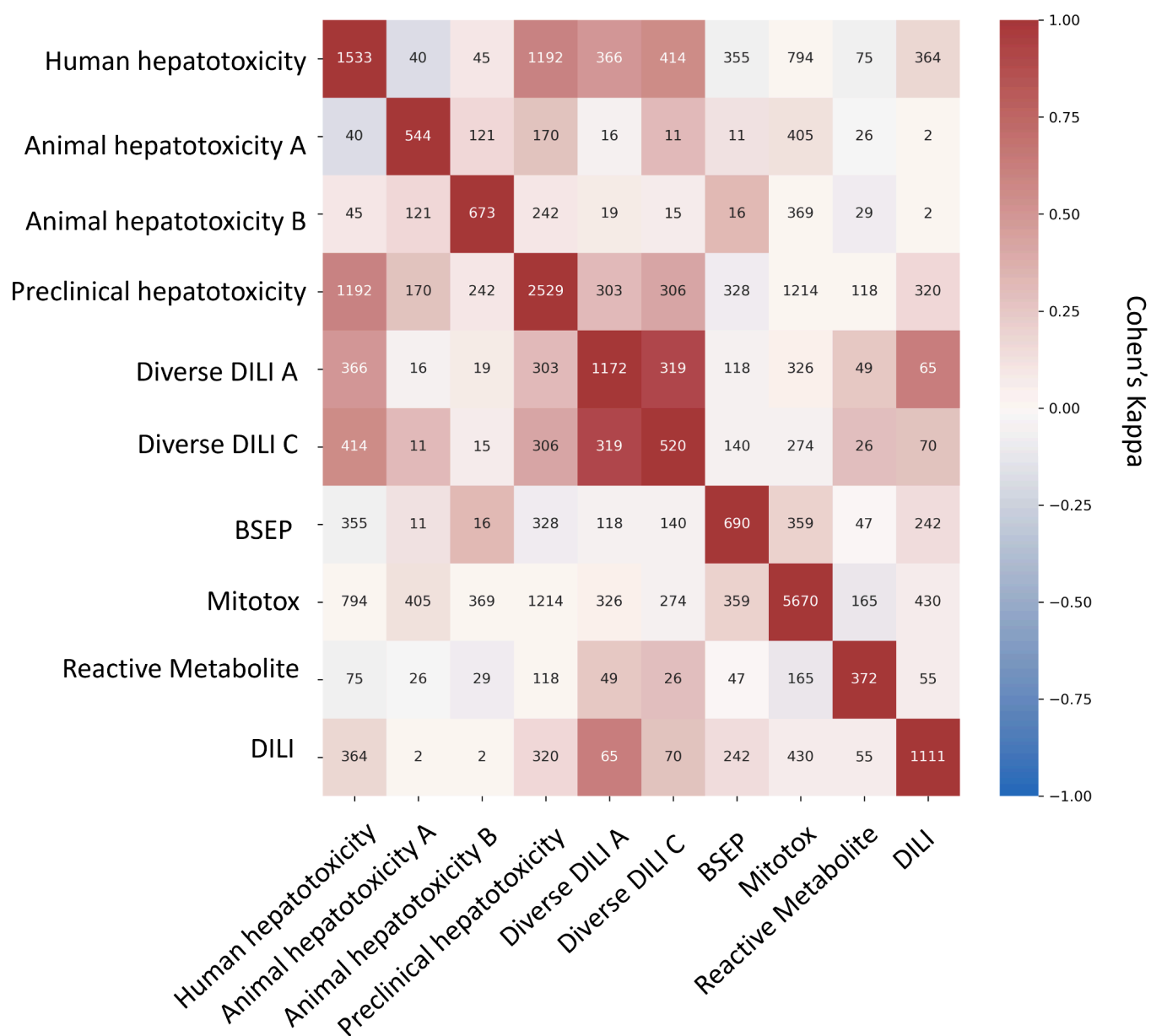


Figure 2. Concordance of compounds overlapping across nine labels in the proxy-DILI data set (13,703 compounds) including compounds that overlapped with DILI data (1,111 compounds). Concordance is given using Cohen's kappa (and the number of overlapping compounds given as annotations). Overall, the human-related proxy-DILI labels and diverse heterogeneous DILI labels showed high concordance with DILI compounds and among each other.

MACCS key substructures that contributed to DILI toxicity/safety. Further, we show how the DILIPredictor can be used to elucidate the causes of DILI, both in chemistry and via mechanisms on the biological level using the importance measures on proxy-DILI labels. We analyzed 4 compounds that were not present in the training data of these models. Two of these compounds, enzalutamide and sitaxentan, are known to cause DILI in humans while two compounds, 2-butoxyethanol and astaxanthin, did not cause DILI in humans. Additionally predicted profiles from another 12 compounds (also present in the training data) are shown in [Supplementary Table S6](#). Several of toxic compounds were related to the study by Chang et al., who compiled compounds causing DILI in patients undergoing chemotherapy.⁷⁵ We also included two pairs of compounds studied by Chen et al. such as doxycycline/minocycline and moxifloxacin/trovafloxacin; these pairs were defined by a similar chemical structure and mechanism of action but differed in their liver toxicity effects.⁷⁶

Statistics and Reproducibility. We have released the data sets used in this proof-of-concept study, which are publicly available at

<https://broad.io/DILIPredictor>. We released the Python code for the models which are publicly available on GitHub at <https://github.com/srijitseal/DILI>.

RESULTS AND DISCUSSION

In this work, we trained models on each of nine proxy-DILI end points related to liver toxicity. We used these models to obtain predicted proxy-DILI labels for 1,111 compounds in the gold standard DILI data set (as defined in Methods), none of which overlapped with the proxy-DILI data set. We then trained new models using those predicted proxy-DILI labels as inputs, which refers to a FeatureNet approach, together with the compounds' structural fingerprints, physicochemical properties, and a combination thereof, for 888 compounds the gold standard DILI data sets. We then evaluated the models on a held-out test set of 223 compounds.

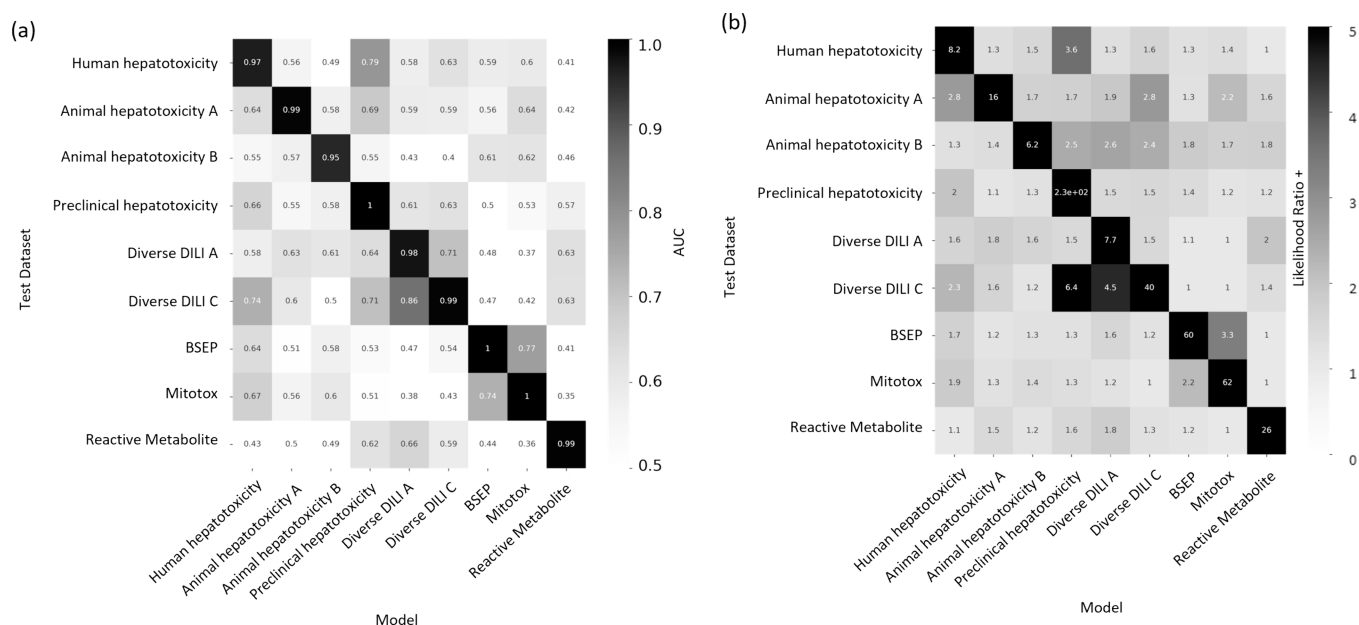


Figure 3. Performance metrics for models built on nine proxy-DILI labels when predicting labels for the other proxy-DILI in the model, evaluated using (a) AUC-ROC and (b) likelihood ratio (LR+).

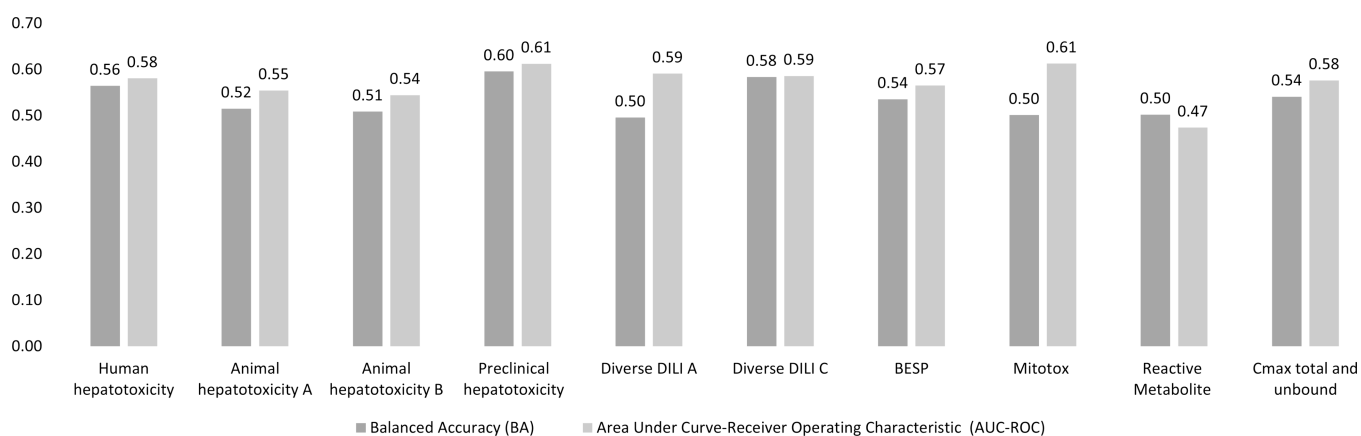


Figure 4. Performance metrics AUC-ROC and balanced accuracy achieved by each of nine individual models built on the proxy-DILI labels and a model built on two pharmacokinetic parameters (C_{\max} total and unbound) when tested on the 223 compounds in the held-out DILI data set.

Comparing Chemical Spaces for the Proxy-DILI and Gold Standard DILI Data Sets. We first examined the diversity and representation of compounds in the proxy-DILI and gold standard DILI data sets, to ensure the evaluation would be reasonable. The distribution of compounds in each of the nine labels of the proxy-DILI data set covers a diverse range of physicochemical parameters as shown in [Supplementary Figure S1](#). Gold standard DILI compounds effectively capture the diversity and representativeness of the compounds in the proxy-DILI data set as shown in [Supplementary Figure S2](#) for the physicochemical space of the 1,111 compounds in the gold standard DILI data set compared to 12,133 compounds in the proxy-DILI data set. Further, the held-out DILI test set (223 compounds) was also representative in the physicochemical parameter space of the training DILI data (888 compounds) as shown in [Supplementary Figure S3](#). The main caveat to consider is that the six characteristic physicochemical descriptors capture the variability of physicochemical space only to a certain extent. Overall, we

conclude that the chemical space covered by the data sets is sufficiently similar for our evaluation to be reliable.

Concordance of Proxy-DILI Data Sets and DILI Compounds. Next, we aimed to evaluate the concordance of labels in the proxy-DILI data set with the gold standard DILI data set. To do so, we compared all 13,703 compounds in the proxy-DILI data set to the 1,111 compounds in the gold standard DILI data set. It is important to note that these compounds (that overlapped between the proxy-DILI and gold standard DILI data set) were only used to analyze concordance in this section and not in training the models, because that would leak information. As depicted in [Figure 2](#), we observed a strong concordance between the data sourced from human hepatotoxicity data set and preclinical data (Cohen's Kappa = 0.60), and the two diverse DILI data sets (0.49 and 0.55) used in this study. The lack of perfect concordance is reasonable given these data sets are primarily derived from human-related data, as opposed to animal data or *in vitro* assays. Note, concordance between DILI and proxy-DILI labels may be affected as the proxy-DILI data set used here includes some of

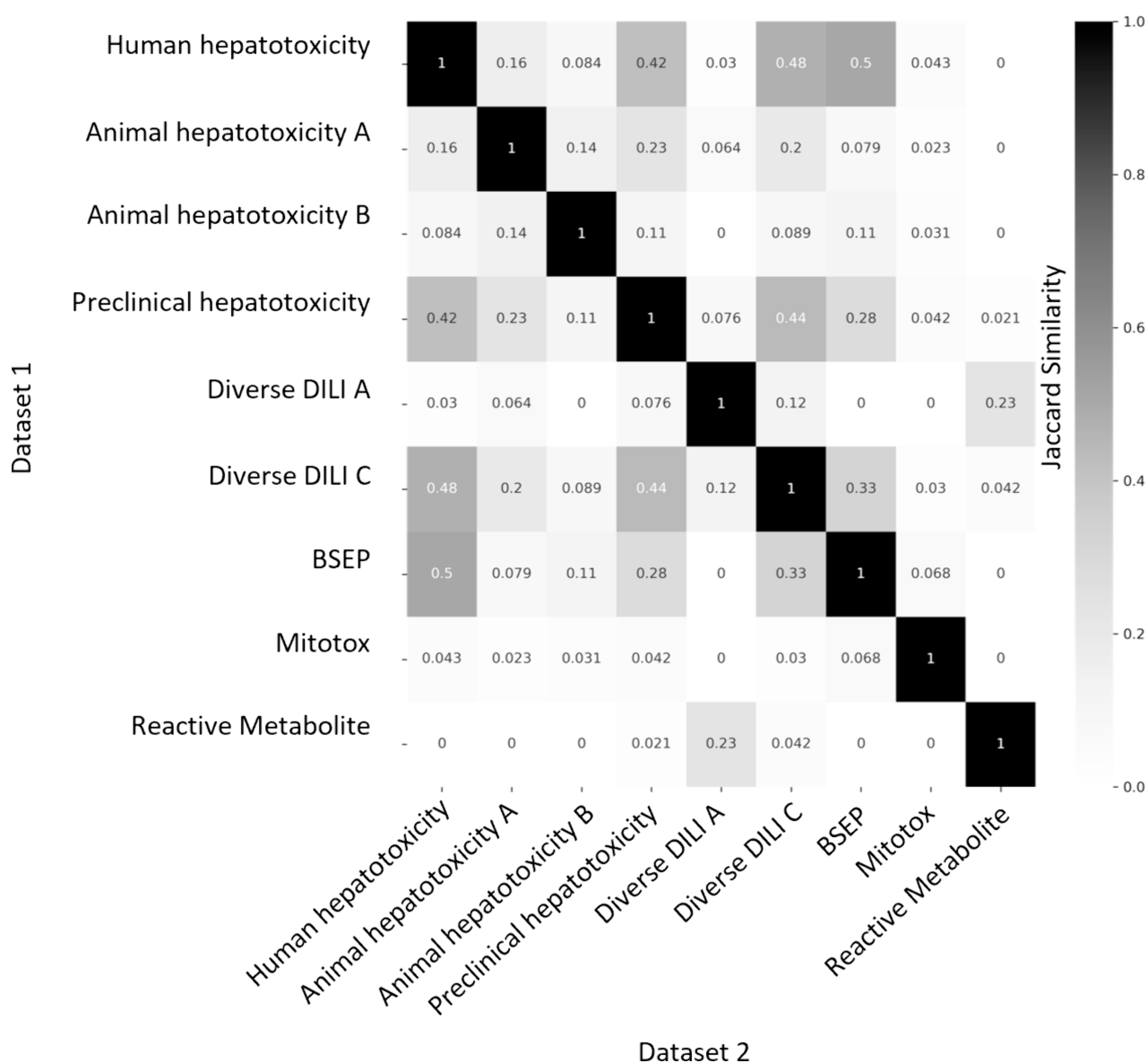


Figure 5. Jaccard similarity of predictions on the held-out DILI data set (223 compounds) for individual models built on nine proxy-DILI labels in the proxy-DILI data set.

the DILI compounds (these overlapped compounds were removed later when training models).

Individual Proxy-DILI Models Are Complementary to Each Other and Distinct in Their Prediction for DILI Compounds. We next used the individual models built on the nine proxy-DILI labels to predict the other proxy-DILI labels (with evaluation metrics as shown in [Supplementary Table S7](#)). As shown in [Figure 3](#), we observed the human hepatotoxicity was well predicted using preclinical hepatotoxicity (LR+ = 3.63, F1 = 0.79). Bile salt export pump inhibition (BSEP) and mitochondrial toxicity were strongly predictive of each other (LR+ = 2.16, F1 = 0.36 when using BSEP to predict mitotox and LR+ = 3.35, F1 = 0.77 when using Mitotox to predict BSEP). Overall, the assays in the proxy-DILI data set can be used to train individual models to generate predicted proxy-DILI labels, which then provide a complementary source of information.

We next analyzed the nine individual proxy-DILI models and a model built on the two PK parameters (C_{max} unbound

and total) for their predictions on the 223 compounds in the held-out compounds of the gold standard DILI data set. As shown in [Figure 4](#) (further details in [Supplementary Table S8](#)), the best-performing models were the model built on the preclinical animal hepatotoxicity (AUC = 0.61, LR+ = 1.63) and the model built on diverse DILI C data set (AUC = 0.59, LR+ = 1.32). Further the proxy-DILI datasets have compounds covering a wider biological and chemical space coverage, which also warrants their inclusion in our study as shown by Jaccard similarity for predictions on the held-out DILI data set. Predictions from models built on animal hepatotoxicity labels were not similar to predictions from models built on human hepatotoxicity labels ([Figure 5](#); mean Jaccard similarity of 0.12). We found that predictions from models built on human-related labels were similar (e.g., predictions from the preclinical hepatotoxicity model have a Jaccard similarity of 0.42). However, predictions from human-related labels were dissimilar to predictions from *in vitro* assays (e.g., predictions from the preclinical hepatotoxicity model had only a 0.04

Table 2. Performance of Various Models from (a) 55 Held Out Test Sets from Repeated Nested Cross Validation (rNCV) and (b) for the 223 Compounds in the Held-Out DILI Dataset^a

| model | features used | evaluation data | balanced accuracy (BA) | Mathew's correlation coefficient (MCC) | area under curve-receiver operating characteristic (AUC-ROC) | sensitivity | specificity | F1 score | likelihood ratio (LR +) | positive predictive value (PPV) | average precision score (AP) |
|---|---|-----------------------|------------------------|--|--|-------------|-------------|----------|-------------------------|---------------------------------|------------------------------|
| proxy-DILI data only | 9 <i>in vitro</i> and <i>in vivo</i> labels, and 2 PK parameters | rNCV (mean) | 0.61 | 0.22 | 0.66 | 0.56 | 0.67 | 0.69 | 1.77 | 0.72 | 0.76 |
| | | held-out DILI dataset | 0.63 | 0.24 | 0.67 | 0.60 | 0.66 | 0.72 | 1.76 | 0.79 | 0.81 |
| chemical structure only | Morgan fingerprints and MACCS keys | rNCV (mean) | 0.63 | 0.25 | 0.68 | 0.60 | 0.65 | 0.69 | 1.78 | 0.74 | 0.76 |
| | | held-out DILI dataset | 0.54 | 0.08 | 0.57 | 0.34 | 0.74 | 0.73 | 1.31 | 0.72 | 0.76 |
| physicochemical properties | Mordred descriptors and physicochemical parameters | rNCV (mean) | 0.64 | 0.27 | 0.68 | 0.61 | 0.66 | 0.70 | 1.83 | 0.75 | 0.76 |
| | | external test | 0.58 | 0.14 | 0.61 | 0.63 | 0.53 | 0.62 | 1.32 | 0.77 | 0.78 |
| chemical structure and physicochemical properties | Morgan fingerprints, MACCS keys, Mordred descriptors, physicochemical parameters | rNCV (mean) | 0.63 | 0.26 | 0.68 | 0.61 | 0.66 | 0.69 | 1.81 | 0.74 | 0.76 |
| | | held-out DILI dataset | 0.57 | 0.13 | 0.62 | 0.43 | 0.71 | 0.72 | 1.47 | 0.74 | 0.79 |
| DILIPredictor | Morgan fingerprints, MACCS keys, Mordred descriptors, physicochemical parameters, 9 <i>in vitro</i> and <i>in vivo</i> labels and 2 PK parameters | rNCV (mean) | 0.64 | 0.28 | 0.68 | 0.64 | 0.64 | 0.69 | 1.84 | 0.76 | 0.76 |
| | | held-out DILI dataset | 0.59 | 0.16 | 0.63 | 0.61 | 0.56 | 0.65 | 1.40 | 0.77 | 0.79 |

^arNCV: repeated nested cross validation.

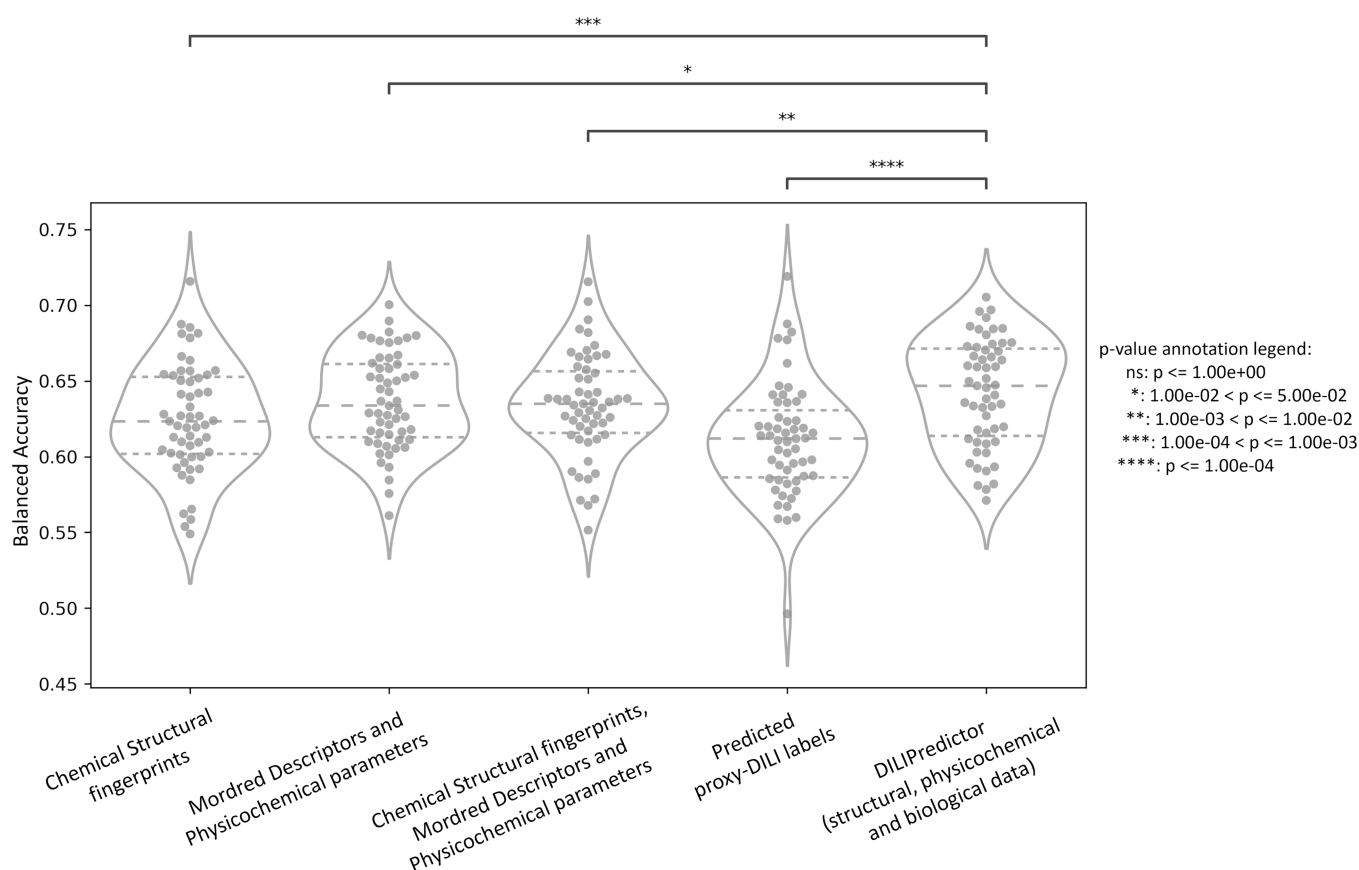


Figure 6. Performance metrics balanced accuracy for combination models from 55 held-out test sets from repeated nested cross-validation using (a) selected 189-bit structural fingerprints, (b) selected 668 molecular descriptors, (c) selected 189-bit structural fingerprints and selected 668 molecular descriptors, (d) predicted nine proxy-DILI labels and 2 PK parameters, and (e) a combination of all three feature spaces, compared with a paired *t* test.

Jaccard similarity to predictions from the Mitotox model and 0.02 Jaccard similarity to predictions from the reactive metabolite formation model). Overall, we conclude that each model built on a proxy-DILI label and the PK parameters was distinctive in its prediction (albeit insufficiently accurate on its own), thus providing complementary information on compounds' potential for DILI.

Models Combining Chemical Structure, Physicochemical Properties, PK Parameters and Predicted Proxy-DILI Data Outperform Individual Models. We next compared models built on combinations of proxy-DILI labels (including PK parameters), chemical structure, and physicochemical properties including Mordred descriptors (Table 2). When comparing results from 55 held-out test sets from the repeated nested cross-validation (as shown in Figure 6 with the comparison of differences in distribution using a paired *t* test), the models combining structural fingerprints, physicochemical properties, Mordred descriptors, PK parameters, and predicted proxy-DILI labels achieved a mean balanced accuracy (BA) of 0.64 (mean LR+ = 1.84), comparable to models using only physicochemical properties and Mordred descriptors with a mean BA of 0.64 (mean LR+ = 1.83) and models using structural fingerprints, physicochemical properties, and Mordred descriptors, which also achieved a mean BA of 0.63 (mean LR+ = 1.81). Models using only structural fingerprints achieved a mean BA of 0.63 (mean LR+ = 1.78) while models using only predicted proxy-DILI labels and PK parameters as features achieved a mean BA of 0.61

(mean LR+ = 1.77) in the nested cross-validation. Supplementary Figure S4 compares the distribution of positive predictive value for all model combinations using all feature sets (predicted proxy-DILI labels and PK parameters, structural features, and Mordred physicochemical descriptors).

We next retrained all hyperparameter-optimized models on the DILI training data (888 compounds) and evaluated the final models on the held-out DILI test set (223 compounds) rather than via cross validation as in the above analysis. The DILIPredictor model (combining all predicted proxy-DILI labels and PK parameters, structural features, and Mordred physicochemical descriptors) achieved an AUC = 0.63 (LR+ = 1.40) on the held-out DILI dataset (Table 2). The model using only proxy-DILI and PK parameters achieved an AUC = 0.67 (LR+ = 1.76). Other models achieved AUC = 0.62 (LR+ = 1.47) using structural, Mordred, and physicochemical descriptors, AUC = 0.54 (LR+ = 1.31) using chemical structural only, and AUC = 0.61 (LR+ = 1.32) for the model using Mordred and physicochemical descriptors.

One metric relevant in predictive safety/toxicology is the positive likelihood ratio⁷³ in the detection of toxic compounds with a lower false-positive rate. Improved detection with lower false positive rates aids in evaluating model performance across various threshold settings, shifting the focus from AUC as a singular statistical value to a more nuanced examination along the AUC-ROC curve from a false positive rate of 0 to 1. When predicting the first 29 true positive compounds (or approximately 13% of the 223 compounds in the held-out

Table 3. Previously Published Models Used in the Evaluation of Hepatotoxicity/Liver Injury (for Test Sets Only)

| model | features | compounds in train set | compounds in test set | test splitting strategy | balanced accuracy | AUC-ROC |
|------------------------|--|------------------------|-----------------------|------------------------------------|---------------------------------|---------------------------------|
| ensemble of RF and SVM | Molecular fingerprints | 1241 | 286 | external data set | 0.82 | 0.9 |
| Random Forests | imaging phenotypes and chemical descriptors | 346 | 41 | literature survey | 0.52 | 0.74 |
| ensemble models | molecular features, physicochemical properties | 1254 | 204 | literature survey | 0.72 | 0.73 |
| Random Forests | 2D molecular descriptors | 996 | 341 | external data set | 0.67 | 0.71 |
| Random Forests | 2D molecular descriptors | 996 | 921 | external data set | 0.57 | 0.59 |
| SVM | Morgan fingerprints | 923 | 49 | external data set | 0.67 | |
| SVM | predicted protein targets | 923 | 49 | external data set | 0.59 | |
| Random Forests | 0–2.5D molecular descriptors | 1075 | 554 | external data set | 0.77 | 0.81 |
| GA-SVM | 2D and 3D molecular descriptors | 3712 | 269 | proprietary data set | 0.64 | 0.68 |
| Random Forests | Morgan fingerprints | 845 | 362 | random splits (repeated 100 times) | | 0.75 |
| naïve Bayes | Morgan fingerprints | 336 | 84 | random split | 0.73 | 0.81 |
| SVM | MACCS keys | 1317 | 88 | external data set | 0.68 | 0.62 |
| rule-based | physicochemical properties and common toxicity mechanisms | | 200 | N/A | 0.32 | |
| average | | | | | 0.64 | 0.73 |
| Random Forests | structural, physicochemical, predicted in vitro, in vivo and PK parameters | 888 | 223 | scaffold-split | 0.59 (scaffold-split) | 0.63 (scaffold-split) |

Table 4. DILI Predictions for 14 Compounds Known to Cause DILI and 2 Compounds That Do Not Cause DILI in Humans (Not Used in Training Models in This Study) and Top 3 Proxy-DILI Labels Positively and Negatively Contributing to the Prediction

| compound name | DILI (literature) | DILI prediction | DILI probability | remarks | most contribution proxy-DILI end points to prediction | | |
|-----------------|-------------------|-----------------|------------------|--|---|----------------------------|-------------------------|
| | | | | | ranked 1 | ranked 2 | ranked 3 |
| 2-butoxyethanol | not toxic | not toxic | 0.56 | known DILI in mice, not in human | Mitotox | human hepatotoxicity | animal hepatotoxicity B |
| astaxanthin | not toxic | not toxic | 0.6 | vitamin A derivative; high structural similarity to retinoid but does not cause DILI | diverse DILI A | preclinical hepatotoxicity | diverse DILI C |
| enzalutamide | toxic | toxic | 0.82 | | preclinical hepatotoxicity | human hepatotoxicity | Mitotox |
| sitaxentan | toxic | toxic | 0.82 | withdrawn | preclinical hepatotoxicity | human hepatotoxicity | Mitotox |

test set), DILIPredictor achieved the highest LR+ score of 2.68 (25 toxic compounds correctly predicted out of 29 compounds, PPV = 0.86) compared to the structural model, which achieved an LR+ score of 1.65 (23 toxic compounds correctly predicted out of 29 compounds, PPV = 0.78). This improvement is mainly from being able to detect compounds at a wider range of structural similarity to training data (as shown in [Supplementary Figure S5](#) using the distribution of the top true positives detected with low false positive rates for each model). Overall, this shows that using all feature types in DILIPredictor allows for the detection of a greater number of toxic compounds with a low false-positive rate.

We subsequently compared our models to those reported in earlier publications. [Table 3](#) presents a selection of recent DILI prediction models that employ chemical features and biological data to predict liver toxicity. Since most previous studies did not emphasize likelihood ratios, and often not scaffold-based splits, it is not possible to compare LR+ scores; therefore, we can only make comparisons within the models developed in this study. It is important to note that the size, source, and consequently the quality of training and test data sets vary across the previous literature, rendering direct comparisons infeasible. In our study, the final DILIPredictor model achieved a AUC-ROC of 0.63 and AUC-PR of 0.79 on the held-out gold

standard DILI data set (223 compounds), which is lower than the average AUC-ROC(0.73) from prior studies. The balanced accuracy of DILIPredictor in this study, standing at 0.59, is lower compared to previous models, averaging 0.64. This difference is likely caused by a stringent scaffold-based split for the external test set used in this study, which decreases the numerical score but better mirrors the real-world drug discovery process. We also adopted a fixed held-out test set that is scaffold-split, in contrast to less rigorous random splits and external data sets used in other studies (which contain some similar compounds to training data). Furthermore, we delved beyond the statistical value of AUC-ROC to examine likelihood ratios and improved detection with low false-positive rates. This approach allows us to evaluate the quality of the AUC-ROC curve rather than reducing it to a single statistical value.

Feature Interpretation. We next used feature interpretation to analyze the chemical and biological mechanisms for compounds known to cause DILI. Four compounds are shown in [Table 4](#), of which two were known for their DILI⁷⁵ (namely, enzalutamide and sitaxentan) and two compounds that do not cause DILI in humans (namely, 2-butoxyethanol and astaxanthin). DILIPredictor could detect structural information relevant to causing DILI (four compounds shown in

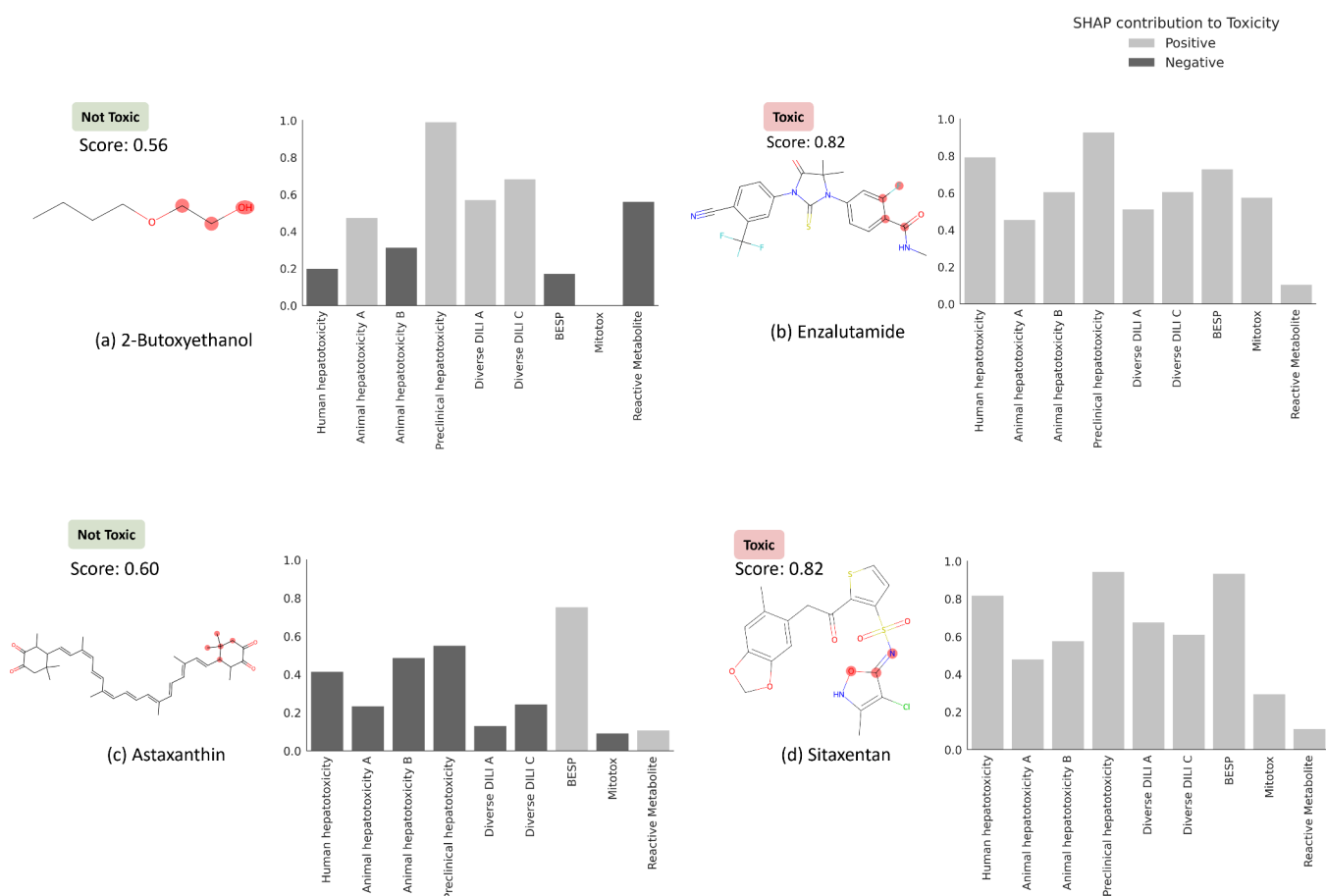


Figure 7. MACCS substructure (highlighted) and proxy-DILI labels contributing to DILI when using DILIPredictor (SHAP values) for two compounds known to cause DILI and for two compounds, which do not cause DILI in humans (further details in Table 4, and for another 12 compounds in Supplementary Figure S6 and Supplementary Table S6). The highest contribution to toxicity/safety from the MACCS substructure is highlighted with the chemical structure.

Figure 7 and Table 4 and Supplementary Figure S6). As shown in Figure 7, sitaxentan (a sulfonamide-based ETA receptor antagonist) was predicted toxic, with a positive contribution from the MACCS substructure near the sulphonamide, which is known to cause human liver injuries.⁷⁷ The MACCS feature most contributing to the toxicity for paclitaxel, docetaxel, and cabazitaxel, (which contain a taxane group as shown in Supplementary Figure S6) was found to be the presence of a taxadiene core. These compounds stabilize microtubules by binding to the β -tubulin and are known to cause mitochondrial toxicity.⁷⁸ Further, DILIPredictor correctly predicted compounds such as 2-butoxyethanol and astaxanthin to be nontoxic in humans even though they cause hepatic injury in animal models^{23,31} (Figure 7). In 2-butoxyethanol, proxy-DILI features associated with either animal hepatotoxicity or preclinical hepatotoxicity contributed to predicting toxicity in humans; however, the proxy-DILI indicators related to human hepatotoxicity ultimately led to the prediction of nontoxicity.

Finally, among structurally similar pairs of compounds, acitretin was correctly predicted as toxic while astaxanthin was correctly predicted to be nontoxic. For acitretin, the preclinical hepatotoxicity label contributed to the toxicity prediction. Conversely, labels associated with human hepatotoxicity contributed to correctly predicting astaxanthin as nontoxic. Among tetracyclines, pairs of compounds doxycycline (prediction scores = 0.66) and minocycline (0.66), and among fluoroquinolones, pairs of compounds moxifloxacin

(0.74) and trovafloxacin (0.84) were correctly predicted toxic. For fluoroquinolones, the prediction scores obtained from DILIPredictor were in agreement with the less-toxic or more-toxic DILI annotations collated by Chen et al.⁷⁶ Among compounds withdrawn from market, sitaxentan and trovafloxacin were flagged with prediction scores above 0.80 threshold; many compounds currently on the market such as docetaxel and paclitaxel were also flagged in the 0.70 to 0.75 threshold as being DILI-toxic. Overall, DILIPredictor combined chemical structures and biological data to correctly predict DILI in humans.

Limitations and Future Directions. The primary focus of this study was the generation of binary classification models for drug-induced liver injury, and this is empirically based on known data sets. Like most previously published models, empirical models can still be very useful for decision making, provided they are well predictive, in particular for novel chemical space.^{79,80} Besides using predicted C_{\max} (unbound and total), we did not incorporate factors such as dose or time point into this study due to its scarcity in available public data. Labeling schemes are not always binary but sometimes include an “ambiguous” class (such as in the DILIRank data set), and these compounds are hence not included in this study. While *in vitro* data can provide valuable insights into drug toxicities, they are still proxy end points for the *in vivo* effects. Toxic compounds detected in *in vitro* assays can often cause corresponding toxicity *in vivo*, but compounds that appear to

be safe in *in vitro* are not necessarily safe in humans.^{29,30} In the future, with the availability of larger relevant -omics data sets, such as from the Omics for Assessing Signatures for Integrated Safety Consortium (OASIS) Consortium,⁸¹ multitask learning models such as Conditional Neural Processes can be implemented, which can train on several predictive tasks simultaneously by shared representation learning.⁸² In the future, with the availability of larger -omics data sets and histopathology readouts, it could be possible to explore mechanistic underpinnings regarding Molecular Initiating Events or Key Events pathways of toxicity.

CONCLUSIONS

In this work, we trained models to predict drug-induced liver injury (DILI) using not only chemical data but also heterogeneous biological *in vivo* (human and animal) and *in vitro* data from various sources. We found a strong concordance in observed data between compounds with the proxy-DILI labels and DILI compounds. The nine proxy-DILI models were not predictive of each other – this complementarity suggests that they could be combined to predict drug-induced liver injury. Random Forest models that combined different types of input data – structural fingerprints, physicochemical properties, PK properties, and proxy-DILI labels – improved predictive performance, especially in detection with low false positive rates, with the highest LR+ score of 2.68 (25 toxic compounds with PPV = 0.93). DILIPredictor accurately predicted the toxicity of various compounds known to cause DILI, including 14 notorious DILI-inducing compounds, by recognizing chemical structure as well as biological mechanisms. DILIPredictor was further able to differentiate between animal and human sensitivity for DILI and exhibited a potential for mechanism evaluation for these compounds. DILIPredictor was trained on existing *in vitro* and *in vivo* data. Using only chemical structures as the input and a FeatureNet approach, it can predict DILI risk more accurately, thus aiding decision-making for new compounds before new *in vivo* toxicology and pharmacokinetic data is collected (which typically involves animal experiments that we aim to reduce). DILIPredictor can also be integrated into Design-Make-Test-Analyze (DMTA) cycles to aid in the selection and modification of compounds before more extensive and expensive testing is conducted.⁸³ Overall, the study demonstrated that incorporating all complementary sources of information can significantly improve the accuracy of DILI prediction models. Furthermore, the availability of larger, high-quality, and standardized data sets for DILI in the public domain can greatly enhance the development of predictive models for drug-induced liver injury such as from the Omics for Assessing Signatures for Integrated Safety Consortium (OASIS).⁸¹ DILIPredictor required only chemical structures as input for prediction. We released our final interpretable models at (with all code available for download at GitHub at <https://github.com/srijitseal/DILI>) and data sets used in this study at <https://broad.io/DILIPredictor>. Further, DILI Predictor is available for direct implementations via <https://pypi.org/project/dilipred/> and installable via “pip install dilipred”.

ASSOCIATED CONTENT

Supporting Information

The Supporting Information is available free of charge at <https://pubs.acs.org/doi/10.1021/acs.chemrestox.4c00015>.

Figure S1: Distribution of 12,133 compounds in proxy-DILI dataset in each of the 9 labels as represented by the t-SNE plot of physicochemical space; Figure S2: Distribution of 12,133 compounds in proxy-DILI dataset and 1,111 in the gold standard DILI dataset as represented by the t-SNE plot of physicochemical space; Figure S3: Distribution of the held-out DILI test set (223 compounds) and the training DILI data (888 compounds) as represented by the t-SNE plot of physicochemical space; Figure S4: Performance metric positive predictive value for combination models from 55 held-out test sets from repeated nested cross-validation; Figure S5: Distribution of 3-nearest neighbour Tanimoto similarity for toxic compounds detected correctly in the early stage (from the top 29 compounds) with a low false positive rate for models using only structural features compared to DILIPredictor models using all feature spaces; Figure S6: MACCS substructure positively contributing to DILI (highlighted) and the proxy-DILI predictions (PDF)

Data sets used in the study (XLSX)

Accession Codes

S.S. designed and performed data analysis, implemented, and trained the models. S.S. analyzed the interpretation of features and the results of models. M.M. assisted with packaging and implementing DILIPred for local implementation. S.S. wrote the manuscript. A.B. and O.S. supervised the project. All the authors (S.S., D.P.W., L.H.G., M.M., A.E.C., O.S., and A.B.) reviewed, edited, contributed to discussions, and approved the final version of the manuscript.

AUTHOR INFORMATION

Corresponding Authors

Ola Spjuth – Department of Pharmaceutical Biosciences and Science for Life Laboratory, Uppsala University, Uppsala SE-75124, Sweden; Email: ola.spjuth@uu.se

Andreas Bender – Yusuf Hamied Department of Chemistry, University of Cambridge, Cambridge CB2 1EW, United Kingdom; orcid.org/0000-0002-6683-7546; Email: ab454@cam.ac.uk

Authors

Srijit Seal – Yusuf Hamied Department of Chemistry, University of Cambridge, Cambridge CB2 1EW, United Kingdom; Imaging Platform, Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02141, United States; orcid.org/0000-0003-2790-8679

Dominic Williams – Safety Innovation, Clinical Pharmacology and Safety Sciences, AstraZeneca, Cambridge CB4 0FZ, United Kingdom; Quantitative Biology, Discovery Sciences, R&D, AstraZeneca, Cambridge CB4 0FZ, United Kingdom

Layla Hosseini-Gerami – Ignota Laboratories, County Hall, London SE1 7PB, United Kingdom

Manas Mahale – Bombay College of Pharmacy Kalina Santacruz (E), Mumbai 400 098, India; orcid.org/0009-0007-3867-996X

Anne E. Carpenter – Imaging Platform, Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02141, United States; orcid.org/0000-0003-1555-8261

Complete contact information is available at: <https://pubs.acs.org/doi/10.1021/acs.chemrestox.4c00015>

Author Contributions

CRedit: **Srijit Seal** conceptualization, data curation, formal analysis, investigation, methodology, resources, software, validation, visualization, writing-original draft, writing-review & editing; **Dominic Williams** methodology, supervision, writing-review & editing; **Layla Hosseini-Gerami** resources, writing-review & editing; **Manas Mahale** software; **Anne E. Carpenter** writing-review & editing; **Ola Spjuth** supervision, writing-review & editing; **Andreas Bender** supervision, writing-review & editing.

Notes

The authors declare no competing financial interest.

All authors have approved the final version of the manuscript.

ACKNOWLEDGMENTS

S.S. acknowledges funding from the Cambridge Commonwealth, European and International Trust, Boak Student Support Fund (Clare Hall), Jawaharlal Nehru Memorial Fund, Allen, Meek and Read Fund, and Trinity Henry Barlow (Trinity College). O.S. acknowledges funding from the Swedish Research Council (grants 2020-03731 and 2020-01865), FORMAS (grant 2022-00940), Swedish Cancer Foundation (22 2412 Pj), and Horizon Europe grant agreement #101057014 (PARC) and #101057442 (REMED4ALL). S.S. and A.E.C. acknowledges funding from the National Institutes of Health (NIH MIRA R35 GM122547 to A.E.C.), the Massachusetts Life Sciences Center Bits to Bytes Capital Call program for funding the data analysis (to Shantanu Singh, Broad Institute of MIT and Harvard), and the OASIS Consortium organised by HESI (OASIS to Shantanu Singh). This work was performed using resources provided by the Cambridge Service for Data-Driven Discovery (CSD3) operated by the University of Cambridge Research Computing Service (www.csd3.cam.ac.uk), provided by Dell EMC and Intel using Tier-2 funding from the Engineering and Physical Sciences Research Council (capital grant EP/P020259/1), and DiRAC funding from the Science and Technology Facilities Council (www.dirac.ac.uk).

REFERENCES

- (1) Andrade, R. Jr; Chalasani, N.; Björnsson, E. S.; Suzuki, A.; Kullak-Ublick, G. A.; Watkins, P. B.; Devarbhavi, H.; Merz, M.; Lucena, M. I.; Kaplowitz, N.; Aithal, G. P. Drug-Induced Liver Injury. *Nat. Rev. Dis. Prim.* **2019**, *5* (1), 1–22.
- (2) Remmer, H. The Role of the Liver in Drug Metabolism. *Am. J. Med.* **1970**, *49* (5), 617–629.
- (3) Licata, A. Adverse Drug Reactions and Organ Damage: The Liver. *Eur. J. Int. Med.* **2016**, *28*, 9–16.
- (4) Ostapowicz, G.; Fontana, R. J.; Schiødt, F. V.; Larson, A.; Davern, T. J.; Han, S. H. B.; McCashland, T. M.; Shakil, A. O.; Hay, J. E.; Hynan, L.; Crippin, J. S.; Blei, A. T.; Samuel, G.; Reich, J.; Lee, W. M.; Santyanarayana, R.; Caldwell, C.; Shick, L.; Bass, N.; Rouillard, S.; Atillasoy, E.; Flamm, S.; Benner, K. G.; Rosen, H. R.; Martin, P.; Stribling, R.; Schiff, E. R.; Torres, M. B.; Navarro, V.; McGuire, B.; Chung, R.; Abraczinskas, D.; Dienstag, J. Results of a Prospective Study of Acute Liver Failure at 17 Tertiary Care Centers in the United States. *Ann. Int. Med.* **2002**, *137* (12), 947–954.
- (5) Onakpoya, I. J.; Heneghan, C. J.; Aronson, J. K. Post-Marketing Withdrawal of 462 Medicinal Products Because of Adverse Drug Reactions: A Systematic Review of the World Literature. *BMC Med.* **2016**, *14* (1), 1–11.
- (6) Raschi, E.; De Ponti, F. Strategies for Early Prediction and Timely Recognition of Drug-Induced Liver Injury: The Case of Cyclin-Dependent Kinase 4/6 Inhibitors. *Front. Pharmacol.* **2019**, *10* (OCT), 1235.

- (7) Weaver, R. J.; Blomme, E. A.; Chadwick, A. E.; Copple, I. M.; Gerets, H. H. J.; Goldring, C. E.; Guillouzo, A.; Hewitt, P. G.; Ingelman-Sundberg, M.; Jensen, K. G.; Juhola, S.; Klingmüller, U.; Labbe, G.; Liguori, M. J.; Lovatt, C. A.; Morgan, P.; Naisbitt, D. J.; Pieters, R. H. H.; Snoeys, J.; van de Water, B.; Williams, D. P.; Park, B. K. Managing the Challenge of Drug-Induced Liver Injury: A Roadmap for the Development and Deployment of Preclinical Predictive Models. *Nat. Rev. Drug Discovery* **2020**, *19* (2), 131–148.
- (8) Mihajlovic, M.; Vinken, M. Mitochondria as the Target of Hepatotoxicity and Drug-Induced Liver Injury: Molecular Mechanisms and Detection Methods. *Int. J. Mol. Sci.* **2022**, *23* (6), 3315.
- (9) Kenna, J. G.; Taskar, K. S.; Battista, C.; Bourdet, D. L.; Brouwer, K. L. R.; Brouwer, K. R.; Dai, D.; Funk, C.; Hafey, M. J.; Lai, Y.; Maher, J.; Pak, Y. A.; Pedersen, J. M.; Polli, J. W.; Rodrigues, A. D.; Watkins, P. B.; Yang, K.; Yucha, R. W. Can Bile Salt Export Pump Inhibition Testing in Drug Discovery and Development Reduce Liver Injury Risk? An International Transporter Consortium Perspective. *Clin. Pharmacol. Ther.* **2018**, *104* (5), 916–932.
- (10) Villanueva-Paz, M.; Morán, L.; López-Alcántara, N.; Freixo, C.; Andrade, R. J.; Lucena, M. I.; Cubero, F. J. Oxidative Stress in Drug-Induced Liver Injury (DILI): From Mechanisms to Biomarkers for Use in Clinical Practice. *Antioxidants* **2021**, *10* (3), 390–35.
- (11) Bao, Y.; Wang, P.; Shao, X.; Zhu, J.; Xiao, J.; Shi, J.; Zhang, L.; Zhu, H. J.; Ma, X.; Manautou, J. E.; Zhong, X. B. Acetaminophen-Induced Liver Injury Alters Expression and Activities of Cytochrome P450 Enzymes in an Age-Dependent Manner in Mouse Liver. *Drug Metab. Dispos.* **2020**, *48* (5), 326–336.
- (12) Johansson, J.; Larsson, M. H.; Hornberg, J. J. Predictive in Vitro Toxicology Screening to Guide Chemical Design in Drug Discovery. *Curr. Opin. Toxicol.* **2019**, *15*, 99–108.
- (13) Hartung, T. The (Misleading) Role of Animal Models in Drug Development. *Front. Drug Discov.* **2024**, *4*, 1355044.
- (14) Biomea Fusion Announces BMF-219 in Diabetes Placed on Clinical Hold <https://www.globenewswire.com/news-release/2024/06/06/2895065/0/en/Biomea-Fusion-Announces-BMF-219-in-Diabetes-Placed-on-Clinical-Hold.html> 2024.
- (15) RAPT Therapeutics Announces Clinical Hold on Studies Evaluating Zel neciron <https://www.globenewswire.com/news-release/2024/02/20/2831721/0/en/RAPT-Therapeutics-Announces-Clinical-Hold-on-Studies-Evaluating-Zel-neciron.html>.
- (16) Tango Therapeutics Announces Discontinuation of TNG348 Program | Business Wire <https://www.businesswire.com/news/home/20240523266195/en/Tango-Therapeutics-Announces-Discontinuation-of-TNG348-Program> 2024.
- (17) Rathman, J.; Yang, C.; Ribeiro, J. V.; Mostrag, A.; Thakkar, S.; Tong, W.; Hobocienski, B.; Sacher, O.; Magdziarz, T.; Bienfait, B. Development of a Battery of in Silico Prediction Tools for Drug-Induced Liver Injury from the Vantage Point of Translational Safety Assessment. *Chem. Res. Toxicol.* **2021**, *34* (2), 601–615.
- (18) Soldatow, V. Y.; Lecluyse, E. L.; Griffith, L. G.; Rusyn, I. In vitro Models for Liver Toxicity Testing. *Toxicol. Res. (Camb.)* **2013**, *2* (1), 23–39.
- (19) Meng, Q. Three-Dimensional Culture of Hepatocytes for Prediction of Drug-Induced Hepatotoxicity. *Expert Opin. Drug Metab. Toxicol.* **2010**, *6* (6), 733–746.
- (20) Sison-Young, R. L. C.; Mitsa, D.; Jenkins, R. E.; Mottram, D.; Alexandre, E.; Richert, L.; Aerts, H.; Weaver, R. J.; Jones, R. P.; Johann, E.; Hewitt, P. G.; Ingelman-Sundberg, M.; Goldring, C. E. P.; Kitteringham, N. R.; Park, B. K. Comparative Proteomic Characterization of 4 Human Liver-Derived Single Cell Culture Models Reveals Significant Variation in the Capacity for Drug Disposition, Bioactivation, and Detoxication. *Toxicol. Sci.* **2015**, *147* (2), 412–424.
- (21) Steger-Hartmann, T.; Raschke, M. Translating in Vitro to in Vivo and Animal to Human. *Curr. Opin. Toxicol.* **2020**, *23-24*, 6–10.
- (22) Kindrat, L.; Dreval, K.; Shpileva, S.; Tryndyak, V.; de Conti, A.; Mudalige, T. K.; Chen, T.; Erstenyuk, A. M.; Beland, F. A.; Pogribny, I. P. Effect of Methapyrilene Hydrochloride on Hepatic Intracellular Iron Metabolism *in vivo* and *in vitro*. *Toxicol. Lett.* **2017**, *281*, 65–73.

- (23) Graham, E. E.; Walsh, R. J.; Hirst, C. M.; Maggs, J. L.; Martin, S.; Wild, M. J.; Wilson, I. D.; Harding, J. R.; Kenna, J. G.; Peter, R. M.; Williams, D. P.; Park, B. K. Identification of the Thiophene Ring of Methapyrilene as a Novel Bioactivation-Dependent Hepatic Toxicophore. *J. Pharmacol. Exp. Ther.* **2008**, *326* (2), 657–671.
- (24) Hamadeh, H. K.; Knight, B. L.; Haugen, A. C.; Sieber, S.; Amin, R. P.; Bushel, P. R.; Stoll, R.; Blanchard, K.; Jayadev, S.; Tennant, R. W.; Cunningham, M. L.; Afshari, C. A.; Paules, R. S. Methapyrilene Toxicity: Anchorage of Pathologic Observations to Gene Expression Alterations. *Toxicol. Pathol.* **2002**, *30* (4), 470–482.
- (25) Mirsalis, J. C. Genotoxicity, Toxicity, and Carcinogenicity of the Antihistamine Methapyrilene (MTR 07216). *Mutat. Res. Genet. Toxicol.* **1987**, *185* (3), 309–317.
- (26) Wright, P. S. R.; Briggs, K. A.; Thomas, R.; Smith, G. F.; Maglennon, G.; Mikulskis, P.; Chapman, M.; Greene, N.; Phillips, B. U.; Bender, A. Statistical Analysis of Preclinical Inter-Species Concordance of Histopathological Findings in the ETOX Database. *Regul. Toxicol. Pharmacol.* **2023**, *138*, No. 105308.
- (27) Olson, H.; Betton, G.; Robinson, D.; Thomas, K.; Monro, A.; Kolaja, G.; Lilly, P.; Sanders, J.; Sipes, G.; Bracken, W.; Dorato, M.; Van Deun, K.; Smith, P.; Berger, B.; Heller, A. Concordance of the Toxicity of Pharmaceuticals in Humans and in Animals. *Regul. Toxicol. Pharmacol.* **2000**, *32* (1), 56–67.
- (28) Fourches, D.; Barnes, J. C.; Day, N. C.; Bradley, P.; Reed, J. Z.; Tropsha, A. Cheminformatics Analysis of Assertions Mined from Literature That Describe Drug-Induced Liver Injury in Different Species. *Chem. Res. Toxicol.* **2010**, *23* (1), 171–183.
- (29) Bender, A.; Cortés-Ciriano, I. *Artificial Intelligence in Drug Discovery: What Is Realistic, What Are Illusions? Part 1: Ways to Make an Impact, and Why We Are Not There Yet.* *Drug Discovery Today*. Elsevier Ltd, 2021; pp 511–524. DOI: 10.1016/j.drudis.2020.12.009.
- (30) van Norman, G. A. Limitations of Animal Studies for Predicting Toxicity in Clinical Trials: Is It Time to Rethink Our Current Approach? *JACC Basic Transl. Sci.* **2019**, *4* (7), 845–854.
- (31) Cunningham, M. L. A Mouse Is Not a Rat Is Not a Human: Species Differences Exist. *Toxicol. Sci.* **2002**, *70* (2), 157–158.
- (32) Thakkar, S.; Li, T.; Liu, Z.; Wu, L.; Roberts, R.; Tong, W. Drug-Induced Liver Injury Severity and Toxicity (DILIst): Binary Classification of 1279 Drugs by Human Hepatotoxicity. *Drug Discovery Today* **2020**, *25* (1), 201–208.
- (33) Chen, M.; Suzuki, A.; Thakkar, S.; Yu, K.; Hu, C.; Tong, W. DILIrank: The Largest Reference Drug List Ranked by the Risk for Developing Drug-Induced Liver Injury in Humans. *Drug Discovery Today* **2016**, *21* (4), 648–653.
- (34) Lawrence, E.; El-Shazly, A.; Seal, S.; Joshi, C. K.; Liò, P.; Singh, S.; Bender, A.; Sormanni, P.; Greenig, M. Understanding Biology in the Age of Artificial Intelligence. *arXiv* **2024**, arXiv:2403.04106.
- (35) Hewitt, M.; Enoch, S. J.; Madden, J. C.; Przybylak, K. R.; Cronin, M. T. D. Hepatotoxicity: A Scheme for Generating Chemical Categories for Read-across, Structural Alerts and Insights into Mechanism(s) of Action. *Crit. Rev. Toxicol.* **2013**, *43* (7), 537–558.
- (36) Ye, L.; Ngan, D. K.; Xu, T.; Liu, Z.; Zhao, J.; Sakamuru, S.; Zhang, L.; Zhao, T.; Xia, M.; Simeonov, A.; Huang, R. Prediction of Drug-Induced Liver Injury and Cardiotoxicity Using Chemical Structure and *in vitro* Assay Data. *Toxicol. Appl. Pharmacol.* **2022**, *454*, No. 116250.
- (37) Liu, A.; Walter, M.; Wright, P.; Bartosik, A.; Dolciemi, D.; Elbasir, A.; Yang, H.; Bender, A. Prediction and Mechanistic Analysis of Drug-Induced Liver Injury (DILI) Based on Chemical Structure. *Biol. Direct* **2021**, *16* (1), 1–15.
- (38) Mora, J. R.; Marrero-Ponce, Y.; García-Jacas, C. R.; Causado, A. S. Ensemble Models Based on QuBiLS-MAS Features and Shallow Learning for the Prediction of Drug-Induced Liver Toxicity: Improving Deep Learning and Traditional Approaches. *Chem. Res. Toxicol.* **2020**, *33* (7), 1855–1873.
- (39) Liu, A.; Seal, S.; Yang, H.; Bender, A. Using Chemical and Biological Data to Predict Drug Toxicity. *SLAS Discovery* **2023**, *28*, 53.
- (40) Seal, S.; Yang, H.; Trapotsi, M.-A.; Singh, S.; Carreras-Puigvert, J.; Spjuth, O.; Bender, A. Merging Bioactivity Predictions from Cell Morphology and Chemical Fingerprint Models Using Similarity to Training Data. *J. Cheminform.* **2023**, *15* (1), 56.
- (41) Kaito, S.; Takeshita, J.; Iwata, M.; Sasaki, T.; Hosaka, T.; Shizu, R.; Yoshinari, K. Utility of Human Cytochrome P450 Inhibition Data in the Assessment of Drug-Induced Liver Injury. *Xenobiotica* **2024**, 1–30.
- (42) Rao, M.; Nassiri, V.; Alhambra, C.; Snoeys, J.; Van Goethem, F.; Irrechukwu, O.; Aleo, M. D.; Geys, H.; Mitra, K.; Will, Y. AI/ML Models to Predict the Severity of Drug-Induced Liver Injury for Small Molecules. *Chem. Res. Toxicol.* **2023**, *36*, 1129–1139.
- (43) Aleo, M. D.; Shah, F.; Allen, S.; Barton, H. A.; Costales, C.; Lazzaro, S.; Leung, L.; Nilson, A.; Obach, R. S.; Rodrigues, A. D.; Will, Y. Moving beyond Binary Predictions of Human Drug-Induced Liver Injury (DILI) toward Contrasting Relative Risk Potential. *Chem. Res. Toxicol.* **2020**, *33* (1), 223–238.
- (44) Chavan, S.; Scherbak, N.; Engwall, M.; Repsilber, D. Predicting Chemical-Induced Liver Toxicity Using High-Content Imaging Phenotypes and Chemical Descriptors: A Random Forest Approach. *Chem. Res. Toxicol.* **2020**, *33* (9), 2261–2275.
- (45) Seal, S.; Carreras-Puigvert, J.; Trapotsi, M. A.; Yang, H.; Spjuth, O.; Bender, A. Integrating Cell Morphology with Gene Expression and Chemical Structure to Aid Mitochondrial Toxicity Detection. *Commun. Biol.* **2022**, *5* (1), 858.
- (46) Seal, S.; Yang, H.; Vollmers, L.; Bender, A. Comparison of Cellular Morphological Descriptors and Molecular Fingerprints for the Prediction of Cytotoxicity- And Proliferation-Related Assays. *Chem. Res. Toxicol.* **2021**, *34* (2), 422–437.
- (47) Seal, S.; Spjuth, O.; Hosseini-Gerami, L.; García-Ortegón, M.; Singh, S.; Bender, A.; Carpenter, A. E. Insights into Drug Cardiotoxicity from Biological and Chemical Data: The First Public Classifiers for FDA Drug-Induced Cardiotoxicity Rank. *J. Chem. Inf. Model.* **2024**, *64*, 1172.
- (48) Mulliner, D.; Schmidt, F.; Stolte, M.; Spirkl, H. P.; Czich, A.; Amberg, A. Computational Models for Human and Animal Hepatotoxicity with a Global Application Scope. *Chem. Res. Toxicol.* **2016**, *29* (5), 757–767.
- (49) Liu, J.; Mansouri, K.; Judson, R. S.; Martin, M. T.; Hong, H.; Chen, M.; Xu, X.; Thomas, R. S.; Shah, I. Predicting Hepatotoxicity Using ToxCast *in vitro* Bioactivity and Chemical Structure. *Chem. Res. Toxicol.* **2015**, *28* (4), 738–751.
- (50) Ambe, K.; Ishihara, K.; Ochibe, T.; Ohya, K.; Tamura, S.; Inoue, K.; Yoshida, M.; Tohkin, M. In Silico Prediction of Chemical-Induced Hepatocellular Hypertrophy Using Molecular Descriptors. *Toxicol. Sci.* **2018**, *162* (2), 667–675.
- (51) He, S.; Ye, T.; Wang, R.; Zhang, C.; Zhang, X.; Sun, G.; Sun, X. An In Silico Model for Predicting Drug-Induced Hepatotoxicity. *Int. J. Mol. Sci.* **2019**, *20* (8), 1897.
- (52) Hemmerich, J.; Troger, F.; Füzi, B.; Ecker, G. F. Using Machine Learning Methods and Structural Alerts for Prediction of Mitochondrial Toxicity. *Mol. Inf.* **2020**, *39* (5), 2000005 DOI: 10.1002/minf.202000005.
- (53) McLoughlin, K. S.; Jeong, C. G.; Sweitzer, T. D.; Minnich, A. J.; Tse, M. J.; Bennion, B. J.; Allen, J. E.; Calad-Thomson, S.; Rush, T. S.; Brase, J. M. Machine Learning Models to Predict Inhibition of the Bile Salt Export Pump. *J. Chem. Inf. Model.* **2021**, *61* (2), 587–602.
- (54) Mazzolari, A.; Vistoli, G.; Testa, B.; Pedretti, A. Prediction of the Formation of Reactive Metabolites by A Novel Classifier Approach Based on Enrichment Factor Optimization (EFO) as Implemented in the VEGA Program. *Molecules* **2018**, *23* (11), 2955.
- (55) Zhang, C.; Zhao, Y.; Yu, M.; Qin, J.; Ye, B.; Wang, Q. Mitochondrial Dysfunction and Chronic Liver Disease. *Curr. Issues Mol. Biol.* **2022**, *44* (7), 3156–3165.
- (56) Kenna, J. G.; Taskar, K. S.; Battista, C.; Bourdet, D. L.; Brouwer, K. L. R.; Brouwer, K. R.; Dai, D.; Funk, C.; Hafey, M. J.; Lai, Y.; Maher, J.; Pak, Y. A.; Pedersen, J. M.; Polli, J. W.; Rodrigues, A. D.; Watkins, P. B.; Yang, K.; Yucha, R. W. Can Bile Salt Export Pump Inhibition Testing in Drug Discovery and Development Reduce Liver

Injury Risk? An International Transporter Consortium Perspective. *Clin. Pharmacol. Ther.* **2018**, *104* (5), 916–932.

(57) Attia, S. M. Deleterious Effects of Reactive Metabolites. *Oxid. Med. Cell. Longevity* **2010**, *3* (4), 238–253.

(58) Seal, S.; Trapotsi, M.-A.; Subramanian, V.; Spjuth, O.; Greene, N.; Bender, A.-D. PKSmart: An Open-Source Computational Model to Predict in Vivo Pharmacokinetics of Small Molecules. *bioRxiv* **2024**, 2024.02.02.578658.

(59) Varnek, A.; Gaudin, C.; Marcou, G.; Baskin, I.; Pandey, A. K.; Tetko, I. V. Inductive Transfer of Knowledge: Application of Multi-Task Learning and Feature Net Approaches to Model Tissue-Air Partition Coefficients. *J. Chem. Inf. Model.* **2009**, *49* (1), 133–144.

(60) U.S. Environmental Protection Agency: Washington, DC, 2010. *Release User-Friendly Web-Based Tool for Mining ToxRefDB*.

(61) Sakuratani, Y.; Zhang, H. Q.; Nishikawa, S.; Yamazaki, K.; Yamada, T.; Yamada, J.; Gerova, K.; Chankov, G.; Mekenyan, O.; Hayashi, M. Hazard Evaluation Support System (HESS) for Predicting Repeated Dose Toxicity Using Toxicological Categories. *SAR QSAR Environ. Res.* **2013**, *24* (5), 351–363.

(62) Williams, D. P.; Lazic, S. E.; Foster, A. J.; Semenova, E.; Morgan, P. Predicting Drug-Induced Liver Injury with Bayesian Machine Learning. *Chem. Res. Toxicol.* **2020**, *33* (1), 239–248.

(63) Horne, R. L.; Wilson-Godber, J.; Diaz, A. G.; Brotzak, Z. F.; Seal, S.; Gregory, R. C.; Possenti, A.; Chia, S.; Vendruscolo, M. Using Generative Modeling to Endow with Potency Initially Inert Compounds with Good Bioavailability and Low Toxicity. *J. Chem. Inf. Model.* **2024**, *64*, 590 DOI: 10.1021/acs.jcim.3c01777.

(64) Smit, I. A.; Afzal, A. M.; Allen, C. H. G.; Svensson, F.; Hanser, T.; Bender, A. Systematic Analysis of Protein Targets Associated with Adverse Events of Drugs from Clinical Trials and Postmarketing Reports. *Chem. Res. Toxicol.* **2021**, *34* (2), 365–384.

(65) RDKit (v2022.09.5) <https://www.rdkit.org/>.

(66) Ropp, P. J.; Kaminsky, J. C.; Yablonski, S.; Durrant, J. D. Dimorphite-DL: An Open-Source Program for Enumerating the Ionization States of Drug-like Small Molecules. *J. Cheminform.* **2019**, *11* (1), 1–8.

(67) scikit-learn: machine learning in Python — scikit-learn 1.2.0 documentation <https://scikit-learn.org/stable/index.html>.

(68) Rogers, D.; Hahn, M. Extended-Connectivity Fingerprints. *J. Chem. Inf. Model.* **2010**, *50* (5), 742–754.

(69) Durant, J. L.; Leland, B. A.; Henry, D. R.; Nourse, J. G. Reoptimization of MDL Keys for Use in Drug Discovery. *J. Chem. Inf. Comput. Sci.* **2002**, *42* (6), 1273–1280.

(70) Moriwaki, H.; Tian, Y. S.; Kawashita, N.; Takagi, T. Mordred: A Molecular Descriptor Calculator. *J. Cheminform.* **2018**, *10* (1), 1–14.

(71) The DeepChem Project — deepchem 2.7.2.dev documentation <https://deepchem.readthedocs.io/en/latest/>.

(72) Ramsudar, B.; Eastman, P.; Walters, P.; Pande, V. *Deep Learning for Life Sciences*. 2019, 222.

(73) Fierz, W.; Bossuyt, X. Likelihood Ratio Approach and Clinical Interpretation of Laboratory Tests. *Front. Immunol.* **2021**, *12*, 1170.

(74) SHAP. Welcome to the SHAP documentation — SHAP latest documentation <https://shap.readthedocs.io/en/latest/index.html>.

(75) Azad, A.; Chang, P.; Devuni, D.; Bichoupan, K.; Kesar, V.; Branch, A. D.; Oh, W. K.; Galsky, M. D.; Ahmad, J.; Odin, J. A. Real World Experience of Drug Induced Liver Injury in Patients Undergoing Chemotherapy. *J. Clin. Gastroenterol. Hepatol.* **2018**, *2* (3), 18 DOI: 10.21767/2575-7733.1000047.

(76) Chen, M.; Borlak, J.; Tong, W. A Model to Predict Severity of Drug-Induced Liver Injury in Humans. *Hepatology* **2016**, *64* (3), 931–940.

(77) Khalili, H.; Soudbakhsh, A.; Talasaz, A. H. Severe Hepatotoxicity and Probable Hepatorenal Syndrome Associated with Sulfadiazine. *Am. J. Heal. Pharm.* **2011**, *68* (10), 888–892.

(78) Galletti, E.; Magnani, M.; Renzulli, M. L.; Botta, M. Paclitaxel and Docetaxel Resistance: Molecular Mechanisms and Development of New Generation Taxanes. *ChemMedChem.* **2007**, *2* (7), 920–942.

(79) Badwan, B. A.; Liaropoulos, G.; Kyrodimos, E.; Skaltsas, D.; Tsigros, A.; Gorgoulis, V. G. Machine Learning Approaches to Predict Drug Efficacy and Toxicity in Oncology. *Cell Reports. Methods* **2023**, *3* (2), No. 100413.

(80) Barelier, S.; Eidam, O.; Fish, I.; Hollander, J.; Figaroa, F.; Nachane, R.; Irwin, J. J.; Shoichet, B. K.; Siegal, G. Increasing Chemical Space Coverage by Combining Empirical and Computational Fragment Screens. *ACS Chem. Biol.* **2014**, *9* (7), 1528–1535.

(81) HESI Insights - February 2023 - HESI - Health and Environmental Sciences Institute <https://hesiglobal.org/hesi-insights-february-2023/>.

(82) Garcia-Ortegon, M.; Singh, S.; Bender, A.; Bacallado, S. *Calibrated Prediction of Scarce Adverse Drug Reaction Labels with Conditional Neural Processes*; Garcia-Ortegon, M. et al. et Al.

(83) Plowright, A. T.; Johnstone, C.; Kihlberg, J.; Pettersson, J.; Robb, G.; Thompson, R. A. Hypothesis Driven Drug Design: Improving Quality and Effectiveness of the Design-Make-Test-Analyse Cycle. *Drug Discovery Today* **2012**, *17* (1–2), 56–62.

NOTE ADDED AFTER ASAP PUBLICATION

This paper was published ASAP on July 9, 2024, with errors in the Abstract paragraph, Tables 2 and 3. The corrected version was reposted July 23, 2024.