



Cite this: DOI: 10.1039/d5dd00483g

# Machine learning models for catalytic asymmetric reactions of simple alkenes: from enantioselectivity predictions to chemical insights

Ajnabiul Hoque,<sup>†a</sup> Nupur Jain,<sup>†a</sup> Divya Chenna<sup>a</sup> and Raghavan B. Sunoj<sup>b</sup>  <sup>\*ab</sup>

The increasing number of applications of machine learning (ML) in chemical catalysis has engendered considerable confidence in predicting reaction outcomes. Despite the successful applications of ML to high-throughput experimentation (HTE) datasets, extension to small real-world datasets prevalent in organic synthesis remained more difficult, primarily due to their imbalanced and sparse distribution. Herein, we present a new chemical reaction dataset curated from published literature that bears class imbalance (CI) with a skewness of  $-1.37$ . The reactions in focus belong to an important class of transition metal-catalysed asymmetric transformations of alkenes such as cyclopropanation, aziridination, and arylation. Such reactions are indispensable for the construction of three-membered structural motifs, a versatile building block found in complex bioactive molecules. In cognizance of the CI in the reaction outcome, measured in terms of enantiomeric excess (% ee), we employ the AttentiveFP-CI model to predict % ee. This class-imbalance aware graph-based model with an attention mechanism exhibits commendable performance, as evidenced by the root mean square error (RMSE) of  $9.80 \pm 1.40$ . Upon evaluation across various molecular representations of these reactions (OHE, fingerprints, SMILES, and graphs) and ML algorithms (DNN, T5Chem, Transformer, and MPNN), AttentiveFP-CI emerged as the best model distinguished by its minimal overfitting (train-test RMSE difference of 3.59, compared to up to 5.40 for other CI-aware models). When extended to other important reaction datasets such as *N,S*-acetylation, asymmetric hydrogenation of alkenes, and USPTO, improved predictions could be obtained by using AttentiveFP-CI. Furthermore, attention visualization identifies key atoms and substructures contributing to high enantioselectivity, offering valuable chemical insights for planning the synthesis of new molecular targets. Harnessing insights derived from ML models could serve as an efficient and cost-effective approach for expedited developments in asymmetric catalysis.

Received 4th November 2025  
Accepted 26th February 2026

DOI: 10.1039/d5dd00483g

rsc.li/digitaldiscovery

## 1 Introduction

In the domain of organic synthesis, efforts toward generating chiral molecules have remained a vibrant activity for several decades. Such developments were motivated by the utility of target compounds as drugs and pharmaceuticals.<sup>1,2</sup> Historically, the discovery and optimization of catalytic reactions have relied heavily on empirically gathered domain knowledge, the intuition of experienced chemists, and an inevitable series of trial-and-error explorations.<sup>3</sup> Such traditional and heuristic approaches might often encounter limitations while predicting the behaviour of specific substrates/catalysts or when fine-tuning reaction conditions to achieve a desired transformation.<sup>4,5</sup> The intricate interactions between various

parameters of a reaction, some obvious and some subtle, could alter the course of the reaction in unpredictable ways.<sup>6,7</sup> A change in reaction outcome, such as the yield or selectivity, could stem from a number of factors, such as the choice of catalysts, ligands, solvents, and additives, besides temperature and other parameters pertaining to the reaction conditions. Since yield and/or selectivity depend on multiple components, it becomes an inherently complex high-dimensional problem. This very complexity associated with reaction outcome predictions makes them an all the more interesting research problem to consider.

There have been a good number of previous efforts aimed at predicting reaction outcomes. The predictive capabilities of quantum chemically derived molecular descriptors have been exploited to build bespoke linear regression models for catalytic reactions.<sup>8,9</sup> Molecular descriptors such as charge, NMR chemical shifts, vibrational frequencies and intensities, Sterimol parameters, buried volumes, *etc.*, bearing electronic and steric features of the participating molecules, have served as useful inputs for modelling reactions.<sup>10,11</sup> The use of such descriptors

<sup>a</sup>Department of Chemistry, Indian Institute of Technology Bombay, Powai, Mumbai 400076, India. E-mail: sunoj@chem.iitb.ac.in

<sup>b</sup>Centre for Machine Intelligence and Data Science, Indian Institute of Technology Bombay, Powai, Mumbai 400076, India

<sup>†</sup> AH and NJ contributed equally.



comes with their own challenges such as higher computational cost, particularly in the case of complex molecular systems, and the requirement of annotation/curation by domain experts.<sup>12,13</sup> Given these challenges, an increasing interest in exploring alternative approaches became more prominent in the current literature. Modern ML algorithms, capable of handling complex and diverse reaction data, can offer promising solutions on this front.<sup>14,15</sup> In practice, when one attempts reaction optimization by way of changing various controllable parameters (as described above), sparsely distributed reaction data become available. It would be of interest to try and see whether such data could possibly present potential opportunities for applying suitable ML algorithms for reaction modelling.<sup>16–18</sup> The key advantage of an early ML intervention would be to help make an informed choice of substrates/catalysts/solvent during the reaction development phase.

The unprecedented growth in computational capabilities has rendered the applications of ML to chemical reactivity problems increasingly more feasible.<sup>19,20</sup> Deployment of incredibly complex language models such as BERT for yield prediction became possible in very recent years.<sup>21,22</sup> The use of hybrid graph neural networks on molecular graphs to derive features for selectivity predictions in the case of chiral phosphoric acid catalysed thiol addition to *N*-acyl imines is now available.<sup>23</sup> These contemporary ML models for reaction outcome prediction have offered robust performance on High Throughput Experimental (HTE) datasets.<sup>24–26</sup> Analyses have shown that HTE datasets, as used in many of the recent studies, exhibited reduced variability in data quality, high internal consistency and high fidelity.<sup>27,28</sup> Another aspect of the HTE settings is that exhaustive permutations of reactants/reagents are affordable under uniform reaction conditions. However, in real-life reaction development only a few combinations between the reactants and the associated conditions can be practically explored. For example, the dataset sourced from the AstraZeneca electronic laboratory notebooks (ELNs) potentially encompassed approximately 470 M possible combinations of reactants. However, in practice, only 1000 reactions were experimentally examined by engaging 340 aryl halides, 260 amines, 24 ligands, 15 bases, and 15 solvents for the Buchwald–Hartwig reaction.<sup>29</sup> In this context, we consider it highly timely to develop accurate ML models for small-sized reaction datasets with different distribution characteristics.

Recent years witnessed several successful applications of ML in predicting yields or enantioselectivities of various catalytic reactions such as Buchwald–Hartwig cross-coupling,<sup>24</sup> Lewis base-catalysed propargylation,<sup>30</sup>  $\beta$ -C–H activation,<sup>31</sup> asymmetric hydrogenation,<sup>32</sup> relay Heck,<sup>33</sup> Negishi cross-coupling,<sup>34</sup> and palladaelectro-catalyzed C–H annulation reactions.<sup>35</sup> Needless to say, most of these studies are early examples of implementing deep learning (DL) methods and were confined to only a few reaction types, leaving out a large family of important asymmetric catalytic reactions. One of the key reasons for such exclusions from ML studies can be traced to the lack of good datasets. One such important catalytic asymmetric reaction that has not received attention is shown in Scheme 1a, which employs simple alkenes as the core substrate. Alkenes are

abundant precursors that can participate in a wide array of reactions to provide valuable products. For example, under suitably chosen Cu/Pd catalytic conditions, alkenes can react with (a) diazoester to form cyclopropane, (b) aryl boronic acid to yield important 1,1-diaryl compounds, and (c) aliphatic or aromatic *N*-tosyloxycarbamates to access key structural motifs such as aziridines. This class of reaction holds promise as it can help synthesize stereochemically well-defined cyclopropanes and aziridines, which are key constituents in medicinal and agrochemical compounds.<sup>36–39</sup> A few representative examples bearing these substructures are shown in Scheme 1b to convey the significance of these small ring containing molecules.<sup>40–43</sup>

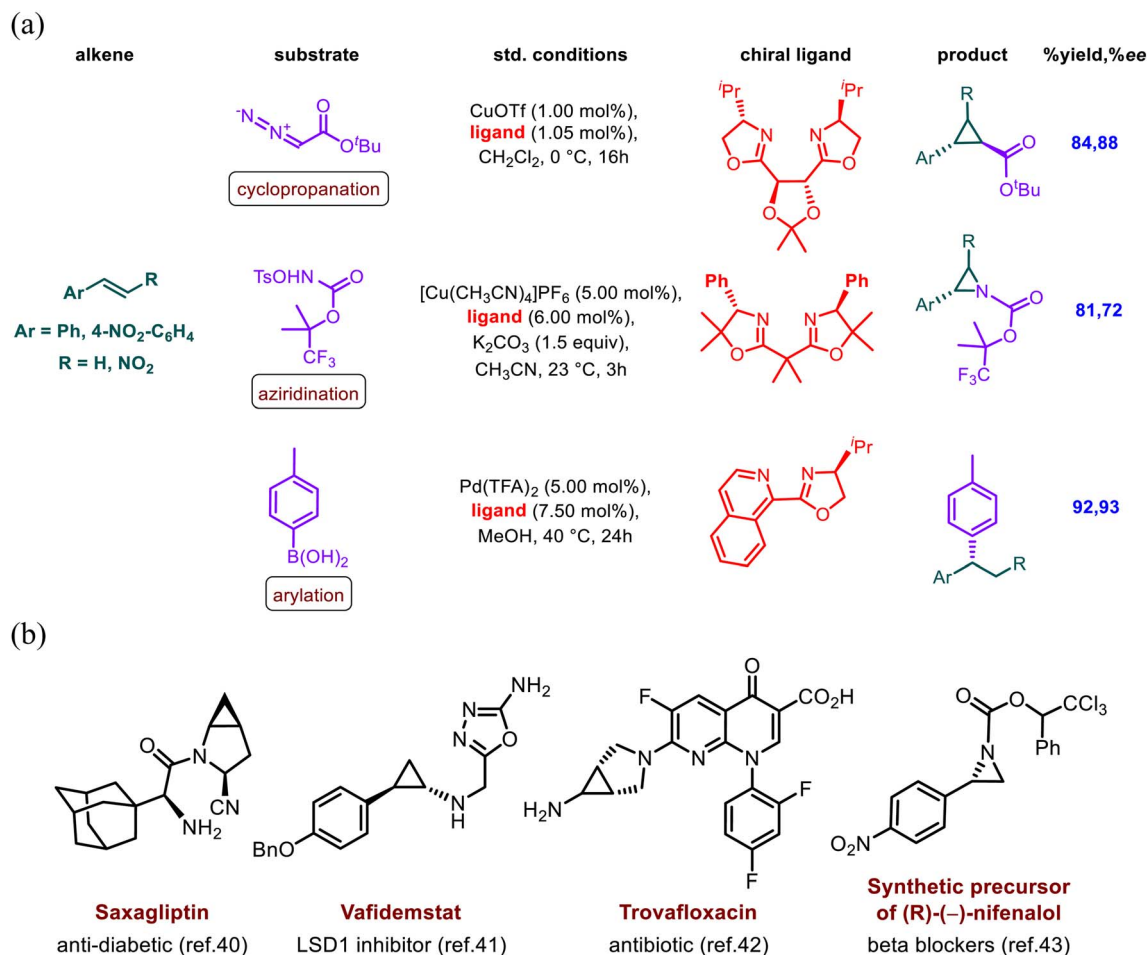
Apart from the synthetic utility of this class of reactions, the use of one of the most widely found ligands, such as chiral bis(oxazoline), is to be taken cognizance of, as the applications of such chiral motifs go well beyond these reactions.<sup>44–46</sup> The conformationally rigid framework of bis(oxazoline) metal chelates, bearing chiral centres close to the donor nitrogen atoms, can provide the desired chiral environment nearer to the catalytic site. The modular architecture of these ligands can help create desirable variations in both steric and electronic attributes, thus allowing fine-tuning of their catalytic activity for specific applications.<sup>47,48</sup> It would therefore be of importance to identify the key regions in the chiral catalyst that impact the stereochemical outcome of such reactions, potentially using ML tools (*vide infra*).

Motivated by recent advancements in machine learning approaches for reaction outcome prediction,<sup>49–53</sup> including contributions from our laboratory,<sup>31–33</sup> we became interested in the catalytic enantioselective reactions of alkenes as shown in Scheme 1a.<sup>54,55</sup> The availability of reliable predictive ML models can help identify optimal reactant triads comprising alkene, chiral ligand, and substrate that are likely to offer higher % *ee*. Such ML models might help reduce the typical timelines involved in reaction discovery. Given these motivations, we set the following major objectives in this work: (a) evaluation of the effectiveness of DL methods in enantioselectivity predictions in transition metal-catalysed asymmetric reactions of alkenes, (b) identification of an optimal featurisation strategy from among One-Hot Encoding (OHE), molecular fingerprints, SMILES, and graph representations, (c) addressing the issues associated with data imbalance consisting of more samples in the high % *ee* region by implementing cost-sensitive training loss, (d) identification of the better combinations of reactants (alkene, chiral ligand, and substrate) that are likely to offer superior reaction outcomes, and (e) examination of learning ability of DL models by using the attention mechanism to identify the critical regions in chiral ligands and substrates that can influence the reaction outcomes. Utilization of a trained DL model can streamline and help expedite the reaction discovery pipeline by identifying and eliminating low selectivity reactions in the initial screening, which would save time and effort.

## 2 Results and discussion

We have arranged the discussion into eight subsections. First, an overview of the reaction dataset is provided to shed light on





Scheme 1 (a) A representative example of catalytic asymmetric cyclopropanation, aziridination, and arylation of alkenes; (b) select examples of pharmaceuticals/drugs containing cyclopropane/aziridine structural motifs.

reaction-specific details such as the type and diversity of reactions considered and the distribution characteristics in the dataset. Next, we describe data encoding modality, followed by an outline of the AttentiveFP model. Subsequently, the model performance across different featurization methods and further refinements to the model to address the class imbalance (CI) issues is provided. Then, we discuss model interpretability by drawing parallels between graph attention and chemical intuition. The next section is on the prospects of exploiting the model interpretability and predictive abilities to identify optimal catalysts potentially useful for the synthesis of important drug molecules. Finally, we extend our AttentiveFP-CI framework to some of the commonly used reaction datasets as additional case studies—an asymmetric hydrogenation,<sup>32</sup> an *N,S*-acylation reaction,<sup>56</sup> and the USPTO<sup>57</sup> to evaluate the applicability of our model across different reaction classes.<sup>58</sup>

## 2.1. Reaction dataset

The dataset was manually curated by collecting reaction details from peer-reviewed publications, which led to a comprehensive collection of 376 reactions.<sup>59–65</sup> The three classes of catalytic asymmetric reactions of alkenes considered in this study are

cyclopropanation, aziridination, and arylation. The dataset is abbreviated as ART that stands for AlkeneReactionTriad to reflect the three classes of alkene reactions.<sup>66</sup> The reactions in the ART dataset differ in terms of ligands bound to the transition metal (Pd or Cu), alkenes, and substrates. Each sample represents a distinct combination of reactants and the corresponding % *ee* as the output value. The diversity of each reaction component can be gleaned from Fig. 1a. In transition metal–ligand complexes, the major difference arises from the decorations in the chiral ligands, ranging from bis(oxazoline) combined with Cu and pyridine-oxazolines bound to Pd or Cu (see Tables S1–S3 in Section 1 in the SI). Either of these two chiral ligand families, when bound to a suitable transition metal, can serve as the catalyst with intriguingly diverse steric and electronic characteristics. They form a sufficiently wide range of catalysts in active use today. Additionally, the reactions exhibit differences in the substrate involved, comprising 73 distinct entities, encompassing electron-rich to electron-deficient aryl boronic acids, aliphatic or aromatic *N*-tosyloxycarbamates, carbenoids, diazo compounds, diaryliodonium salts, and other species. The alkenes used in these reactions span a wide array of compounds, including mono- and di-substituted alkenes, alkyl/aryl-



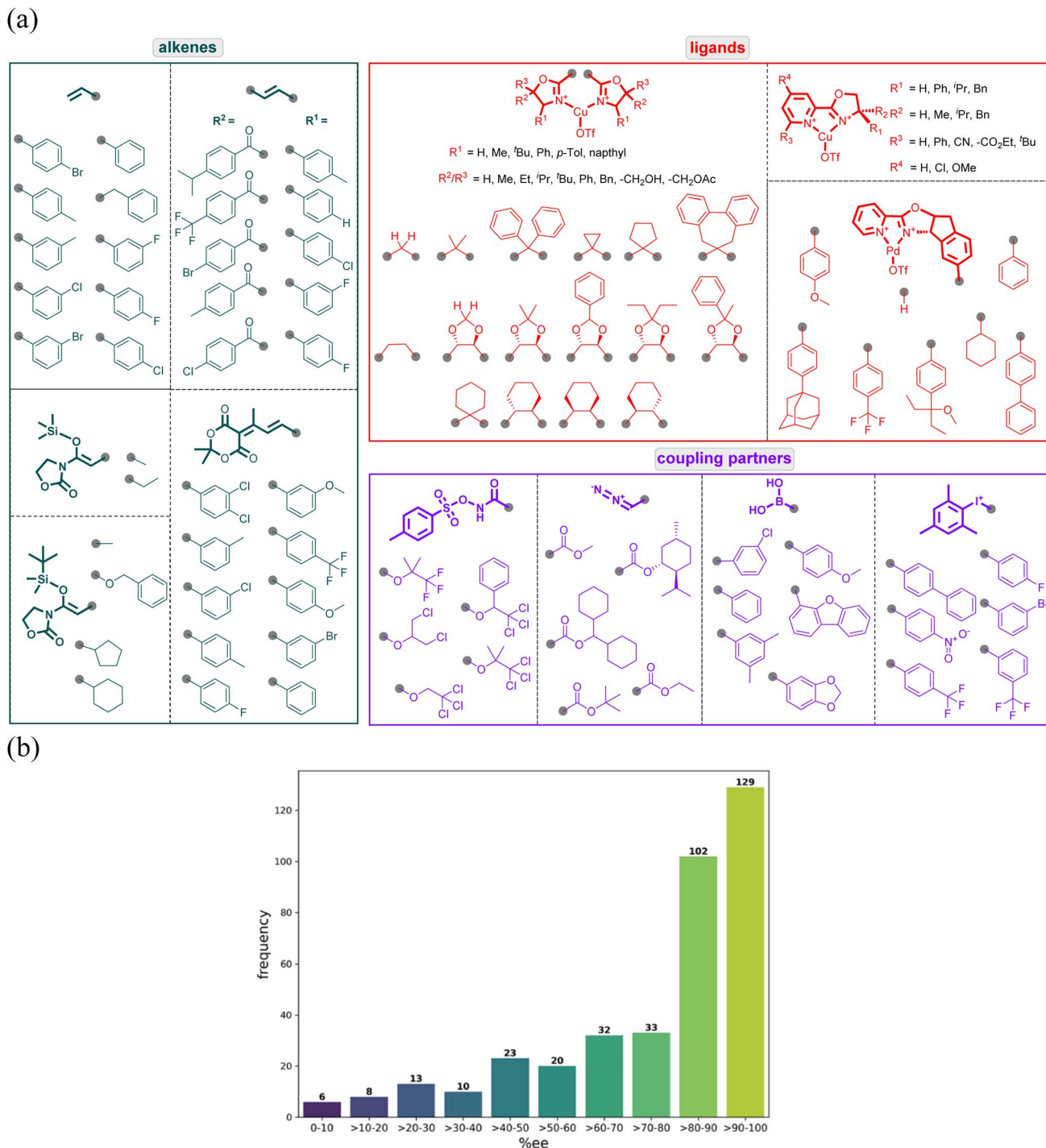


Fig. 1 (a) Details of various substituents present in the individual reacting partner; (b) yield distribution in the ART dataset.

substituted silylketenimides, and Meldrum's acid-derived dienes. All these diverse combinations across each reaction partner together make the ART dataset richly diverse and, hence, can be considered a representative of ground truth situations one would typically encounter in reaction discovery. Thus, the utility of the ML approach, designed for such a family of reactions with sparse and imbalanced data distribution (Fig. 1b), should remain good for other reactions as well.

The chemical space spanned by the ART dataset is very sparse compared to the combinatorial possibilities arising from the number of reactions between the compatible partners. For instance, in the cyclopropanation subset there are 130 examples (67 catalysts, 10 substrates, and 23 alkenes) while the aziridination reaction class comprises only 91 reactions (14 catalysts, 13 substrates, and 44 alkenes), and the arylation reaction class contains as many as 155 reactions (19 catalysts, 50 substrates,



and 55 alkenes). Given the possible combinations between the reactants, the theoretically likely reactions for cyclopropanation, aziridination, and arylation are 15 410, 52 250, and 8008, respectively, totaling 75 668 possibilities. However, the ART dataset contains only 376 experimentally known reactions from among these combinations, indicating a highly sparse distribution. Furthermore, distribution of the % *ee* values is also skewed towards the high *ee* regions (Fig. 1b).<sup>67,68</sup> The diversity of chemical structures, skewed distribution of reaction outcomes, and sparsity in the dataset can together make the ML model building rather challenging.

## 2.2. Data pre-processing and training protocol

Since each reaction variable can play a role in influencing the outcome to different extents, their meaningful representation is central to the input data. To achieve this, we encoded all the molecular partners as linear strings using the Simplified Molecular Input Line-Entry System (SMILES) notation.<sup>69</sup> These individual molecules are concatenated to form a complete reaction. Since the target % *ee* values are continuous, the ML task can be considered a regression problem. Within the ART dataset, each catalyst scaffold consistently delivers a preferred enantiomeric outcome, allowing the model to focus on predicting numerical % *ee* values. The dataset is randomly divided into training, validation, and test sets in a 70 : 10 : 20 ratio for training purposes. We have conducted hyperparameter tuning on the validation set based on the criterion of achieving the lowest mean validation loss and then employ such optimal

hyperparameters for prediction on the test set. To mitigate potential bias due to sample distribution, while creating the test-train splits (70 : 10 : 20), 30 independent runs with randomised splits were considered. The model performance is evaluated using root mean squared error (RMSE) on the test sets as the average RMSE over these runs.

## 2.3. AttentiveFP model

In this study, we use the AttentiveFP-based DL model as the primary model to predict the output of the reaction (*i.e.*, the enantioselectivity expressed in % *ee*).<sup>70</sup> This model uses molecular graphs as the input, where nodes and edges of the graphs respectively represent atoms and bonds. For a full representation of the reaction, the three major reaction entities *i.e.*, alkene, chiral ligand, and coupling partner, are concatenated to form a single composite graph. These molecular graphs remain disconnected within this composite representation, implying no explicit edges linking different reactants involved in the reaction. The bond features and atom features are obtained from the RDKit and encoded as described in Section 2.5 in the SI. The model utilizes a graph-attention mechanism to focus on the most relevant regions in the composite graph of the concatenated 'reactants'. Fig. 2 outlines the AttentiveFP network with an attention mechanism.

Since the model focuses on individual atoms, each atom includes features from its neighbouring atoms and the bonds connecting them to form the respective initial state vectors  $h_i^0$  for each atom. These initial vectors are further embedded

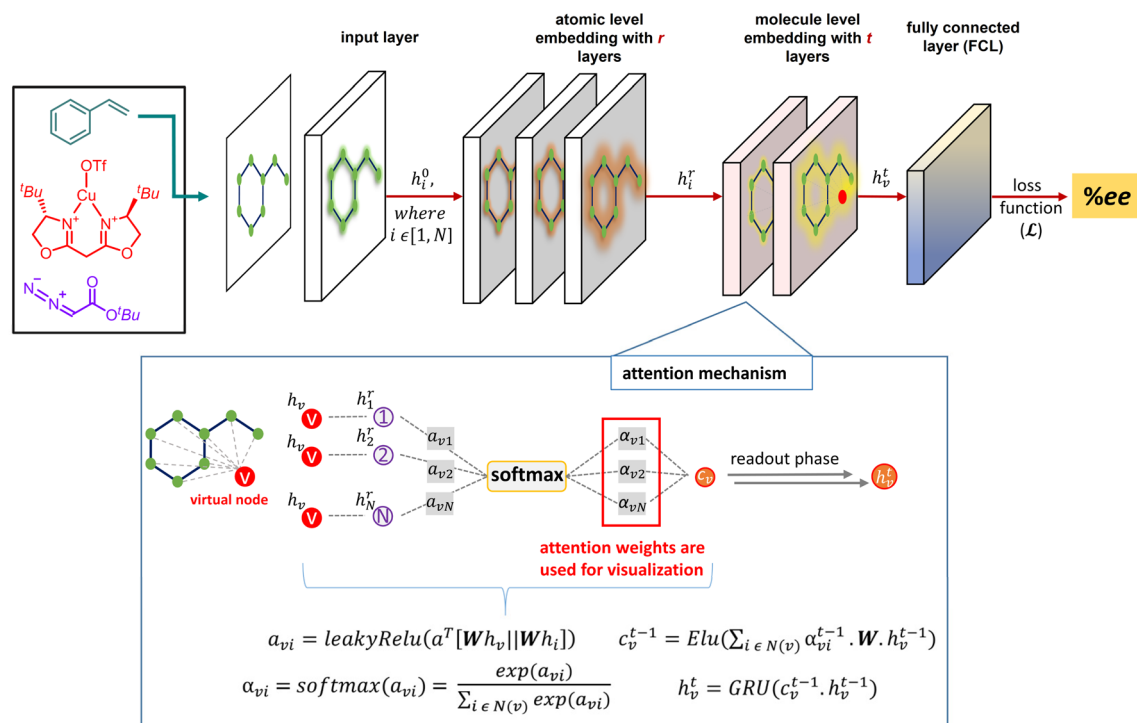


Fig. 2 Major components of the AttentiveFP model architecture. While the input sample is a composite molecular graph formed by concatenating all three reaction components (alkene, chiral ligand, and coupling partner), only a representative alkene (styrene) is shown here for brevity.



with an  $r$  number of stacked attentive layers, allowing atoms to aggregate relevant “messages” from their neighborhood. This step is expected to capture the nuances of atomic local environments by effectively propagating node information over various distances. For molecule-level embedding, the entire molecule is treated as a super-virtual node ( $V$ ), connecting every atom, and embedded using the attention mechanism as shown in Fig. 2. This process, over  $t$  stacked layers, produces a state vector  $h_v^t$  for the whole molecule. In this mechanism, the first step is to concatenate the state vectors of the virtual node ( $h_v$ ) with all connected nodes ( $h_1^t$ ), followed by a linear transformation ( $W$ ) and nonlinear activation (leakyRelu) to produce  $a_{vi}$ . This  $a_{vi}$  is then normalised using a softmax function over the neighbour nodes, resulting in  $\alpha_{vi}$  that captures the importance (weight) of each neighbour node to the virtual node. These attention weights are then passed through the message and readout functions to obtain the final state vector  $h_v^t$ , which encodes structural information about the molecular graph. Finally,  $h_v^t$  is fed through the fully connected layer (FCL) for the regression task.

We have optimised the model using the Adam optimization algorithm, in conjunction with Bayesian optimization (BO) for model-specific parameters such as the number of graph layers for atom embedding, the number of time steps for molecule embedding (denoted respectively using  $r$  and  $t$  in Fig. 2), graph feature size, dropout rate, and optimizer parameters such as the learning rate. Tuning hyperparameters in the DNN can be challenging as they are used for model parameter estimation rather than the model directly assessing them. Hence, BO, as implemented in the Optuna Python package,<sup>71</sup> is used to optimize model-specific parameters. The optimal sets of validation hyperparameters for all 30 runs are provided in Section 2.5 in the SI.

#### 2.4. Effect of different featurization methods and DL models on enantioselectivity prediction

Since ML model performance can depend on the nature of chemical featurisation, we examined the influence of different featurization methods such as one-hot encoding (OHE), molecular fingerprints (FPs), SMILES, and graphs in this work. This aspect has received relatively limited attention in previous studies on % *ee* prediction, prompting us to assess the three traditional featurisation modalities to determine their efficacy on the ART dataset. We have encoded the samples using one-hot encoded vectors (OHEs) to capture the presence or absence of chemical entities such as catalysts, ligands, and substrates in a given reaction. This approach therefore treats each reaction component as an independent categorical variable without explicitly encoding the corresponding chemical structures. While such representations may lack chemical relevance, they might help unveil statistical patterns within complex datasets.<sup>72</sup> The second approach utilizes molecular fingerprints (FPs), a well-established molecular encoding technique in the form of fixed-length binary strings.<sup>73</sup> These FPs capture molecular topology by delineating atom neighbourhoods within a specified radius, providing a compact yet

informative representation of the molecular structure. We used a variety of FPs, such as circular, atom pair, and layered fingerprints, each furnishing unique structural details. For instance, circular FPs encode sub-structural patterns around each atom, atom pair FPs capture pairwise interactions, and layered FPs provide structural information at multiple abstraction levels. Leveraging the open-source cheminformatics package RDKit, we generated molecular FPs with specific parameters tailored to our experimental setup. Additionally, we explored the utility of molecular SMILES, a string-based representation that not only encapsulates structural details like atoms, bonds, and connectivity but also captures nuanced features such as stereochemistry and unique bonding patterns. Our fourth approach involved molecular graphs, offering a robust representation of molecules where atoms correspond to nodes and bond to edges in an undirected graph. By augmenting nodes and edges respectively with atomic identity and bond order as features, these attributed molecular graphs can exhibit versatile predictive capabilities in DL endeavours.<sup>74</sup>

In addition to examining the effects of featurisation on the ART dataset, we aim to compare the performance of the AttentiveFP model with that of some of the state-of-the-art (SOTA) DL models commonly used in enantioselectivity predictions. We have conducted comprehensive evaluations of various DNN models employing both OHE and molecular FPs.<sup>75</sup> Additionally, we explored transformer-based language models such as Transformer,<sup>76</sup> ULMFiT,<sup>77</sup> and T5Chem,<sup>78</sup> which utilize reaction SMILES, and graph neural networks like MPNN<sup>79</sup> and AttentiveFP, which leverage graphs for reactant featurisation. Each combination between a DL model and featurization is then evaluated on the basis of the corresponding training and test RMSEs (see Section 2 in the SI for more details).

The model performances compiled in Table 1 highlight the effect of different featurization techniques on the ART dataset. With OHE, the test RMSE of the DNN is found to be  $14.43 \pm 3.05$  (the details of the DNN architecture are given in Section 2.2 in the SI). These OHE-based models serve as a statistical probe, offering an internal baseline performance for other models built using chemically meaningful descriptors. In contrast to OHE, models that utilize FPs offered improved performance.<sup>71,80</sup> It is worth noting that although the fingerprint-based DNN model exhibited a relatively lower test RMSE ( $9.55 \pm 1.31$ )

Table 1 Performance comparison in terms of RMSE (in % *ee*) of different models and various featurization techniques obtained as the average over 30 independent runs<sup>a</sup>

Featurization	Model	Training	Test
OHE	DNN	$3.67 \pm 1.94$	$14.43 \pm 3.05$
Fingerprint	DNN	$5.54 \pm 1.50$	$9.55 \pm 1.31$
SMILES	T5Chem	$6.74 \pm 0.39$	$10.83 \pm 1.73$
	ULMFiT	$10.94 \pm 0.51$	$11.30 \pm 1.30$
	Transformer	$5.28 \pm 1.18$	$12.26 \pm 2.02$
Graph	<b>AttentiveFP</b>	<b><math>7.41 \pm 1.77</math></b>	<b><math>10.56 \pm 1.86</math></b>
	MPNN	$8.01 \pm 1.18$	$11.00 \pm 2.22$

<sup>a</sup> The datasets are randomly divided into 70 : 10 : 20 training, validation, and test sets.



compared to AttentiveFP, a larger gap with the training RMSE ( $5.54 \pm 1.50$ ) suggests overfitting; hence it should be considered with caution when applied to out-of-bag situations. The use of SMILES representations in conjunction with advanced DL architectures, including T5Chem, ULMFiT, and Transformer, yielded slightly higher test RMSEs of  $10.83 \pm 1.73$ ,  $11.30 \pm 1.30$ , and  $12.26 \pm 2.02$ , respectively.<sup>81</sup> The graph-based MPNN model showed a high test RMSE of  $11.00 \pm 2.00$ .<sup>82</sup> Thus, despite not being the top performer in terms of the lowest test RMSE, the balanced performance of AttentiveFP suggests that it is a robust model with a lower susceptibility to overfitting compared to the other models considered here.<sup>83</sup> In addition to the good performance, the graph attention mechanism inherent to the AttentiveFP model that allows for chemically meaningful interpretability (*vide infra*) has made AttentiveFP our primary framework for the *ee* prediction task.<sup>84</sup>

### 2.5. Class imbalance aware reweighting strategies

While the AttentiveFP model built on molecular graphs showed good predictive performance (Table 1), it does not address the class imbalance issue, given the inherently skewed distribution toward the high *ee* values found in several real-world cases, including that in the ART dataset. Although class imbalance mitigation methods are commonly applied to classification problems, regression tasks like the *ee* prediction have received little attention toward building ML models that address distribution imbalances. Previous research has often relied on the SMOTE technique to generate synthetic data in the minority class for certain regression models.<sup>31,33</sup> The ART dataset bears a skewed distribution with a majority of samples in the higher end of the % *ee* values (Fig. 1b), thus necessitating the development of a customized regression model aware of this class imbalance.

Herein, we propose a customized model for class imbalance, namely, AttentiveFP-CI. Unlike conventional Mean Squared Error (MSE) loss, our model incorporates a class imbalance loss, assigning different weights to training samples based on their actual *ee* values (Fig. 3). The idea is to reduce the influence of the majority class samples while prioritising the more challenging minority class instances during training.<sup>85</sup> We have examined the effect of using different class boundaries, from 30 to 60, by placing the boundary at statistically important points of the dataset such as the mean ( $\mu$ ) value of 76 and  $\mu - \sigma$  of 54. In most cases, the AttentiveFP-CI models performed better than the AttentiveFP model, as evident from the corresponding test RMSE (see Tables S61–S69 in the SI). The model with a class boundary of 30 achieved the best test RMSE of  $9.80 \pm 1.40$  as compared to other class boundaries considered.<sup>86</sup> Moreover, the

*t*-test resulted in a *p*-value  $< 0.05$ , indicating that the gain in performance is statistically significant as compared to the model without the CI (test RMSE for AttentiveFP is  $10.56 \pm 1.86$ ). Incorporation of the CI-aware loss into other DL models also improved the respective test RMSEs, except in the case of ULMFiT.<sup>87</sup> A comparison between different deep learning architectures reveals that AttentiveFP-CI outperforms Transformer-CI and ULMFiT, with *p*-values  $< 0.05$  endorsing their statistical significance.<sup>88</sup> However, most of these models tend to exhibit overfitting issues, evident from the train-test RMSE differences as follows: MPNN-CI (4.59), T5Chem-CI (5.4), and DNN-CI (3.93) as opposed to AttentiveFP-CI (3.59).<sup>89</sup> Additionally, the number of model parameters in AttentiveFP-CI is in the order of 1.93 M, which is much fewer than those in T5Chem-CI (14.71 M), assuring us of better computational scalability. Given the lower RMSE, reduced overfitting, and computational efficiency, AttentiveFP-CI stands out as the optimal choice from among all the SOTA models for % *ee* predictions on the ART dataset.<sup>90</sup>

Efforts were also expended to assess whether the performance issues could be traced to the sparse and imbalanced distribution in the ART dataset. We have compared the model performances on more balanced and denser datasets such as Buchwald–Hartwig Amination (BHA), which is a catalytic transformation of high practical utility. The high throughput experimental (HTE) dataset of the BHA reaction, denoted as BHA-HTE, is a commonly used dataset for baseline comparisons for yield prediction tasks.<sup>24</sup> BHA-HTE comprises 3955 labeled reactions and their corresponding experimentally measured yields. The AttentiveFP model offered a good test performance, with an RMSE of  $6.49 \pm 0.33$  and a coefficient of determination ( $R^2$ ) of  $0.94 \pm 0.01$  (see Table S72 in the SI), surpassing the previously reported  $R^2$  of 0.92 obtained using physical-organic descriptors.<sup>29</sup>

In the present context, we have done random sampling of the full BHA-HTE dataset to create a few sparser subsets, denoted as BHA-LTE (low throughput), each containing about 500 reactions. The idea is to induce skewness to produce an imbalance in labels such that the distribution ( $\mu$  and  $\sigma$ ) in the BHA-LTE subsets resembles that of the ART dataset. These subsets are then employed for evaluating the baseline performance of various deep learning models considered in this study.<sup>91</sup> In general, the BHA-LTE subsets have a  $\mu$  of 75 and a  $\sigma$  of 14 (see Fig. S1 in the SI), similar to those in our ART dataset ( $\mu = 76$  and  $\sigma = 22$ ). When the AttentiveFP model was trained using these subsets bearing an induced sparse distribution, the test RMSE dropped from  $6.49 \pm 0.33$  with the original BHA-HTE dataset (see Table S72 in the SI) to  $9.14 \pm 0.80$  (or higher, depending on

AttentiveFP	AttentiveFP-CI	
$\mathcal{L} = \frac{1}{M} \sum_{j=1}^M (y_j - \hat{y}_j)^2; \quad \hat{y}_j = \text{FCL}(h_v^t)$	$\mathcal{L}_{\text{class\_imb}} = \frac{1}{M} \sum_{j=1}^M \begin{cases} (y_j - \hat{y}_j)^2 & \text{if } y_j < 50 \\ 0.5 \times (y_j - \hat{y}_j)^2 & \text{if } y_j \geq 50 \end{cases}$	boundary set at 50, ≥ 50 majority samples and <50 minority samples

Fig. 3 Different types of loss used in the AttentiveFP and AttentiveFP-CI models.



the BHA-LTE subset used). The lower performance of the same AttentiveFP model can be attributed to the induced sparse distribution and CI in the BHA-LTE subsets. Interestingly, inclusion of the CI loss with a class boundary of 50 improved the test RMSE to  $8.70 \pm 0.52$ . Such test performances are analogous in quality of predictions by the same model on the ART dataset bearing comparable distribution characteristics ( $\mu$  of 76 and a  $\sigma$  of 22). Additional details on the model performance with varying class boundaries, spanning 30 to 70, and with a  $\mu$  of 75 are provided in Section 3 in the SI. An improvement in RMSE ( $9.14 \pm 0.80$  for AttentiveFP *versus*  $8.70 \pm 0.52$  for AttentiveFP-CI) could possibly be due to the use of a customized loss function to mitigate CI issues.

A similar performance trend is conspicuous in our ART dataset as well. For instance, the test RMSE of the AttentiveFP-CI model (with a class boundary of 30) is found to be 9.80 over 10.56 obtained with the AttentiveFP model without the CI loss. On the basis of the model performance, with and without the inclusion of the CI loss, on the ART and BHT-LTE datasets, we could conclude that the data sparsity is primarily responsible for the higher RMSEs. These insights would be valuable in developing suitable deep learning models with customized loss functions for chemical reaction datasets bearing skewed distribution.

An alternative for imbalanced and sparsely distributed chemical datasets is to consider a two-step model,<sup>92</sup> wherein a classification of samples is done first on the basis of a pre-defined class label. Subsequently, separate regressors are developed for the major and minor classes. This approach, termed classification-followed-by-regression (CFR), is likely of relevance to the ART dataset.<sup>93</sup> In the first step, reactions are classified as 'major' or 'minor' using a statistically meaningful class boundary set at a ( $\mu - \sigma$ ) of 54 % *ee*.<sup>94</sup> We found that a hyperparameter optimized custom built DNN classifier could achieve a very good accuracy of  $0.98 \pm 0.003$ .<sup>95</sup> The reactions thus classified are employed in training two separate AttentiveFP regression models, one for the major class and the other for the minor class. The AttentiveFP regressor achieved test RMSEs of  $10.74 \pm 1.98$  for the major class and  $8.73 \pm 0.90$  for the minor class, outperforming our direct regression in the case of minor class reactions (*i.e.*, reactions with less than 54 % *ee* as their true label).<sup>96</sup> To ensure a balanced assessment of the overall model performance, we have also considered the use of weighted RMSE, which accounts for class imbalance by combining error contributions proportionate to the sample size.<sup>97</sup> For the AttentiveFP model, the weighted RMSE is  $8.97 \pm 1.28$ , which is poorer than our AttentiveFP-CI with a class boundary of 30 (RMSE =  $9.80 \pm 1.40$ ;  $R^2 = 0.80 \pm 0.05$ ) (see Table S64 in the SI). Similarly, an ULMFiT regression model showed a test RMSE of  $10.15 \pm 2.26$  and  $8.48 \pm 1.26$  respectively for the major and minor classes with a corresponding weighted RMSE of  $8.75 \pm 1.05$  ( $R^2 = 0.40 \pm 0.29$ ).<sup>98</sup> Since no significant improvement is found with CFR-major and CFR-minor classes, our original AttentiveFP-CI, with its interpretable characteristics, can be considered a more appropriate model for the ART dataset.

AttentiveFP-CI showed good performance in predicting % *ee*, achieving a test RMSE of  $9.80 \pm 1.40$ . Importantly, the difference between the training and validation RMSEs suggests

a lower overfitting, which is good for model generalizability when predicting on unseen samples. In the 30 independent runs, the model predicts % *ee* thousands of times for the 76 reactions present in the test set. Furthermore, every reaction gets predicted multiple times whenever it appears in the test set. A comparison of the predicted % *ee* with the experimentally known ground truth values revealed good correlation as shown in Fig. 4. In fact, ~87% of the predictions remain within 15 units of the actual values (Fig. 4a). In the optimal run with an RMSE of 8.2 % *ee*, as many as 70 out of 76 test samples are predicted well within an error limit of 15 units with respect to the corresponding true values (Fig. 4b). In a typical run (RMSE = 10.1), only 12 out of the 76 samples incurred prediction errors in excess of 15 units (Fig. 4c). The parity plot also conveys a good correlation between the % *ee* predicted by the AttentiveFP-CI model and the corresponding experimental values with an  $R^2$  of 0.84 (Fig. 4d).<sup>99</sup> These assuring findings highlight the efficacy of AttentiveFP-CI in learning from the sparse ART dataset for catalytic asymmetric reactions of alkenes.

To evaluate the learning ability of an ML model and to examine its robustness, control experiments are required. For this purpose, the dependence of the model performance on the quality of the input data is assessed using techniques such as Y-scrambling. We created a straw model of AttentiveFP-CI, which intentionally breaks the potential connections between the input features and the output variable. Here, each sample is incorrectly mapped to an output value belonging to some other sample within the dataset. The considerably worse test RMSE of  $25.2 \pm 2.1$  obtained with the Y-scrambling run shows that the model learns from the true features it was provided with in the correct training campaigns. The inferior performance also highlights the effectiveness of AttentiveFP-CI in learning the chemically meaningful aspects of the catalytic reaction investigated in this work (*vide infra*).

## 2.6. Interpretability and graph attention visualization

With these promising findings on % *ee* prediction, it becomes all the more relevant to gather additional insights into what the model would have learned from the molecular representations as provided. We sought details such as the key atoms that are likely to exert more impact on the reaction output expressed in % *ee*. Leveraging the interpretability of attention weights (see  $\alpha_{vi}$  in Fig. 2) learned by the model, we probed the latent connections between various molecules present in a given reaction sample and the predicted outcome. Here, we aim to visualize the atomic attention weights using the similarity map as implemented in the RDKit program.<sup>100</sup> In this mode of visualization shown in Fig. 5, the regions with a green glow indicate a positive influence on % *ee*, a pink glow suggests a negative influence, and a grey colour conveys no significant overall effect. As a representative example, we have compared the attention weights for two of the chosen reactions involving the same reactants but different chiral ligands (pyridine-oxazoline and bis(oxazoline)) in Fig. 5a. It should be noted that in both these reactions, the predicted % *ee* is comparable to the experimentally reported % *ee*.



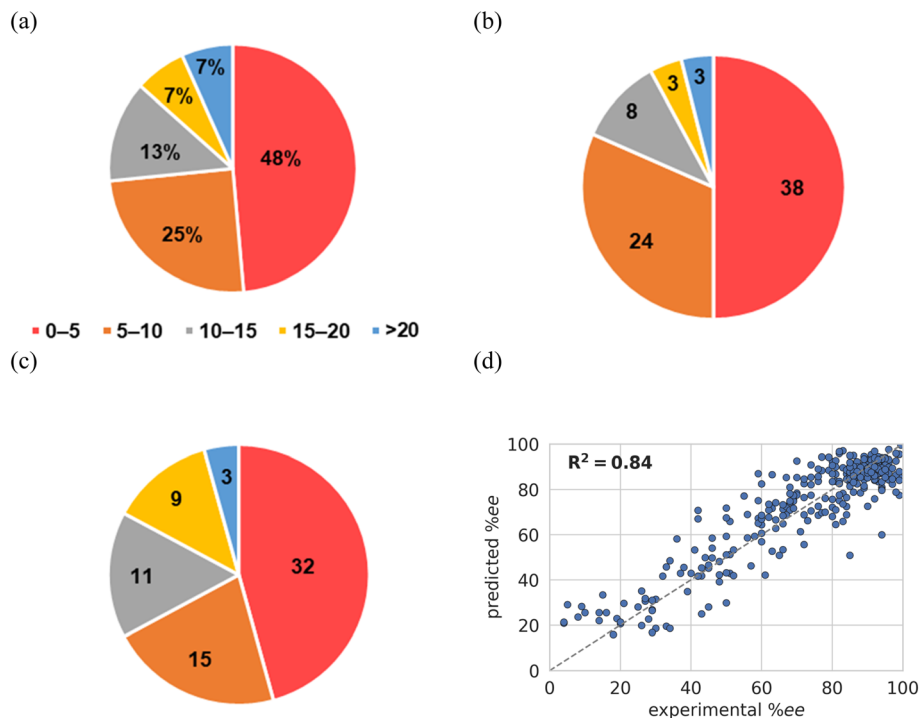


Fig. 4 Pie charts summarize the prediction errors (% ee differences) obtained using the AttentiveFP-CI model across 76 test samples: (a) aggregated over 30 independent runs, (b) from the best run, and (c) in a typical run whose performance is similar to the average RMSE over 30 runs. (d) Parity plot showing the correlation between experimental and predicted % ee values for all test samples.

From perusal of the attention map shown in the inset of Fig. 5a and the bar plot in Fig. 5b for the reaction involving catalyst-1 (pyridine-oxazoline catalyst), it can be learned that the Cu centre and pyridine N with two nearest C atoms make positive contributions while all other atoms or substructures have relatively lower negative or negligible contributions to the % ee. In the case of catalyst-2, the side arm (SA) on the bridge carbon of the bis(oxazoline) ligands positively contributes to high % ee, consistent with the trends observed with this family of ligands.<sup>101</sup> It is further evident that the other positive contributors to the reaction outcome are (i) the SA on the chiral ligand, (ii) the transition metal-bound triflate ligand, (iii) the styrenyl double bond, and (iv) one of the carbon atoms of the cyclohexyl diazo compound. These positive attention values are suggestive of their synergistic role in the enantioselectivity of the cyclopropanation reaction. One of the significances of this analysis is that installation of suitable substituents on the SA group could be key to achieving enhanced enantioselectivity. This prediction by the model is chemically meaningful and intuitively appealing, as it aligns with the fact that most variants of reported bis(oxazoline) ligands rely on modifications of the SA.<sup>101</sup>

After visualizing attention for two representative examples, we have analysed the global effects of the critical regions/atoms that likely exert a significant contribution toward the quality of the % ee prediction. To accomplish this, it is essential to identify a common region present in each reaction partner across all the samples. Fig. 5b highlights such shared regions in all the reaction components, along with their atom numbering.

The steps involved in estimating the effect of each atom are as follows: first, attention values for each atom in the shared region are extracted using the corresponding SMART pattern. Second, the variance of these attention values is plotted, since the variance is crucial in assessing the feature importance as it captures the most significant and informative variations in the data.<sup>102</sup> Interestingly, the bar diagram shown in Fig. 5b indicates a higher importance of the chiral ligand (atom numbers are given with L in parentheses to denote the chiral ligand) as compared to the reactants such as the alkene and other substrates. It is gratifying to see that our attention-based model deciphered chemically intuitive patterns present in the chiral ligand as the most relevant contributor to asymmetric induction. The variances in the attention values exhibited by the atoms in the chiral ligands are found to be much higher than those of the alkene and other substrates. This observation is in line with one's chemical intuition that chiral ligands play a pivotal role in transferring the chiral information to the developing product.<sup>103</sup> Notably, the chiral carbon centre, denoted as [C\*8(L)], in the ligand exhibits the highest variance in the attention, corroborating the domain knowledge that the substituents at this centre are largely known to influence enantioselectivity. Additionally, the carbon atom near the bridge or SA [C1(L)] shows the second-highest variance, suggesting that modifications to this atom, where branching from the bridge carbon begins, could help in fine-tuning the reaction outcome. The identification of these atoms as important ones indicates that AttentiveFP-CI accurately captures the relationship between the molecular factors and the desired outcome.



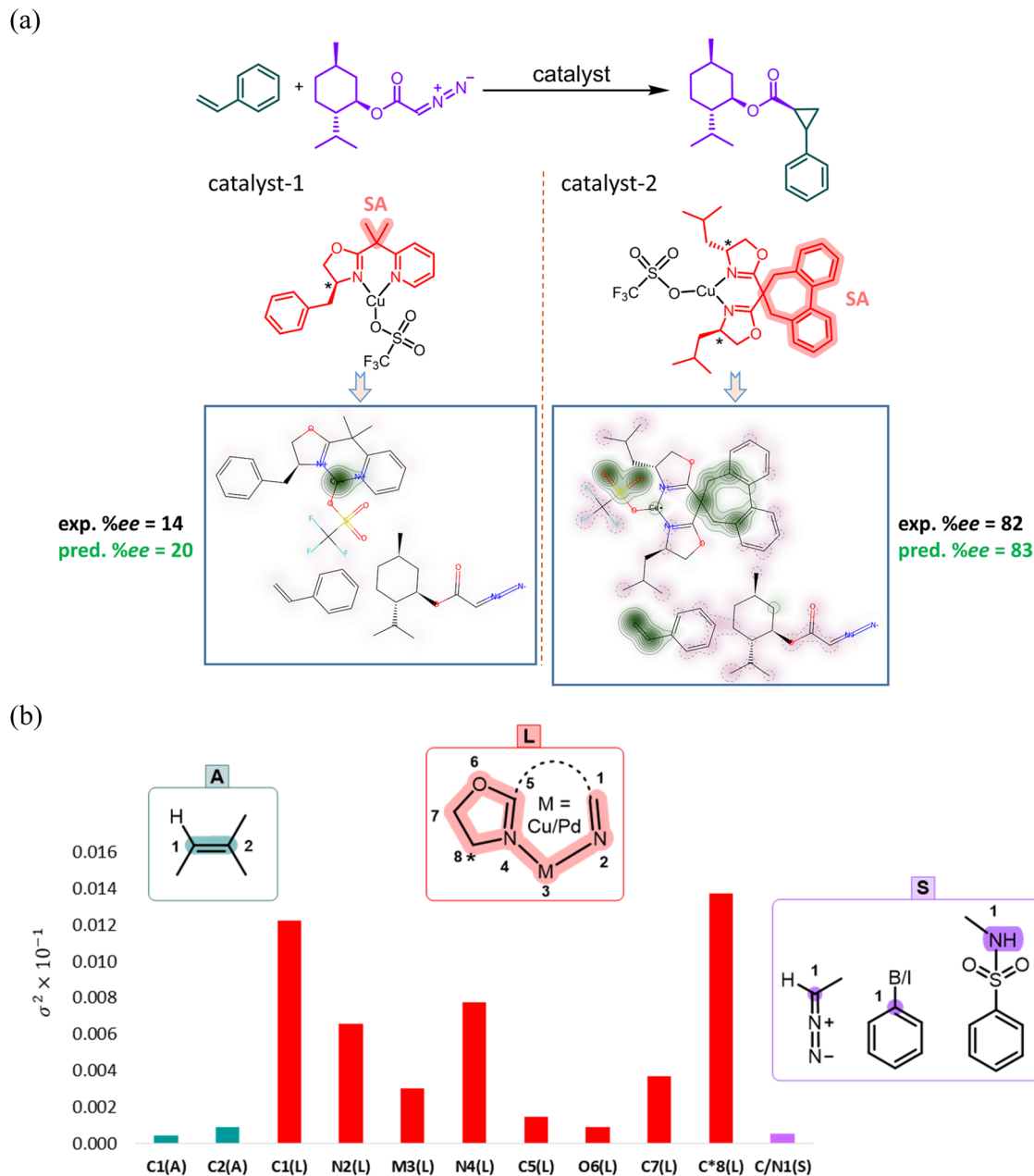


Fig. 5 Visualization of the attention weights of two representative reactions involving the same substrates and different catalysts. (a) In the predicted probability maps, atoms with positive contributions to the reaction outcome appear in green, while red indicates that the corresponding attention weight is negative. The larger the absolute value, the darker the colour in the attention map. (b) Atoms within the shared regions in each reaction partner involved in the overall reaction are shown highlighted, along with the variance ( $\sigma^2$ ) of the attention value ( $y$ -axis) for the highlighted atoms ( $x$ -axis). Different colours such as cyan for the alkene (A), red for the ligand (L), and purple for substrates (and diazo compounds, boronic acid, or *N*-tosyloxycarbamate) (S) are used in the bar plots.

Exploiting this protocol by fine-tuning these key features, particularly during reaction development, or while expanding the scope of this reaction family, could prove advantageous.

### 2.7. Application of the AttentiveFP-CI regressor for identifying potential catalysts

We wish to demonstrate one of the potential applications of the attention-based regression model here in terms of identifying suitable substrates or catalysts to expedite the synthesis of

target molecules of interest, such as a drug molecule. To illustrate this aspect, we have used the AttentiveFP-CI model to identify an efficient catalyst for synthesizing (*S,S*)-naproxen, which is oral non-steroidal anti-inflammatory drug. One of the experimental reports suggests the use of (*S,S*)-PhBox as the chiral ligand in an enantioselective Cu-catalyzed arylation reaction as shown in Fig. 6a.<sup>104,105</sup> The observed *ee* was 94%, which is very close to the predicted *ee* of 92% by the AttentiveFP-CI model. Given this encouraging agreement, we have screened



a library of 35 chiral Box ligands from the PubChem database using the AttentiveFP-CI model and identified several promising chiral ligands with high predicted % *ee* values, as shown using a heatmap representation in Fig. 6b (see Table S88 in the SI for more details).<sup>106,107</sup> These candidates prioritized by the model can be utilized to streamline future screening efforts. However, their activity remains to be established through wet-lab validation.

It is worth noting that although these ligands were previously reported, their use in copper-catalyzed arylation reactions, involving silylketenes and diaryliodonium salts, remains unexplored (Fig. 6a). Thus, the potential of these ligands for such alkene–substrate pairs is novel, even when the reaction conditions are retained.<sup>63–65</sup> A heatmap representation of the predicted % *ee* for different chiral ligands shown in Fig. 6b conveys the significance of fast (virtual) screening of chiral ligands by using our regression model. It can be seen (top left) that a ligand with only one of the oxazoline rings chiral is predicted to exhibit a lower % *ee* of 65. Interestingly, the attention analysis identifies the groups on the side arm at the bridge carbon (SA) and the chiral carbon of the bis(oxazoline) ligand as the dominant contributors to % *ee*. In light of the attention as noted, we considered two representative variants of the bis(oxazoline) family of ligands such as the (*S,S*)-Ph-Box for

further illustration as shown in Fig. 6b. One of these ligands is obtained by replacing the Ph group at the chiral carbon with 4-*t*-Bu-Ph, and the other is obtained by replacing the 1,1-dimethyl on the side arm with a 1,1-diisopropyl group. Both of these ligands are predicted to show high % *ee*. More importantly, a higher attention value noted for 1,1-diisopropyl and in the Ph regions (green color contours) indicates their positive contribution to enantioselectivity. These can be considered as indicative of how an attention-based approach could be utilized in catalyst design for asymmetric reactions.

## 2.8. Extension of the AttentiveFP model to asymmetric hydrogenation, *N,S*-acetylation, and USPTO datasets

Motivated by the promising interpretability of the AttentiveFP model, we have considered two more important reactions from the domain of asymmetric catalysis involving axially chiral ligands/phosphoric acid (CPAs). These are *N,S*-acetylation<sup>56</sup> and asymmetric hydrogenation of alkenes,<sup>32</sup> both known for their high industrial importance.<sup>108–110</sup> *N,S*-acetylation involves enantioselective coupling between an imine and a thiol, while the hydrogenation reaction encompasses the reduction of imines and alkenes using BINOL- and BINAP-derived chiral ligands. The original *N,S*-acetylation dataset contains 1075

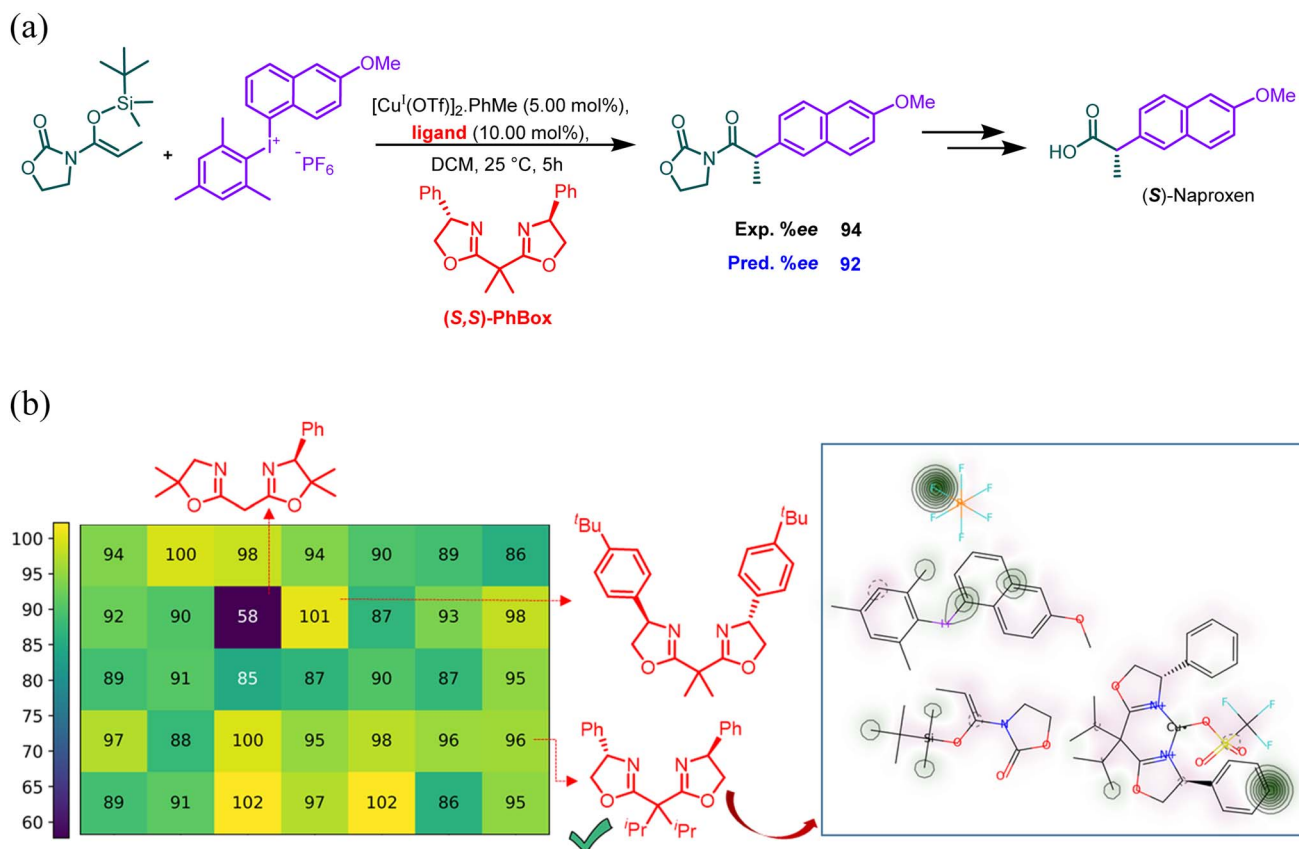


Fig. 6 (a) Experimental conditions for making a precursor involved in the synthesis of (*S*)-naproxen using an asymmetric Cu-catalyzed arylation reaction of alkenes, (b) representative examples of some chiral bis(oxazoline) ligands with their predicted % *ee* obtained using the AttentiveFP-CI model. Shown in the inset is attention weight visualization for the reaction catalyzed by [Cu-(*S,S*)-Ph-Box] with 1,1-diisopropyl groups at the bridging carbon.



reactions, wherein the enantioselectivity is expressed using the corresponding  $\Delta\Delta G^\ddagger$  values, ranging from positive to negative values, or a signed *ee* value is used to denote the experimentally observed major enantiomer (*R* or *S*) in the reaction. Hence, it should be noted that in the present work, we formulate the problem as a regression task to predict % *ee* values on a 0–100 scale, focusing on the magnitude of enantioselectivity in line with the general practice followed in the organic chemistry literature. Consequently, 48 reactions bearing negative *ee* values (used to indicate the opposite enantiomer) are excluded during data curation to maintain a consistent target value, making it 1027 reactions. For consistency in benchmarking and comparison with prior studies, the AttentiveFP-CI model is also evaluated on the full 1075 *N,S*-acetylation reaction dataset using the same 80:10:10 training, validation, and test split as used before besides a  $\mu$ -based class boundary of 0.98. Our model achieved a test  $R^2$  of  $0.90 \pm 0.02$  and an RMSE of  $0.21 \pm 0.02$ .<sup>111</sup> These results are comparable to those of the SEMG-MIGNN ( $R^2 = 0.915$ ; RMSE = 0.197)<sup>23</sup> and ChemAH ( $R^2 = 0.918$ ; RMSE = 0.209).<sup>112</sup> In the regression setting with 1027 reactions, AttentiveFP-CI achieves a good predictive performance of  $8.06 \pm 1.00$  ( $R^2$  of  $0.92 \pm 0.10$ ) on the 1027 *N,S*-acetylation reactions, which is found to be a statistically significant improvement over the corresponding AttentiveFP model devoid of CI loss.<sup>113a</sup> In the case of asymmetric hydrogenation, although the performance of AttentiveFP-CI as indicated by an RMSE of  $10.48 \pm 1.10$  ( $R^2$  of  $0.60 \pm 0.17$ ) is good, the gain as compared to the AttentiveFP model is not statistically significant.<sup>113b</sup> Notably, for the *N,S*-acetylation reaction, AttentiveFP-CI significantly outperformed an often used baseline such as the ULMFIT model ( $p = 0.0036$ ), with a test RMSE of  $8.88 \pm 1.03$ .<sup>114</sup> Although in the hydrogenation case, baseline ULMFIT provides better predictive accuracy (test RMSE of  $8.56 \pm 1.46$ ), it lacks the advantage of interpretability as afforded by AttentiveFP-CI.

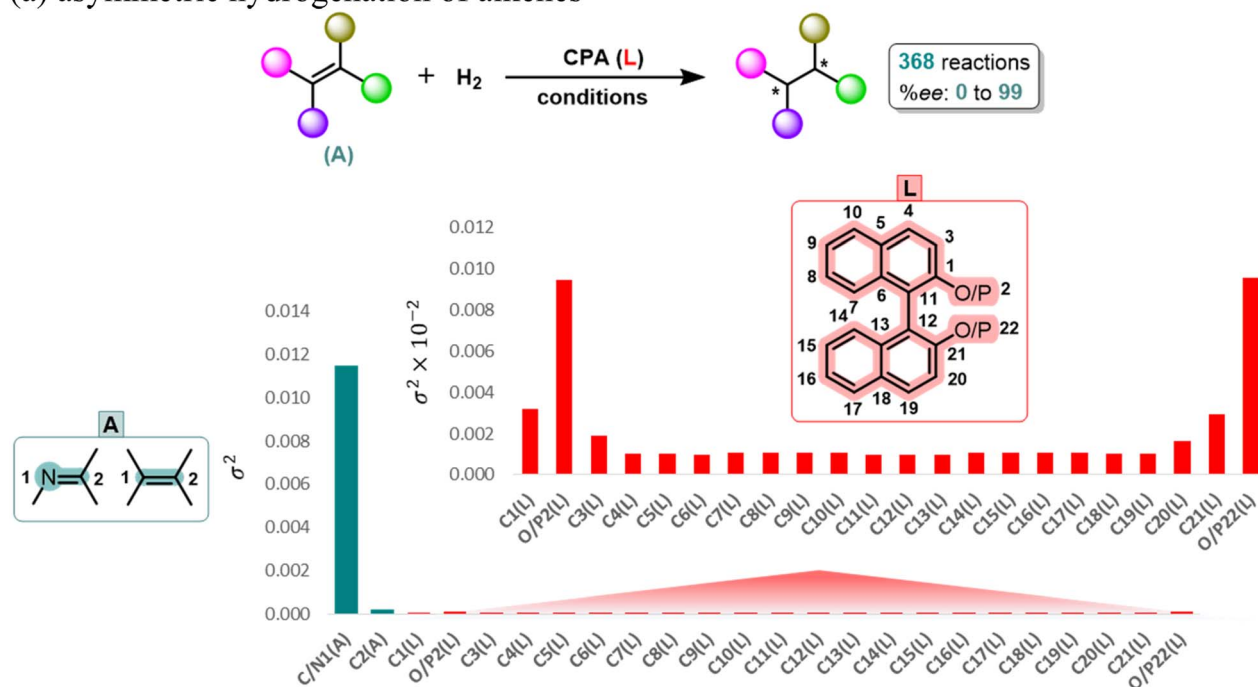
In addition to the enantioselectivity prediction on three important chemical reaction datasets, the utility of the AttentiveFP-CI framework is also evaluated for yield prediction tasks on the USPTO (grams) reaction dataset.<sup>21</sup> This dataset comprises  $1.9 \times 10^5$  reactions, each annotated with the corresponding yield values. The yield distribution exhibits a skewness of  $-0.86$ , indicating the presence of CI in the USPTO dataset. Furthermore, the yield values in this dataset are reported as scaled values to fit in the interval of 0 to 1. In light of this skewed distribution, we used the AttentiveFP-CI model with a statistically relevant class boundary of  $\mu + \sigma$  (0.94) to note a test RMSE of  $0.20 \pm 0.01$  and a marginally better  $R^2$  of  $0.08 \pm 0.00$ , as compared to AttentiveFP without the CI consideration (test RMSE =  $0.21 \pm 0.01$ ;  $R^2 = 0.04 \pm 0.01$ ).<sup>115</sup> However, the *t*-test resulted in a *p*-value  $> 0.05$ , indicating that the numerical improvement is not statistically significant in the case of the USPTO dataset. Notably, the performance of AttentiveFP-CI is even comparable to the previously reported RMSE of 0.195 obtained using a more complex transformer-based model on the same dataset.<sup>116</sup> Overall, these results indicate that AttentiveFP-CI could be useful in addressing CI issues in chemical datasets, even if its performance does not always surpass state-of-the-art benchmark performances.

Similar to the approach employed earlier in this manuscript for global attention analysis, we have visualized the attention weights by using the computed variance in the attention values of atoms within the shared region as shown in Fig. 7, for both asymmetric *N,S*-acetylation and asymmetric hydrogenation reactions.<sup>100,117–119</sup> The variance in the attention values of the shared region atoms of the catalysts and substrates across different samples (*i.e.*, reactions in the dataset) can be considered a measure of sensitivity of the reaction outcome to the environment of such atoms. Therefore, such analysis might help decipher how the changes in the local substituents are likely to influence the enantioselectivity. A relatively larger change in variance implies a higher attention on such atoms, which might stem from the changes in their substituents or local environment. A modest change indicates that the atom concerned consistently gets similar attention weights. Interestingly, in the context of the reactions in our ART dataset, we notice that reactive positions on the alkene and the vital regions around the chiral centre exhibit high variances in their attention. These variances both in the ART reactions and in the case of asymmetric hydrogenation are in line with our chemical intuition, where alkene is the key substrate undergoing the reaction. Since this analysis collectively reveals mechanistically valuable insights consistent with the domain knowledge on the origin of enantioselectivity catalyzed by axially chiral ligands,<sup>101,120,121</sup> we consider that the AttentiveFP-CI model is meaningfully interpretable.

In the asymmetric hydrogenation reaction, the atoms belonging to the imine or alkene (*e.g.*, C/N1(S)) exhibited consistently higher variance in attention, conveying that they might play a key role in enantioselectivity. While the atoms in the chiral ligand showed relatively lower variance compared to those of the substrate, the P/O center in BINOL- and BINAP-derived ligands (*e.g.*, O/P2(L) or O/P22(L)) has maximal attention variance. This is an interesting aspect, which is in line with the chemical intuition that these positions in the chiral BINOL/BINAP frameworks are expected to influence effectiveness of enantioinduction. Thus, suitable substitution at the *ortho* positions in the biphenyl ring can potentially impact enantioselectivity values. For the *N,S*-acetylation reaction, high attention variance is found across all three components: the thiol (C3(S)), imine (C1(A)), and ligand (C5(L)), reflecting the impact of local substituent changes on enantioselectivity. Notably, within the ligand, attention values corresponding to the atoms C5(L) and C6(L), which bridge the biphenyl units, showed substantial variance, linked to the dihedral angle fluctuations (*e.g.*, C4–C5–C6–C7) that can modulate the chiral environment. Additionally, variance in attention scores suggests that the *meta*-positions on the biphenyl rings (C19(L) and C10(L)) are important where a change in substituents is likely to influence the reaction outcome. Thus, the attention variance analysis suggests that introducing suitable substituents at the high-variance sites (*e.g.*, P/O centers and dihedral-sensitive atoms) might shift the enantioselectivity. These findings offer a firm basis for the interpretability of AttentiveFP-CI by way of identifying hotspots for enantioselectivity tuning. The results could become useful for rational catalyst design and for making an



## (a) asymmetric hydrogenation of alkenes



## (b) asymmetric N,S-acetylation

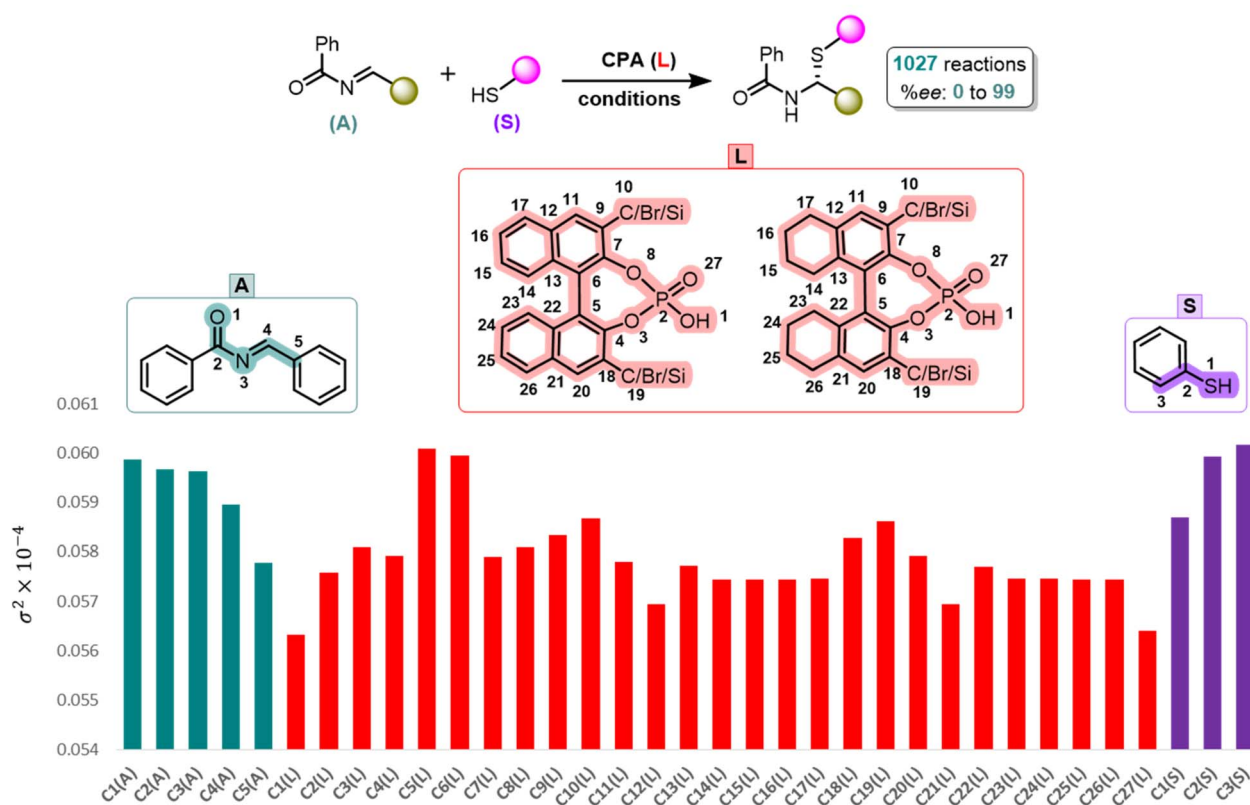


Fig. 7 Highlighted atoms in the alkene (A), chiral ligand (L), and substrate (S) from the shared region across different reactions/samples are identified as exerting a positive/negative impact on the reaction outcome. The bar plot shows the variance in attention values ( $\sigma^2$ ) obtained from the AttentiveFP-CI model ( $y$ -axis) for each atom as shown labeled in the  $x$ -axis. Colors indicate different reaction components: cyan for alkenes, red for ligands, and purple for substrates for better visualization.



informed choice of substrates during reaction scope investigations in asymmetric catalysis.

### 3 Conclusions

We have curated a comprehensive dataset, named ART (AlkeneReactionTriad) that contains 376 catalytic asymmetric reactions of alkenes, as a representative of the real-world chemical reaction dataset suitable for machine learning model building. We have employed an attention-based class imbalance aware GNN model (AttentiveFP-CI) on the ART dataset for predicting enantioselectivity in terms of % *ee* as the label. Our ML model is found to be effective in addressing one of the known challenges in reaction outcome prediction, arising due to the sparse and imbalanced distributions often found in chemical reaction datasets. Through a comparative analysis of various featurization techniques, including one-hot encoding, molecular fingerprints, SMILES, and molecular graphs, we found that graph-based AttentiveFP-CI is the most suitable model offering a test RMSE of  $9.80 \pm 1.40$  on the ART dataset. Importantly, our findings indicate that an RMSE in this range is comparable to the performance of the same AttentiveFP-CI on a sparsity-induced Buchwald–Hartwig amination dataset ( $10.86 \pm 1.32$  test RMSE as opposed to  $6.49 \pm 0.33$  with the full BHA dataset). Lower performance on the sparse subset of the BHA dataset can be considered to stem from the sparse distribution and CI issues, similar to that in the ART dataset.

Visualization of the atomic attention weights could identify the pivotal regions in the reaction partners, such as the chiral centre, as a high attention spot in the chiral catalyst. Similarly, critical atoms/substructures in the reactant(s) are identified as an important contributor to high enantioselectivity. Thus, AttentiveFP-CI not only serves as a good predictive model, but it also offers chemically meaningful insights for reaction optimization. This method can therefore pave the way for informing ligand design and reaction development, as exemplified by the identification of the (*S,S*)-PhBox ligand variant, featuring 1,1-diiisopropyl groups on the side arm as a potentially effective catalyst relevant to the synthesis of (*S*)-naproxen. When extended to an important enantioselective reaction, such as the axially chiral phosphoric acid (CPA) catalyzed *N,S*-acetylation, AttentiveFP-CI offered a very good RMSE of  $8.06 \pm 1.00$ . The interpretability of our model sheds light on the factors governing enantioselectivity in the form of identifying the reactive olefinic sites in imines and alkenes in asymmetric hydrogenation reactions as the key contributors and the binaphthyl axis of the axially chiral ligands in the case of asymmetric acetylation reactions. Overall, the AttentiveFP-CI model not only serves as a robust predictive framework but also as a chemically interpretable tool that complements intuition. Interpretable models can therefore be exploited in data-driven discovery of chiral ligands and substrates in asymmetric catalysis.

### Author contributions

DC prepared the ART dataset. AH, NJ and RBS wrote the manuscript. RBS supervised the research and analysis.

### Conflicts of interest

There are no conflicts to declare.

### Data availability

Data and codes related to this work are publicly available through our GitHub repository at <https://github.com/alhqlearn/ART-AttentiveFP-CI>. A citable, versioned snapshot of the codebase is archived on Zenodo with an assigned DOI: <https://doi.org/10.5281/zenodo.18256995>.

Supplementary information (SI): details of the machine learning setups with their hyperparameter tuning, chemical reaction datasets, various control experiments, and other relevant information are provided. See DOI: <https://doi.org/10.1039/d5dd00483g>.

### Acknowledgements

We are thankful to the Institution of Eminence (IoE) Data and Information Science computing facility for generous computational resources. N. J. acknowledges the Prime Minister Research Fellowship (PMRF).

### References

- 1 K. R. Campos, P. J. Coleman, J. C. Alvarez, S. D. Dreher, R. M. Garbaccio, N. K. Terrett, R. D. Tillyer, M. D. Truppo and E. R. Parmee, *Science*, 2019, **363**, 6424.
- 2 G. M. Whitesides, *Angew. Chem., Int. Ed.*, 2015, **54**, 3196–3209.
- 3 I. W. Davies, *Nature*, 2019, **570**, 175–181.
- 4 P. Raccuglia, K. C. Elbert, P. D. F. Adler, C. Falk, M. B. Wenny, A. Mollo, M. Zeller, S. A. Friedler, J. Schrier and A. J. Norquist, *Nature*, 2016, **533**, 73–76.
- 5 S. Lin, S. Dikler, W. D. Blincoe, R. D. Ferguson, R. P. Sheridan, Z. Peng, D. V. Conway, K. Zawatzky, H. Wang, T. Cernak, I. W. Davies, D. A. DiRocco, H. Sheng, C. J. Welch and S. D. Dreher, *Science*, 2018, **361**, 6402.
- 6 S. Stocker, G. Csányi, K. Reuter and J. T. Margraf, *Nat. Commun.*, 2020, **11**, 5505.
- 7 P. Schwaller, A. C. Vaucher, R. Laplaza, C. Bunne, A. Krause, C. Corminboeuf and T. Laino, *Wiley Interdiscip. Rev.: Comput. Mol. Sci.*, 2022, **12**, 5.
- 8 J. Werth and M. S. Sigman, *J. Am. Chem. Soc.*, 2020, **142**, 16382–16391.
- 9 J. J. Henle, A. F. Zahrt, B. T. Rose, W. T. Darrow, Y. Wang and S. E. Denmark, *J. Am. Chem. Soc.*, 2020, **142**, 11578–11592.
- 10 A. Milo, A. J. Neel, F. D. Toste and M. S. Sigman, *Science*, 2015, **347**, 737–743.
- 11 S. Zhao, T. Gensch, B. Murray, Z. L. Niemeyer, M. S. Sigman and M. R. Biscoe, *Science*, 2018, **362**, 670–674.
- 12 K. H. Hopmann, *Int. J. Quantum Chem.*, 2015, **115**, 1232–1249.



- 13 A. S. K. Tsang, I. A. Sanhueza and F. Schoenebeck, *Chem.–Eur. J.*, 2014, **20**, 16432–16441.
- 14 W. Beker, R. Roszak, A. Wołos, N. H. Angello, V. Rathore, M. D. Burke and B. A. Grzybowski, *J. Am. Chem. Soc.*, 2022, **144**, 4819–4827.
- 15 M. Meuwly, *Chem. Rev.*, 2021, **121**, 10218–10239.
- 16 S. Singh and R. B. Sunoj, *Acc. Chem. Res.*, 2023, **56**, 402–412.
- 17 A. Hoque, M. Das, M. Baranwal and R. B. Sunoj, *arXiv*, 2024, preprint, arXiv:2407.10090, DOI: [10.48550/arXiv.2407.10090](https://doi.org/10.48550/arXiv.2407.10090).
- 18 P.-X. Hua, Z. Huang, Z.-Y. Xu, Q. Zhao, C.-Y. Ye, Y.-F. Wang, Y.-H. Xu, Y. Fu and H. Ding, *Commun. Chem.*, 2025, **8**, 42.
- 19 S. Shilpa, G. Kashyap and R. B. Sunoj, *J. Phys. Chem. A*, 2023, **127**, 8253–8271.
- 20 A. Hoque, M. Surve, S. Kalyanakrishnan and R. B. Sunoj, *J. Am. Chem. Soc.*, 2024, **146**, 28250–28267.
- 21 P. Schwaller, A. C. Vaucher, T. Laino and J.-L. Reymond, *Machine Learning: Science and Technology*, 2021, **2**, 015016.
- 22 M. Das, A. Ghosh and R. B. Sunoj, *J. Comput. Chem.*, 2024, **45**, 1160–1176.
- 23 S.-W. Li, L.-C. Xu, C. Zhang, S.-Q. Zhang and X. Hong, *Nat. Commun.*, 2023, **14**, 1.
- 24 D. T. Ahneman, J. G. Estrada, S. Lin, S. D. Dreher and A. G. Doyle, *Science*, 2018, **360**, 186–190.
- 25 D. Perera, J. W. Tucker, S. Brahmabhatt, C. J. Helal, A. Chong, W. Farrell, P. Richardson and N. W. Sach, *Science*, 2018, **359**, 429–434.
- 26 J. P. Reid and M. S. Sigman, *Nature*, 2019, **571**, 343–348.
- 27 V. Voinarovska, M. Kabeshov, D. Dudenko, S. Genheden and I. V. Tetko, *J. Chem. Inf. Model.*, 2023, **64**, 42–56.
- 28 M. J. Gaunt, J. M. Janey, D. M. Schultz and T. Cernak, *Chem*, 2021, **7**, 2259–2260.
- 29 M. Saebi, B. Nan, J. E. Herr, J. Wahlers, Z. Guo, A. M. Zurański, T. Kogej, P.-O. Norrby, A. G. Doyle, N. V. Chawla and O. Wiest, *Chem. Sci.*, 2023, **14**, 4997–5005.
- 30 S. Gallarati, R. Fabregat, R. Laplaza, S. Bhattacharjee, M. D. Wodrich and C. Corminboeuf, *Chem. Sci.*, 2021, **12**, 6879–6889.
- 31 A. Hoque and R. B. Sunoj, *Digital Discovery*, 2022, **1**, 926–940.
- 32 S. Singh, M. Pareek, A. Changotra, S. Banerjee, B. Bhaskararao, P. Balamurugan and R. B. Sunoj, *Proc. Natl. Acad. Sci. U. S. A.*, 2020, **117**, 1339–1345.
- 33 M. Das, P. Sharma and R. B. Sunoj, *J. Chem. Phys.*, 2022, **156**, 114303.
- 34 A. E. Cuomo, S. Ibarraran, S. Sreekumar, H. Li, J. Eun, J. P. Menzel, P. Zhang, F. Buono, J. J. Song, R. H. Crabtree, V. S. Batista and T. R. Newhouse, *ACS Cent. Sci.*, 2023, **9**, 1768–1774.
- 35 X. Hou, S. Li, J. Frey, X. Hong and L. Ackermann, *Chem*, 2024, **10**, 1–12.
- 36 Y. Imai, W. Zhang, T. Kida, Y. Nakatsuji and I. Ikeda, *Tetrahedron Lett.*, 1997, **38**, 2681–2684.
- 37 A. V. Bedekar and P. G. Andersson, *Tetrahedron Lett.*, 1996, **37**, 4073–4076.
- 38 A. V. Bedekar, E. B. Koroleva and P. G. Andersson, *J. Org. Chem.*, 1997, **62**, 2518–2526.
- 39 K. Alexander, S. Cook and C. L. Gibson, *Tetrahedron Lett.*, 2000, **41**, 7135–7138.
- 40 A. Ramirez, V. C. Truc, M. Lawler, Y. K. Ye, J. Wang, C. Wang, S. Chen, T. Laporte, N. Liu, S. Kolotuchin, S. Jones, S. Bordawekar, S. Tummala, R. E. Waltermire and D. Kronenthal, *J. Org. Chem.*, 2014, **79**, 6233–6243.
- 41 T. Nakamura, T. Yoshihara, C. Tanegashima, M. Kadota, Y. Kobayashi, K. Honda, M. Ishiwata, J. Ueda, T. Hara, M. Nakanishi, T. Takumi, S. Itoharu, S. Kuraku, M. Asano, T. Kasahara, K. Nakajima, T. Tsuboi, A. Takata and T. Kato, *Mol. Psychiatry*, 2024, **29**, 2888–2904.
- 42 S. Kurosawa, F. Hasebe, H. Okamura, A. Yoshida, K. Matsuda, Y. Sone, T. Tomita, T. Shinada, H. Takikawa, T. Kuzuyama, S. Kosono and M. Nishiyama, *J. Am. Chem. Soc.*, 2022, **144**, 16164–16170.
- 43 T. Norris, T. F. Braish, M. Butters, K. M. DeVries, J. M. Hawkins, S. S. Massett, P. R. Rose, D. Santafianos and C. Sklavounos, *J. Chem. Soc., Perkin Trans. 1*, 2000, 1615–1622.
- 44 K. Kikushima, J. C. Holder, M. Gatti and B. M. Stoltz, *J. Am. Chem. Soc.*, 2011, **133**, 6902–6905.
- 45 G. Chelucci, M. G. Sanna and S. Gladiali, *Tetrahedron*, 2000, **56**, 2889–2893.
- 46 J. I. Garcia, B. Lopez-Sanchez, J. A. Mayoral, E. Pires and I. Villalba, *J. Catal.*, 2008, **258**, 378–385.
- 47 J. I. Garcia, G. Jimenez-Oses, B. Lopez-Sanchez, J. A. Mayoral and A. Velez, *Dalton Trans.*, 2010, **39**, 2098–2107.
- 48 C. Mazet, V. Köhler and A. Pfaltz, *Angew. Chem., Int. Ed.*, 2005, **44**, 4888–4891.
- 49 A. Button, D. Merk, J. A. Hiss and G. Schneider, *Nat. Mach. Intell.*, 2019, **1**, 307–315.
- 50 C. W. Coley, R. Barzilay, T. S. Jaakkola, W. H. Green and K. F. Jensen, *ACS Cent. Sci.*, 2017, **3**, 434–443.
- 51 A. M. Żurański, J. I. Martinez Alvarado, B. J. Shields and A. G. Doyle, *Acc. Chem. Res.*, 2021, **54**, 1856–1865.
- 52 Y. Ma, X. Zhang, L. Zhu, X. Feng, J. A. H. Kowah, J. Jiang, L. Wang, L. Jiang and X. Liu, *Molecules*, 2023, **29**, 5995.
- 53 B. J. Shields, J. Stevens, J. Li, M. Parasram, F. Damani, J. I. M. Alvarado, J. M. Janey, R. P. Adams and A. G. Doyle, *Nature*, 2021, **590**, 89–96.
- 54 J. Li, S. H. Liao, H. Xiong, Y. Y. Zhou, X. L. Sun, Y. Zhang, X. G. Zhou and Y. Tang, *Angew. Chem., Int. Ed.*, 2012, **51**, 8838–8841.
- 55 D. A. Evans, K. A. Woerpel, M. M. Hinman and M. M. Faul, *J. Am. Chem. Soc.*, 1991, **113**, 726–728.
- 56 A. F. Zahrt, J. J. Henle, B. T. Rose, Y. Wang, W. T. Darrow and S. E. Denmark, *Science*, 2019, **363**, eaau5631.
- 57 D. M. Lowe, Extraction of Chemical Structures and Reactions from the Literature, PhD thesis, University of Cambridge, 2012.
- 58 Additional details about each of the reaction datasets used in this study are provided in Section 1.4 in the SI.
- 59 R. E. Lowenthal, A. Abiko and S. Masamune, *Tetrahedron Lett.*, 1990, **31**, 6005–6008.
- 60 R. E. Lowenthal and S. Masamune, *Tetrahedron Lett.*, 1991, **32**, 7373–7376.



- 61 M. B. France, A. K. Milojević, T. A. Stitt and A. J. Kim, *Tetrahedron Lett.*, 2003, **44**, 9287–9290.
- 62 M. Itagaki, K. Masumoto, K. Suenobu and Y. Yamamoto, *Org. Process Res. Dev.*, 2006, **10**, 245–250.
- 63 A. Ebinger, T. Heinz, G. Umbricht and A. Pfaltz, *Tetrahedron*, 1998, **54**, 10469–10480.
- 64 H. Lebel, M. Parmentier, O. Leogane, K. Ross and C. Spitz, *Tetrahedron*, 2012, **68**, 3396–3409.
- 65 S. Chen, L. Wu, Q. Shao, G. Yang and W. Zhang, *Chem. Commun.*, 2018, **54**, 2522–2525.
- 66 As the ART dataset aggregates reactions reported across different peer reviewed literature, it inherently reflects the reporting practices, which may introduce biases such as an overrepresentation of successful or high-yielding reactions. Since the experimental data were not generated within this work, such biases cannot be fully eliminated.
- 67 W. Gao, P. Raghavan, R. Shprints and C. W. Coley, *J. Am. Chem. Soc.*, 2025, **147**, 8959–8968.
- 68 To meaningfully capture the distribution of the label (% ee) in our ART dataset, skewness is calculated, which is a statistical measure of how the data deviate from symmetry about the mean. The fact that a skewness of zero implies a symmetric distribution and a negative value indicates left-skewness, the skewness of  $-1.37$  found for the ART dataset indicates a pronounced deviation toward the lower end.
- 69 D. Weininger, *J. Chem. Inf. Comput. Sci.*, 1988, **28**, 31–36.
- 70 Z. Xiong, D. Wang, X. Liu, F. Zhong, X. Wan, X. Li, Z. Li, X. Luo, K. Chen, H. Jiang and M. Zheng, *J. Med. Chem.*, 2019, **63**, 8749–8760.
- 71 T. Akiba, S. Sano, T. Yanase, T. Ohta and M. Koyama, in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, ACM, New York, NY, USA, 2019, pp. 2623–2631.
- 72 K. V. Chuang and M. J. Keiser, *Science*, 2018, **362**, 6416.
- 73 F. Sandfort, F. Strieth-Kalthoff, M. Kühnemund, C. Beecks and F. Glorius, *Chem*, 2020, **6**, 1379–1390.
- 74 Additional details on the featurization strategies are provided in Section 1.5 in the SI.
- 75 In addition to DNN model building, we built various tree-based models (random forest, decision tree, and gradient boosting) and other models like support vector machines to predict % ee on our ART dataset. Although performances were comparable among these models, they are found to be less effective than the DNN. Each of these model performances with their hyperparameters are compiled in Tables S4–S11 in the SI.
- 76 A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser and I. Polosukhin, *Adv. Neural Inf. Process. Syst.*, 2017, **30**, 5998–6008.
- 77 S. Singh and R. B. Sunoj, *Digital Discovery*, 2022, **1**, 303–312.
- 78 J. Lu and Y. Zhang, *J. Chem. Inf. Model.*, 2022, **62**, 1376–1387.
- 79 Y. Kwon, D. Lee, Y.-S. Choi and S. Kang, *J. Cheminf.*, 2022, **14**, 2.
- 80 For each of the 30 splits, hyperparameter optimization of the DNN architecture is performed using the Optuna framework (ref. 71). The optimal hyperparameters are selected based on the lowest validation RMSE. Details of the optimal hyperparameters and corresponding model performance are provided in Tables S12–S14 in the SI.
- 81 Details of the hyperparameters and corresponding model performances for T5Chem, Transformer encoder, and ULMFiT are respectively provided in Tables S23, S24, S31, S32 and S39–S46 in the SI.
- 82 Details of the hyperparameters and corresponding model performances for Yield-BERT and MPNN are respectively provided in Tables S21, S22, S53 and S54 in the SI.
- 83 To evaluate alternative target values, we computed  $\Delta\Delta G^\ddagger$  from the reported % ee values and trained the AttentiveFP model accordingly. We note that training on % ee offers better predictive performance (RMSE =  $10.56 \pm 1.86$ ;  $R^2 = 0.77 \pm 0.08$ ) than on  $\Delta\Delta G^\ddagger$  (RMSE =  $0.43 \pm 0.05$ ;  $R^2 = 0.58 \pm 0.10$ ). This highlights that the model captures enantioselectivity trends more effectively when learning from % ee directly. See Section 5 in the SI for more details.
- 84 Additional justification for the use of AttentiveFP as the primary framework is provided in the SI (Section 2.7).
- 85 In this case, if a true ee < 50, the standard squared error (SE) between the predicted and true ee is taken into consideration just as in the case of AttentiveFP implementation. For ee > 50, on the other hand, the SE is halved, to enhance the model sensitivity to minority class samples.
- 86 The AttentiveFP-CI model maintains strong predictive performance (MAE  $\approx 7.10$ , RMSE  $\approx 9.80$ , RSE  $\approx 0.16$ , and  $R^2 \approx 0.84$ ) with low relative error across all class boundaries, confirming that it is free from significant systematic bias (see Table S71 in the SI).
- 87 In the ULMFiT regressor, each SMILES in the reaction dataset is augmented by 125 non-canonicalized/randomized forms (hyperparameter tuning using various degrees of SMILES augmentation (from 25, 50, 75, 100, and 125 to 150) helped us choose 125 as the optimal choice). For further details on hyperparameter tuning and model performance, see Section 2.4 in the SI. Additionally, Gaussian noise is added to the corresponding output of each randomized SMILES. This augmentation, with slight modifications to their ee values, effectively serves to increase samples in the low ee region. This is likely the reason why the addition of CI loss did not enhance performance, as the augmentation inherently addresses CI. The performance of the ULMFiT-CI algorithm with various class boundaries is provided in Tables S47–S52 in the SI.
- 88 The *t*-test results resulted in *p*-values of  $9.0 \times 10^{-6}$  for the transformer model and 0.0002 for ULMFiT.
- 89 Details and performance of all other CI-aware models considered in this study are provided in Section 2 in the SI. The results with various class boundaries for DNN-CI are provided in Tables S15–S20, T5Chem-CI are provided in Tables S25–S30, Transformer encoder-CI are provided in Tables S33–S38 and MPNN-CI are provided in Tables S55–S60.



- 90 To assess model robustness evaluated using random partitioning, we have additionally implemented a scaffold-based splitting using the Murcko scaffolds. The AttentiveFP-CI model with a class boundary of 30 exhibits an increase in test RMSE from  $9.80 \pm 1.40$  under random splitting to  $13.85 \pm 3.41$  under scaffold splitting (see Table S70 in the SI). Importantly, these results confirm that the AttentiveFP-CI algorithm is able to generalise to reactions with unseen core scaffolds, rather than relying solely on close structural analogues.
- 91 (a) Additional details about dataset creation and performance of the AttentiveFP model are provided in Section 3 in the SI.; (b) The AttentiveFP-CI model on different sparsity induced subsets of the BHA dataset, denoted as BHA-LTE, revealed a similar performance to that obtained using the same model on our ART dataset (see Table S73 in the SI).
- 92 (a) A conceptually related two-step strategy by Chung *et al.* first partitioned the reactions on the basis of their input chemical diversity and used cluster-specific regressors to improve accuracy across varied reaction types. See J. Chung, J. Li, A. I. Saimon, P. Hong and Z. Kong, *Sci. Rep.*, 2024, **14**, 12131; (b) In contrast, we employ the classification followed by regression (CFR) framework, which partitions reactions according to their output % *ee* labels to directly address class imbalance and data sparsity in the ART dataset. Thus, although both approaches share a two-step structure, they target distinct modelling challenges.
- 93 S. Ghosh, N. Jain and R. B. Sunoj, *ACS Catal.*, 2025, **15**, 20251–20269.
- 94 Using alternative boundaries as  $\mu$  and  $(\mu + \sigma)$ , the maximum attainable classification accuracies are found to be inferior to that obtained with  $(\mu - \sigma)$ . Further details are provided in Table S74 in the SI.
- 95 See Section 4.1 in the SI for more details on the DNN architecture and their performances in Tables S75 and S76.
- 96 See Tables S81–S83 in the SI for details of performance of the AttentiveFP model for the CFR-major and -minor classes.
- 97 Details of calculation of the weighted average of RMSE and  $R^2$  are provided in Table S84 in Section 4.4 in the SI.
- 98 See Tables S77–S80 in the SI for more details about the ULMFiT model performances for the CFR-major and -minor classes.
- 99 See Fig. S4 and Section 6 in the SI for more details.
- 100 S. Riniker and G. A. Landrum, *J. Cheminf.*, 2013, **5**, 1.
- 101 S. Liao, X.-L. Sun and Y. Tang, *Acc. Chem. Res.*, 2014, **47**, 2260–2272.
- 102 S. Sadeghyan, *arXiv*, 2018, preprint, arXiv:1804.05092, DOI: [10.48550/arXiv.1804.05092](https://doi.org/10.48550/arXiv.1804.05092).
- 103 R. B. Sunoj, *Acc. Chem. Res.*, 2016, **49**, 1019–1028.
- 104 J. S. Harvey, S. P. Simonovich, C. R. Jamison and D. W. C. MacMillan, *J. Am. Chem. Soc.*, 2011, **133**, 13782–13785.
- 105 R. Connon, B. Roche, B. V. Rokade and P. J. Guiry, *Chem. Rev.*, 2021, **121**, 6373–6521.
- 106 Additional details on the extraction of chiral Box ligands from the PubChem library and the computational filtering protocol are provided in Section 7 in the SI.
- 107 The predicted % *ee* exceeding 99% is likely due to data sparsity and the use of linear output layers, which, in most of the DL regression tasks, inherently allow extrapolation beyond the training distribution. Mathematically, a linear output layer applies a transformation of the kind  $y = Wx + b$ , where  $W$  and  $b$  are the parameters learned by the model. Without constraints, if  $W = 0.98$  and  $b = 4.98$ , an input ( $x = 99$ ) could yield ( $y = 102$ ), exceeding the expected range of 0 to 99. Thus, the model can return overestimated predictions, especially when trained on imbalanced datasets with sparse data near the extrema. See H. Shimakawa, A. Kumada and M. Sato, *npj Comput. Mater.*, 2024, **10**, 11.
- 108 L. Simón and J. M. Goodman, *J. Am. Chem. Soc.*, 2008, **130**, 8741–8747.
- 109 D. J. Ager, A. H. M. de Vries and J. G. de Vries, *Chem. Soc. Rev.*, 2012, **41**, 3340–3380.
- 110 J. F. Teichert and B. L. Feringa, *Angew Chem. Int. Ed. Engl.*, 2010, **49**, 2486–2528.
- 111 The details of the benchmarking on the full 1075 *N,S*-acetylation reaction dataset are provided in Section 5 in the SI.
- 112 L. Cheng, P.-L. Shao, J. Lv, H. Xiao, Y. Sun, J. Yang, Z. Xu, M. Lv, G. Wang, S. Zhao, J. Li, Z. Jin, X. Tan, G. Xing and B. Zhang, *Nat. Comput. Sci.*, 2026, **6**, 145–155.
- 113 (a) Statistical significance tests are carried out for an effective comparison of performance obtained with AttentiveFP-CI and AttentiveFP models for both *N,S*-acetylation and asymmetric hydrogenation reactions. For *N,S*-acetylation, the AttentiveFP model yielded RMSE =  $8.57 \pm 0.82$  and  $R^2 = 0.91 \pm 0.02$ .; (b) In the case of hydrogenation, the difference in RMSE between AttentiveFP-CI ( $11.00 \pm 0.52$ ) and AttentiveFP ( $10.48 \pm 1.10$ ) is not statistically significant ( $p > 0.05$ ) (see Table S90 in the SI). However, the change in  $R^2$  (AttentiveFP-CI:  $0.52 \pm 0.03$ ,  $p = 0.0355$ ; AttentiveFP:  $0.60 \pm 0.17$ ) indicates a meaningful gain in performance.
- 114 The dataset is randomly divided into 70 : 10 : 20 training, validation, and test sets. The hyperparameter tuning for the ULMFiT model is performed on the validation set using the Optuna framework (ref. 71), and the resulting optimal hyperparameters are applied to the model for predictions on the test set. The model performance is reported in terms of RMSE and  $R^2$  obtained as the average over 30 different runs each using a randomly created train-test split.
- 115 The dataset is randomly divided into training, validation, and test sets in a 70 : 10 : 20 ratio and trained on five such independent random splits.
- 116 X. Yin, C.-Y. Hsieh, X. Wang, Z. Wu, Q. Ye, H. Bao, Y. Deng, H. Chen, P. Luo, H. Liu, T. Hou and X. Yao, *Research*, 2024, **7**, 0292.



- 117 X. Wang, C.-Y. Hsieh, X. Yin, J. Wang, Y. Li, Y. Deng, D. Jiang, Z. Wu, H. Du, H. Chen, Y. Li, H. Liu, Y. Wang, P. Luo, T. Hou and X. Yao, *Research*, 2023, **6**, 0231.
- 118 E. Mathew, K. A. Emmitte and J. Liu, *ACS Omega*, 2025, **10**, 32968–32986.
- 119 C. Liu, Y. Sun, R. Davis, S. T. Cardona and P. Hu, *J. Cheminf.*, 2023, **15**, 29.
- 120 K. Ermanis, A. C. Colgan, R. S. J. Proctor, B. W. Hadrys, R. J. Phipps and J. M. Goodman, *J. Am. Chem. Soc.*, 2020, **142**, 21091–21101.
- 121 M. N. Grayson, *J. Org. Chem.*, 2021, **86**, 13631–13635.

