



## Pharmaceutical Biotechnology

## A machine learning approach to identify active polysorbate 20 degrading hydrolases in biopharmaceutical formulations

Melanie Maier<sup>a,b</sup>, Simon Kluters<sup>a</sup>, Joey Studts<sup>a</sup>, Matthias Franzreb<sup>b</sup>, Patrick Garidel<sup>c,\*</sup>, Viktor Groß<sup>c</sup><sup>a</sup> Drug Substance Development Biologicals, Boehringer Ingelheim Pharma GmbH & Co., KG, Biberach an der Riss, Germany<sup>b</sup> Institute of Functional Interfaces, Karlsruhe Institute of Technology, Karlsruhe, Germany<sup>c</sup> Drug Product Development Biologicals, DPB-TIP, Boehringer Ingelheim Pharma GmbH & Co., KG, Biberach an der Riss, Germany

## ARTICLE INFO

## Article history:

Received 20 February 2026

Revised 22 April 2026

Accepted 10 May 2026

Available online 4 June 2026

## Keywords:

Polysorbate degradation

Hydrolases

Host cell proteins

CHO

RP-UPLC-MS

Machine learning

Classification models

Biopharmaceutical formulations

Degradation fingerprints

## ABSTRACT

Polysorbate degradation by host cell-derived hydrolases presents a critical challenge in biopharmaceutical formulations. It can lead to fatty acid release, particle formation and reduced product stability. Mass spectrometry-based host cell protein (HCP) analysis is widely used for HCP identification, but detection becomes challenging in formulations where monoclonal antibodies are present in large excess. In such cases, hydrolases can remain undetected, despite being enzymatically active at trace levels.

In this study, we demonstrate that individual CHO-derived hydrolases generate distinct polysorbate degradation fingerprints that can be detected by reverse phase ultra performance liquid chromatography coupled to mass spectrometry (RP-UPLC-MS) and classified using supervised machine learning. Models were trained on single time point fingerprints comprising approximately 50 measurements for five hydrolases (CES1F, CES2C, LPLA2, PPT1 and PAF-AH). Evaluated algorithms included Logistic Regression, Random Forest, Gradient Boosting, Support Vector Classifier, AdaBoost, and Artificial Neural Networks. Seven out of eight models achieved 100 % accuracy on the test set, confirming that enzyme-specific information is preserved in single measurements in the presence of individual enzymes, independent of enzyme concentration or degradation time.

External validation using an independently prepared hydrolase spike sample confirmed the robustness of the models. Prediction confidence was high at early degradation stages and decreased at late stages, as enzyme-specific degradation fingerprints became more similar. This work presents an activity-based classification framework for the functional identification of polysorbate degrading hydrolases. The approach supports downstream monitoring and risk-based mitigation strategies by identifying the enzymes that drive polysorbate hydrolysis under formulation conditions.

© 2026 The Authors. Published by Elsevier Inc. on behalf of American Pharmacists Association®. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>)

## Introduction

Polysorbates, particularly polysorbate 20 (PS20) and polysorbate 80 (PS80), are widely used as non-ionic surfactants in monoclonal antibody formulations.<sup>1,2</sup> Their primary function is to stabilise therapeutic proteins in the final drug product by preventing interfacial stress, aggregation, and particle formation.<sup>3</sup> Polysorbates are chemically highly heterogeneous mixtures consisting of a sorbitan or isosorbide core, polyoxyethylene (POE) chains and esterified fatty acids, resulting in hundreds of distinct molecular species.<sup>4,5</sup>

Despite their stabilising properties, polysorbates are susceptible to chemical and enzymatic degradation.<sup>6–11</sup> Of particular concern is the hydrolysis of polysorbates by host cell proteins (HCPs), which can

lead to the release of free fatty acids and the formation of visible and subvisible particles.<sup>12</sup> Notably, enzymatic degradation can occur even at trace concentrations, highlighting the need for sensitive detection methods.<sup>13,14</sup>

The functional identification of hydrolases in complex formulations remains a significant analytical challenge.<sup>13</sup> Mass spectrometry (MS)-based proteomics is widely used for HCP profiling in biopharmaceutical development. However, its application to formulated drug products is limited. High antibody concentrations (>100 g/L) dominate the MS signal and can leave trace level impurities undetected. Furthermore, not all HCPs are enzymatically active under formulation conditions, and only a small subset contributes to polysorbate degradation.<sup>15–17</sup> The lack of correlation between enzyme abundance and functional activity further complicates risk assessment and targeted mitigation.

\* Corresponding author.

E-mail address: [Patrick.Garidel@boehringer-ingelheim.com](mailto:Patrick.Garidel@boehringer-ingelheim.com) (P. Garidel).

Several activity-based assays for assessing polysorbate degradation have been reported previously, including fluorescence-based esterase assays and activity-based protein profiling probes.<sup>18–20</sup> These approaches offer high sensitivity and throughput and are valuable for screening overall hydrolytic potential. However, because they rely on indirect readouts, their results do not necessarily correlate quantitatively with actual polysorbate degradation under formulation conditions.<sup>21</sup>

Identifying which specific hydrolases are responsible for polysorbate degradation is therefore critical for risk-based process development. Such knowledge can guide interventions during upstream and downstream processing, for example by implementing targeted gene knockouts in production cell lines or by targeted optimisation of purification steps.<sup>22–25</sup>

Previous studies have shown that individual hydrolases exhibit distinct substrate preferences, resulting in characteristic degradation patterns or ‘fingerprints’.<sup>21,26</sup> These fingerprints reflect the enzyme’s specificity towards different polysorbate species and can be captured using reverse phase ultra performance liquid chromatography coupled to mass spectrometry (RP-UPLC-MS). The RP-UPLC-MS fingerprinting approach directly measures polysorbate degradation using polysorbate itself as the substrate. This allows enzyme-specific degradation patterns and substrate selectivity to be resolved rather than providing a generic activity readout.

Because these fingerprints are highly complex and contain a lot of co-elution of PS species, manual integration is not feasible. Recent advances in artificial intelligence and machine learning have shown that algorithmic pattern recognition approaches significantly outperform manual evaluation in settings involving highly multidimensional or co-eluting chromatographic and MS signals. AI-based feature extraction methods have become widely used in analytical chemistry to simplify complex datasets, reduce integration bias, and enable consistent feature extraction.<sup>27,28</sup> Therefore, an in-house developed generative model was employed to analyse the MS fingerprints and extract the intensities of the individual PS species that are the features later used for classification.<sup>29</sup> Since there are 35 PS species extracted, manual interpretation is cumbersome and not reliable.

As mentioned before, at a feature space of 35-dimension algorithmic pattern recognition approaches outperform manual evaluation by far. To address this challenge a comprehensive comparison across different model types, with diverse learning paradigms, such as Logistic Regression (logReg), Random Forest Classifier (RFC), and Artificial Neural Networks (ANNs) to differentiate hydrolases based on their degradation fingerprints. Machine learning algorithms have been successfully applied in biopharmaceutical analytics. For example, they can be used to predict product quality from LC-MS data, select excipient for drug product formulations,<sup>30</sup> or interpret UV-based measurements for process monitoring and sub-visible particle classification.<sup>31,32</sup> They can also be used to determine whether a process is running under normal or critical conditions.<sup>31,32</sup> Their ability to detect subtle, multidimensional patterns makes them suitable for differentiating hydrolases based on degradation fingerprints.

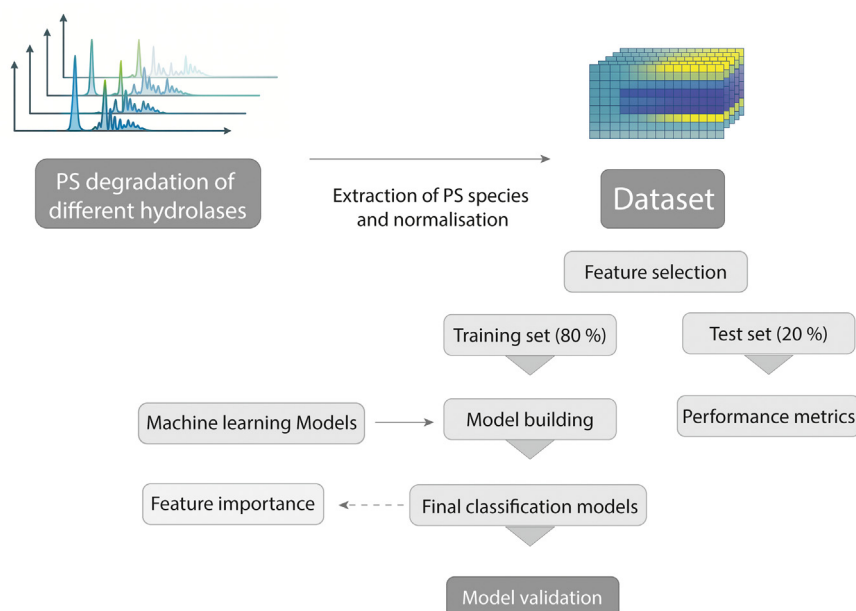
In this work, we demonstrate an activity-based classification approach that combines RP-UPLC-MS fingerprinting with supervised machine learning algorithms to identify polysorbate degrading hydrolases from single measurements.

## Materials and methods

An overview of the experimental and computational workflow from enzymatic degradation to model validation is illustrated in Fig. 1. It comprises (i) enzymatic degradation of polysorbate-containing formulations, (ii) acquisition of degradation fingerprints by RP-UPLC-MS, (iii) data preprocessing and normalisation, and (iv) supervised machine learning for hydrolase classification. The following sections describe each step in detail.

### Materials

Polysorbate 20 high purity (PS20 HP) was obtained from Croda GmbH (Nettetal, Germany). Polysorbate hydrolases were produced in-house at Boehringer Ingelheim Pharma GmbH & Co. KG (Biberach an der Riss, Germany). Histidine and sodium chloride of parenteral grade were purchased from Sigma-Aldrich (St. Louis, USA).



**Fig. 1.** Overview of the experimental and computational workflow. CHO-derived hydrolases were incubated with polysorbate 20 HP-containing formulations and analysed by RP-UPLC-MS over hydrolase-specific time courses. Species were extracted via generative modelling, and intensities were normalised to the initial timepoint, with the data compiled into a fingerprint dataset comprising five hydrolases, each with up to 50 timepoints and 23 polysorbate species. The dataset was split into training (80%) and test (20%) sets, and different classifiers were trained. The final models were then applied to external validation samples.

**Table 1**

Expressed and purified polysorbate hydrolases with abbreviations and their Uniprot accession numbers.

Uniprot accession no.	Protein name	Gene name
AOA06117 × 9	CES1F (Liver carboxylesterase B-1-like protein)	<i>Ces1f</i>
G3IIG1	CES2C (Carboxylic ester hydrolase 2C)	<i>Ces2c</i>
G3HKV9	LPLA2 (Group XV phospholipase A2)	<i>Pla2g15</i>
AOA8C2M2A7	PAF-AH (1-alkyl-2-acetylglucosphosphocholine esterase)	<i>Pla2g7</i>
G3HN89	PPT1 (Palmitoyl-proteinthioesterase 1)	<i>Ppt1</i>

### Recombinant expression and purification of hydrolases

Recombinant hydrolases (Table 1) were recombinantly expressed and purified as previously described.<sup>17</sup> Each protein was produced using a stably expressing Chinese Hamster Ovary CHO cell line, followed by orthogonal affinity chromatography and polishing steps.

In addition to the five hydrolases included in this study (CES1F, CES2C, LPLA2, PPT1 and PAF-AH), several other CHO-derived hydrolases previously reported to exhibit polysorbate-degrading activity were initially evaluated under identical formulation-relevant conditions. However, within the experimental time window, no measurable polysorbate degradation was observed for these candidates. Prolonged incubation led to the onset of oxidative polysorbate degradation in the corresponding blanks, which precluded a reliable assessment of low enzymatic activity. Consequently, these hydrolases were not included in the subsequent fingerprinting analyses.

### Sample generation and degradation fingerprint acquisition

Each purified hydrolase was incubated individually in 10 mM histidine pH 6 containing 0.4 g·L<sup>-1</sup> PS20 HP at room temperature in Eppendorf tubes (Eppendorf, Hamburg, Germany). To closely reflect actual conditions, a low hydrolase concentration of 20 ng·mL<sup>-1</sup> was chosen (protein concentration determined by UV absorbance at 280 nm). Samples were collected over a total duration of up to 50 days. To capture both early and late degradation stages, a non-uniform sampling scheme was applied, with denser sampling at the beginning of the incubation and wider intervals at later stages, resulting in up to 50 sampling time points per hydrolase.

Polysorbate blanks (no enzyme added) were monitored in parallel under identical conditions. Oxidative degradation in the blank became apparent after approximately two months of incubation, and time points beyond the onset of measurable oxidative degradation were therefore excluded for hydrolases that did not show enzymatic activity within the initial monitoring period. At each time point, samples were analysed by RP-UPLC-MS and PS subspecies were extracted via generative modelling to quantify the relative abundance of different PS species accordingly to Roelants et al. 2025.<sup>29</sup> Peak intensities were normalised to the initial time point ( $t_0$ ) to generate degradation fingerprints for each hydrolase.

### Fingerprint dataset preprocessing and feature selection

The complete dataset comprised normalised intensity values for up to 50 time points across 35 polysorbate (PS20) subspecies per hydrolase. During initial data evaluation, particularly in strongly degraded samples, outliers in PS subspecies were identified, suggesting potential inaccuracies in the extraction process. The extraction process was performed accordingly to Roelants et al. 2025.<sup>29</sup> To ensure data reliability, all polysorbate subspecies were manually re-evaluated for the presence of outliers. Given the relatively large feature space (35 PS20 subspecies) compared with the dataset size, subspecies showing even a small number of outlier events were

excluded to obtain a clean input dataset and minimise potential bias in the evaluated models.

It has been established that PS20 and PS80 hydrolysis invariably yields free sorbitan, isosorbide, and POE species. Consequently, these subspecies are consistently present at high levels. As they lack informative properties for distinguishing among hydrolases, they were thus excluded from the classification. Furthermore, the feature importance score of the Random Forest Classifier was also considered during the reduction of the number of features.

This refinement reduced the number of subspecies from 35 to 23, a step taken to minimise bias and enhance the performance and interpretability of subsequent classification models.

To focus the classification exclusively on degradation patterns, all sample profiles were normalised to the initial reference sample. This normalisation strategy ensures that classification is based solely on the relative ratio among the 23 selected subspecies at each sampling time point. Consequently, metadata such as time, enzyme concentration, enzymatic activity, and overall PS concentration were excluded from the classification input, as they are not required for capturing the intrinsic degradation signature.

### Model training and evaluation

All computational analysis were conducted in Python using Scikit-Learn, Numpy, Pandas and Tensorflow. The processed fingerprint dataset was randomly divided into training (80 %) and test (20 %) subsets using stratified sampling to preserve class balance among the five hydrolases (CES1F, CES2C, LPLA2, PPT1, and PAF-AH).

Eight supervised classification models were evaluated for their ability to assign the correct hydrolase based on their individual degradation fingerprint profile. Therefore, multiple model classes were evaluated, including linear models (Logistic Regression), tree-based ensemble models (Random Forest, Gradient Boosting, AdaBoost), kernel-based classifiers (Support Vector Classifier), and Artificial Neural Networks of varying depth (Tables 2 and 3). These model types were selected to represent complementary learning paradigms and to assess whether hydrolase-specific information is consistently captured across conceptually different algorithms.

The models were tested for the following parameters listed in Tables 2 and 3. The models and the best performing parameters are as follows: Logistic Regression (multi class: ovr, penalty: L2 regularization, solver: lbfgs), Random Forest (estimators: 120, bootstrap: False, max features: 3), Gradient Boosting (learning rate: 1.5, max depth = 3, 100 estimators), Support Vector Classifier (RBF kernel, C = 1.0). Additionally, three Artificial Neural Networks were implemented. Since the dataset is relatively small, drop out rates were

**Table 2**

Hyperparameters evaluated for supervised classification models during model training and optimization.

Model	Hyperparameter	Values
Logistic Regression (logReg)	solver	['saga', 'lbfgs']
	multi_class	['ovr', 'multinomial']
	penalty	['l2', 'l1']
	C	['logspace(0.4,10)', 'logspace(0.4,15)', 'logspace(0.4,20)']
Random Forest Classifier (RFC)	n_estimators	[60, 90, 120, 150, 200]
	max_features	[2, 3, 4, 5, 10]
	bootstrap	[True, False]
AdaBoost (AB)	n_estimators	[50, 100, 150, 200]
	learning_rate	[0.01, 0.1, 0.2, 0.5]
Gradient Boosting (GB)	n_estimators	[5, 20, 40, 100]
	max_depth	[3, 5, 10, 20]
	learning_rate	[0.01, 0.1, 0.3, 0.6, 1.0, 1.5, 2.0]
Support Vector Classifier (SVC)	C	[0.01, 0.1, 1]
	kernel	['rbf', 'sigmoid']
	degree	[2, 3, 4]

**Table 3**  
Architecture and configuration details of Artificial Neural Networks (ANNs) for hydrolase classification.

Model	Hidden layers	Dropout	Output activation	Loss function	Optimizer	Classes	Mutually exclusive
ANN-1	23 → 16	None	Softmax	Sparse categorical cross-entropy	Adam	5	Yes
ANN-2	23 → 64 → 32 → 16	30 %, 30 %, 20 %	Softmax	Sparse categorical cross-entropy	Adam	5	Yes
ANN-3	23 → 64 → 32 → 16	30 %, 30 %, 20 %	Sigmoid	Sparse categorical cross-entropy	Adam	5	No

chosen with 20–30 % for the multi-layer neural networks (ANN-2 and ANN-3) to reduce the possibility of overfitting.<sup>33</sup> More detailed specifications of the models are shown in Table 3. Hyperparameter tuning was performed using five-fold cross-validation on the training set. Performance metrics included accuracy, macro-averaged F1-score, and class-wise precision and recall.

Feature importance scores were extracted from tree-based models, highlighting which polysorbate species contributed the most to hydrolase differentiation.

#### External validation and application to unknown samples

To evaluate the performance of the trained classification models, two types of external samples were analysed. For the first scenario that was used for validation, a formulation containing polysorbate that was spiked with a known CHO-derived hydrolase, PAF-AH. Importantly, the PAF-AH used for this validation experiment originated from an independent expression and purification campaign and was therefore distinct from the enzyme material used to generate the training dataset. This design ensured that the validation samples were fully independent at the enzyme batch level and had not been seen by the model during training. PAF-AH was spiked to a final concentration of 500 ng·mL<sup>-1</sup>.

The spiked formulation was incubated in 10 mM Histidine at pH 6 and samples were collected at four incubation time points. Each sample was analysed by RP-UPLC-MS to generate a degradation fingerprint that was then submitted to the trained models. This approach allowed assessment of whether the models could correctly identify the responsible hydrolase from a single measurement. For the application to an unknown sample, a purified monoclonal antibody formulation containing polysorbate 20 was analysed. The antibody concentration was 50 mg·mL<sup>-1</sup> and the initial PS20 concentration was 0.4 mg·mL<sup>-1</sup>.

The samples were subjected to RP-UPLC-MS analysis to generate a degradation fingerprint. The host cell protein composition and

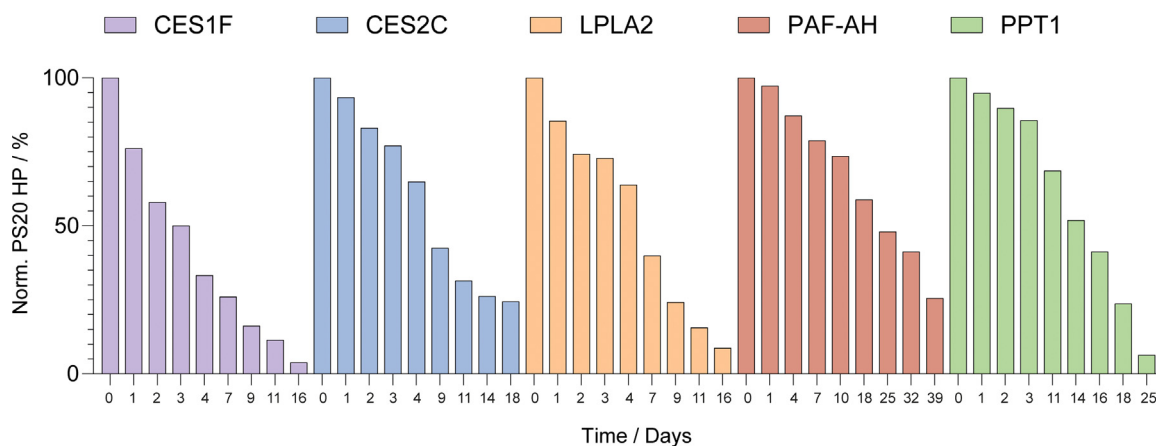
hydrolase content of this formulation were unknown. The fingerprint was submitted to the trained model, and the model was used to assign the most likely hydrolase class based on the learned degradation patterns.

## Results and discussion

### Overview of polysorbate hydrolase activity

To assess the activity of CHO-derived hydrolases under formulation relevant conditions, enzymes were incubated in a 10 mM histidine formulation (pH 6) at 0.4 mg·mL<sup>-1</sup> PS20 HP. To reflect realistic levels that may occur in drug product formulations,<sup>14</sup> 20 ng·mL<sup>-1</sup> of the respective hydrolases were added to the formulations. Degradation was monitored over several weeks at 25 °C, and the residual PS20 content was quantified by RP-UPLC-MS (Fig. 2). Among the tested enzymes, CES1F, CES2C, LPLA2, PAF-AH and PPT1 exhibited clear hydrolytic activity. CES1F and LPLA2 showed the most rapid degradation, while PAF-AH was the least active among the five enzymes. These findings confirm previous reports that these hydrolases are capable of degrading polysorbates.<sup>15,17</sup> Importantly, the observed activity occurred at very low enzyme concentrations, highlighting their potential relevance for product stability.

It has been demonstrated that a number of hydrolases, including LPL, LIPA, CES1 and IAH1, are active against polysorbate.<sup>15,17</sup> These enzymes have also been tested here; however, at a concentration of 20 ng·mL<sup>-1</sup>, no measurable PS degradation was observed within the first weeks. Extended incubation of these samples led to oxidative degradation in the polysorbate blank. This indicates that their activity is lower compared to the five enzymes highlighted above. These enzymes were therefore not spiked at artificially high concentrations (in the μg·mL<sup>-1</sup> range) and excluded from further fingerprinting and classification analysis. Nevertheless, the presented framework can be extended to incorporate additional enzymes as more data becomes available.



**Fig. 2.** Degradation of PS20 HP (0.4 mg·mL<sup>-1</sup>) by CHO-derived hydrolases at 20 ng·mL<sup>-1</sup>. PS content was quantified via RP-UPLC-MS and normalised to the initial PS concentration. Each colour represents a different hydrolase: CES1F (purple), CES2C (blue), LPLA2 (orange), PAF-AH (red) and PPT1 (green). Please note that only selected time points are shown for clarity.

### Distinct fingerprint patterns across hydrolases

To assess the separability of hydrolase-specific fingerprints, the normalised intensities of polysorbate species were visualised using pair plots. Each fingerprint corresponds to a single time point of a hydrolase incubation and represents normalised intensities of individual polysorbate species relative to the initial time point ( $t_0$ ). Pair plots were generated using a subset of selected polysorbate species: S08, S12, I08, S12/12/14, S12/12/16, S10 and POE10 (Fig. 2). These species were chosen based on their high feature importance scores (Fig. 5) and chemical diversity, including both short- and long-chain esters. Each point represents the normalised intensities of two polysorbate species at one time point, resulting from the activity of one hydrolase. These visualisations allow the comparison of relative species abundances between enzymes. The diagonal histograms show the distribution of each species individually by hydrolase class. Notably, several species combinations exhibit well-separated clusters, indicating class-specific degradation patterns.

CES2C and CES1F form compact and for some of the species non-overlapping clusters. This suggests that, despite their shared classification as carboxylesterases, the two enzymes generate distinct degradation patterns for some of the PS subspecies. Especially in plots like S12 vs I08, CES1F and CES2C show well-defined separable clusters. This supports previous biochemical findings that CES1F and CES2C differ in their substrate preferences.<sup>34</sup>

PAF-AH forms a distinct cluster in the S08 vs S12 plot, indicating it can be differentiated from other enzymes using specific species combinations.

LPLA2 shows broader or even bimodal profiles for several species such as POE10 and S12. This may be due to non-linear degradation kinetics, where initially a rapid degradation (high affinity towards certain subspecies) can be observed that slows down over time (subspecies with little affinity are hydrolysed). In some scatter plots, the data points, which represent the paired normalised intensities of two polysorbate species, are spread across the full range, showing high variability but still with a recognizable pattern.

PPT1 shows a preference for triester species, especially S12/12/14 and S12/12/16. This is reflected in the diagonal histograms, where the intensities are close to zero, indicating rapid and consistent degradation of these species. In contrast, monoesters such as S10 and S12 retain high intensities, suggesting low activity on these substrates. In scatter plots that include triesters, PPT1 forms a distinct cluster, making it clearly distinguishable from the other hydrolases due to its pronounced triester selectivity.

Overall, this visualisation illustrates the differences between hydrolase degradation fingerprints. As the data is normalised to the initial time point, the observed differences are independent of the absolute degradation time or enzyme concentration, allowing direct comparison between fingerprints even when the time point is unknown. The class-specific clustering across *multiple species* suggests that polysorbate degradation fingerprints contain sufficient discriminatory information for hydrolase classification and provide the basis for supervised model training.

### Classification model performance

To evaluate whether the degradation fingerprints are sufficient to predict the responsible hydrolase, different classification approaches were compared (Tables 2 and 3). Models were trained on normalised fingerprints from five CHO-derived hydrolases, using 80 % of the dataset for training and 20 % for testing. Each fingerprint represented a single time point and contained normalised intensities of 23 polysorbate species. Due to the normalisation only the ratios at a specific time point of the used PS20 species is evaluated, therefore the classification is independent of the PS concentration itself. Furthermore, no information on incubation time or enzyme concentration was

included, ensuring that predictions relied solely on the degradation pattern.

The models evaluated in this study were Logistic Regression (log-Reg), Random Forest Classifier (RFC), Gradient Boosting (GB), AdaBoost (AB), Support Vector Classifier (SVC), and three Artificial Neural Networks (ANNs) of either increasing depth or designed for multi-class classification (Tables 2 and 3).

Fig. 4 shows the confusion matrices for the results of the seven top-performing models and for AdaBoost. Logistic Regression, Random Forest, Gradient Boosting, Support Vector Classifier and all three ANNs achieved perfect classification (accuracy = 1), producing identical confusion matrices. These results confirmed that even closely related fingerprint patterns (e.g. CES1F and CES2C) could be reliably separated without misclassifications. AdaBoost showed slightly lower performance (accuracy = 0.92), with occasional misclassifications between CES1F and CES2C and one misclassification of LPLA2 as CES1F. Although pair plots of selected species (Fig. 3) reveal distinct clusters for CES1F and LPLA2, this model still misclassified these enzymes.

Overall, the results demonstrate that all evaluated models, except AdaBoost, can accurately identify the enzyme from a single time point fingerprint, confirming that the relevant discriminatory features are embedded in the degradation pattern itself rather than requiring kinetic information.

### Feature importance and contribution of PS species

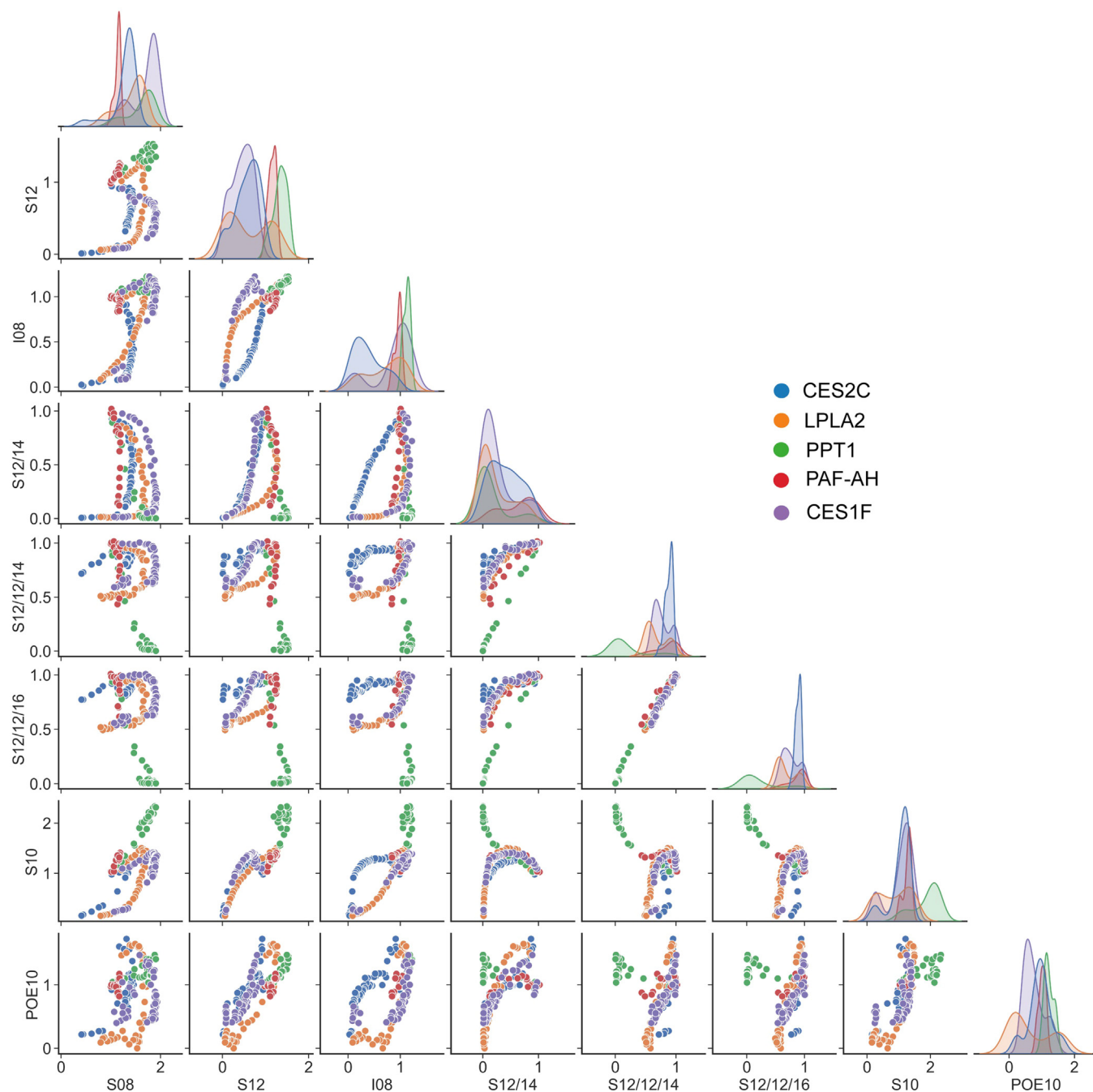
To understand which species contributed most to classification performance, feature importance scores were extracted from the trained Random Forest model (Fig. 5). These scores reflect how often and how effectively each species was used to split the data during training. A higher importance score indicates that the species was frequently used in decision splits and significantly improved class separation. Among these, the monoester S08 emerged as the most informative feature, followed by the triester species S12/12/14 and S12/12/16, and the isosorbide-based species I08. Additional contributors included S12 and S12/18, although their influence was comparatively lower.

This matches previous observations in the pair plots, where enzymes showed specific patterns in these species. For example, PPT1 consistently led to a rapid and complete depletion of triesters, resulting in near-zero intensities across time points. This makes these species particularly useful for distinguishing between hydrolases. PAF-AH showed a pronounced degradation of the monoester S08 while CES2C and CES1F showed distinct patterns for combinations with I08 species.

From a practical perspective, the feature importance patterns suggest that only a subset of polysorbate species is required for robust enzyme classification. This indicates that future models could be built on a reduced set of species, potentially decreasing data complexity while maintaining classification accuracy.

### External validation with spiked samples

To assess the robustness of the classification models beyond the internal test set, an external validation was performed using a formulation sample spiked with the hydrolase PAF-AH. This sample was generated independently of the training dataset and was analysed at four individual time points ( $t_1$ - $t_4$ ). The goal was to assess whether the trained models could correctly assign the enzyme class based on the polysorbate 20 fingerprint. Table 3 summarizes the classification probabilities for each model across the four time points. The results demonstrated that most models correctly identified PAF-AH at early degradation stages with high confidence, while prediction certainty decreases at later time points (Fig. 6).



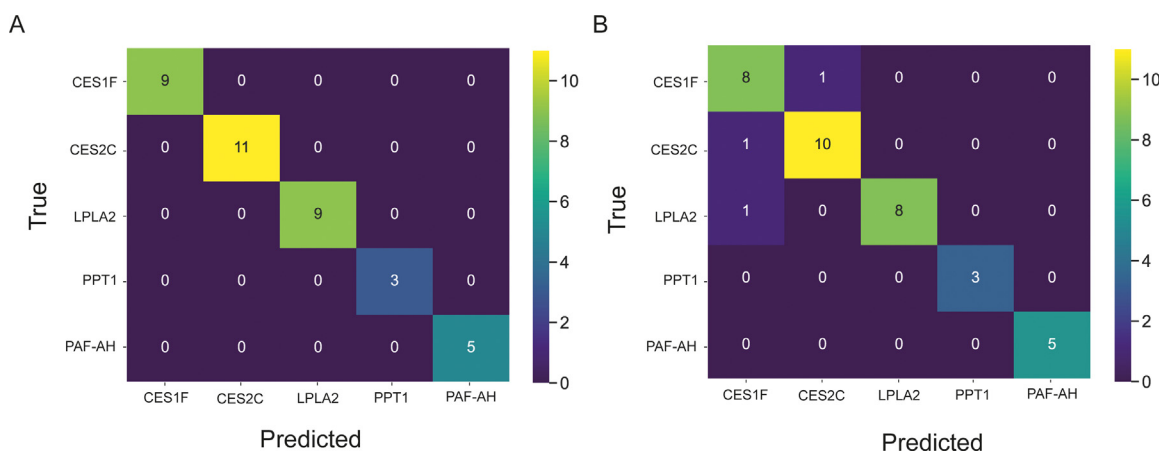
**Fig. 3.** Pair plots of selected polysorbate species. Each point represents a single fingerprint for a single timepoint and a hydrolase. The selected species were chosen based on their contribution to classification performance (feature importance) and chemical representativeness. Colours indicate the responsible hydrolase. Axes show normalised intensities.

This observation indicates that fingerprint-based classification has inherent limitations at advanced degradation stages. Initially, degradation patterns are enzyme-specific and well-defined, allowing models to confidently assign classes. However, as degradation progresses, many species are already degraded, and the fingerprints become less distinctive. Consequently, specific fingerprint markers disappear and become non-detectable, reducing classification accuracy and increasing the likelihood of misclassifications. This highlights the importance of considering degradation stage when applying fingerprint-based approaches.

Gradient Boosting and Logistic Regression achieved perfect or near-perfect confidence for the first two samples, whereas confidence dropped substantially for the last sample, which was misclassified by Gradient Boosting as PPT1 with full certainty. Random Forest

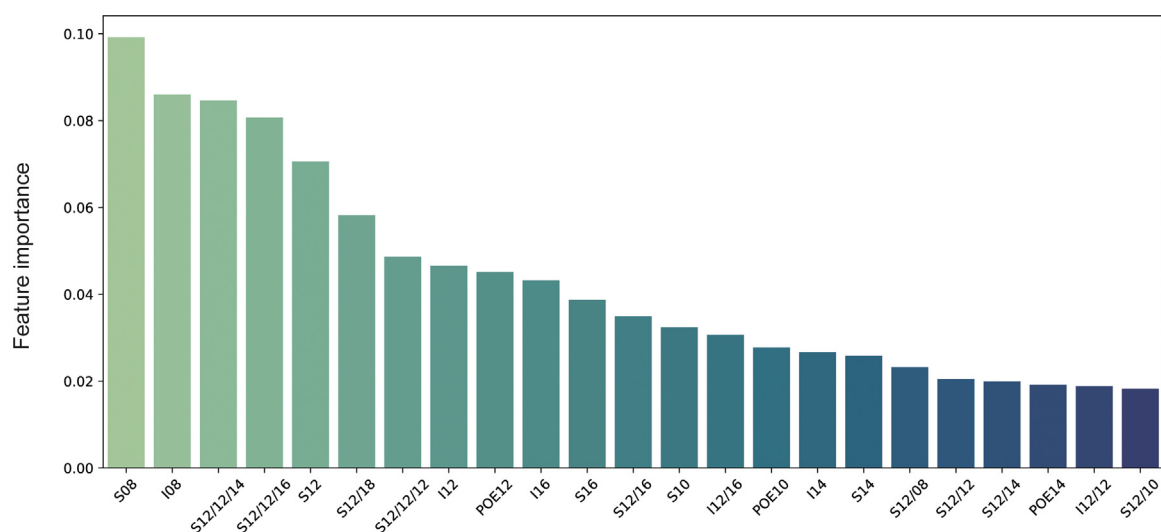
and Support Vector Classifier also showed declining confidence over time, although Support Vector Classifier maintained correct predictions across all four samples. AdaBoost exhibited moderate performance, correctly classifying early samples but misclassifying the latest time point.

ANN-2 and ANN-3 maintained high confidence for PAF-AH across all time points, outperforming several tree-based models, while ANN-1 showed inconsistent predictions and even misclassified late-stage samples. A likely reason is that ANN-1 has only one hidden layer and is therefore too shallow to capture the more complex patterns in the data. ANN-2 and ANN-3, which each contain three hidden layers, seem to extract these features more effectively. This allows them to correctly classify PAF-AH even when PS20 is already heavily degraded.

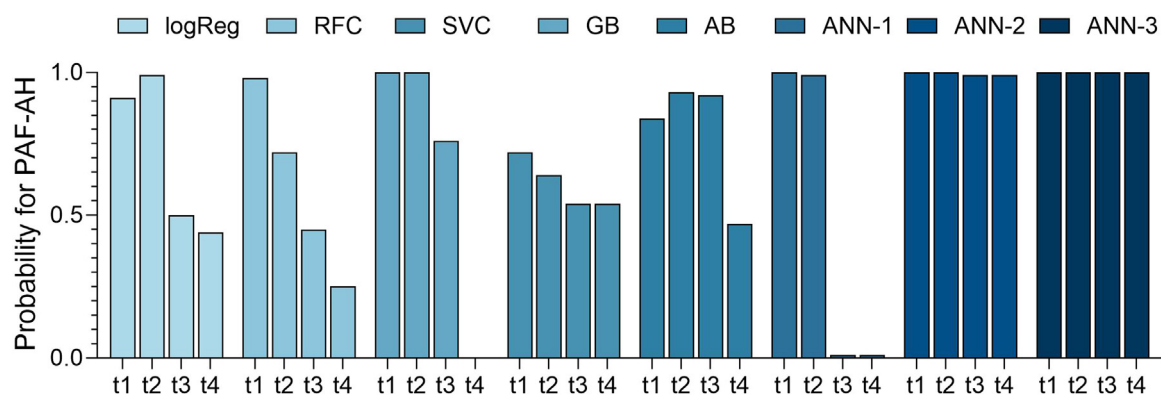


**Fig. 4.** Confusion matrices for classification models trained on single time point degradation fingerprints of five CHO-derived polysorbate hydrolases. A: Logistic Regression, Random Forest, Gradient Boosting, Support Vector Classifier, and three Artificial Neural Networks (ANN-1 to ANN-3). B: AdaBoost.

All models were evaluated on an internal test set. Confusion matrices display the class-wise prediction results for all five hydrolases: CES1F, CES2C, LPLA2, PAF-AH and PPT1.



**Fig. 5.** Feature importance scores of individual polysorbate species derived from the Random Forest model. Bar plots display the relative importance of the selected polysorbate species, as determined by the trained Random Forest classifier. Feature importance reflects how frequently and effectively each species was used to split decision trees during model training. Higher values indicate a stronger contribution to hydrolase classification.



**Fig. 6.** Predicted probability for an external PAF-AH spike sample. Predictions are shown for four samples taken at increasing degradation time points (sample t1 = earliest, sample t4 = latest) for eight models: Logistic Regression (logReg), Random Forest (RFC), Gradient Boosting (GB), Support Vector Classifier (SVC), AdaBoost (AB), and three Artificial Neural Networks (ANN-1, ANN-2, ANN-3). Probability values range from 0 to 1, with 1.0 indicating the highest level of confidence in a correct prediction. Bars are colour-coded by model. Most models show high confidence at the earliest time points, while confidence decreases at later stages, reflecting reduced fingerprint specificity due to extensive degradation.

These findings have important implications for more complex scenarios, such as mixtures of hydrolases in formulated drug products. At advanced degradation stages, fingerprints may reflect overlapping activity patterns from multiple enzymes, reducing the distinctiveness

of any single class. Models that rely predominantly on a small set of dominant features, such as tree-based classifiers, tend to perform poorly under these conditions. In contrast, deeper neural networks may benefit from their ability to integrate information across a large

**Table 4**  
Model-based classification probabilities for external test data using a PAF-AH spiked sample.

Model	Sample	CES1F	CES2C	LPLA2	PPT1	PAF-AH
Logistic Regression	t1	0.00	0.00	0.00	0.00	0.91
	t2	0.00	0.00	0.00	0.01	0.99
	t3	0.00	0.00	0.00	0.50	0.50
	t4	0.00	0.00	0.00	0.56	0.44
Random Forest Classifier	t1	0.01	0.00	0.01	0.00	0.98
	t2	0.00	0.00	0.08	0.20	0.72
	t3	0.01	0.06	0.06	0.47	0.45
	t4	0.02	0.21	0.21	0.48	0.25
Gradient Boosting	t1	0.00	0.00	0.00	0.00	1.00
	t2	0.00	0.00	0.00	0.00	1.00
	t3	0.00	0.00	0.00	0.24	0.76
	t4	0.00	0.00	0.00	1.0	0.00
Support Vector Classifier	t1	0.03	0.11	0.08	0.07	0.72
	t2	0.03	0.07	0.05	0.22	0.64
	t3	0.04	0.06	0.05	0.31	0.54
	t4	0.04	0.07	0.07	0.28	0.54
AdaBoost	t1	0.00	0.16	0.00	0.00	0.84
	t2	0.00	0.00	0.00	0.07	0.93
	t3	0.00	0.00	0.00	0.08	0.92
	t4	0.00	0.01	0.00	0.52	0.47
ANN-1	t1	0.00	0.00	0.00	0.00	1.00
	t2	0.00	0.00	0.00	0.01	0.99
	t3	0.00	0.00	0.00	0.99	0.01
	t4	0.00	0.00	0.00	0.99	0.01
ANN-2	t1	0.00	0.00	0.00	0.00	1.00
	t2	0.00	0.00	0.00	0.00	1.00
	t3	0.00	0.00	0.00	0.01	0.99
	t4	0.00	0.00	0.00	0.01	0.99
ANN-3	t1	0.00	0.49	0.01	0.56	1.00
	t2	0.00	0.13	0.01	0.83	1.00
	t3	0.00	0.01	0.00	1.00	1.00
	t4	0.00	0.02	0.13	0.95	1.00

number of features, enabling them to better approximate mixed signal patterns. To address this challenge, ANN-3 was trained as a multi-label classifier, allowing it to predict multiple independent labels (hydrolases) for a single input instance and to treat each label as an independent binary decision.

As shown in Table 4, ANN-3 also assigns high likelihood values to PPT1, especially at later sampling time points. This occurs because both hydrolases have stronger preferences for di- and triesters compared to the other investigated hydrolases. A more general, but also more complex solution would be to incorporate training data that mimic mixed-enzyme conditions into the current framework to tackle the problem when having more than one hydrolase in the same sample. This could be achieved for example by generating synthetic fingerprints that combine profiles from different hydrolases. Such an approach would help improve classification robustness when multiple enzymes contribute simultaneously to polysorbate degradation.

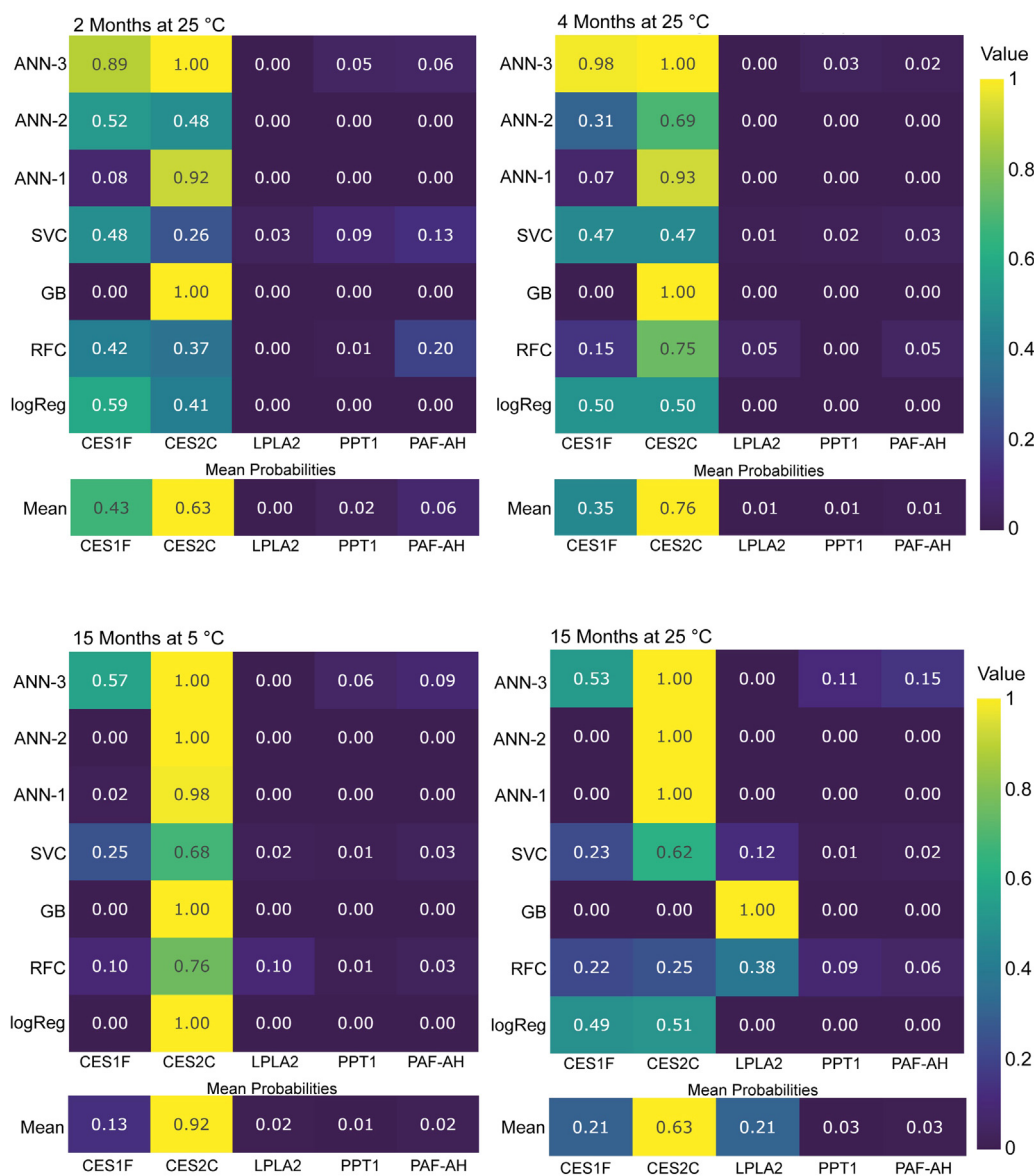
#### Prediction of unknown samples

To evaluate the applicability of the classification models beyond controlled conditions, we applied them to degradation fingerprints derived from a monoclonal antibody formulation. These model antibodies were purified under a typical downstream purification process and contained polysorbate 20 HP as part of the formulation. The composition and abundance of polysorbate hydrolases in these samples are unknown, making them a realistic test case for fingerprint-based

classification. Degradation fingerprints were generated at four incubation time points (2 months and 4 months at 25 °C, and 15 months at 5 °C and 25 °C). Following data extraction and normalisation, the fingerprints were submitted to the trained classification models.

The classification results are shown in Fig. 7. CES2C had the highest mean probability at all time points (0.63–0.92), indicating that it is the most likely contributor to polysorbate degradation in the tested mAb formulation. The second highest probability was observed for CES1F.

At the first time point (two months), CES2C was not clearly dominant, which may indicate that degradation was still in the early stages and that the fingerprint was not yet specific enough. This is supported by the fact that CES1F and CES2C are both carboxylesterases with a preference for monoester species, making them difficult to distinguish when only minor changes have occurred. In the 15-month sampling time point at 25 °C storage condition, the confidence for CES2C decreased again. This is expected when a subset of polysorbate subspecies has already been fully hydrolysed, because fewer informative species remain and the fingerprint provides less differentiation between enzymes, making the signal less pronounced. In general, classification is most challenging when degradation is either minimal or advanced, because many values are close to one or zero, which reduces class separability. For samples with very little change, only a few of the 23 polysorbate species differ from the reference. In this case, Random Forest trees may split on features that remain almost unchanged, which can lead to less reliable predictions. If only six of the 23 species show any deviation, the model may select



**Fig. 7.** Model-based classification probabilities for a monoclonal antibody formulation at four sampling points. Panels show heatmaps of predicted class probabilities for five CHO-derived polysorbate hydrolases (CES1F, CES2C, LPLA2, PPT1, PAF-AH), generated by seven classifiers at the following time points: 2 months (25 °C), 4 months (25 °C), 15 months (5 °C) and 15 months (25 °C). Each cell represents the probability assigned by one model to one hydrolase at the given time point. The bottom row in each panel displays the mean probability across all models. Colours encode probability from low (dark) to high (yellow).

features with values still close to 1, yielding near-random assignments for individual trees. Aggregation across 100 estimators stabilises the overall prediction, but the confidence remains lower. This makes the Random Forest less reliable at these stages compared to other applied algorithms here. Based on the probabilities for LPLA2, PPT1 and PAF-AH (Fig. 7), these enzymes are unlikely to play a critical role in this sample. Overall, the results suggest that CES2C is the most likely contributor to polysorbate degradation in this sample, with CES1F possibly contributing as well. A key advantage of the models is their ability to clearly exclude hydrolases that are not relevant for polysorbate degradation based on their mismatch with the observed fingerprints.

The fingerprint-based classification approach presented in this study offers several applications to improve biopharmaceutical development. Fingerprint data enables the identification of hydrolases that consistently contribute to polysorbate degradation across formulations. If certain enzymes emerge as dominant contributors while others show negligible activity, this supports risk-based

prioritization of hydrolases. These findings can directly guide downstream process development. For example, polishing steps such as ion exchange or hydrophobic interaction chromatography can be systematically optimised to selectively remove high-risk hydrolases.<sup>24</sup> Furthermore, the method allows comparison of different purification trains.<sup>35</sup> By analysing degradation fingerprints from the same monoclonal antibody purified through different downstream processes, it can be systematically assessed whether certain hydrolases are consistently removed or retained. This allows the evaluation which purification strategy is most effective at mitigating specific enzymes.

The approach also supports cell line engineering. Once high-risk hydrolases are identified, targeted gene knockouts can be implemented to eliminate their expression in CHO production cell lines. This strategy has already been applied successfully to reduce polysorbate degradation and improve product stability.<sup>23,25,36</sup>

Beyond process optimisation, fingerprint analysis provides insights into hitchhiking phenomena where specific hydrolases co-purify with certain mAbs due to molecular interactions.<sup>37</sup> If a

hydrolase consistently appears in fingerprints across multiple purification trains for a given antibody, this may be indicative of a hydrolyase-antibody binding interaction. Conversely, if the same purification strategy yields different fingerprint profiles for different antibodies, this suggests molecule-specific effects that cannot be addressed by standardized purification approaches.

Recent studies indicate that certain host cell-derived hydrolases can also associate with protein aggregates in formulated drug products.<sup>38,39</sup> This observation highlights a critical limitation of conventional HCP monitoring approaches that measure protein abundance but do not provide information on functional activity. Enzymes embedded in aggregates may be detected by proteomics, yet their actual contribution to polysorbate degradation remains uncertain. Conversely, low abundance hydrolases with high enzymatic activity may not only remain undetected but are often underestimated.

Importantly, fingerprint-based classification can be combined with LC-MS-based HCP profiling to strengthen enzyme identification. While proteomics provides broad coverage, it may miss low-abundance enzymes that are functionally active. This can be due to the extreme dynamic range in monoclonal antibody drug products.<sup>40,41</sup> Integrating both methods could provide a more complete picture of degradation risk and support systematic mitigation strategies.

A current limitation of the framework is that models were trained exclusively on single-enzyme profiles. In real formulation samples, however, multiple hydrolases may be present simultaneously, either in additive or overlapping activity. This is particularly relevant at late degradation stages, where multiple hydrolases may have contributed to the degradation and where model confidence tends to decline. In such cases, the fingerprint may not match one enzyme perfectly but instead represent a hybrid pattern, resulting in lower confidence scores or different class assignments at different time points within the sample. While the framework performs well for isolated enzymes, its interpretability decreases in the presence of complex mixtures. In these scenarios, ANN-3 (multi-label classifier) is better suited than strictly exclusive models because it does not force a single class and can indicate partially overlapping patterns. However, since it was trained only in single-enzyme profiles, such indications should not be interpreted as definitive evidence of mixtures. Moreover, when high-confidence assignment is not possible, predictions can still narrow down the set of likely candidates. For example, if the model consistently excludes certain enzymes across multiple time points. This enables a targeted focus on a smaller subset of hydrolases for further investigation.

One possible extension of the current approach would be to train the model on simulated mixtures of enzyme fingerprints. For instance, by linearly combining normalised degradation profiles from two or more hydrolases. This could help the model to learn intermediate patterns and help improve classification robustness in mixed samples.

In future work, predictions on unknown samples may be complemented by orthogonal methods such as LC-MS-based HCP profiling, allowing experimental confirmation of predicted enzyme identities and refinement of classification in complex mixtures.

## Conclusion and outlook

This work establishes degradation fingerprinting as a robust and functional approach for classifying polysorbate-degrading hydrolases. The approach enables enzyme identification from single RP-UPLC-MS measurements, independent of enzyme concentration or degradation time, demonstrating that the degradation pattern itself contains sufficient discriminatory information.

Comparable performance across multiple model classes indicates that hydrolase identification is primarily driven by the degradation fingerprints themselves rather than by the choice of algorithm. However, at sampling time points exhibiting increased PS20 degradation

(Fig. 6), ANN-2 and ANN-3 outperformed the other tested algorithms. In this context, non-mutually exclusive models such as ANN-3 offer particular advantages under conditions of reduced fingerprint specificity.

The current study focused on five highly active CHO hydrolases that represent the main contributors to polysorbate hydrolysis at formulation relevant conditions. Additional hydrolases with lower activity were excluded but can readily be integrated into the classification framework as more data become available.

Future work should extend the current framework to mixed-enzyme samples, as multiple hydrolases may act simultaneously in real formulations. Generating in-silico mixed fingerprints could improve model robustness for overlapping degradation profiles. In addition, ensemble learning strategies combining multiple base learners via a meta-classifier could further improve robustness and confidence of hydrolase identification and represent a promising direction for future work. Combining fingerprint-based classification with orthogonal analytical methods such as HCP proteomics may further strengthen enzyme identification in complex samples.

Overall, the method presented here offers a valuable tool for risk-based prioritization of polysorbate degrading hydrolases. Fingerprint-based classification can help identify enzymes that are functionally active under formulation conditions and thus guide targeted mitigation of polysorbate degradation in therapeutic formulations.

## Funding statement

This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors. This research was wholly funded by Boehringer Ingelheim Pharma GmbH & Co.KG.

## Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgements

We would like to thank Holger Thie for excellent project management and support throughout the study.

## Supplementary materials

Supplementary material associated with this article can be found, in the online version, at [doi:10.1016/j.xphs.2026.104355](https://doi.org/10.1016/j.xphs.2026.104355).

## References

1. Khan TA, Mahler HC, Kishore RSK. Key interactions of surfactants in therapeutic protein formulations: a review. *Eur J Pharm Biopharm.* 2015;97(Pt A):60–67. <https://doi.org/10.1016/j.ejpb.2015.09.016>.
2. Jones MT, Mahler HC, Yadav S, et al. Considerations for the use of polysorbates in biopharmaceuticals. *Pharm Res.* 2018;35(8):148. <https://doi.org/10.1007/s11095-018-2430-5>.
3. Singh SM, Bandi S, Jones DNM, Mallela KMG. Effect of polysorbate 20 and polysorbate 80 on the higher-order structure of a monoclonal antibody and its fab and fc fragments probed using 2D nuclear Magnetic resonance spectroscopy. *J Pharm Sci.* 2017;106(12):3486–3498. <https://doi.org/10.1016/j.xphs.2017.08.011>.
4. Evers DH, Schultz-Fademrecht T, Garidel P, Buske J. Development and validation of a selective marker-based quantification of polysorbate 20 in biopharmaceutical formulations using UPLC QDa detection. *J Chromatogr B.* 2020;1157:122287. <https://doi.org/10.1016/j.jchromb.2020.122287>.
5. Carle S, Evers DH, Hagelskamp E, Garidel P, Buske J. All-in-one stability indicating polysorbate 20 degradation root-cause analytics via UPLC-QDa. *J Chromatogr B.* 2024;1232:123955. <https://doi.org/10.1016/j.jchromb.2023.123955>.
6. Kishore RSK, Kiese S, Fischer S, Pappenberger A, Grauschopf U, Mahler HC. The degradation of polysorbates 20 and 80 and its potential impact on the stability of

- biotherapeutics. *Pharm Res.* 2011;28(5):1194–1210. <https://doi.org/10.1007/s11095-011-0385-x>.
7. Dwivedi M, Blech M, Presser I, Garidel P. Polysorbate degradation in biotherapeutic formulations: identification and discussion of current root causes. *Int J Pharm.* 2018;552(1–2):422–436. <https://doi.org/10.1016/j.ijpharm.2018.10.008>.
  8. Schultz-Fademrecht T, Schmid K, Scholz Y, Garidel P. Comprehensive investigation of factors influencing the degradation of polysorbate 20: insights into mechanisms and interactions using a design of experiments approach. *J Drug Deliv Sci Technol.* 2024;101:106156. <https://doi.org/10.1016/j.jddst.2024.106156>.
  9. Kozuch B, Weber J, Buske J, Mäder K, Garidel P, Diederichs T. Comparative stability study of polysorbate 20 and polysorbate 80 related to oxidative degradation. *Pharmaceutics.* 2023;15(9):2332. <https://doi.org/10.3390/pharmaceutics15092332>.
  10. Weber J, Pedri L, Peters LP, et al. Micellar solvent accessibility of esterified polyoxyethylene chains as crucial element of polysorbate oxidation: a density functional theory, molecular dynamics simulation and liquid chromatography/mass spectrometry investigation. *Mol Pharm.* 2025;22(3):1348–1364. <https://doi.org/10.1021/acs.molpharmaceut.4c01015>.
  11. Dwivedi M, Buske J, Haemmerling F, Blech M, Garidel P. Acidic and alkaline hydrolysis of polysorbates under aqueous conditions: towards understanding polysorbate degradation in biopharmaceutical formulations. *Eur J Pharm Sci.* 2020;144:105211. <https://doi.org/10.1016/j.ejps.2019.105211>.
  12. Groß V, Hornberger A, Schultz-Fademrecht T, Garidel P, Diederichs T. Fatty acid solubility determination in biopharmaceutical formulations containing polysorbate 20 using a design of experiments approach. *J Drug Deliv Sci Technol.* 2025;111:107153. <https://doi.org/10.1016/j.jddst.2025.107153>.
  13. Felix MN, Waerner T, Lakatos D, Reisinger B, Fischer S, Garidel P. Polysorbates degrading enzymes in biotherapeutics – a current status and future perspectives. *Front Bioeng Biotechnol.* 2025;12:1490276. <https://doi.org/10.3389/fbioe.2024.1490276>.
  14. Hall T, Sandefur SL, Frye CC, Tuley TL, Huang L. Polysorbates 20 and 80 degradation by group XV lysosomal phospholipase A2 isomer X1 in monoclonal antibody formulations. *J Pharm Sci.* 2016;105(5):1633–1642. <https://doi.org/10.1016/j.xphs.2016.02.022>.
  15. Kovner D, Yuk IH, Shen A, et al. Characterization of recombinantly-expressed hydrolytic enzymes from Chinese hamster ovary cells: identification of host cell proteins that degrade polysorbate. *J Pharm Sci.* 2023;112(5):1351–1363. <https://doi.org/10.1016/j.xphs.2023.01.003>.
  16. Graf T, Tomlinson A, Yuk IH, et al. Identification and characterization of polysorbate-degrading enzymes in a monoclonal antibody formulation. *J Pharm Sci.* 2021;110(11):3558–3567. <https://doi.org/10.1016/j.xphs.2021.06.033>.
  17. Maier M, Weiß L, Zeh N, et al. Illuminating a biologics development challenge: systematic characterization of CHO cell-derived hydrolases identified in monoclonal antibody formulations. *mAbs.* 2024;16(1):2375798. <https://doi.org/10.1080/19420862.2024.2375798>.
  18. Bhargava AC, Mains K, Siu A, et al. High-throughput, fluorescence-based esterase activity assay for assessing polysorbate degradation risk during biopharmaceutical development. *Pharmaceut Res.* 2021;38(3):397–413. <https://doi.org/10.1007/s11095-021-03011-1>.
  19. Liu GY, Nie S, Zheng X, Li N. Activity-based protein profiling probe for the detection of enzymes catalyzing polysorbate degradation. *Anal Chem.* 2022;94(24):8625–8632. <https://doi.org/10.1021/acs.analchem.2c00059>.
  20. Gupta SK, Graf T, Edelmann FT, et al. A fast and sensitive high-throughput assay to assess polysorbate-degrading hydrolytic activity in biopharmaceuticals. *Eur J Pharm Biopharm.* 2023;187:120–129. <https://doi.org/10.1016/j.ejpb.2023.04.021>.
  21. Maier M, Gross V, Weiss L, et al. Specific polysorbate fingerprints of CHO hydrolases and implications for indirect assays. *J Pharm Sci.* 2025;115(2):104126. <https://doi.org/10.1016/j.xphs.2025.104126>.
  22. Weiß L, Schmieder-Todtenhaupt V, Haemmerling F, Lakatos D, Schulz P, Fischer S. Multi-lipase gene knockdown in Chinese hamster ovary cells using artificial microRNAs to reduce host cell protein mediated polysorbate degradation. *Biotechnol Bioeng.* 2024;121(1):329–340. <https://doi.org/10.1002/bit.28563>.
  23. Chiu J, Valente KN, Levy NE, Min L, Lenhoff AM, Lee KH. Knockout of a difficult-to-remove CHO host cell protein, lipoprotein lipase, for improved polysorbate stability in monoclonal antibody formulations. *Biotechnol Bioeng.* 2017;114(5):1006–1015. <https://doi.org/10.1002/bit.26237>.
  24. Maier M, Schneider S, Weiss L, et al. Tailoring polishing steps for effective removal of polysorbate-degrading host cell proteins in antibody purification. *Biotechnol Bioeng.* 2024;121(10):3181–3195. <https://doi.org/10.1002/bit.28767>.
  25. Weiß L, Zeh N, Maier M, Lakatos D, Otte K, Fischer S. Without a trace: multiple knockout of CHO host cell hydrolases to prevent polysorbate degradation in biologics. *Trends Biotechnol.* 2025;43(8):1982–2002. <https://doi.org/10.1016/j.tibtech.2025.04.016>.
  26. Glücklich N, Carle S, Buske J, Mäder K, Garidel P. Assessing the polysorbate degradation fingerprints and kinetics of lipases – how the activity of polysorbate degrading hydrolases is influenced by the assay and assay conditions. *Eur J Pharm Sci.* 2021;166:105980. <https://doi.org/10.1016/j.ejps.2021.105980>.
  27. Bilbao A, Munoz N, Kim J, et al. PeakDecoder enables machine learning-based metabolite annotation and accurate profiling in multidimensional mass spectrometry measurements. *Nat Commun.* 2023;14(1):2461. <https://doi.org/10.1038/s41467-023-37031-9>.
  28. Beck AG, Muhoberac M, Randolph CE, et al. Recent developments in machine learning for mass spectrometry. *ACS Meas Sci Au.* 2024;4(3):233–246. <https://doi.org/10.1021/acsmesuresci.3c00060>.
  29. Roelants P, Choubh RR, Verbeeck N, et al. Extraction of the polysorbate 20 and 80 fingerprint via generative modeling. *Int J Pharm: X.* 2025;10:100433. <https://doi.org/10.1016/j.ijpx.2025.100433>.
  30. Vidal-Henriquez E, Holder T, Lee NF, Pompe C, Teese MG. Machine learning driven acceleration of biopharmaceutical formulation development using Excipient Prediction Software (ExPreSo). *bioRxiv.* 2025. <https://doi.org/10.1101/2025.02.12.637685>. Published online 2025.02.12.637685.
  31. del Rio AL, Pacios-Michelena A, Picart-Armada S, Garidel P, Nikels F, Kube S. Sub-visible particle classification and label consistency analysis for flow-imaging microscopy via machine learning methods. *J Pharm Sci.* 2024;113(4):880–890. <https://doi.org/10.1016/j.xphs.2023.10.041>.
  32. Rathore AS, Nikita S, Thakur G, Mishra S. Artificial intelligence and machine learning applications in biopharmaceutical manufacturing. *Trends Biotechnol.* 2023;41(4):497–510. <https://doi.org/10.1016/j.tibtech.2022.08.007>.
  33. Labach A, Salehinejad H, Valaee S. Survey of dropout methods for deep neural networks. *arXiv.* 2019. <https://doi.org/10.48550/arxiv.1904.13310>. Published online.
  34. Wang D, Zou L, Jin Q, Hou J, Ge G, Yang L. Human carboxylesterases: a comprehensive review. *Acta Pharm Sin B.* 2018;8(5):699–712. <https://doi.org/10.1016/j.apsb.2018.05.005>.
  35. Lakatos D, Idler M, Stibitzky S, et al. Buffer system improves the removal of host cell protein impurities in monoclonal antibody purification. *Biotechnol Bioeng.* 2024;121(12):3869–3880. <https://doi.org/10.1002/bit.28844>.
  36. Yuk IH, Ko P, Ahyow P, et al. Engineering Chinese hamster ovary cells to mitigate polysorbate degradation in biotherapeutics. *Biotechnol Bioeng.* 2025;122(11):3139–3159. <https://doi.org/10.1002/bit.70037>.
  37. Hecht ES, Mehta S, Weckler AT, et al. Insights into ultra-low affinity lipase-antibody noncovalent complex binding mechanisms. *Mabs.* 2022;14(1):2135183. <https://doi.org/10.1080/19420862.2022.2135183>.
  38. Herman CE, Min L, Choe LH, et al. Analytical characterization of host-cell-protein-rich aggregates in monoclonal antibody solutions. *Biotechnol Prog.* 2023;39(4):e3343. <https://doi.org/10.1002/btpr.3343>.
  39. Zhao B, Abdubek P, Zhang S, Xiao H, Li N. Analysis of host cell proteins in monoclonal antibody therapeutics through size exclusion chromatography. *Pharmaceut Res.* 2022;39(11):3029–3037. <https://doi.org/10.1007/s11095-022-03381-0>.
  40. Guo J, Kufer R, Li D, Wohrlab S, Greenwood-Goodwin M, Yang F. Technical advancement and practical considerations of LC-MS/MS-based methods for host cell protein identification and quantitation to support process development. *mAbs.* 2023;15(1):2113365. <https://doi.org/10.1080/19420862.2023.2213365>.
  41. Yang Y, Wu J, Wang F, Lefers M. Improved identification of host cell proteins in monoclonal antibodies by combining filter-aided sample preparation and native digestion. *J Pharm Sci.* 2025;114(8):103875. <https://doi.org/10.1016/j.xphs.2025.103875>.