

---

# Explainable Deep Learning Framework for SERS Bio-quantification

---

**Jihan K. Zaki**

Melville Laboratory for Polymer Synthesis  
Yusuf Hamied Department of Chemistry  
University of Cambridge  
Lensfield Rd, CB2 1EW

**Jakub Tomasik**

Department of Chemical Engineering and Biotechnology  
University of Cambridge  
Philippa Fawcett Drive, CB3 0AS

**Jade A. McCune**

Melville Laboratory for Polymer Synthesis  
Yusuf Hamied Department of Chemistry  
University of Cambridge  
Lensfield Rd, CB2 1EW

**Sabine Bahn**

Department of Chemical Engineering and Biotechnology  
University of Cambridge  
Philippa Fawcett Drive, CB3 0AS

**Pietro Liò \***

Department of Computer Science and Technology  
University of Cambridge  
15 JJ Thomson Ave, CB3 0FD  
p1219@cam.ac.uk

**Oren A. Scherman \***

Melville Laboratory for Polymer Synthesis  
Yusuf Hamied Department of Chemistry  
University of Cambridge  
Lensfield Rd, CB2 1EW  
oas23@cam.ac.uk

## Abstract

Surface-enhanced Raman spectroscopy (SERS) is rapidly gaining attention as a potential fast and inexpensive method for biomarker quantification, which can be combined with deep learning methodology to elucidate complex biomarker-disease relationships. Current standard practices in computational SERS analysis are substantially behind the state-of-the-art machine learning approaches; however, present

---

\*Corresponding author

challenges of SERS analysis could be effectively addressed with a robust computational framework. Additionally, there is a particular need for improved model explainability for SERS analysis. While current methods are capable of providing global explainability, they are insufficient in assessing the lower level contexts where other factors, such as confounding factors, outliers and measurement errors, or different predictors of the same outcome variable, could affect prediction outcomes. This study presents a novel framework for SERS bio-quantification rooted in a three-step process, including spectral processing, analyte quantification, and model explainability. To this end, a serotonin quantification task in a urine medium was assessed as a model task with 682 SERS spectra measured in a micromolar range using gold nanoparticles and cucurbit[8]uril chemical spacers. A denoising autoencoder was developed for spectral enhancement, convolutional neural networks (CNN), and vision transformers were utilized for biomarker quantification. Lastly, a novel context representative interpretable model explanations (CRIME) method was developed to suit the current needs of SERS mixture analysis explainability. Serotonin quantification was most efficient in denoised spectra analysed using a convolutional neural network with a three-parameter logistic output layer (Validation: mean absolute error (MAE) = 0.24  $\mu\text{M}$ , mean percentage error (MPE) = 15.00%, Test: MAE = 0.15  $\mu\text{M}$ , MPE = 4.67%). Subsequently, the CRIME method revealed the CNN model to present six unique prediction contexts, of which three were associated with serotonin. The proposed framework could unlock a novel, untargeted hypothesis generating method of biomarker discovery considering the rapid and inexpensive nature of SERS measurements, and the potential to identify biomarkers from CRIME contexts, which should be validated in a clinical setting.

## 1 Introduction

Deep learning methods are increasingly being used in biomarker research, as yet-to-be discovered relationships between biomarkers and disease outcomes increase in complexity with expanding numbers of biomarkers investigated[23]. Surface-enhanced Raman spectroscopy (SERS) presents a novel molecular assaying method and allows for the delivery of consistent, accurate, and sizable data that can be utilized with deep learning methods. The technique capitalizes on ‘hot-spots’, localized regions of intense optical fields, created by the aggregation of noble metal nanoparticles[16, 24]. These nanoparticles offer a robust platform for in situ analysis within liquid media, rendering SERS a practical choice for broad applications. There are a set number of challenges in SERS-based analyte detection, which if solved could unlock the significant potential of the method. These primarily relate to the reproducibility, and readability of spectra, particularly in mixtures[32]. Namely, inherent variability in SERS affects the signal intensities of spectra in repeat measurements, and the complexity of biological media measured, alongside potential *intra*- and cross-individual variation in molecular patterns cloud the spectra through biological noise[34]. While experimental developments are seeing improvements in the applicability of SERS, computational frameworks must be developed to supplement and enable the application of SERS in clinical practice.

The SERS domain is far behind current state-of-the-art practices developed in the field of machine learning, and there are several promising methodologies yet to be integrated to SERS analyses. SERS analysis has relied on traditional dimensionality reduction methods to reduce the variations of the spectra and to account for the high noise levels of SERS spectra particularly in biological samples. Principal component analysis and discriminant analysis (PCA-DA) and partial least squares regression (PLSR) have been the *de facto* industry standard within the SERS domain[14, 22, 15, 7, 2, 1, 26], however while these methods can be useful for feature extraction, the increased complexity of biological media can hinder predictions. Over the last 5 years, the application of convolutional neural networks (CNN) [17] for SERS analysis has become more common[27, 21, 12, 28]. Nevertheless, the application of CNNs has been predominantly applied with established model architectures with limited exploration of domain specific modified layers. Moreover, while transformers have revolutionized various domains of machine learning by enabling models to handle sequential data with long-range dependencies effectively, their application in SERS spectral analysis remains underexplored. To date, we were able to source only two studies where vision transformer-models (ViT) [8] were applied in SERS based applications[18, 29]. Similarly to transformer models, the application of autoencoders in

the field of SERS analysis is scarce and is primarily deployed for improved feature extraction[3, 4], despite their strong promise as robust denoising models. Most notably, a majority of SERS studies fail to show adequate model explainability. Without significant exploration of the models selected features, it cannot be determined without reasonable doubt that the predictions are due to the intended signal, or confounders or other sources of sample bias.

This study aims to develop computational methods to mitigate the described primary challenges of SERS, namely the variability between spectra, and the effects of biological noise on measurements. To this end, the present study proposes a complete and up-to-date SERS analysis framework enabling robust bio-quantification, and explainability. The assessment of urinary biomarkers for mental health disorders was selected as the model system to assess the strengths of SERS. Urine is non-invasive, and easy to obtain in large quantities, and shows significant potential as a biomarker source with over 5000 analytes identified to date[31]. In turn, serotonin plays a crucial role in regulating mood, emotion, and sleep, among other physiological functions, and imbalances in serotonin levels have been closely associated with various mental health disorders and has been long hypothesized to be a causal factor in major depressive disorder (MDD), anxiety, and schizophrenia[19, 6, 13, 11, 9]. In brief, this study aims to extend the clinical applicability of SERS in three main directions. First, it seeks to computationally mitigate inherent biological and method-based variations in SERS. Second, it aims to explore deep learning models for more accurate targeted analyte quantification. Lastly, it aims to propose a framework for explaining the decision making process of the developed models, which could identify all contexts in which a model uses different predictors to reach the desired target outcome.

## 2 Methods

### 2.1 Dataset Preparation

The study design is summarized in **Figure 1**. The dataset assessed consisted of 318 SERS spectra measured in a lyophilized urine medium and 364 SERS spectra measured in a water medium. Samples were measured using a 785 nm laser and cucurbit[8]uril (CB[8]) spacers (0.9 nm) with 60 nm gold nanoparticles, as visualised in **Figure 1A**. Both urine- and water-medium samples were spiked with three key neurotransmitters: epinephrine, dopamine, and serotonin, with concentrations ranging from 0 to 9  $\mu\text{M}$ . Serotonin was used as the target analyte. This was due to the lyophilized urine medium containing varying endogenous concentrations of epinephrine and dopamine resulting in their unknown absolute concentrations. The SERS spectra were shortened to feature a relevant range of Raman shifts from 300 to 2000  $\text{cm}^{-1}$ , resulting in a total of 842 datapoints per spectrum. The specific concentrations of the neurotransmitters in each sample are presented in **Table 1**. Prior to assessment, spectra were processed using the asymmetric least squares (ALS) algorithm [10] for baseline correction with pre-specified parameters ( $\lambda=1000$ ,  $p=0.1$ ,  $n=10$ ). Following the correction, the data was normalized to intensities between 0 and 1.

### 2.2 Denoising Autoencoder

Following baseline correction and normalization, the spectra were denoised using a denoising autoencoder. For conventional autoencoders, an encoder neural network is trained to convolute input data into a latent space, and simultaneously a decoder neural network is trained to restructure the original data from the latent transformation. A denoising autoencoder differs by attempting to reconstruct clean outputs from a latent space formed by encoding noisy data [30], which could prove useful in SERS applications with significant biological noise. A simple denoising autoencoder was trained using full water-medium spectra consisting of 364 spectra with 937 datapoints using an 80:20 train-test split, with the urine background samples incorporated as noise. Noisy data were generated using urine background data (Table 1: sample U), where a randomly selected measurement was overlaid to each clean spectra following baseline correction but before normalization. The denoising autoencoder was implemented in TensorFlow. The encoder comprises two dense layers. The first layer has  $2 \times 200$  units and utilizes a ReLU activation function. The second layer further compresses the data into an encoding space of dimension 200, with ReLU activation. Symmetric to the encoder, the decoder also consists of two dense layers. The first layer expands the encoded data back to  $2 \times 200$  dimensions using ReLU activation. Finally, the second layer reconstructs the data to its original dimension (937) using a sigmoid activation function. The model was compiled

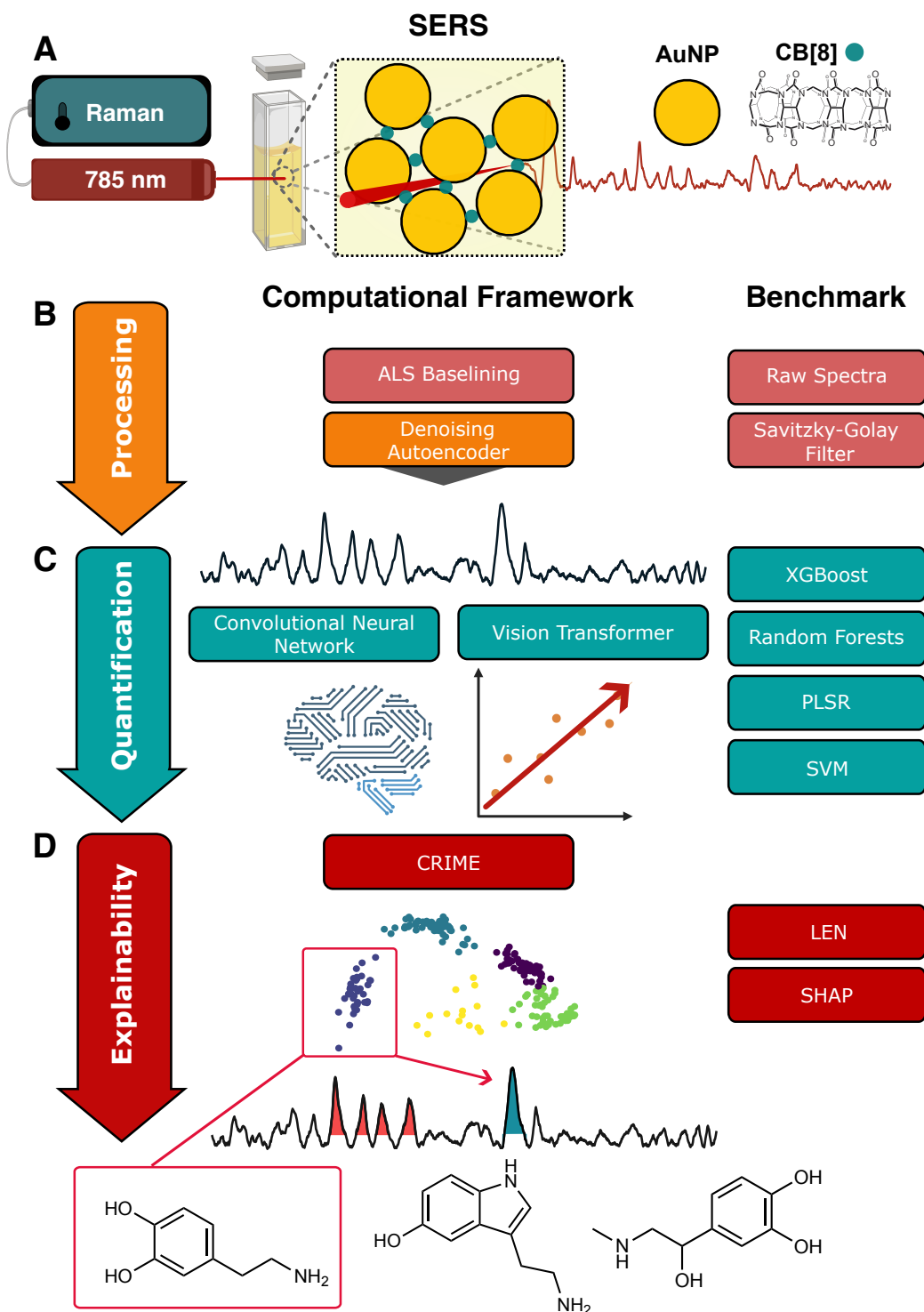


Figure 1: **SERS deep learning framework development pipeline.** Illustrated are the SERS measurement process applied (A), and the computational framework pipeline. Benchmark comparisons of alternative methodology are presented on the right. Preprocessing methods (B) are marked in orange and light red, quantification methods (C) in blue, and explainability methods (D) in dark red. Asymmetric least squares (ALS) baselining is applied to all spectra prior to assessing the framework or the benchmarks. SERS = surface-enhanced raman spectroscopy, AuNP = gold nanoparticle, CB[8] = cucurbit[8]uril, CNN = convolutional neural network, XGBoost = extreme gradient boosting trees, PLSR = partial least squares regression, SVM = support vector machines, CRIME = context representative interpretable model explanations, LEN = logic explained networks, SHAP = Shapley additive explanations.

Table 1: **Added concentrations and number of spectra for all three neurotransmitters in both water and urine backgrounds.** \*Samples D, E, and F were not present in the water dataset, and sample M was not present in the urine dataset. Sample U represents a baseline measurement with no added or measured concentrations of the neurotransmitters. EPI = Epinephrine, DA = Dopamine, 5-HT = Serotonin, n = number of spectra.

Sample	EPI	DA	5-HT	Water	Urine
	( $\mu\text{M}$ )	( $\mu\text{M}$ )	( $\mu\text{M}$ )	(n)	(n)
A	2	0	0	25	22
B	0	2	0	17	21
C	0	0	2	22	22
D	3	0	7	0*	26
E	0	8	3	0*	28
F	7	3	0	0*	21
G	3	2	7	99	26
H	1	1	9	38	22
I	2	8	3	37	22
J	6	9	1	93	23
K	7	3	2	11	23
L	9	6	4	11	21
M	3	3	3	11	0*
U	0	0	0	58	41

using mean squared error (MSE) as the loss function and optimized using the Nesterov-accelerated Adaptive Moment Estimation (Nadam) optimizer. Training was conducted for 128 epochs with a batch size of 32. Both training and testing was performed on noisy data to facilitate the denoising objective. The quality and utility of denoised spectra was subsequently evaluated through effects on performance in quantification models.

### 2.3 Quantification models

The quantification of serotonin was primarily evaluated using state-of-the-art neural network models, as shown in **Figure 1C**. The two model types applied to analyse the spectra were the CNN and the ViT, with custom SERS-specific layers evaluated for the CNN. Both the CNNs and the ViT models were implemented in TensorFlow and designed to adapt to SERS spectral data. A core CNN architecture comprised of ReLu and Tanh-ReLu paired 1-D convolution layers. The core architectures of the quantification models are described in more detail in the **Supplementary Material section B**. Three variants of the CNN were evaluated with varying additional custom layers. These included the base CNN with a linear final activation output layer, a CNN with a modified custom Three-Parameter Logistic (3PL) activation function, and a CNN model with inherent scale-adjusting capabilities through scaling layers.

The scale-adjusting CNN model was developed with two unique scaling layers implemented. These were a multi-scale assessing layer, and a local scaling layer. Both layers were utilized prior to the half-peak ReLu layer in the core CNN architecture. The multi-scale layer captures features from the input  $X$ , with three layers sized 8, 25, and 50, each with 8 filters. Each convolutional operation is defined as  $C_i(X) = W_i * X + b_i$ , where  $*$  denotes a convolutional operation,  $W$  the weight, and  $b$  the bias. To assess the spectra at different scales simultaneously, the output of each convolutional layer is combined following the convolutional operations along the feature dimension. The local peak scaling layer in turn was developed to scale regions of the spectra which were assessed not to be relevant to the outcome variable, identified from the reference spectra of the pure compounds in water. The layer applies a set of scaling factors  $s_j$  unique to the number of pre-registered regions of interest (or non-interest), which are defined by start and end indices  $a_j, b_j$ , in the spectra. The scaling operation for each region is expressed as:  $S_j(X) = X_{a_j:b_j} \odot s_j$ , where  $\odot$  denotes the element-wise multiplication. The output of the scaling operation is concatenated to reconstruct the spectra with scaled regions. The modified output layer assessed in both custom layer CNN models utilizes a Three-Parameter Logistic (3PL) activation function. The ViT architecture in turn consisted of an embedding layer with a patch size of 25 matching the CNN architectures initial layer, with a hidden size of 64 and a dropout rate of 0.1. Subsequently, the architecture consisted of 6 transformer blocks

with 6 multi-head perceptron’s. The transformer blocks each employed Gaussian Error Linear Unit (GELU) activation functions.

Each model was evaluated in both the raw and denoised data incorporating unseen spectra as well as repeat spectra, with spectra defined as repeat if separate measurements of the specific sample were used in training either the denoising autoencoder or the quantification models. Repeat spectra were split into training and validation sets with a 90:10 split, and furthermore measurements taken from an unseen serotonin free sample (sample F) were included in the validation set. The remaining unseen samples (D, and E) were included exclusively in the test set. Final spectra counts for both datasets consisted of 218 training spectra, with a validation set of 46, and a test set of 54 spectra. Hyperparameter tuning and architecture search for both the CNN variants and the ViT was conducted iteratively, guided by the model’s performance on the validation set. Each model variant for both the CNNs and the ViT was trained 100 times, with an ensemble average used for evaluation. Both model types were optimized using the Adaptive Moment Estimation (Adam) algorithm with a learning rate of 0.001, a batch size of 64, and 256 epochs, and compiled with a mean absolute error (MAE) loss function. Additional evaluation metrics included mean squared error (MSE) and mean percentage error (MPE). Early stopping with a patience of 64 was employed to mitigate overfitting, and model checkpoints were saved for epochs that minimized validation loss. Reproducibility was ensured by setting random seeds for TensorFlow, NumPy, and the train-test split. Of the 100 trained models in the ensemble, the model with the lowest MAE in the validation set was selected as the final model, which was assessed in the holdout test set.

## 2.4 Context representative interpretable model explanations

The reliability and explainability of the final quantification model were assessed using the Context Representative Interpretable Model Explanations (CRIME) framework, developed in this study for machine learning interpretations of data with expected contextual prediction clusters. The CRIME framework expands on the widely applied local interpretable model agnostic explanations (LIME) framework[25], by assessing model explanations through contexts. Contexts can be defined within this framework as prominent and consistent explanation outcomes across a number of prediction instances. While contexts can have numerous explanations as to why they are prominent from other explanation contexts, they can be broadly categorized to stem from differing sources of prediction reasoning, such as confounding factors, outliers and measurement errors, or different predictors of the same outcome variable. The framework is summarized in **Algorithm 1**.

---

### Algorithm 1 Context Representative Interpretable Model Explanations

---

```

1: Input: Dataset  $X$ 
2: Initialize: Explained model  $\mathcal{M}$ 
3: Initialize: LIME model  $\mathcal{L}$ 
4: Output: CRIME context explanations  $C$ 
5: Extract local explanations  $\epsilon$  using LIME model  $\mathcal{L}$  on model  $M$  predictions on dataset  $X$ 
6: Initialize CRIME variational autoencoder model  $V$  using explanations  $\epsilon$ 
7: Apply CRIME encoder  $e$  to explanations  $\epsilon$  to assess the initialized latent space  $L$ 
8: Apply K-means clustering algorithm  $\mathcal{K}$  to initialized latent space  $L$  with  $n$  visually plausible contexts
9: Assign context labels  $c$  to latent space  $L$  points based on clustering
10: Initialize CRIME context explanation array  $C$  of size  $n$ 
11: for  $c = 1$  to  $n$  do
12:   Compute mean LIME explanation  $\epsilon_c$  from explanation instances  $\epsilon_c$ 
13:    $C[c] \leftarrow \epsilon_c$ 
14: end for
15: return  $C$ 

```

---

The CRIME framework attempts to identify all prediction contexts of the input data space through the latent space of a variational autoencoder (VAE) trained on the LIME predictions of all instances in the available data. The LIME predictions are flattened with regards to perturbation limits and weights prior to input and are subsequently projected to the two-dimensional latent space. The VAE architecture used in this study consisted of a simple encoder, sampler, and decoder. Notably, the architecture can be fine-tuned depending on the individual requirements for the CRIME framework in future

applications. Details regarding the VAE of the CRIME method are described in the **Supplementary Material section C.1**. Following training of the VAE, the latent space is utilized to identify context clusters representing all the possible ways in which the quantification model interprets the input data. The latent space instances are clustered into the final contexts using K-means clustering, and the latent space is visually inspected for selecting the number of clusters. Finally, a mean LIME explanation is assessed through averaging all instances in each cluster to represent the contexts. To identify the defining features of each context representation, normalized LIME feature weights are combined with mean feature values representing the spectral intensities within the context clusters. They are then set in a three-dimensional space, together with normalized feature positions, which are then further clustered into 15 clusters using K-means clustering. Following the clustering of mean spectral feature values, position z-scores, and LIME weights, the clusters are ordered according to the product of their LIME weights and spectral intensities. The five clusters with the highest score are selected to represent the regions of the spectra which contribute most to the contextual predictions.

Following the identification of the most relevant context prediction regions, the highlighted regions of the mean context spectra are assessed against measured clean spectra of the neurotransmitters known to be present in the mixture. To emphasize the explanation weights in the spectra, both the reference clean spectra, and the mean context spectra are scaled according to the explanation weights in the specific feature location. To determine the cause or identity of the recognized context clusters, the final mean context indicators are compared to the weighted reference spectra using cosine similarity  $S_c$ , which is defined in **Equation 1**.

$$S_{\cos} = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} \quad (1)$$

## 2.5 Benchmarking

Each segment of the framework presented in this study is carefully benchmarked against alternative models or methods, previously established benchmarks, or common practices in the SERS domain.

The utility of the denoising autoencoder was assessed by measuring performance in the raw and denoised data, and additionally, comparing it with fifth order polynomial second Savitzky-Golay derivative processing with a window length of 33[14], which was previously assessed as a reference standard. The developed CNN architecture was benchmarked against simpler architectures without Tanh-ReLu pairing layers. Furthermore, the best-performing quantification model was assessed against methods used previously to quantify neurotransmitter concentrations, as well as other machine learning methods, including partial least squares regression (PLSR), random forests, support vector machines (SVM), and extreme gradient boosting (XGBoost). The primary benchmark for serotonin quantification accuracy is the previous study by Kaser *et al.* 2014. Hyperparameter tuning was determined using grid search and 3-fold cross validation. The searched parameters are included in the **Supplementary Material section D**. Final model robustness was tested using perturbation testing. Gaussian noise was added in two ways across the input spectra, first through applying across the entire spectra to assess overall noise tolerance, and second through applying noise to identified relevant spectral regions to assess how the model can adapt to noise by assessing contextual cues from neighboring regions. In total, four noise levels were assessed at 5% noise, 10% noise, 20% noise, and 30% noise for both tested methods. Robustness of predictions was assessed using a 0.5  $\mu\text{M}$  MAE cutoff.

For comparison with CRIME, feature importance and model explainability was assessed using Logic Explained Networks (LEN)[5], and Shapley Additive Explanations (SHAP)[20]. The LEN was implemented using PyTorch in python. In order to apply the LEN, the spectral input data was sectioned to discrete categories, and concept mapping was done through taking the mean of the min-max scaled feature map activations across each layer of the model. Each concept corresponded to approximately 25 x-axis points across the SERS spectrum corresponding to approximately half a peak. Further details of the LEN implementation are described in **Supplementary Material section C.2**. SHAP calculations were done using the above-mentioned sectioned categories separately using Gradient Explainer.

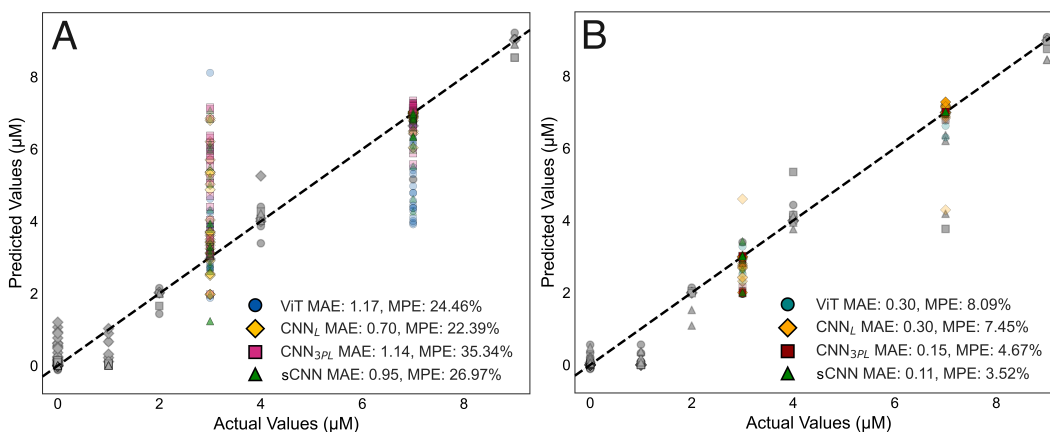
### 3 Results

#### 3.1 Denoising autoencoder

In the present work, the utility of deep learning methods was assessed for identifying target serotonin concentrations from SERS measurements with CB[8] additives as chemical spacers. To this end, a denoising autoencoder and neural network quantification models were developed using 682 spectral measurements taken from water and artificial urine-based samples, with concentrations of serotonin ranging from 0 to 9  $\mu\text{M}$ . Following training, the denoising autoencoder was able to robustly reconstruct the clean data from noisy inputs in the test-set (MSE=0.025). Examples of input noisy data and subsequently denoised spectra are presented in **Supplementary Figure 1A**. Similar trends are observed in denoising experimental urine measurement spectra, which have been presented in **Supplementary Figure 1B**.

#### 3.2 Quantification models

Four different neural network models were evaluated in both raw and denoised datasets. These included the ViT model, the linear output layer CNN model ( $\text{CNN}_L$ ), the three-parameter logistic output layer CNN model ( $\text{CNN}_{3PL}$ ), and the scale-adjusting three-parameter logistic output layer CNN model (sCNN). All ensemble models followed a similar trend in predictions in the validation set and are presented in **Supplementary Figure 2**. Final selected best model validation set predictions for all three model types are presented in **Figure 2** and **Supplementary Table 2**, and for both datasets, all four models showed strong performance in the validation set. The models were then applied on the test set which comprised of unseen concentration combinations by either the denoising autoencoder or the neural network quantification models. The performance of the selected best models is presented in **Figure 2B**. None of the models were able to reach satisfactory differentiation of serotonin from the other neurotransmitters in the raw urine dataset (ViT: MAE = 1.17  $\mu\text{M}$ , MPE = 24.46%,  $\text{CNN}_L$ : MAE = 0.70  $\mu\text{M}$ , MPE = 22.39%, sCNN: MAE = 0.95  $\mu\text{M}$ , MPE = 26.97%,  $\text{CNN}_{3PL}$ : MAE = 1.14  $\mu\text{M}$ , MPE = 35.34%). However, in the denoised dataset, all models were capable of robust quantification of serotonin, with the  $\text{CNN}_{3PL}$  model (MAE = 0.15  $\mu\text{M}$ , MPE = 4.67%) and the sCNN model (MAE = 0.11  $\mu\text{M}$ , MPE = 3.52%) outperforming both the ViT model (MAE = 0.30  $\mu\text{M}$ , MPE = 8.09%) and the  $\text{CNN}_L$  model (MAE = 0.30  $\mu\text{M}$ , MPE = 7.45%).



**Figure 2: Results of the final models in the validation and test sets for the four model types in both raw (A) and denoised datasets (B).** Validation set results are shown in grey, and test set results are shown in color: the linear CNN model is shown in yellow (diamond), the vision transformer model in blue (circle), the scale-adjusting CNN in green (triangle), and the three-parameter logistic output layer CNN model in red (square). The shown values were obtained from the final test set. Validation set results are presented in Table 2 in the appendix. MAE = mean absolute error, MPE = mean percentage error.



### 3.3 CRIME framework

The CRIME framework was fit on a VAE using the LIME explanations of the  $\text{CNN}_{3PL}$  predictions. This setting was selected due to its strong performance across both the validation and test data. Following the K-means clustering, the latent space of the VAE was clustered into six contexts. The latent space is visualized in **Figure 3**. Among them, four distinct contexts were identified (contexts A, B, C, and F), as well as one intermediate context (context E), and one outlier context (context D). The mean LIME explanations for each CRIME context cluster are presented in **Figures 3A to 3F**. Peak regions were selected further from the most prominent clusters in a three-dimensional representation of the spectral intensities, explanation weights, and position z-scores. The peak-region cluster plots are visualised in **Supplementary Figures 3A to 3F**. Contexts A ( $S_{cos}=0.87$ ), E ( $S_{cos}=0.46$ ), and F ( $S_{cos}=0.54$ ) were correctly associated with serotonin, while contexts B and C were associated with dopamine ( $S_{cos}=0.98$ ) and epinephrine ( $S_{cos}=0.97$ ) respectively. Complete cosine similarity values between mean CRIME context spectra and reference neurotransmitter spectra are presented for each CRIME context in **Supplementary Table 3**.

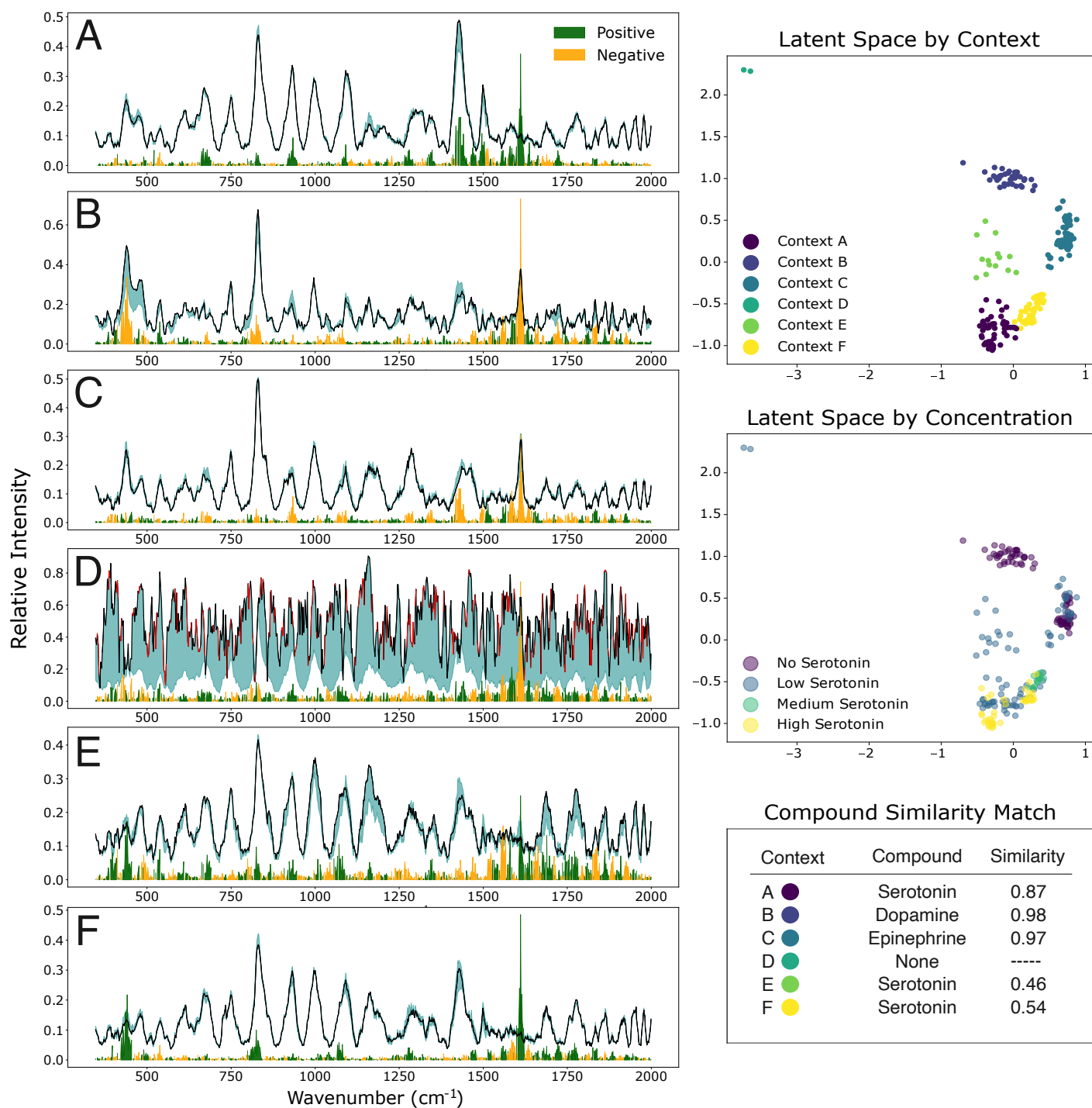
### 3.4 Benchmarking

The neural network models were benchmarked between different architectures and against other machine learning models, with the results summarized in **Table 2** and **Supplementary Figure 4**. The benchmarking results of the CNN architectures are presented in **Supplementary Table 1**. Of all the benchmarked models, the PLSR model trained with autoencoder-denoised spectra showed the best performance amongst non-neural network-based models with a MAE of  $0.70 \mu\text{M}$ , a sixfold error compared to the sCNN. Models trained with denoised spectra near universally showed superior performance when compared to raw spectra trained models, or second derivative Savitzky-Golay denoised spectra trained models with the exception of the random forests model. Lastly, perturbation testing revealed the final neural network quantification model to be robust to noise, with model predictions crossing the set threshold ( $\text{MAE} < 0.5 \mu\text{M}$ ) at 30% localized noise, and 10% universal noise.

Table 2: **Comparison of test-set mean absolute errors across different machine learning models for serotonin quantification from SERS spectra.** Best performing model MAEs within each dataset has been bolded, and the two best models overall have been marked with an asterisk (\*). Additionally, for baseline comparison, the mean absolute error of the previously published PLSR model has been presented. sCNN = convolutional neural network with scaling layers and three parameter logistic (3PL) output layer,  $\text{CNN}_{3PL}$  = convolutional neural network with 3PL output layer,  $\text{CNN}_L$  = convolutional neural network with linear output layer, SVM = support vector machine, RF = random forests, PLSR = partial least squares regression, XGB = extreme gradient boosting, ViT = vision transformer.

Dataset	XGB	PLSR	RF	SVM	$\text{CNN}_{3PL}$	$\text{CNN}_L$	sCNN	ViT
Denoising Autoencoder	0.78	<b>0.70</b>	0.93	0.96	0.15*	0.30	<b>0.11*</b>	0.30
Raw Spectra	0.82	2.05	0.88	1.26	1.14	<b>0.70</b>	0.95	1.17
Savitzky-Golay	1.15	2.25	1.46	1.37	-	-	-	-
Kasera et al. 2014	-	<b>0.52</b>	-	-	-	-	-	-

To compare the context explanations to methods of global explainability, LEN and SHAP frameworks were evaluated as a reference standard. The mean feature activations of the  $\text{CNN}_{3PL}$  model across all layers is presented in **Supplementary Figure 5**. The LEN identified logic statements explaining the four categories of serotonin concentrations with fair (0.69, no serotonin) to excellent (0.98, medium serotonin) explanation accuracy. The logic statements are visualized in **Supplementary Figures 6-9**. Peak regions near wavenumbers of 800, 1000, 1200, and  $1450 \text{ cm}^{-1}$  were consistently selected within the first-order-logic (FOL) statements for all concentration ranges and were deemed to be relevant for serotonin concentration prediction. SHAP values were assessed for all concentration ranges separately and have been visualized on an averaged spectra in **Supplementary Figure 10**.



**Figure 3: Results for Context Representative Interpretable Model Explanations (CRIME) analysis.** Six distinct contexts were identified, which are visualized across mean spectra in subfigures A - F. Positive prediction weights are presented in green, negative prediction weights in yellow, and perturbation limits have been shaded in teal. Red regions in the mean spectra correspond to average perturbation limits at either the top or bottom of the feature weight range for the simplicity of the plot. Latent spaces are visualized by context and concentrations, and compound similarity matching was done using cosine similarity. The highest similarity score is presented alongside the matched compound.

## 4 Discussion

Within the present study, a comprehensive framework of spectral quantification from complex biological media was developed, consisting of data preprocessing and denoising, bio-quantification, and model explanation through the CRIME framework. To this end, data from 318 spectra from lyophilized urine media, as well as 364 spectra from water media were utilized for the development of neural network models for denoising of urine backgrounds and quantification of serotonin. The trained denoising autoencoder improved prediction outcomes near-universally across all model types and enabled robust quantification. The assessed state-of-the-art neural network models substantially outperformed traditional machine learning methods commonly used in the SERS domain, with the  $\text{CNN}_{3PL}$  and the scale-adjusting sCNN models yielding the lowest prediction errors. Notably, the models developed in the present study substantially outperformed existing methodologies[14]. Additionally, the developed CRIME explainability framework identified the spectral contexts in which the model was reliably assessing the relevant serotonin peaks, as well as contexts representing confounding factors or other sample artefacts, which were not readily observable from the outputs of the LEN or the SHAP explanations.

It can be assessed that the custom layers developed in the present study for the sCNN and  $\text{CNN}_{3PL}$  models significantly improved quantification performance. The final output layer of a model acts as a type of calibration curve to the final transformations from the input data. In this sense, a regression task of biomolecular quantification is effectively a calibration task. The use of logistic calibration curves can often yield a better fit on data compared to linear curves, as near the limit of quantification an assay can become saturated, or as the values approach the limit of detection, the linearity of the signal can deteriorate. The scaling layers similarly were able to improve the predictions, most likely due to their added capability to handle variation in scale. This was surprisingly relevant despite present samples all being measured using an identical setup, spectrometer, chemicals, and location which would suggest limited reasons for significant changes in intensity scaling. These layers should be further assessed in different SERS tasks to confirm their utility in subsequent studies.

The CRIME framework within this study was able to effectively explain the CNN model exhibiting a significant improvement in the understanding of the model decisions. The found contexts were associated with the neurotransmitters present in the mixture, and it could be assessed that the two largest contexts were accurately representing serotonin, while two contexts were associated with unwanted signals, and one context was ambiguous in its associations. The dopamine- and epinephrine-associated contexts reveal the imbalances present within the dataset, as the misidentified contexts were primarily in samples with low or absent serotonin concentrations. Following this observation, it can be concluded that within the present dataset, a correlation existed between a lack of serotonin and the presence of other neurotransmitters. Therefore, expanding the dataset to include samples with low concentrations of all neurotransmitters could remedy the apparent inability of the model to generalize to lower serotonin concentrations. When compared directly to the explanations provided by the trained LEN and SHAP explainers, it can be seen that there is a robust association of peaks to serotonin in all three methods. However, SHAP explanations could not effectively communicate the potential presence of confounding factors. Similarly, LEN statements, while effective in identifying potential confounders, were complex enough to make their presentation unintuitive. However, this could be remedied through a similar context-seeking variation of LEN. Within applications where the identification of all prediction reasoning is crucial, the application of the CRIME framework could see benefit over current best practices for global explainability.

It must be highlighted, that the CRIME framework combined with SERS could see clinically relevant use through acting as the first step in biomarker discovery trials. Instead of assessing individual biomarkers of disease through established hypotheses, a biomarker discovery study could be initiated in a non-targeted, hypothesis generating fashion. Applying a machine learning model on raw spectra presently is not advisable due to the lack of confidence in the model assessing true biomarkers as opposed to confounding factors. The exact identification of the signalling biomarkers is challenging when global explainability methods are used for peak detection, as the spectral signals could be a result of multiple overlapping compounds. However, were the CRIME framework applied, individual target biomarkers could be identified through contexts uniquely, and subsequently assigned to the likely biomarkers through a complete library of present compounds, as well as hypothesized biomarkers. With the advent of computational hypothesis generating methodologies such as Mendelian randomization[33], future biomarker discovery trials could see a significant

change in early-stage methodology and approach. Similarly, effects of potential comorbidities or medication on spectral signals could be identified through association of contexts to such groupings, and therefore their effects could be mitigated within a biomarker analysis. For example, in a diagnostic task predicting schizophrenia, there would be a significant risk of misidentifying the effects of comorbidities and antipsychotic medications as disease biomarkers. However, following a CRIME analysis, different contexts could be identified corresponding to these effects or their known biomarkers. Various strategies could then be applied to remove these effects from the spectra. To validate the CRIME framework's explainability and clinical utility, ongoing efforts are being directed towards a biomarker study, which aims to establish the framework's effectiveness in relevant clinical scenarios.

There are several limitations to consider in this study. The development of a denoising autoencoder in patient urine samples as opposed to artificial urine samples could prove to be more challenging. While it could be explained by the CRIME framework, the neural network models were not always able to assess the association of the peaks to the target serotonin compound directly, and instead assessed the presence of dopamine or epinephrine as a predictor of the absence of serotonin. The CRIME framework in turn presents limitations inherent in LIME explanations, as the explanations are dependent on a simpler model fit to complex model predictions and as such do not completely represent an explanation of the actual model. Additionally, it could prove challenging to identify the true reasoning behind CRIME contexts if there were more potential confounding compounds or effects present. Lastly, while in the present task it was feasible to calculate the LIME explanations for all instances, this might not be applicable in larger datasets.

In conclusion, the present study set out to develop a machine learning predictive approach which was capable of achieving three main aims: to perform enhanced single analyte detection of serotonin from a mixture of varying neurotransmitter concentrations in urine, to recognize and adjust for method-based variation in SERS measurements, and lastly to identify all prediction contexts applied by the quantification model. A denoising autoencoder was developed to improve the targeting of relevant neurotransmitter peaks. Additionally, to assess serotonin quantification in raw and denoised spectra, three different state-of-the-art neural network models were developed: a CNN with a three parameter logistic output layer, a scale-adjusting CNN, and a ViT model. In addition, all models were compared to other machine learning methods. Finally, a novel model explainability framework, CRIME, was built around LIME explanations through the assessment of prediction contexts using a combination of VAE and clustering algorithms. The model interpretability was compared between the novel framework and global prediction methods: LEN, and SHAP. Within this study, it was shown that the denoising autoencoder substantially improved predictive capabilities of applied machine learning models, of which the developed three-parameter logistic output layer CNN outperformed other models assessed. Moreover, model explainability was strongly enhanced through the CRIME framework. To our knowledge, this marks the first instance where an autoencoder has been successfully applied to biological 'noise' within the SERS domain. Within the chemical spectral domain, the CRIME framework promises to enable deployment of spectral quantification methods to directly identify disease features in biological fluids, which could be further refined into specific biomarkers through the identification of relevant contexts.

## **Acknowledgements**

The authors would like to thank Dr Setu Kasera for their detailed and elaborate data collection enabling the repurposing of previously measured SERS spectra for the present study.

## **Availability of Data and Materials**

The developed code for the CRIME framework can be found in the following GitHub repository: <https://github.com/jkz22/CRIME>

## **Funding**

This work was supported by the Stanley Medical Research Institute (grant number: O7R-1888) by grants to Sabine Bahn, and by the Oskar Huttunen Foundation grant to Jihan K. Zaki.

## Conflicts of Interest

The authors have no conflicts to declare.

## Author Contribution Statement

Conceptualization: JKZ, PL, OAS;  
Methodology: JKZ, PL;  
Data analysis: JKZ;  
Resources: SB, OAS;  
Writing - Original Draft: JKZ;  
Writing - Review & Editing: All co-authors;  
Supervision: PL, SB, OAS, JT;  
Funding acquisition: OAS, JKZ

## References

- [1] Saba Akbar, Muhammad Irfan Majeed, Haq Nawaz, Nosheen Rashid, Ayesha Tariq, Wajeeha Hameed, Samra Shakeel, Ghulam Dastgir, Rana Zaki Abdul Bari, Maham Iqbal, Amna Nawaz, and Maria Akram. Surface-enhanced raman spectroscopic (sers) characterization of low molecular weight fraction of the serum of breast cancer patients with principal component analysis (pca) and partial least square-discriminant analysis (pls-da). *Analytical Letters*, 55:1588–1604, 2022.
- [2] Dawei Cao, Hechuan Lin, Ziyang Liu, Jiayi Qiu, Shengjie Ge, Weiwei Hua, Xiaowei Cao, Yayun Qian, Huiying Xu, and Xinzhong Zhu. Pca-tlnn-based sers analysis platform for label-free detection and identification of cisplatin-treated gastric cancer. *Sensors and Actuators B: Chemical*, 375:132903, 1 2023.
- [3] Fatma Uysal Ciloglu, Abdullah Caliskan, Ayse Mine Saridag, Ibrahim Halil Kilic, Mahmut Tokmakci, Mehmet Kahraman, and Omer Aydin. Drug-resistant staphylococcus aureus bacteria detection by combining surface-enhanced raman spectroscopy (sers) and deep learning techniques. *Scientific Reports*, 11:1–12, 9 2021.
- [4] Fatma Uysal Ciloglu, Mehmet Hora, Aycan Gundogdu, Mehmet Kahraman, Mahmut Tokmakci, and Omer Aydin. Sers-based sensor with a machine learning based effective feature extraction technique for fast detection of colistin-resistant klebsiella pneumoniae. *Analytica Chimica Acta*, 1221:340094, 8 2022.
- [5] Gabriele Ciravegna, Pietro Barbiero, Francesco Giannini, Marco Gori, Pietro Lió, Marco Maggini, and Stefano Melacci. Logic explained networks. *Artificial Intelligence*, 314, 8 2021.
- [6] A. Coppen. The biochemistry of affective disorders. *The British journal of psychiatry : the journal of mental science*, 113:1237–1264, 1967.
- [7] M. Czaplicka, A. A. Kowalska, A. B. Nowicka, D. Kurzydłowski, Z. Gronkiewicz, A. Machulak, W. Kukwa, and A. Kamińska. Raman spectroscopy and surface-enhanced raman spectroscopy (sers) spectra of salivary glands carcinoma, tumor and healthy tissues and their homogenates analyzed by chemometry: Towards development of the novel tool for clinical diagnosis. *Analytica Chimica Acta*, 1177:338784, 9 2021.
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR 2021 - 9th International Conference on Learning Representations*, 10 2020.
- [9] Arnold E. Eggers. A serotonin hypothesis of schizophrenia. *Medical hypotheses*, 80:791–794, 6 2013.
- [10] Paul H. C. Eilers. Parametric time warping. *Analytical Chemistry*, 76(2):404–411, 2004.

- [11] Joshua A. Gordon and Rene Hen. The serotonergic system and anxiety. *Neuromolecular medicine*, 5:27–40, 2004.
- [12] T. Y. Huang and Jorn Chi Chung Yu. Development of crime scene intelligence using a hand-held raman spectrometer and transfer learning. *Analytical Chemistry*, 28:18, 2021.
- [13] Sameer Jauhar, Philip J. Cowen, and Michael Browning. Fifty years on: Serotonin and depression. *Journal of Psychopharmacology (Oxford, England)*, 37:237, 3 2023.
- [14] Setu Kasera, Lars O. Herrmann, Jesús Del Barrio, Jeremy J. Baumberg, and Oren A. Scherman. Quantitative multiplexing with nano-self-assemblies in sers. *Scientific Reports*, 4:1–6, 10 2014.
- [15] Soogeun Kim, Tae Gi Kim, Soo Hyun Lee, Wansun Kim, Ayoung Bang, Sang Woong Moon, Jeongyoon Song, Jae Ho Shin, Jae Su Yu, and Samjin Choi. Label-free surface-enhanced raman spectroscopy biosensor for on-site breast cancer detection using human tears. *ACS Applied Materials and Interfaces*, 12:7897–7904, 2 2020.
- [16] Judith Langer, Dorleta Jimenez de Aberasturi, Javier Aizpurua, Ramon A. Alvarez-Puebla, Baptiste Auguie, Jeremy J. Baumberg, Guillermo C. Bazan, Steven E.J. Bell, Anja Boisen, Alexandre G. Brolo, Jaebum Choo, Dana Cialla-May, Volker Deckert, Laura Fabris, Karen Faulds, F. Javier García de Abajo, Royston Goodacre, Duncan Graham, Amanda J. Haes, Christy L. Haynes, Christian Huck, Tamitake Itoh, Mikael Käll, Janina Kneipp, Nicholas A. Kotov, Hua Kuang, Eric C. Le Ru, Hiang Kwee Lee, Jian Feng Li, Xing Yi Ling, Stefan A. Maier, Thomas Mayerhöfer, Martin Moskovits, Kei Murakoshi, Jwa Min Nam, Shuming Nie, Yukihiko Ozaki, Isabel Pastoriza-Santos, Jorge Perez-Juste, Juergen Popp, Annemarie Pucci, Stephanie Reich, Bin Ren, George C. Schatz, Timur Shegai, Sebastian Schlücker, Li Lin Tay, K. George Thomas, Zhong Qun Tian, Richard P. van Duyne, Tuan Vo-Dinh, Yue Wang, Katherine A. Willets, Chuanlai Xu, Hongxing Xu, Yikai Xu, Yuko S. Yamamoto, Bing Zhao, and Luis M. Liz-Marzán. Present and future of surface-enhanced raman scattering. *ACS Nano*, 14:28–117, 1 2020.
- [17] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86:2278–2323, 1998.
- [18] Bo Li, Giulia Zappalá, Elodie Dumont, Anja Boisen, Tomas Rindzevicius, Mikkel N. Schmidt, and Tommy S. Alstrøm. Nitroaromatic explosives’ detection and quantification using an attention-based transformer on surface-enhanced raman spectroscopy maps. *Analyst*, 148:4787–4798, 9 2023.
- [19] Shih Hsien Lin, Lan Ting Lee, and Yen Kuang Yang. Serotonin and mental disorders: A concise review on molecular neuroimaging evidence. *Clinical Psychopharmacology and Neuroscience*, 12:196, 12 2014.
- [20] Scott M. Lundberg and Su In Lee. A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 2017-December:4766–4775, 5 2017.
- [21] Félix Lussier, Dimitris Missirlis, Joachim P. Spatz, and Jean François Masson. Machine-learning-driven surface-enhanced raman scattering optophysiology reveals multiplexed metabolite gradients near cells. *ACS Nano*, 2019.
- [22] Félix Lussier, Vincent Thibault, Benjamin Charron, Gregory Q. Wallace, and Jean Francois Masson. Deep learning and artificial intelligence methods for raman and surface-enhanced raman scattering. *TrAC Trends in Analytical Chemistry*, 124:115796, 3 2020.
- [23] Vivek Bhakta Mathema, Partho Sen, Santosh Lamichhane, Matej Orešič, and Sakda Khoomrung. Deep learning facilitates multi-data type analysis and predictive biomarker discovery in cancer precision medicine. *Computational and Structural Biotechnology Journal*, 21:1372–1382, 1 2023.
- [24] Pilot, Signorini, Durante, Orian, Bhamidipati, and Fabris. A review on surface-enhanced raman scattering. *Biosensors*, 9:57, 4 2019.

- [25] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the predictions of any classifier. *NAACL-HLT 2016 - 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Proceedings of the Demonstrations Session*, pages 97–101, 2 2016.
- [26] Fatima Tahir, Ali Kamran, Muhammad Irfan Majeed, Abeer Ahmed Alghamdi, Muhammad Rizwan Javed, Haq Nawaz, Muhammad Adnan Iqbal, Muhammad Tahir, Anam Tariq, Nosheen Rashid, Urwa Shahid, Ahmad Hassan, and Umar Sohail Shoukat. Surface-enhanced raman scattering (sers) in combination with pca and pls-da for the evaluation of antibacterial activity of 1-isopentyl-3-pentyl-1h-imidazole-3-ium bromide against bacillus subtilis. *ACS Omega*, 9:6861–6872, 2 2024.
- [27] Jia Wei Tang, Qing Hua Liu, Xiao Cong Yin, Ya Cheng Pan, Peng Bo Wen, Xin Liu, Xing Xing Kang, Bing Gu, Zuo Bin Zhu, and Liang Wang. Comparative analysis of machine learning algorithms on surface enhanced raman spectra of clinical staphylococcus species. *Frontiers in Microbiology*, 12:696921, 8 2021.
- [28] William John Thrift and Regina Ragan. Quantification of analyte concentration in the single molecule regime using convolutional neural networks. *Analytical Chemistry*, 91:13337–13342, 11 2019.
- [29] Yi Ming Tseng, Ko Lun Chen, Po Hsuan Chao, Yin Yi Han, and Nien Tsu Huang. Deep learning-assisted surface-enhanced raman scattering for rapid bacterial identification. *ACS Applied Materials and Interfaces*, 15:26398–26406, 6 2023.
- [30] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. *Proceedings of the 25th International Conference on Machine Learning*, 2008.
- [31] David S. Wishart, An Chi Guo, Eponine Oler, Fei Wang, Afia Anjum, Harrison Peters, Raynard Dizon, Zinat Sayeeda, Siyang Tian, Brian L. Lee, Mark Berjanskii, Robert Mah, Mai Yamamoto, Juan Jovel, Claudia Torres-Calzada, Mickel Hiebert-Giesbrecht, Vicki W. Lui, Dorna Varshavi, Dorsa Varshavi, Dana Allen, David Arndt, Nitya Khetarpal, Aadhavya Sivakumaran, Karxena Harford, Selena Sanford, Kristen Yee, Xuan Cao, Zachary Budinski, Jaanus Liigand, Lun Zhang, Jiamin Zheng, Rupasri Mandal, Naama Karu, Maija Dambrova, Helgi B. Schiöth, Russell Greiner, and Vasuk Gautam. Hmdb 5.0: the human metabolome database for 2022. *Nucleic acids research*, 50:D622–D631, 1 2022.
- [32] Min Xiong and Jian Ye. Reproducibility in surface-enhanced raman spectroscopy. *Journal of Shanghai Jiaotong University (Science)*, 19:681–690, 12 2014.
- [33] Jihan K. Zaki, Jakub Tomasik, Jade McCune, Oren A. Scherman, and Sabine Bahn. Discovery of urinary metabolite biomarkers of psychiatric disorders using two-sample mendelian randomization. *medRxiv*, page 2023.09.26.23296078, 9 2023.
- [34] Cheng Zong, Mengxi Xu, Li Jia Xu, Ting Wei, Xin Ma, Xiao Shan Zheng, Ren Hu, and Bin Ren. Surface-enhanced raman spectroscopy for bioanalysis: Reliability and challenges. *Chemical Reviews*, 118:4946–4980, 5 2018.

## Supplementary Material

### A Experimental Methods

All initial reagents were sourced from Alfa Aesar and Merck and were utilized in their received state unless otherwise specified. Cucurbit[8]uril was prepared following established literature protocols. Millipore water with a resistivity of 18 M  $\Omega$ -cm was employed in all experiments, unless otherwise indicated. Fresh standard stock solutions of neurotransmitters, specifically dopamine, epinephrine, and serotonin, were prepared at varying concentrations to simulate potential interfering background analytes. Gold nanoparticles (AuNP) with a diameter of 60 nm stabilized by citrate were procured from British Biocell International (BBI). Lyophilized urine samples designated as Calibrator Lot No. 150 and Control Level II Lot No. 230 were obtained from RECIPE ClinChek-Control and were reconstituted in dilute hydrochloric acid as per the supplier’s guidelines.

Spectra for both Raman and SERS were collected with a 785 nm laser operating at 17.5 mW, using an Ocean Optics QE65000 Spectrometer. Each spectrum was acquired over a 10 s interval. AuNPs with a 60 nm diameter were first centrifuged at 12,000 rpm for 45 s, repeated twice, and 900  $\mu$ L of the supernatant were removed. Subsequently, a sample preparation sequence was followed: neurotransmitters (dopamine, epinephrine, and serotonin) were first added, followed by 50  $\mu$ L of the centrifuged AuNPs, then 20  $\mu$ L of CB[8] at a final concentration of 20  $\mu$ M, and finally, 50  $\mu$ L of thawed urine, which had been initially stored on ice. An identical procedure was replicated, replacing urine with water for control experiments.

### B Neural network architectures

Both the CNNs and the ViT models were implemented in TensorFlow and designed to adapt to SERS spectral data. The CNN architecture comprised sequential layers optimized for 1D convolution operations, and the core CNN architecture was used in all trained CNN models, where the initial layer is a convolutional layer featuring 8 filters and a large kernel size approximately the width of half a peak (25 datapoints), aimed to capture broader features in the spectrum. This initial layer employs a Rectified Linear Unit (ReLU) activation function and reduces sequence length through striding. Intermediate layers employ paired hyperbolic tangent (Tanh) and ReLU activation functions, designed to capture complex patterns while maintaining non-linearity. The combination of consequential Tanh and ReLU layers is to direct the model to assess the upper half of the identified general peaks from the sweeper layer. These layers maintain the same padding to avoid changes in sequence length. The model contains two Tanh-ReLU paired layers with  $2 \times 16$  and  $2 \times 32$  filters respectively, with a filter size of 9. The final convolutional stage employs 64 filters with a smaller kernel size of 2 using ReLU activation, aimed to capture fine-grained details in the data. Subsequently, the data is flattened and passed through two fully connected layers employing the same Tanh to ReLU structure with 32 and 16 nodes respectively, to serve the regression task. The core architecture of the CNN models was benchmarked against similar architectures with Tanh-ReLU pairs replaced with ReLU pairs, inversed ReLU-Tanh pairs, or single ReLU layers. Each of the benchmark architectures were trained and validated as described in the methods section, and the results are summarized in **Supplementary Table 1**.

Table 1: **Comparison of core CNN architectures.** All architectures were trained on denoised urine medium datasets, and validated using the holdout test set. ReLU = Rectified linear unit, Tanh = Hyperbolic tangent, MAE = Mean absolute error, MPE = Mean percentage error.

Error	Tanh-ReLU	ReLU-Tanh	2xReLU	ReLU
MAE	0.30	0.56	0.88	0.78
MPE	7.45	17.68	20.71	16.61



Table 2: **Validation set results for neural network models.** sCNN = convolutional neural network with scaling layers and three parameter logistic (3PL) output layer,  $\text{CNN}_{3PL}$  = convolutional neural network with 3PL output layer,  $\text{CNN}_L$  = convolutional neural network with linear output layer, ViT = vision transformer.

Dataset	$\text{CNN}_{3PL}$	$\text{CNN}_L$	sCNN	ViT
Denosing Autoencoder	0.24	0.13	0.33	0.20
Raw Spectra	0.19	0.31	0.14	0.25

## C AI explainability

### C.1 CRIME variational autoencoder

Within the encoder of the CRIME VAE, the input data  $X$ , is transformed into the mean ( $\mu$ ) and logarithm of the variance ( $\log(\sigma^2)$ ) in a proposed Gaussian distribution in the latent space through a fully connected ReLu layer with 256 nodes. During training, the outputs of the encoder are passed to a sampling layer which generates a random noise variable  $\epsilon$  generated from a standard gaussian distribution, which is then transformed using the encoder outputs to draw samples  $z$  as such:  $z = \mu + \sigma * \epsilon$ . The decoder then applies a mirrored dense layer network to the encoder with a ReLu layer with 256 nodes, and a final output sigmoid layer with  $3 \times 842$  nodes. The model was trained for 128 epochs and a batch size of 32, with the Adam optimizer with a learning rate of 0.001 and using the sum of the mean squared error, and Kullback–Leibler divergence as the loss function.

Table 3: **Cosine similarity values across explanation weighted reference spectra and explanation weighted mean context spectra.** Highest similarity values within a context cluster are bolded. X = context.

Reference compound	A	B	C	D	E	F
Serotonin	<b>0.87</b>	0.85	0.60	-0.79	<b>0.46</b>	<b>0.54</b>
Dopamine	0.05	<b>0.98</b>	0.91	-0.80	0.29	0.06
Epinephrine	0.0	0.81	<b>0.97</b>	-0.79	0.12	0.08

### C.2 Logic explained network architecture

The LEN architecture was modified to be similar to the original model architecture, consisting of an entropy layer with 164 input nodes, a leaky ReLu layer with 32 nodes, a Tanh layer with 16 nodes, a ReLu layer with 4 nodes, and a final linear output layer. The LEN was trained using weight decay Adam as the optimizer, using binary cross entropy with logits loss as the loss function with a scaled auxiliary entropy loss at a 0.000001 multiplier. The model was trained with 5001 epochs using a learning rate of 0.0001.

## D Benchmark hyperparameter search

The grids used for the hyperparameter search of the benchmark machine learning models are presented below with the final hyperparameters for the denoised models in bold.

### XGBoost

Hyperparameter	Values
colsample_bytree	<b>0.5</b> , 0.7, 0.8
learning_rate	<b>0.01</b> , 0.1, 0.2, 0.3
max_depth	3, <b>6</b> , 9, 12
alpha	1, <b>3</b> , 5
n_estimators	100, 300, 600, 900, <b>1200</b>

## Random Forests

Hyperparameter	Values
max_depth	1, 2, 3, 6, 7, <b>8</b> , 10
n_estimators	100, 300, 600, 900, <b>1200</b>

## PLSR

Hyperparameter	Values
n_components	5, 8, <b>12</b>

## SVM

Hyperparameter	Values
C	0.1, 1, 10, <b>50</b> , 100
epsilon	<b>0.01</b> , 0.1, 1
gamma	<b>scale</b> , auto

## Supplementary Figures

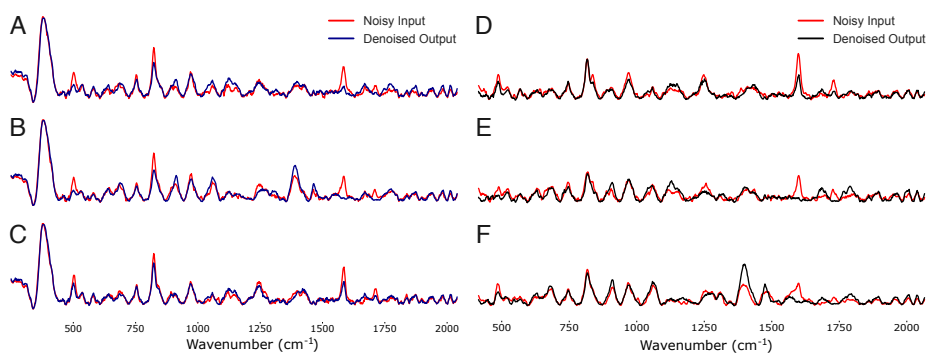


Figure 1: **Examples of denoised spectra in the artificial training data (A-C), and in lyophilized urine spectra (D-F).** Y-axis (relative intensity) is omitted for clarity.

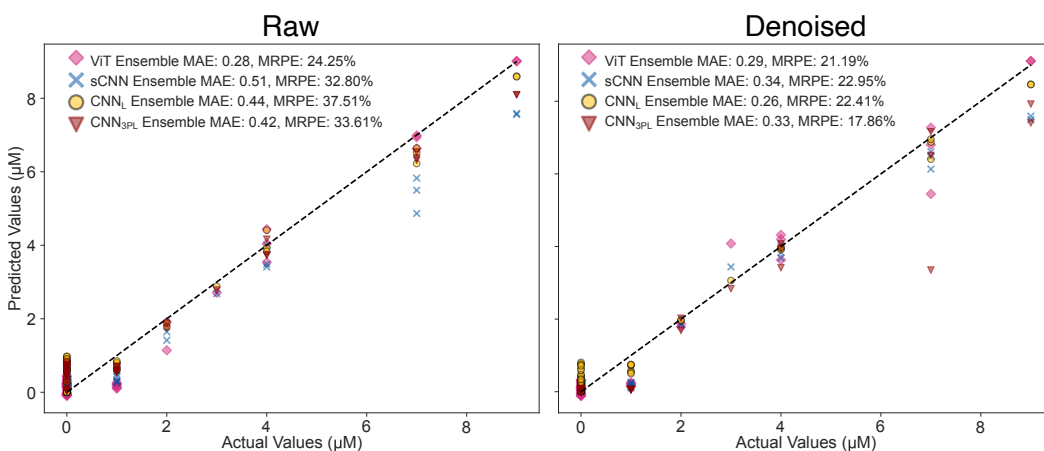


Figure 2: **Predictions of the trained ensembles on the validation set for all types of neural networks on both raw spectra and denoised spectra.**

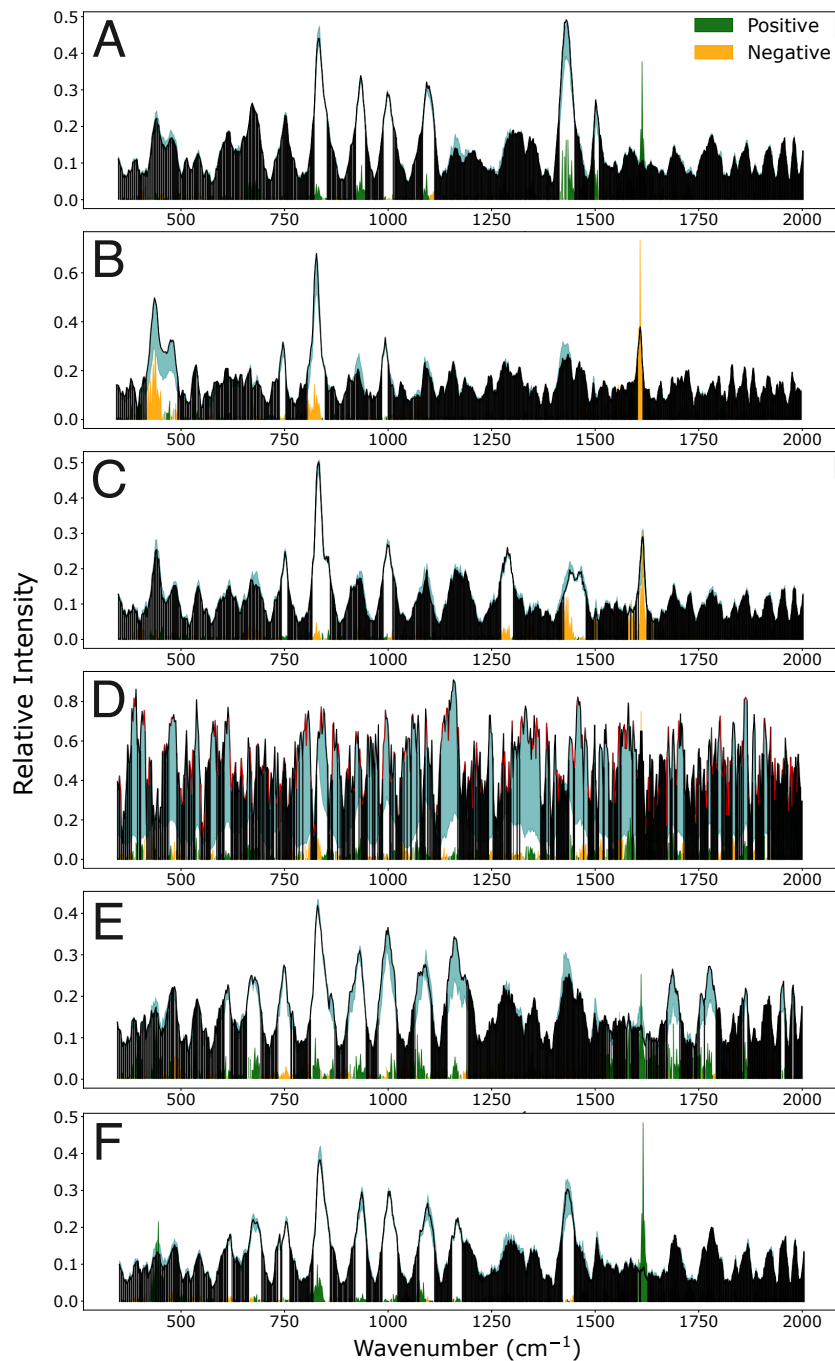


Figure 3: **Peak-region clusters of high relevance extracted from CRIME contexts for compound identification.** Context labels correspond to labels in Figure 3. Positive prediction weights are presented in green, negative prediction weights in yellow, and perturbation limits have been shaded in teal. Red regions in the mean spectra correspond to average perturbation limits at either the top or bottom of the feature weight range for the simplicity of the plot. Areas not relevant are marked in black. High-relevance clusters were obtained with K-means clustering of the product of peak height and LIME weights, with the top 5 largest clusters selected.

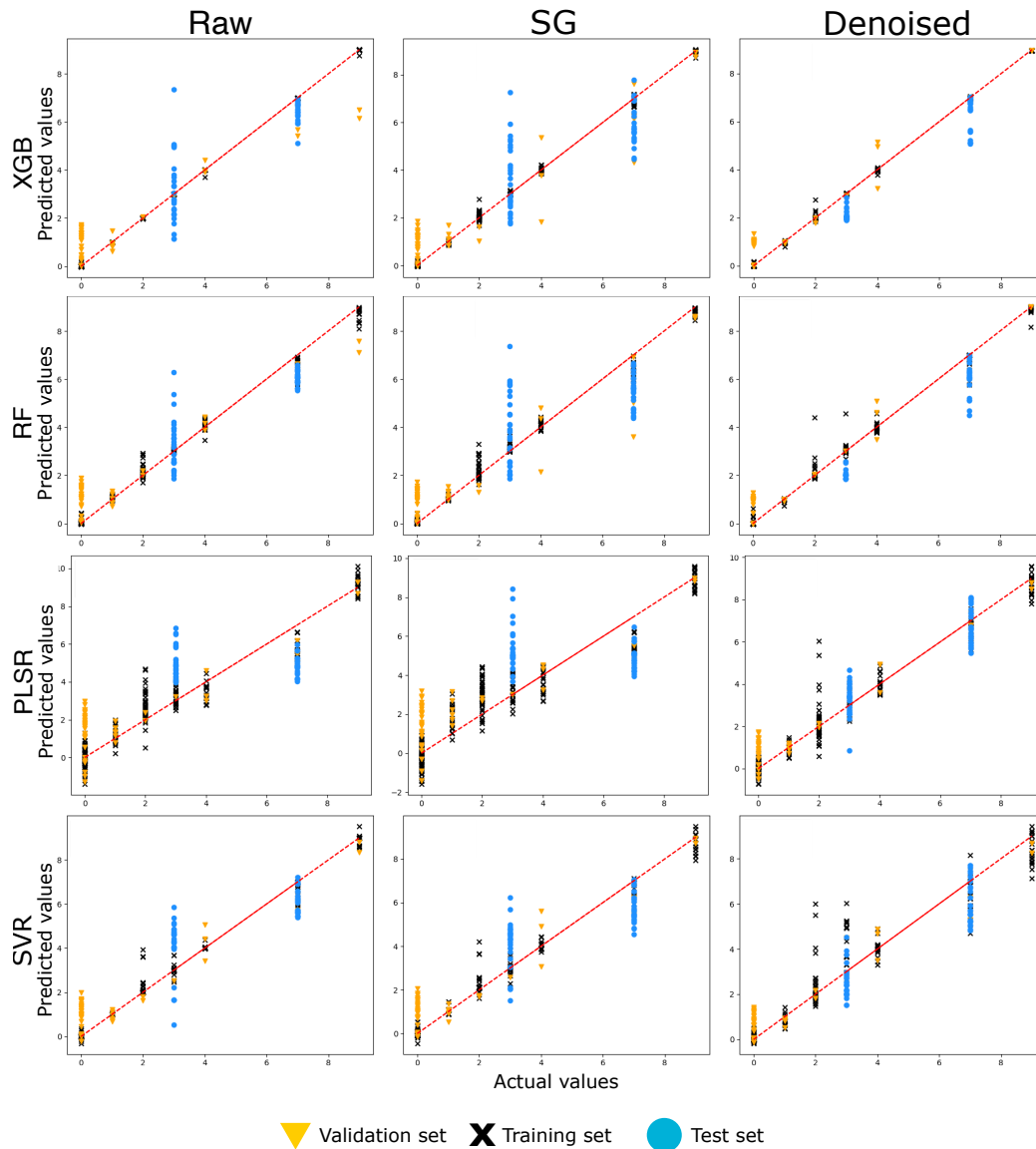


Figure 4: **Results for quantification model benchmarking.** Training set predictions are marked in black (cross), validation set predictions in orange (triangle), and test set predictions in blue (circle). XGB = extreme gradient boosting, RF = random forests, PLSR = partial least squares regression, SVM = support vector machine regression, SG = Savitzky-Golay filter.

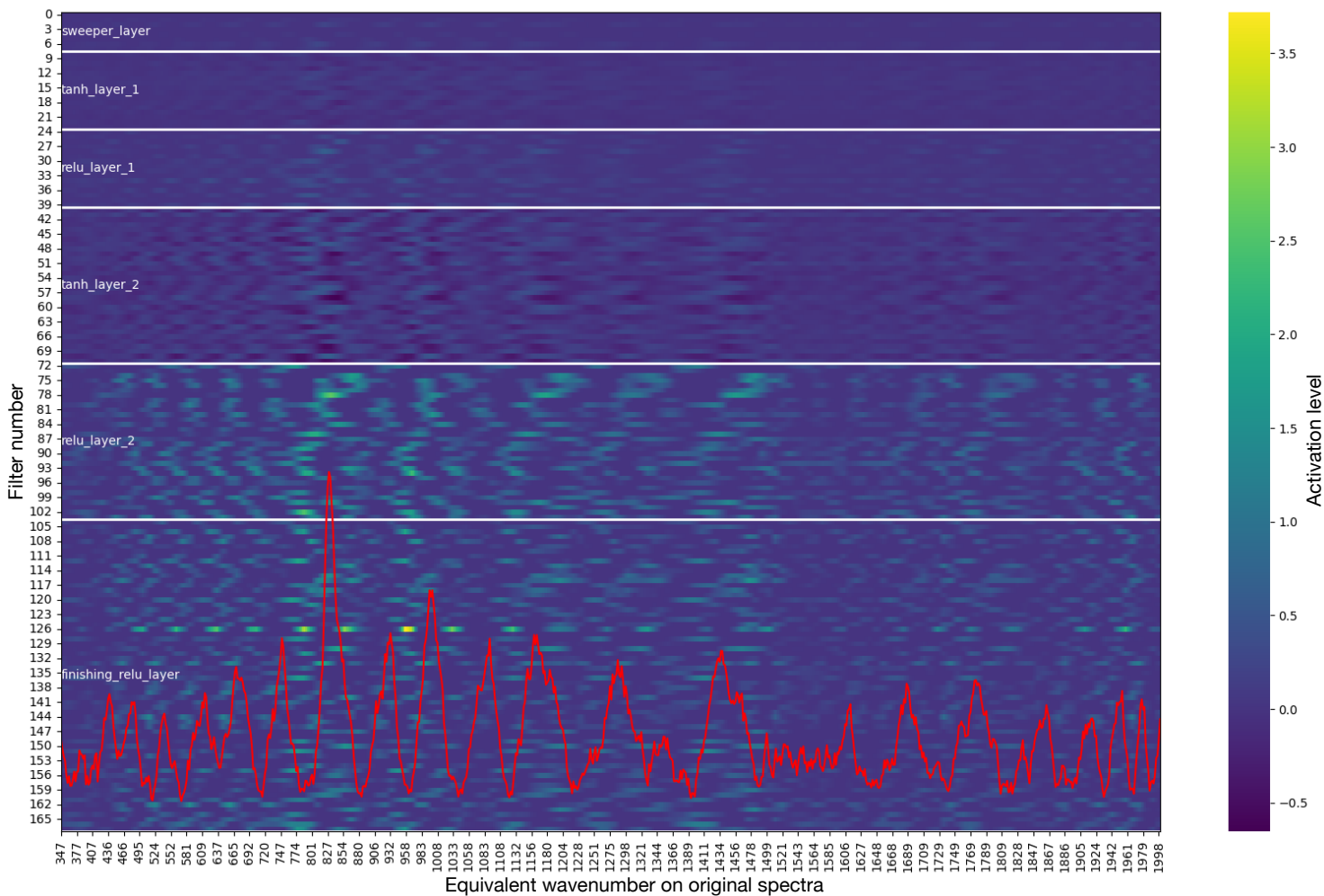


Figure 5: Feature activation map for each convolutional layer in the CNN model overlaid with an example spectra. SERS spectra is shown in red, and higher activations are marked with a yellow hue, with lower activations marked in blue.

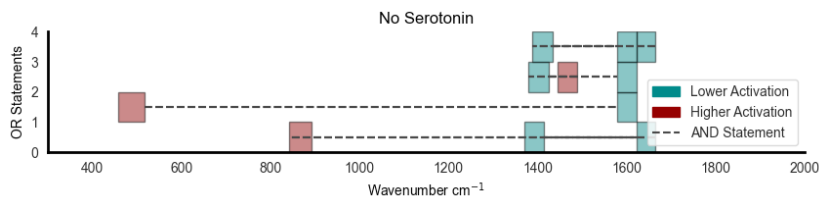


Figure 6: LEN results visualized for samples with no serotonin concentration. OR statements are separated vertically, and AND statements are presented level with a dashed line connecting the statements. Blue squares denote lower activation and red squares higher activation.

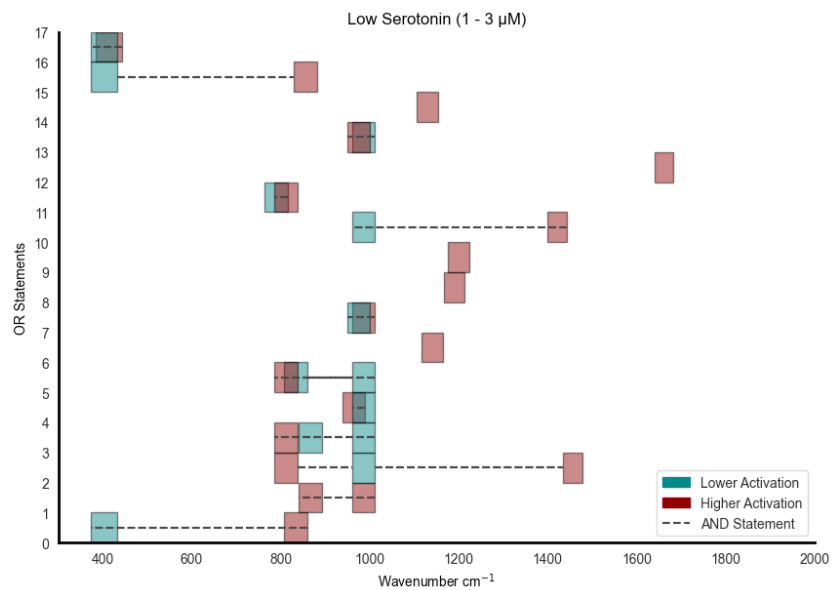


Figure 7: **LEN results visualized for low serotonin concentrations.** OR statements are separated vertically, and AND statements are presented level with a dashed line connecting the statements. Blue squares denote lower activation and red squares higher activation.

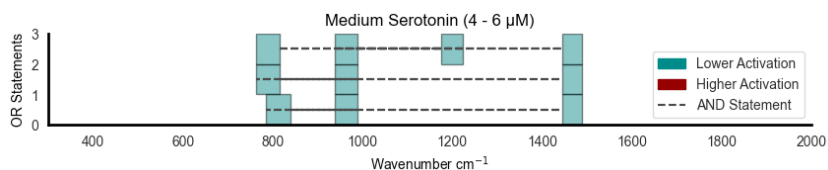


Figure 8: **LEN results visualized for medium serotonin concentrations.** OR statements are separated vertically, and AND statements are presented level with a dashed line connecting the statements. Blue squares denote lower activation and red squares higher activation.

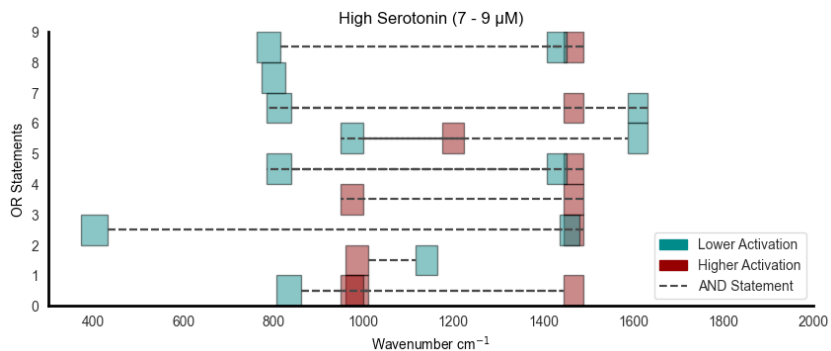


Figure 9: **LEN results visualized for high serotonin concentrations.** OR statements are separated vertically, and AND statements are presented level with a dashed line connecting the statements. Blue squares denote lower activation and red squares higher activation.

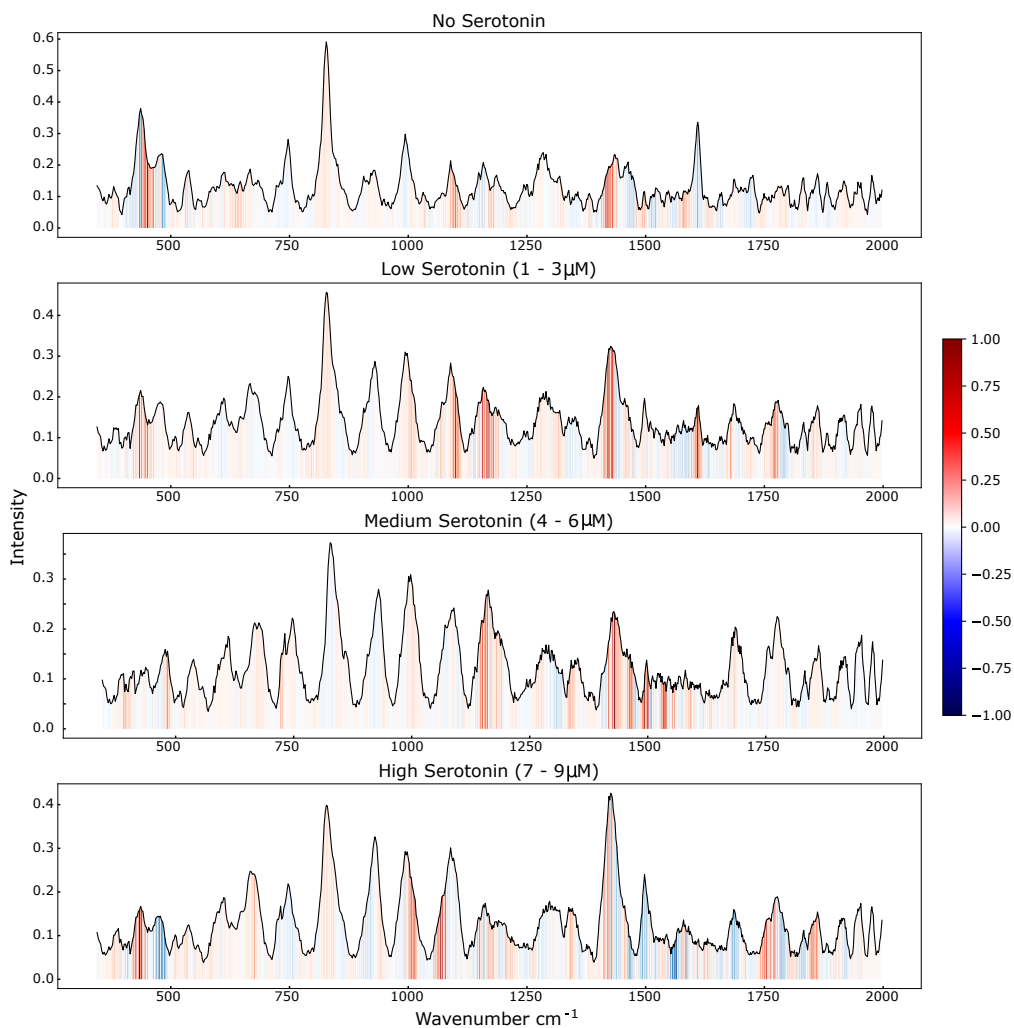


Figure 10: **Shapley additive explanations (SHAP) visualized for all serotonin concentration ranges.** Spectra shown are mean spectra across the respective concentration ranges. SHAP values were obtained using Gradient Explainer, and red areas correspond to positive SHAP values and blue areas to negative SHAP values.