

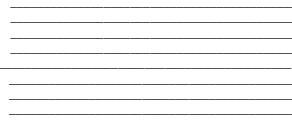
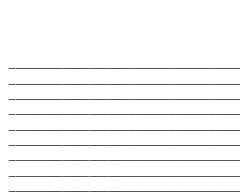
Encoding Strategies for Voltammetry Data and their Machine
Learning Applications in Self-Driving Laboratories

By

Rainey Fu

Supervisor: Alán Aspuru-Guzik
Advisor: Sergio Pablo-García Carrillo
April 2024

B.A.Sc. Thesis



Division of Engineering Science
UNIVERSITY OF TORONTO

Abstract

Self-driving laboratories (SDLs) represent a cutting-edge concept in scientific research and experimentation. SDLs utilize automated instruments, recommendation algorithms, and an orchestration device to conduct experiments and analyze data without human intervention. Among the array of experiments conducted by SDLs, cyclic voltammetry (CV) and differential pulse voltammetry (DPV) are prominent, offering insights into electrochemical processes. However, efficiently extracting crucial information, such as overall shape and peaks, from CV and DPV data remains challenging. This thesis presents a novel encoding technique tailored for CV and DPV data to enhance SDLs' understanding of chemical environments. With this encoding method, SDLs can discern intricate patterns and relationships within the data more effectively. Experiments consisting of various machine learning tasks, such as clustering, classification, denoising, and synthetic data generation, that an SDL may encounter showed excellent results. Beyond SDLs, the utility of this encoding technique extends to any 2-dimensional data. Its versatility opens avenues for broader scientific and industrial applications, empowering researchers and practitioners to glean valuable insights from complex datasets. As SDLs continue to evolve, incorporating innovative methodologies such as this encoding technique promises to accelerate scientific discovery and advance technological frontiers.

Acknowledgment

I am deeply thankful to my advisor, Sergio Pablo-García Carrillo, for his invaluable guidance, without which this project would not have been achievable. I also extend my gratitude to Prof. Alán Aspuru-Guzik for granting me the opportunity to collaborate with the Matter Lab research group at the University of Toronto.

Additionally, I express my heartfelt appreciation to Lianjiang Fu, Yuchuan Hu, Cindy Fu, and Christina Men for their unwavering support throughout this journey. Their encouragement has played a significant role in shaping who I am today.

Table of Contents

Acknowledgment	ii
List of Tables	v
List of Figures	vii
Chapter 1 Introduction	1
Chapter 2 Background and Motivation	4
2.1 Electrochemistry	4
2.2 Cyclic Voltammetry	9
2.3 Differential Pulse Voltammetry	11
Chapter 3 Clustering	13
3.1 Introduction	13
3.2 Data Collection	14
3.3 Curse of Dimensionality	15
3.4 Ramer–Douglas–Peucker Algorithm	16
3.5 Data Preparation and Encoding	17
3.6 K-Means	19
3.7 Density-Based Spatial Clustering of Applications with Noise	21
3.8 t-Distributed Stochastic Neighbour Embedding	22
3.9 Uniform Manifold Approximation and Projection	23
3.10 Results and Discussion	25
Chapter 4 Classification	38
4.1 Introduction	38
4.2 Variational Autoencoders	39
4.3 Conditional Variational Autoencoders	40
4.4 Classifier Model Architecture	40
4.5 Results and Discussion	42

Table of Contents

Chapter 5 Denoising	56
5.1 Introduction	56
5.2 Autoencoder	56
5.3 Results and Discussion	57
Chapter 6 Conclusion	60
6.1 Data and Code Availability	61
Bibliographic references	62
Appendix A	66
A.1 CV K-Means Cluster Results	66
A.2 Metals and Ligands	71

List of Tables

Table 4.1	Classification Accuracy	42
Table 4.2	Classification Accuracy with Synthetic Data	42
Table 4.3	CV Metals Classification Report	44
Table 4.4	CV Ligands Classification Report	45
Table 4.5	DPV Ligands Classification Report	46
Table 4.6	DPV Metals Classification Report	47
Table A.1	Table of Metals	71
Table A.2	Table of Ligands	72

List of Figures

Figure 2.1	Schematic of Electrochemical Cell [7]	6
Figure 2.2	Potentiostat Circuit Diagram	7
Figure 2.3	Cyclic Voltammogram	9
Figure 2.4	Differential Pulse Voltammogram	11
Figure 3.1	RDP Algorithm	16
Figure 3.2	Raw Data and Processed Data	18
Figure 3.3	K-Means Clustering Visualization with CV data	19
Figure 3.4	DBSCAN Clustering Visualization with CV Data	21
Figure 3.5	K-Means Elbow Method	25
Figure 3.6	Silhouette Method for Promising K Values	27
Figure 3.7	DBSCAN Clusters	31
Figure 3.8	UMAP Visualization of Normal CV Data and Incorrectly Shaped Data	32
Figure 3.9	Cyclic Voltammetry t-SNE Projection	33
Figure 3.10	Differential Pulse Voltammetry t-SNE Projection	34
Figure 3.11	Cyclic Voltammetry UMAP Projection	35
Figure 3.12	Differential Pulse Voltammetry UMAP Projection	36
Figure 3.13	Bokeh Interactive Plot with CV Data	37
Figure 4.1	Autoencoder Diagram	39
Figure 4.2	Classification Model Architecture	41
Figure 4.3	CV Ligand ROC Curves	48
Figure 4.4	CV Metal ROC Curves	48
Figure 4.5	DPV Ligand ROC Curves	49
Figure 4.6	DPV Metal ROC Curves	50
Figure 4.7	CV Ligand Confusion Matrix	51
Figure 4.8	Ligand Six and Ligand Seven Voltammogram Comparison . .	52
Figure 4.9	CV Metal Confusion Matrix	53
Figure 4.10	DPV Ligand Confusion Matrix	54

List of Figures

Figure 4.11	DPV Metal Confusion Matrix	55
Figure 5.1	AutoEncoder Results	58

Chapter 1

Introduction

Amidst urgent global challenges like climate change, energy sustainability, and health-care crises, there is a growing need for efficient solutions to address the needs of a growing population and increasing resource demands. Accelerating advancements in materials, technology, and scientific knowledge offer potential avenues for tackling these challenges. However, conventional research methods, marked by gradual progress and limited efficiency, may fall short of meeting the urgency posed by these issues. Self-driving laboratories (SDLs), which integrate laboratory automation and data-driven decision-making, emerge as promising tools to expedite and streamline the exploration of solutions while presenting several advantages over traditional scientific approaches [1]. Developing a fully autonomous self-driving laboratory is a complex endeavor that combines various research disciplines. Machine learning and modeling techniques are utilized to forecast materials properties and propose new experiments. SDLs typically use optimization techniques to guide their decision-making algorithm. An example of this is Atlas, a Python library offering access to different optimization algorithms, which has been used to identify the voltage peak in CV experiments to optimize the oxidation potential of a set of metal complexes [2]. Concurrently, robotics, computer vision, and automated characterization methods are employed to conduct experiments and analyze outcomes. Integrating these disparate technologies into a cohesive platform is central to the design of autonomous labs, facilitating seamless interaction between experiments and computational modeling [3].

Chapter 1. Introduction

SDLs can conduct experiments autonomously, performing tasks quicker and more precisely than manual processes. Moreover, they utilize data-driven algorithms to navigate through experimental spaces, enabling efficient exploration based on feedback from existing data, a process known as "closed-loop" experimentation. Additionally, SDLs address issues such as reproducibility challenges and the underrepresentation of negative results in scientific literature by promoting the digitization of research processes. Through automated systems, experimental protocols are meticulously documented, enhancing repeatability and reproducibility. Furthermore, digitization facilitates comprehensive data recording and sharing, emphasizing the importance of negative or null results, thus providing a more accurate depiction of scientific endeavors. The wealth of high-quality data generated by autonomous experimentation serves as a valuable resource for the development of artificial intelligence (AI) in materials science and chemistry. By improving machine learning (ML) and deep learning (DL) models, this data enhances the decision-making capabilities of SDLs, furthering their effectiveness in optimizing materials or processes and facilitating novel discoveries.

SDLs in chemistry and materials science are characterized by two critical dimensions: software autonomy and hardware autonomy. Regarding software autonomy, which governs experiment selection, SDLs are categorized into three types: (1) single iterations of automated experimentation with data-driven methods for selecting subsequent experiments, (2) multiple iterations within closed-loop systems where experimental results inform subsequent rounds of automated experiments, and (3) generative approaches involving numerous iterations of closed-loop optimization within algorithmically generated search or chemical spaces. By automating high-throughput

Chapter 1. Introduction

experimentation and streamlining experiment planning and execution, SDLs possess the potential to substantially accelerate research in chemistry and materials discovery. SDLs have played a pivotal role and made noteworthy advancements in various fields, including drug discovery, genomics, chemistry, and materials science [1]. Since SDLs can provide solutions to many different problems, SDLs should be widely adopted. However, the accessibility of SDLs is impeded by the substantial financial investment required for high-precision commercial platforms. Thus, there is a pressing demand for affordable alternatives to democratize SDLs [4].

Chapter 2

Background and Motivation

2.1 Electrochemistry

Given the pivotal role of reduction-oxidation (redox) reactions in materials chemistry and industrial applications, electrochemistry stands as a primary beneficiary of advancements in SDLs. According to the broad definition commonly accepted among researchers, electrochemistry encompasses the study of both the physical and chemical characteristics of ionic conductors, along with phenomena taking place at the interfaces between these ionic conductors and electronic conductors, semiconductors, other ionic conductors, and even insulating materials (such as gases and vacuum) [5]. The flow of electrons only occurs between two species, but the transfer of charge can also occur through an oxidation-reduction reaction. When a substance loses an electron, its oxidation state increases, indicating oxidation. When a substance acquires an electron, its oxidation state decreases, indicating reduction. For example, consider the following redox reaction, which has oxidation and reduction components:



Chapter 2. Background and Motivation

A redox reaction is balanced when the number of electrons gained by the oxidant is equal to the number of electrons lost by the reductant. Like any balanced chemical equation, the entire process is electrically neutral, meaning that the net charge remains consistent on both sides of the equation. With redox reactions, it is possible to separate the oxidation and reduction half-reactions physically in space, provided a complete circuit exists using an external electrical link, such as a wire, connecting the two halves. Electrons migrate from the reductant to the oxidant as the reaction progresses through this electrical connection, generating an electric current.

Devices that use redox reactions to generate electricity or use electricity to drive non-spontaneous redox reactions are called electrochemical cells. This device effectively transforms chemical energy into electrical energy or vice-versa. In an electrochemical cell, reduction and oxidation reactions occur at the electrodes. The electrode where reduction occurs is termed the cathode, while oxidation occurs at the anode. An electrode serves as a stable electrical conductor, facilitating the flow of electrical current within non-metallic solids, liquids, gases, plasmas, or even vacuums. Electrodes are typically fabricated from highly conductive materials, including but limited to metals and graphite [6]. In a battery, redox reactions create a flow of electrical current that can be used to power electronic devices.

Electrode potential is the voltage of an electrochemical cell composed of a reference electrode and another electrode to be characterized. Figure 2.1 shows a three-electrode setup typical for electrochemical experiments such as cyclic voltammetry. During the flow of current between the working and counter electrodes, the reference electrode is used to precisely measure the applied potential in relation to a stable

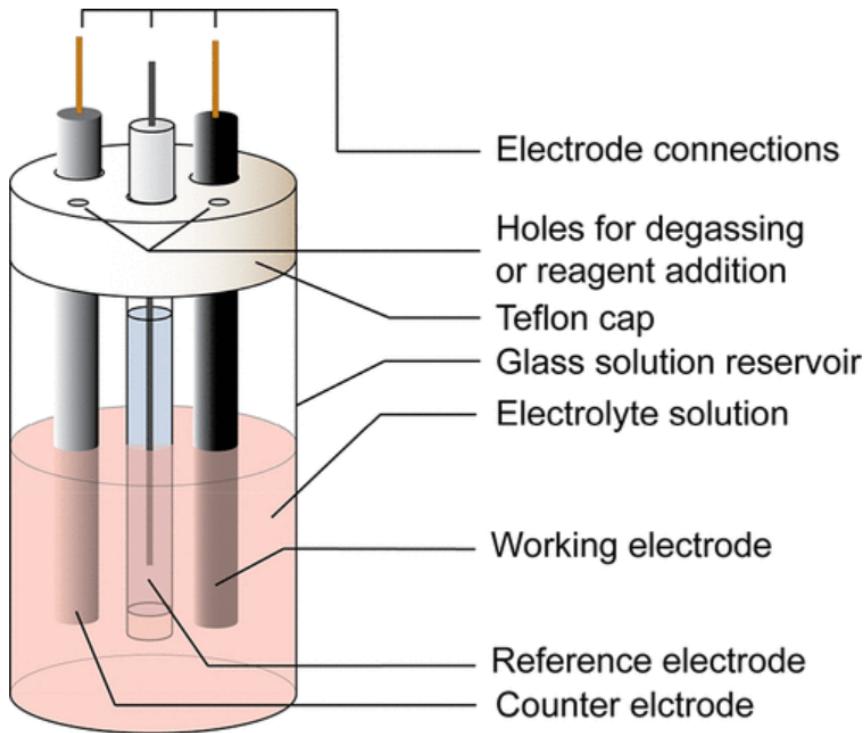


Figure 2.1: Schematic of Electrochemical Cell [7]

reference reaction. A potentiostat, as shown in Figure 2.2, is an analytical instrument designed to control the potential between the working electrode and counter electrode within a multi-electrode cell [8]. The potentiostat contains various internal circuits tailored to fulfil this role, facilitating the generation and measurement of potentials and currents. External wires within a cell cable establish connections between the potentiostat circuit and the electrodes within the electrochemical cell. In a three-electrode configuration, the cell cable links the working, counter, and reference electrodes on one terminal and the potentiostat cell cable connector on the opposite end. The potentiostat's internal circuitry governs the applied signal.

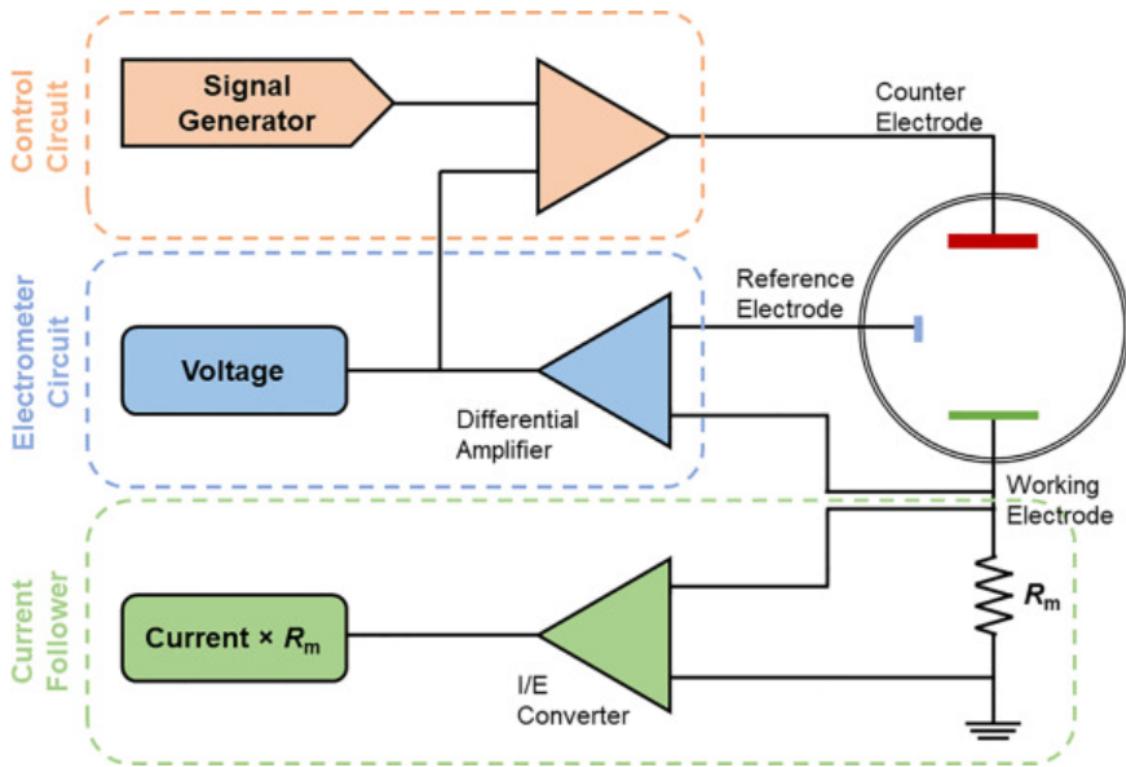


Figure 2.2: Potentiostat Circuit Diagram

The working electrode performs the electrochemical event of interest. Since reactions occur at the cathode and anode surfaces, it is crucial that the surface is spotless and that the surface area is well-defined. The working electrodes should be immediately polished after use to ensure there are no surface contaminants that inhibit electron transfer. Even a few hours of air exposure will degrade the electrode surface. Detecting when surface contamination affects data quality is one of the questions this work addresses. This detection can trigger automatic polishing or replacement with a new disposable electrode [9].

Commercial vendors commonly provide potentiostats that are governed by propri-

Chapter 2. Background and Motivation

etary software, employ graphical user interfaces (GUI), and produce already curated data. These devices are widely used for electroanalytical experiments such as cyclic voltammetry and differential pulse voltammetry. Commercial potentiostats can vary in design, but a typical potentiostat is shown in Figure 2.2 and consists of three component circuits: a control circuit, an electrometer, and a current follower [10]. The electrometer circuit utilizes a differential amplifier to measure the difference in potential between the working and reference electrodes. Subsequently, the measured potential feeds into the control circuit, which administers a current through the counter electrode, altering the relative potential of the working electrode to align with the user-defined parameters. A single generator ensures this potential adheres to a predefined periodic waveform. The current flowing through the working electrode is then assessed by a current follower circuit, commonly in the form of a current-to-voltage converter. This circuit measures the drop in potential across a grounded resistor, allowing the current to be determined using Ohm's law.

However, commercial potentiostats present challenges for integration into automated systems due to their reliance on proprietary software and GUIs. Furthermore, their high cost poses a significant barrier for groups seeking to perform high-throughput analysis.

To address these issues, Pablo-García et al. recently introduced an open-source, low-cost potentiostat [11]. This innovative device aims to democratize electrochemical analysis by reducing the financial barrier to entry and improving integration with automation systems. By providing an affordable and accessible alternative to traditional commercial potentiostats, this open-source solution empowers researchers to

conduct electrochemical experiments with greater flexibility and efficiency. Despite its remarkable advancements, the device's precision falls short of commercial standards. As such, we later explore various machine-learning methodologies to enhance data quality.

2.2 Cyclic Voltammetry

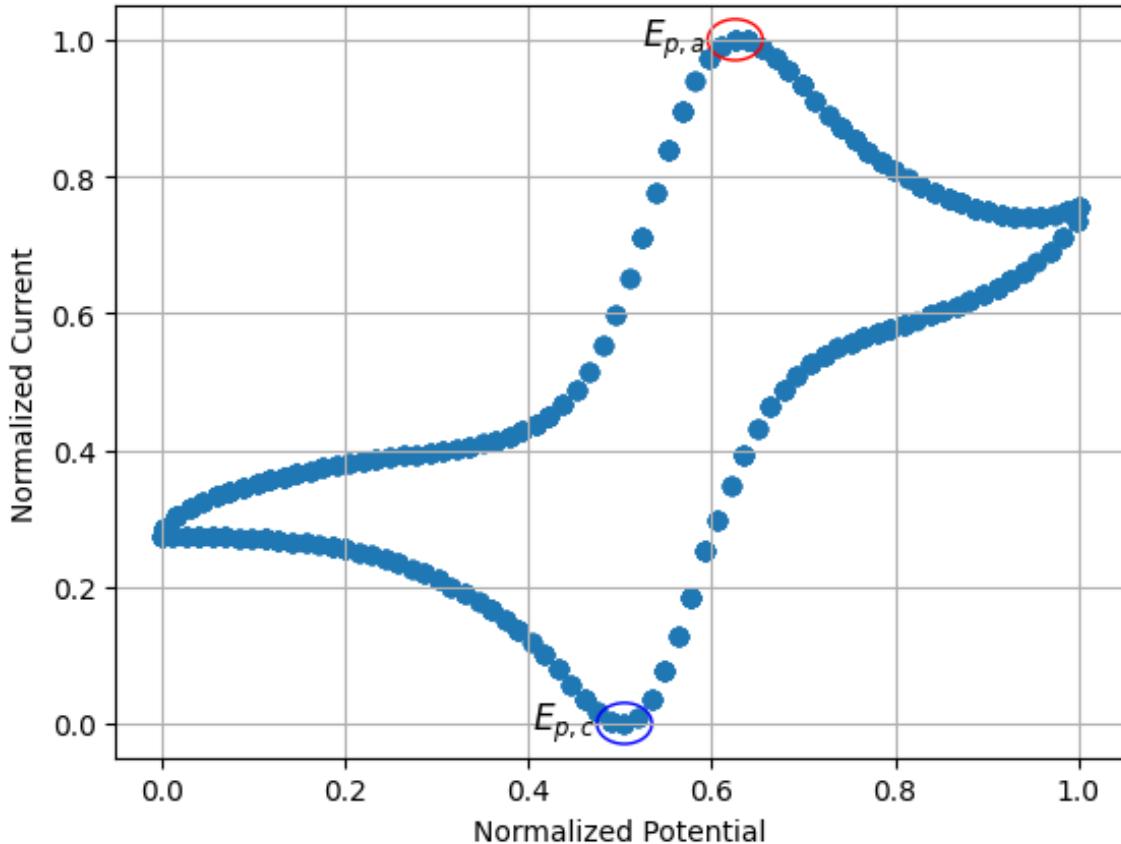


Figure 2.3: Cyclic Voltammogram

Cyclic voltammetry (CV) is a common electrochemical characterization that extracts

Chapter 2. Background and Motivation

important reduction and oxidation information about molecules [12]. Typically, the working electrode potential increases linearly with time. After reaching a predetermined limit, the potential decreases to return to the starting voltage. These cycles can be repeated as frequently as needed to bolster confidence in the obtained data. The rate of voltage change over time is known as the experiment's scan rate (Voltage/Time) and affects how many data points are gathered throughout the experiment [13].

CV is valuable for studying qualitative information about electrochemical processes across diverse conditions. It enables the examination of intermediates in oxidation-reduction reactions and the assessment of reaction reversibility. Other use cases include the determination of electron stoichiometry, analyte diffusion coefficients, and formal reduction potentials, aiding in identification processes [14]. Additionally, in reversible Nernstian systems, the proportional relationship between concentration and current allows for determining unknown solution concentrations by constructing calibration curves correlating current and concentration [15].

In a typical cyclic voltammogram shown in Figure 2.3, peaks represent electrochemical processes occurring at the electrode surface. The anodic peak ($E_{p,a}$) is observed during the scan where oxidation of the electroactive species occurs at the electrode and corresponds to the potential at which oxidation is most favourable. The current increases as the potential applied to the electrode becomes more positive, reaching a maximum at the peak potential. The cathodic peak ($E_{p,c}$) is observed during the reverse scan where reduction of the electroactive species occurs at the working electrode and corresponds to the potential at which reduction is most favorable. The

current increases as the potential becomes more negative, reaching a maximum at the peak potential [16]. Typically, chemists are especially interested in these peaks as they condense the redox behavior of the analyzed compound [17].

2.3 Differential Pulse Voltammetry

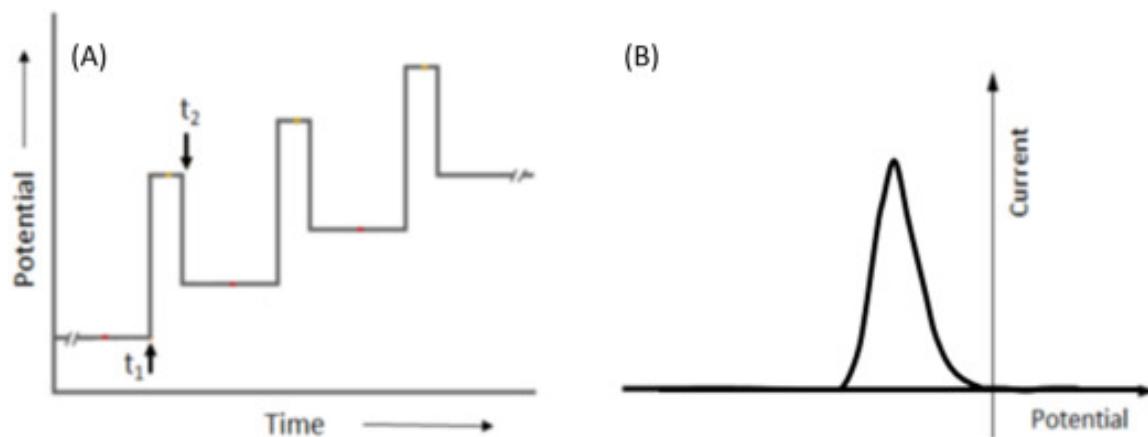


Figure 2.4: Differential Pulse Voltammogram

Differential Pulse Voltammetry (DPV) is a more sophisticated electrochemical measurement technique where a series of increasing pulses are applied across the electrodes in an electrochemical cell [18]. The current I_1 is measured right before applying the pulse at time t_1 , and I_2 is measured again at the end at time t_2 . The difference in current ($\Delta I_2 - I_1$) is plotted against the potential and results in a peak-like shape. This method helps reduce the impact of charging current by sampling the current just before the potential change. DPV is well suited for measurements with extremely low concentrations of chemicals. This is because the effect of the charging current can be minimized to achieve high sensitivity, and only the faradaic current, the electric

Chapter 2. Background and Motivation

current generated by the redox of a chemical at an electrode, is extracted so that electrode reactions can be measured precisely [19].

Furthermore, DPV is a versatile tool for the qualitatively characterizing chemical compounds and their electrochemical properties [18]. By analyzing the shape, position, and area of the peaks in the DPV curve, chemists can glean insights into the nature of the electroactive species present, their concentration, kinetics of electron transfer processes, and other relevant electrochemical parameters. This capability makes DPV invaluable in various fields such as analytical chemistry, environmental monitoring, and pharmaceutical research, where understanding the behaviour of chemical compounds at the molecular level is crucial [18].

Chapter 3

Clustering

3.1 Introduction

Given the capabilities and limitations of SDLs, quick and accurate characterization of the electrical compounds produced is needed. Clustering experimental results becomes crucial for several reasons. Clustering identifies patterns and similarities among experimental results. This aids in the discovery of underlying trends or relationships between different compounds or experimental conditions. It also facilitates quality control by pinpointing outliers or anomalies in experimental data, ensuring the reliability of SDL data.

Moreover, clustering allows researchers to optimize processes by providing insights into the effects of various parameters, such as metal/ligand ratio, on the formation of redox or electrochemical compounds. This optimization can significantly enhance the efficiency of voltammetry automation in SDLs. Additionally, by classifying different types of electrical compounds based on their properties or characteristics, clustering supports classification and prediction tasks, enabling researchers and SDLs to predict the behaviour of new compounds or classify unknown compounds based on their similarities to known clusters. Categorizing the automated characterizations typically done by SDLs can improve the automation capabilities of the laboratory. This is done by improving the behaviour mapping of a certain compound under certain conditions

into a chemical space. Finally, clustering provides a structured way to organize and interpret large volumes of experimental data, facilitating decision-making processes related to the selection of compounds for further analysis or the design of future experiments. In essence, clustering experimental results in the context of SDLs used for voltammetry is indispensable for gaining insights, ensuring data quality, optimizing processes, classifying compounds, and facilitating decision-making processes, which are crucial in an SDL environment.

3.2 Data Collection

To analyze how data gathered from SDLs can be clustered, this work uses an open dataset published by the Aspuru-Guzik group [11]. The data was collected through autonomous electrochemistry experimentation that operates through an iterative workflow [11]. The workflow was used to synthesize ten distinct metals and ten distinct ligands, with specific details available in Appendix A.1 and Appendix A.2, resulting in 100 unique complexes. Each complex was synthesized using a metal/ligand concentration ratio of 1:7 to ensure complete complexation. The synthesis process employed 1.0 M NaCl in water as the electrolyte/solvent and a buffer solution consisting of a 1:1 ratio of HOAc/NaOAc. Following synthesis, comprehensive characterizations were conducted using CV and DPV techniques. The experimentation was done using a low-cost electrochemistry platform designed as an alternative to commercial options. As they were measured with a custom potentiostat, the number of points in each sample can vary due to different scan rates. Higher scan rates lead

to more data points being collected during the experiment and can provide finer resolution of the electrochemical processes occurring. Additionally, it's worth noting that these samples may be duplicated as CV and DPV analyses can be conducted multiple times on the same sample to ensure robustness. Notably, the workflow is adaptable, with the potential to encompass a broader range of parameters, including additional ligands, varying metal/ligand ratios, and reaction times. Mixed ligands and different buffer pH levels can also be configured but was not done so for this dataset. The accumulation of data points is ongoing, contributing to the continuous expansion and refinement of our understanding. The final dataset consists of 800 CV and 200 DPV data points. The dataset used in this work can be found in the article preprint [11].

3.3 Curse of Dimensionality

The curse of dimensionality refers to the phenomena that cause various challenges and complications when analyzing data in high-dimensional spaces. As the number of features in a dataset increases, the amount of data needed to generalize accurately grows exponentially [20]. As the number of dimensions increases, the data becomes increasingly sparse. This makes tasks like clustering and classification more challenging. In higher dimensions, the difference between distances between data points starts to become negligible, making measurements like Euclidean distance negligible. As such, algorithms that rely on distance measurements will experience a drop in performance. Furthermore, more dimensions will require more computational resources and time to process the data. It is good practice to aim to have the data in as low a

dimension as possible, provided that relevant information is maintained.

3.4 Ramer–Douglas–Peucker Algorithm

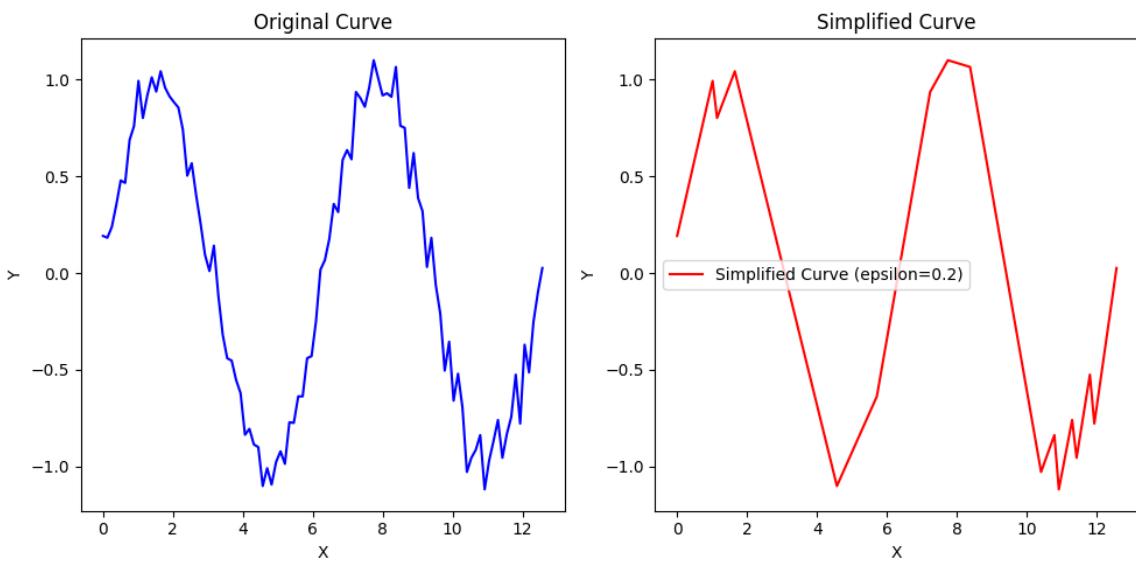


Figure 3.1: RDP Algorithm

The Ramer–Douglas–Peucker (RDP) algorithm is employed to reduce the number of points in a curve approximated by a series of points. It operates by conceptualizing a line between the initial and terminal points within the curve's points. Subsequently, it identifies the point furthest from this line among the intermediary points. If this point, termed the outlier point, and consequently all intervening points, lie within a specified distance epsilon from the line, they are removed. Conversely, suppose the outlier point exceeds the epsilon distance. In that case, the curve is segmented into two parts: from the initial point to the outlier point, inclusive and the outlier and the remaining points. The algorithm is then recursively applied to both resulting segments, and the

reduced forms of the curve are reassembled. Figure 3.1 shows curve simplification done with the RDP algorithm. Since CV and DPV results can be represented as a curve, RDP can be used to remove unnecessary points while maintaining the overall shape of the voltammogram. This will reduce the dimensionality and improve data analysis results.

3.5 Data Preparation and Encoding

Many parameters can be set during CV and DPV analysis, affecting the characterization outcomes. Notably, the experiment's scan rate affects the sampling frequency and the number of points collected within a specific time interval, leading to a variable number of point densities depending on the analyzed compound. Heterogeneity among samples becomes challenging for many ML algorithms, as they often require input data to be the same shape. Similarly, the potential limit at which the potential begins to return to its initial point will affect the overall shape of the cyclic voltammogram. To address these issues, the following steps are used to encode the data:

1. Split experiment cycles into separate data points
2. Normalize values to fit between [0, 1]
3. Reduce points using the Ramer-Douglas-Peucker algorithm
4. Duplicate data points until the total length reaches the longest cycle's length

5. Order data points based on angular position relative to the center

Due to the curse of dimensionality, the RDP algorithm is used to reduce the number of dimensions. Since the RDP algorithm takes only a variable ϵ , the final length after reduction will differ for each data set. Data points are randomly selected and duplicated to ensure the data is the same size as the longest after RDP reduction. Finally, the data is ordered based on its angular position relative to the center for consistency. Plots of the raw and processed data can be seen in Figure 3.2, with the

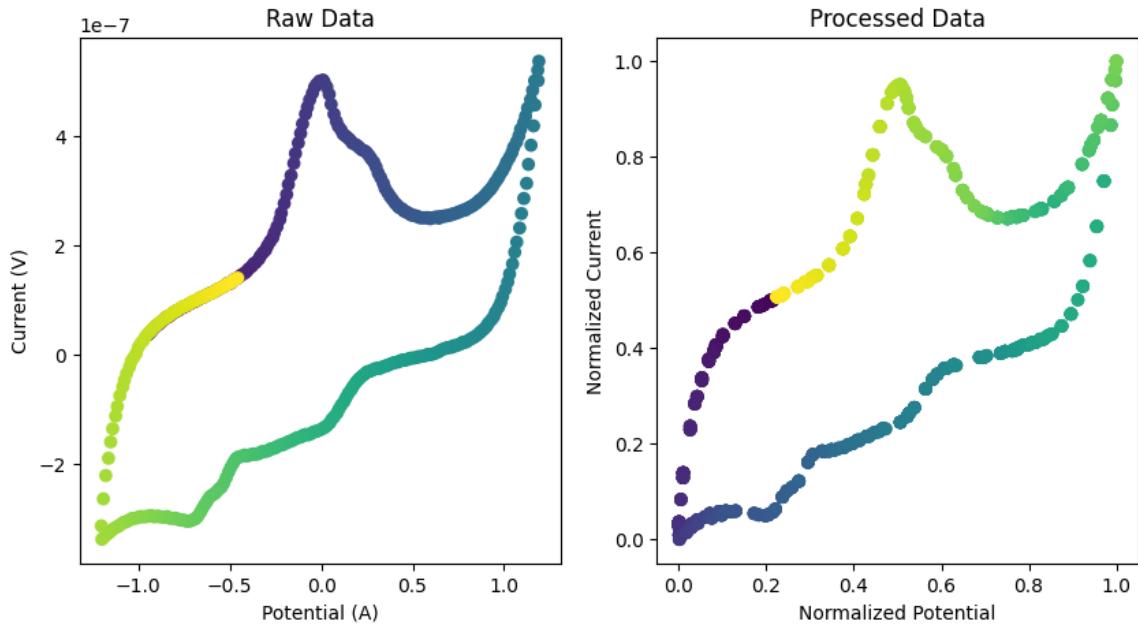


Figure 3.2: Raw Data and Processed Data

colour of the scatter plot representing the order in which the points appear. The starting point and end point varies across different voltammetry experiments. As such, it is important to order the data points so that comparisons can be informative. A significant reduction in dimensionality by $1/3$ can also be seen in the plots. Despite

this, the important characteristics of the voltammogram, such as the overall shape and peaks, are maintained.

3.6 K-Means

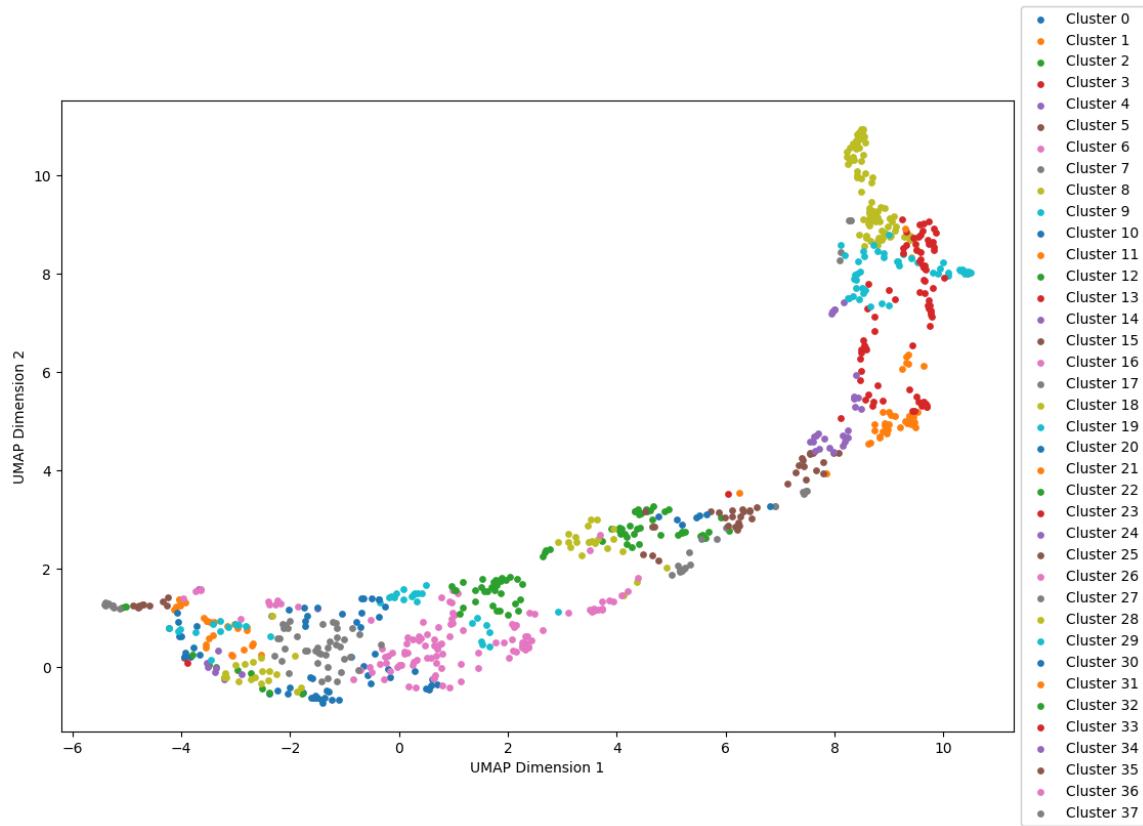


Figure 3.3: K-Means Clustering Visualization with CV data
Due to the dimensionality of the data, UMAP (see subsection 3.9) was applied to reduce the dimensions of the data.

K-Means clustering is an unsupervised machine learning algorithm aimed to divide data points into clusters so that the data points within each cluster are similar and different from the data points in other clusters [21]. The K-Means clustering result

Chapter 3. Clustering

for CV data can be seen in Figure 3.3. The algorithm is explained below, with K representing the desired number of clusters:

1. Initially, K points are selected randomly as the cluster centroids
2. Each data point is assigned to the closest mean, quantified by the Euclidean distance.
3. Each cluster centroid is updated to reflect the average of data points currently assigned to that cluster
4. This process is repeated for a specified number of iterations

One of the questions that needs to be answered is the choice of K. This means finding a balance between the number of clusters represented by K and the average variance of the clusters while minimizing both. There is no approach for determining K that works better than all others. For this problem of clustering CV and DPV data, a combination of the Elbow Method [22] and Silhouette method is used ROUSSEEUW198753. The Elbow Method is performed by plotting the within-cluster sum of squares (WCSS) for a range of K and choosing the value K where adding more clusters does not significantly decrease the WCSS. While the Elbow Method can quickly eliminate many values of K, it also has drawbacks regarding the shape of the WCSS curve. Determining the exact location of the "elbow" can be subjective and depends on the analyst's interpretation. Different individuals may identify different elbows, leading to inconsistency in results. In cases where the relationship between the number of clusters and WCSS is not distinctly elbow-shaped, the Elbow Method may not provide clear guidance for choosing the appropriate number of clusters.

The Silhouette Method addresses some of these drawbacks by providing a more quantitative measure of cluster quality. Instead of relying on subjective interpretation, the Silhouette Method calculates the silhouette coefficient for each data point, quantifying how similar an object is to its cluster compared to others. This provides a more objective measure of cluster cohesion and separation. The process for selecting K for this work includes determining a set of candidate K values using the Elbow Method by eliminating suboptimal values and then using the Silhouette method to find optimal K among the potential candidates.

3.7 Density-Based Spatial Clustering of Applications with Noise

Density-Based Spatial Clustering of Applications with Noise (DBSCAN) is another clustering algorithm that works by partitioning the data into dense regions of points separated by less dense areas [23]. It defines clusters as areas of the dataset with many points close to each other, while the points far from any cluster are considered outliers or noise. In DBSCAN, $\text{eps} (\epsilon)$ represents the maximum distance between two points for them to be considered neighbours, and minimum samples represents the number of points required for a point to be considered a core point. Points that have fewer than minimum samples points are labelled as noise. The key differentiator for DBSCAN is that the number of clusters does not need to be determined beforehand.

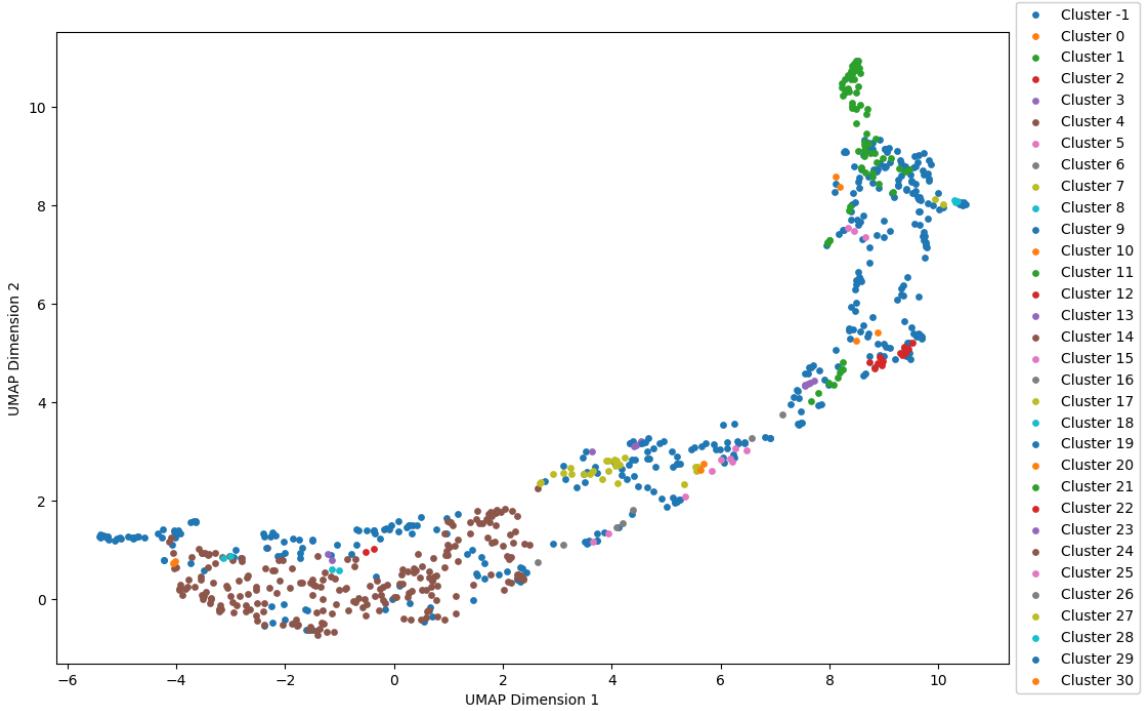


Figure 3.4: DBSCAN Clustering Visualization with CV Data
 UMAP (see subsection 3.9) was applied to reduce the dimensions of the data.

3.8 t-Distributed Stochastic Neighbour Embedding

Dimensionality techniques like t-Distributed Stochastic Neighbour Embedding (t-SNE) are used for visualizing high-dimensional data in a low-dimensional space [24]. This visualization can aid in the clustering process by providing insights into the underlying structure of the data and help in understanding the results of the clustering algorithm. The first step of the algorithm is to create a probability distribution that represents the similarity between neighbours. The similarity between the two data points is represented by their Euclidean distance. Each data point is placed in the middle of the Gaussian curve, and the rest of the data is placed along the curve. This

is represented by the following equation where $j \neq i$ and $p_{i|i} = 1$:

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)} \quad (3.1)$$

"The similarity of datapoint x_j to datapoint x_i is the conditional probability, $p_{j|i}$, that x_i would pick x_j if neighbours were picked in proportion to their probability density under a Gaussian centred at x_i " [24]. The remaining variable, sigma, is not chosen directly but rather by selecting a value for perplexity. Perplexity is defined as:

$$Perp(p) := 2^{-\sum_x p(x) \log_2(p(x))} \quad (3.2)$$

Perplexity represents the density of data and how many neighbours the central point should have with higher values relating to higher variance. After choosing the perplexity value, the corresponding sigma values are found using binary search. Next, the similarities between data points for low-dimensional representations must also be found to ensure that similar data are close after projection.

3.9 Uniform Manifold Approximation and Projection

Uniform Manifold Approximation and Projection (UMAP) is a dimension reduction technique that can be used for visualization similarly to t-SNE [25]. It achieves this by leveraging concepts from algebraic topology and Riemannian geometry. Here's a simplified breakdown of how UMAP works:

Chapter 3. Clustering

1. Constructing a Topological Representation: UMAP starts by creating a fuzzy topological representation of the data. This involves building a simplicial complex, which is a way to represent topological spaces using simple geometric shapes called simplices. The algorithm constructs these simplices based on the proximity of data points.
2. Optimizing Low-Dimensional Representation: Once the topological representation is established, UMAP optimizes a low-dimensional representation of the data to match this topological structure as closely as possible. It does this by minimizing a measure called cross-entropy, which quantifies the difference between the fuzzy topological structures of the high-dimensional and low-dimensional data.
3. Efficient Computations: UMAP employs several strategies to make computations efficient. It focuses on computing only the nearest neighbours of each point and uses algorithms like Nearest-Neighbour-Descent for this purpose. Additionally, it utilizes stochastic gradient descent for optimization and smooth approximations of the membership strength function to ensure differentiability.
4. Preserving Topological Structure: The goal of UMAP is to ensure that the low-dimensional representation maintains the essential topological properties of the original data. It achieves this by balancing attractive forces that pull similar points together and repulsive forces that push dissimilar points apart based on the weights of edges in the topological representation. The farther away the two points are, the more dissimilar they are.

3.10 Results and Discussion

The K-Means clustering algorithm was used to categorize the entire set of experimental voltammetry data after encoding. With K-Means, a value of K will need to be selected. This is done using the elbow method. Figure 3.5 shows the results of the elbow method applied to the dataset. It can be seen that this methodology identifies multiple potential candidates for K-values, necessitating a more comprehensive analysis to select the most appropriate option. The Silhouette method is used to

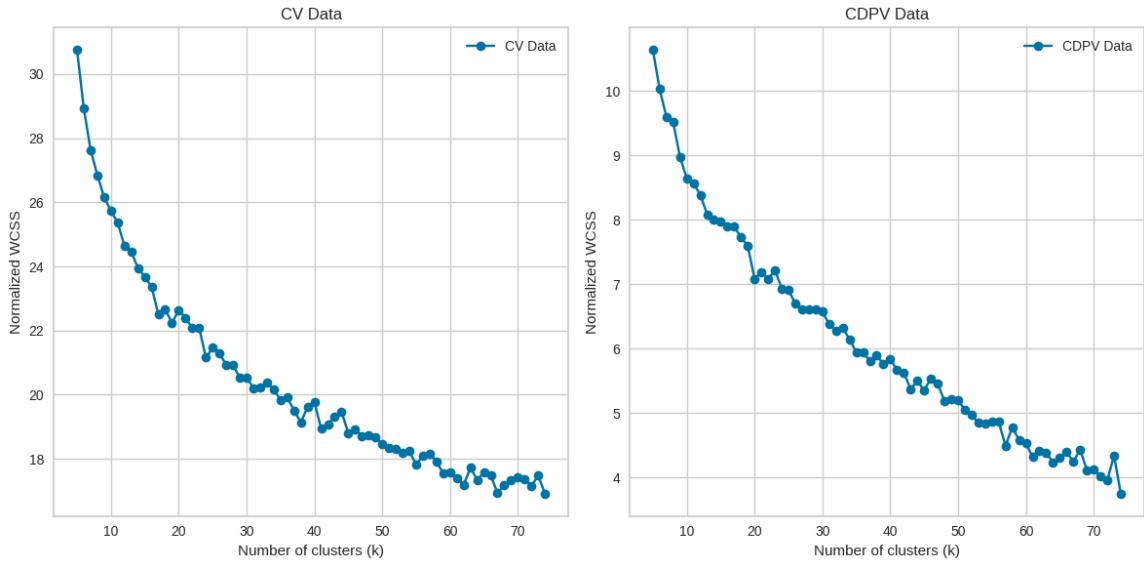


Figure 3.5: K-Means Elbow Method
Plots show no obvious choice for K

analyze promising values to aid the decision-making process. A cluster with a value of one means points are perfectly assigned in a cluster, and clusters are easily distinguishable. Zero means clusters overlap; negative one means points are assigned to the wrong cluster. The K value should be chosen based on which value produces the

Chapter 3. Clustering

most clusters with Silhouette scores greater than the average score of the dataset, represented by the red-dotted line seen in Figure 3.6a and Figure 3.6b. Furthermore, there should not be wide fluctuations in the size of the clusters. The width of the clusters represents the number of data points belonging to the cluster. In Figure 3.6a, which showcases the application of the Silhouette method for CV cross-validation, $K = 38$ yields the highest number of clusters with a score surpassing the mean of the dataset. This configuration reduces the number of clusters scoring below zero and minimizes the variance in cluster sizes. Similarly, in Figure 3.6b showcasing the Silhouette method for DPV, $K = 42$ results in the best quality of clusters. A subset of the cluster results is available in the appendix. Despite having 100 different combinations of metals and ligands, using a relatively small K value still shows promising results, as the data points within each cluster have similar overall shapes, which is crucial for compound identification.

DBSCAN, as an alternative clustering method, demonstrated significant promise in identifying anomalous data points. With the appropriate parameters, DBSCAN efficiently grouped cycles from the same experiment. As Figure 3.7 depicts, DBSCAN defines a cluster comprising cycles solely from a single experiment. This capability could be seamlessly incorporated into SDLs as an error validation mechanism. Any cycle not assigned to the same cluster as others from the identical experiment could trigger an error notification, prompting intervention and investigation. Another method is identifying the erroneous clusters and investigating the assigned data points, as seen in Figure 3.8.

To further demonstrate the efficacy of the encoding, t-SNE and UMAP projections are

Chapter 3. Clustering

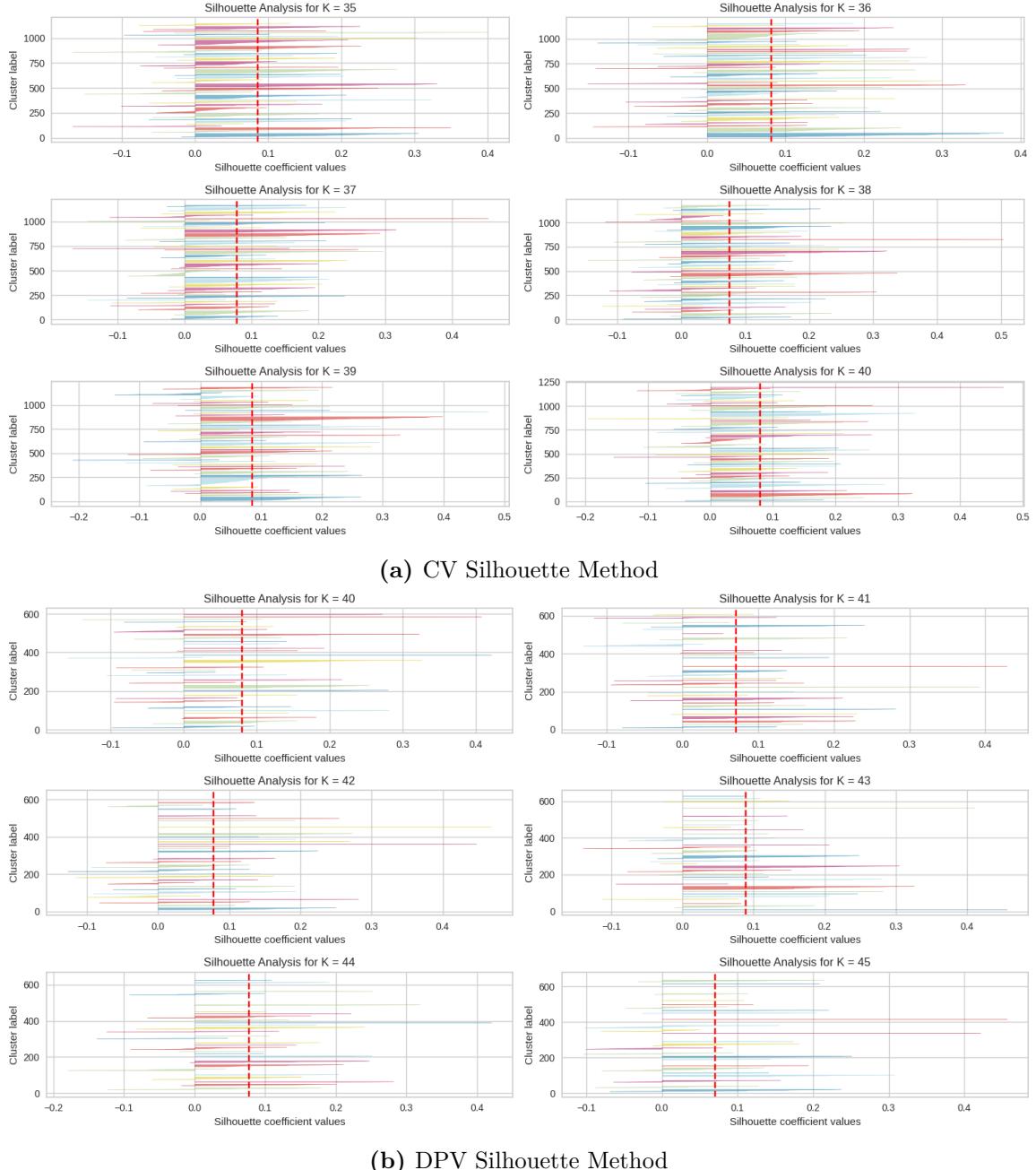


Figure 3.6: Silhouette Method for Promising K Values

created to visualize the data in 2-D and show how the shapes, metals, and ligands are distributed. As seen in Figure 3.9 and Figure 3.10, t-SNE emphasizes local structure

Chapter 3. Clustering

and tends to agglomerate similar data points into tight clusters. As a result, t-SNE plots often show clearer separation between clusters but may not preserve the global structure as effectively. t-SNE primarily preserves local neighbourhoods, which leads to tighter clusters of similar points. However, it may not always capture the global structure accurately, especially for complex datasets. t-SNE embeddings can vary significantly with different random initializations and parameter choices, making it less stable and potentially more sensitive to noise in the data. UMAP tends to focus more on preserving global structure and maintaining relative distance between clusters. Therefore, clusters in the UMAP plot are usually well-separated and evenly distributed. UMAP tries to preserve local and global neighbourhoods, resulting in more evenly spaced clusters and better representation of local and global structures. UMAP embeddings are generally more stable across different runs and parameter settings than t-SNE. Figures 3.11 and 3.12 illustrate these characteristics, especially when contrasted with the projections generated by t-SNE. For example, Figure 3.11 clearly shows more evenly spaced clusters than Figure 3.9.

Interactive plots, as seen in Figure 3.13, made with Bokeh, are available on [Github](#). The corresponding voltammetry plot is shown when hovering over a point with the mouse. Pan, zoom, and rotation tools are also available. These plots are extremely useful for chemists as SDLs can automatically generate them.

Utilizing machine learning techniques to classify voltammetry data based on the overall shape presents numerous benefits over merely employing a script to identify the number of peaks, as has traditionally been done. Machine learning models can be trained to recognize patterns and variations regarding the overall shape and number

Chapter 3. Clustering

of peaks. They can adapt to experimental conditions, electrode materials, and analytes without needing manual adjustment of parameters. Voltammetry data can often be noisy, especially at low concentrations. ML models can be trained to distinguish true peaks from noise more effectively than simple peak-finding algorithms. Voltammograms can vary in characteristics due to electrode deterioration, surface roughness, and solution composition. ML models can learn to handle this variability and provide more reliable peak classification across different experimental conditions. Additionally, ML models can learn when the electrode deteriorates and automatically polish it. ML models can automatically extract relevant features from voltammogram data, such as peak heights, peak widths, peak potential, and overall shape. This allows for more comprehensive analysis beyond locating peaks. Once trained, ML models can be integrated into larger data analysis pipelines to classify cyclic voltammetry data rapidly and efficiently, potentially saving time and effort compared to manual analysis or parameter tuning for peak-finding algorithms. These automated analysis techniques can be integrated into an SDL, updating them with newly generated data to increase accuracy as more experiments are performed.

In summary, clustering techniques are crucial in analyzing and interpreting experimental voltammetry results obtained from SDLs. By organizing data into meaningful clusters, clustering techniques like K-Means and DBSCAN and dimensionality reduction techniques like t-SNE and UMAP uncover patterns, similarities, and trends that enhance our understanding of the electrochemical compounds of interest. Choosing appropriate clustering algorithms and parameter selection methods, such as the Elbow and Silhouette Method, have been discussed to ensure meaningful and reliable

Chapter 3. Clustering

clustering results. The results obtained from clustering algorithms and dimensionality reduction techniques have provided valuable insights into the underlying structure of the experimental data, facilitating compound identification, error detection, and decision-making processes in SDLs.

Chapter 3. Clustering

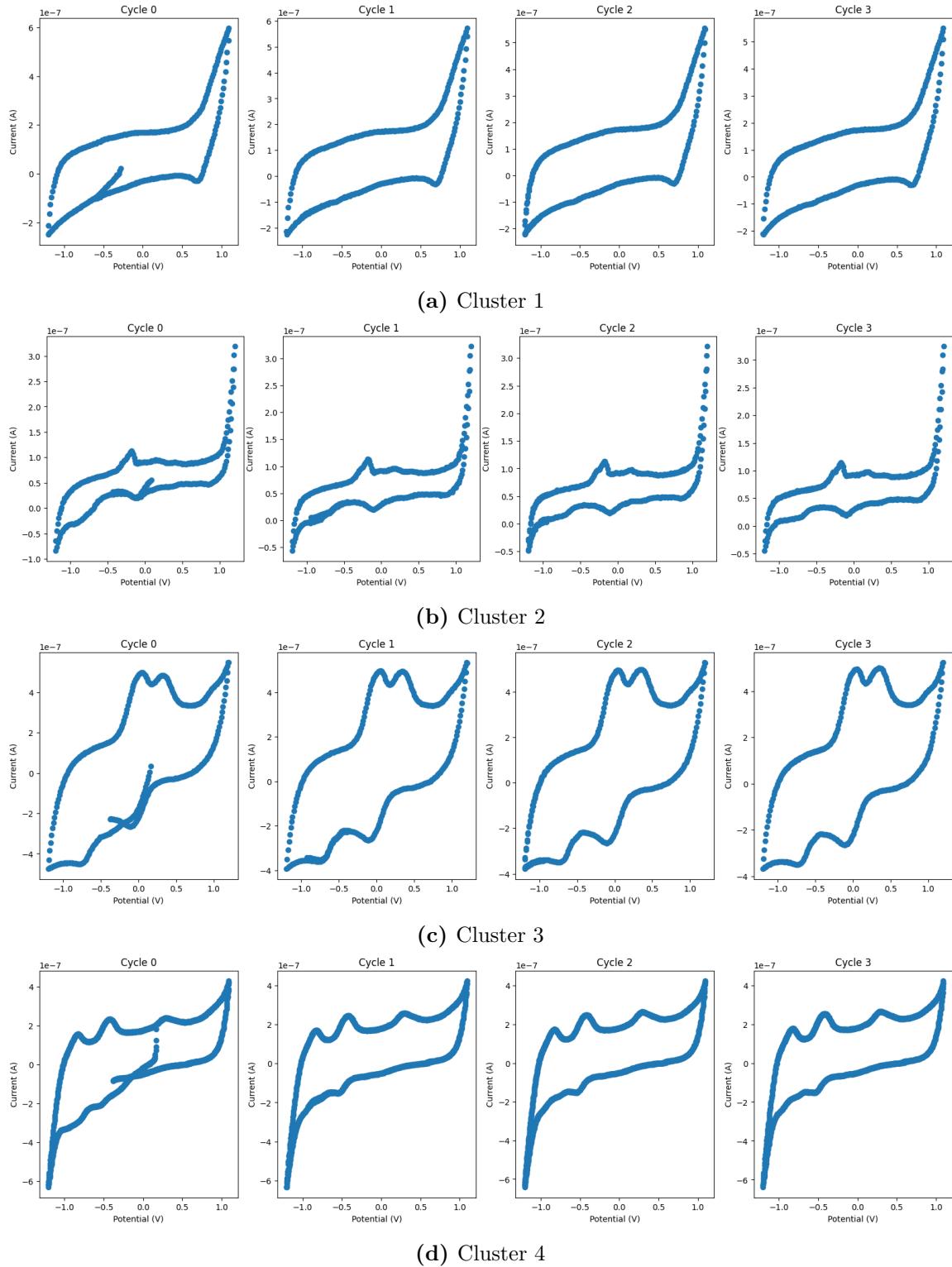


Figure 3.7: DBSCAN Clusters
31

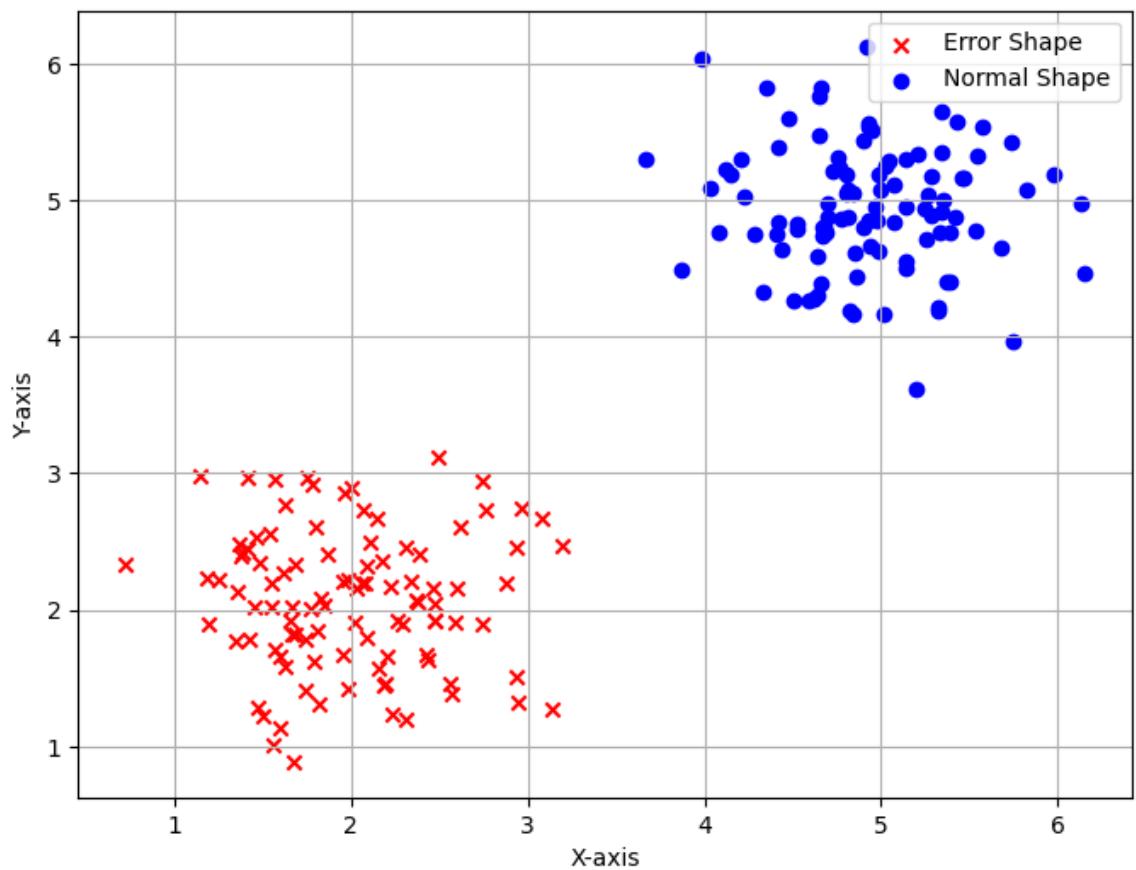


Figure 3.8: UMAP Visualization of Normal CV Data and Incorrectly Shaped Data

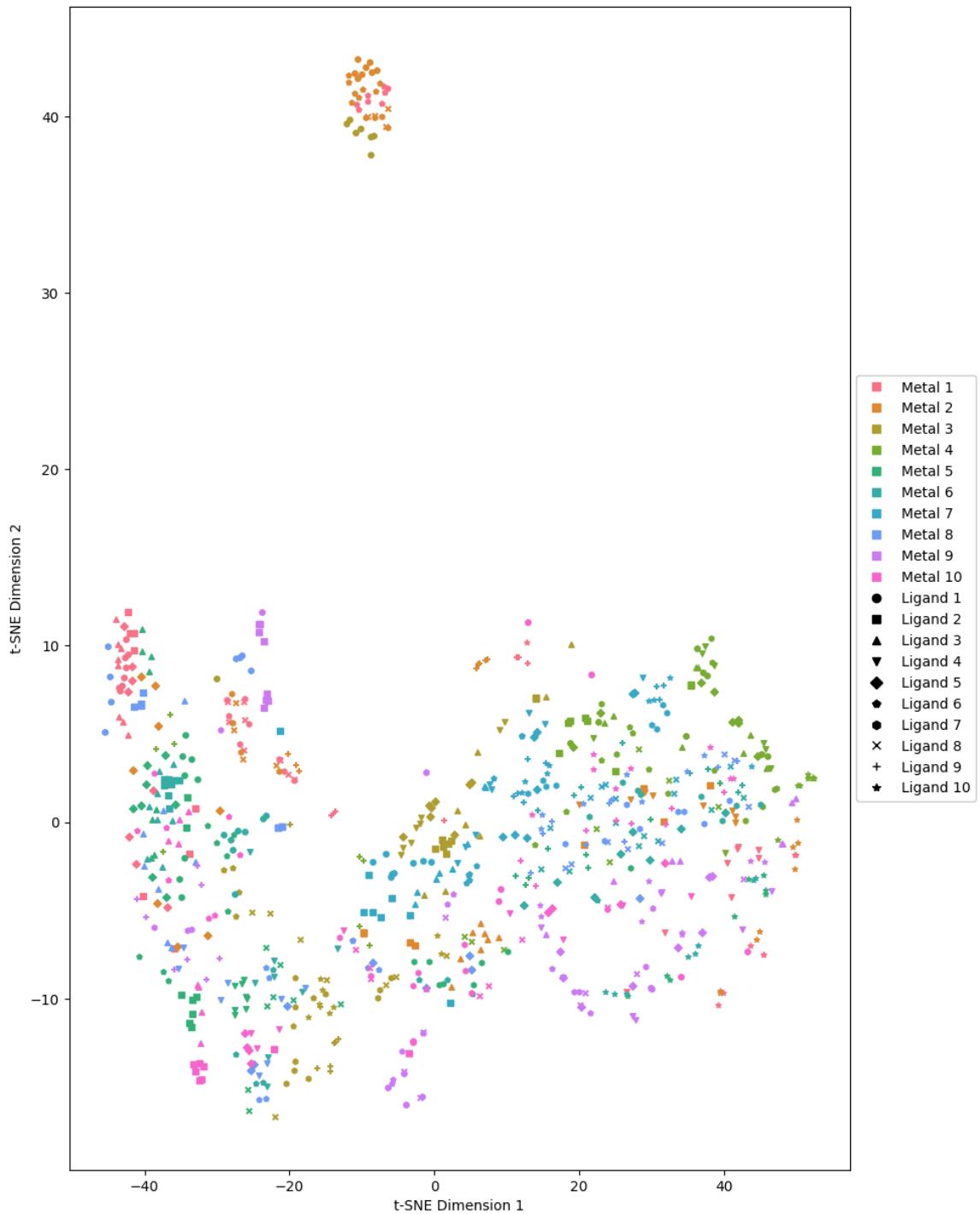


Figure 3.9: Cyclic Voltammetry t-SNE Projection

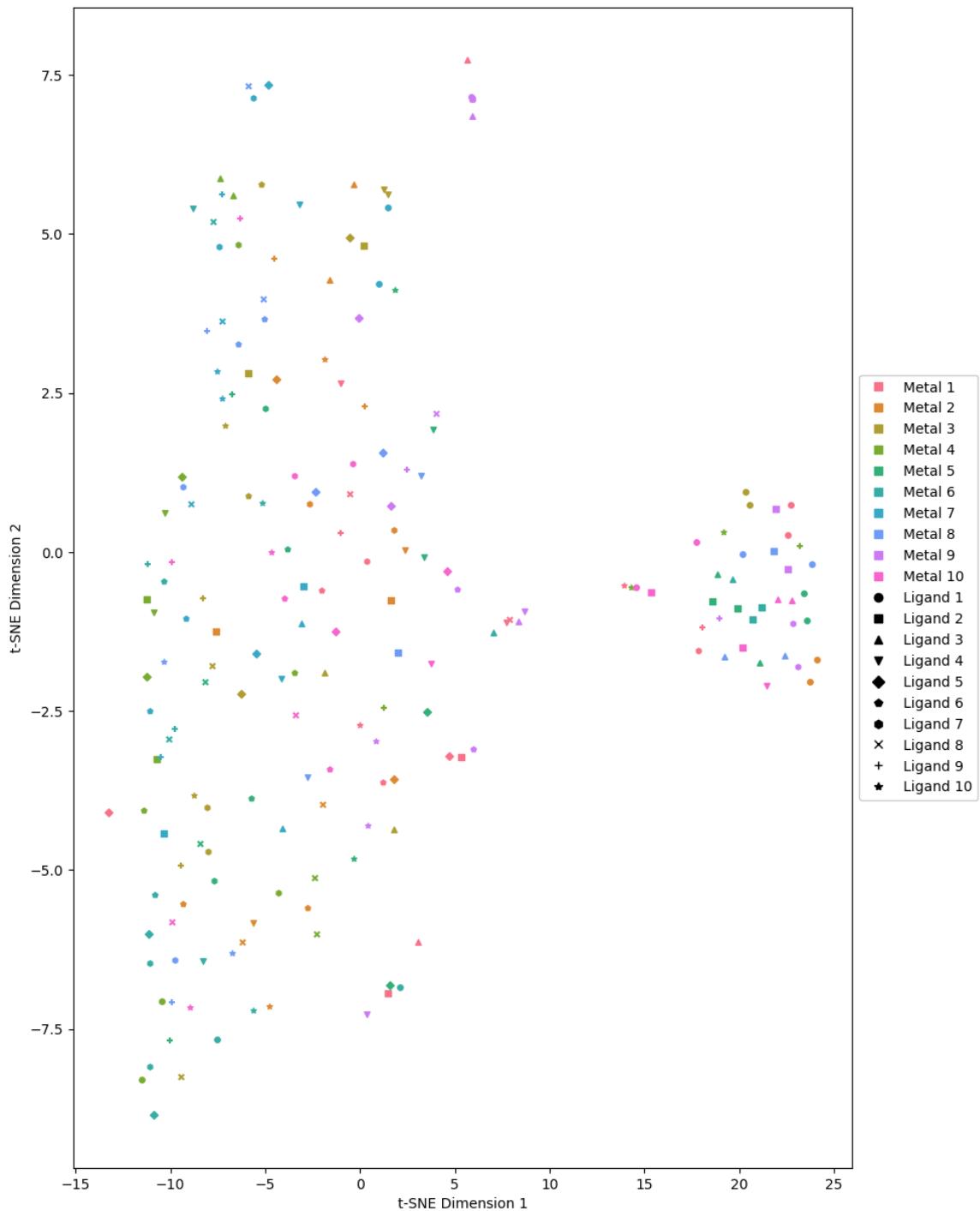


Figure 3.10: Differential Pulse Voltammetry t-SNE Projection

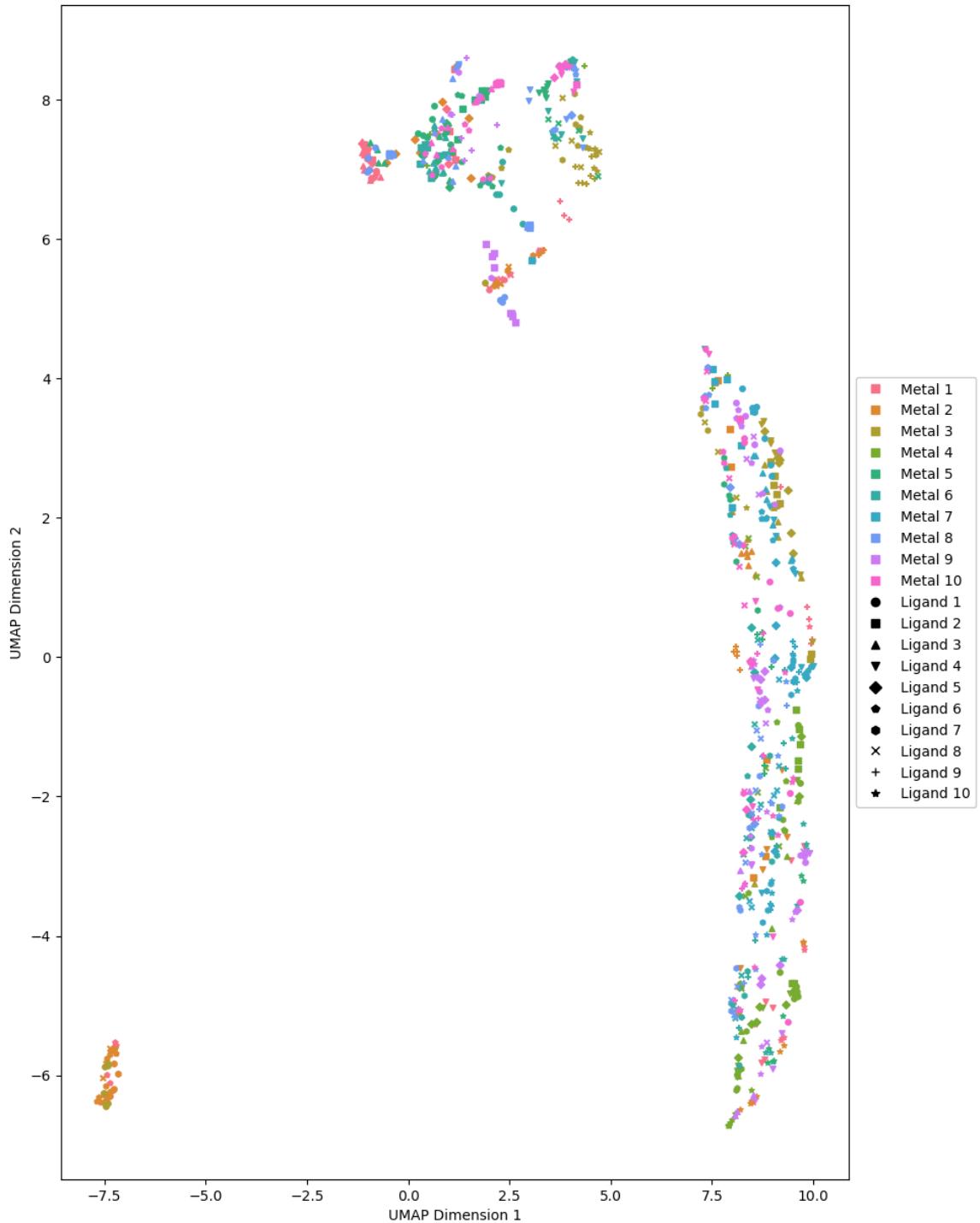


Figure 3.11: Cyclic Voltammetry UMAP Projection

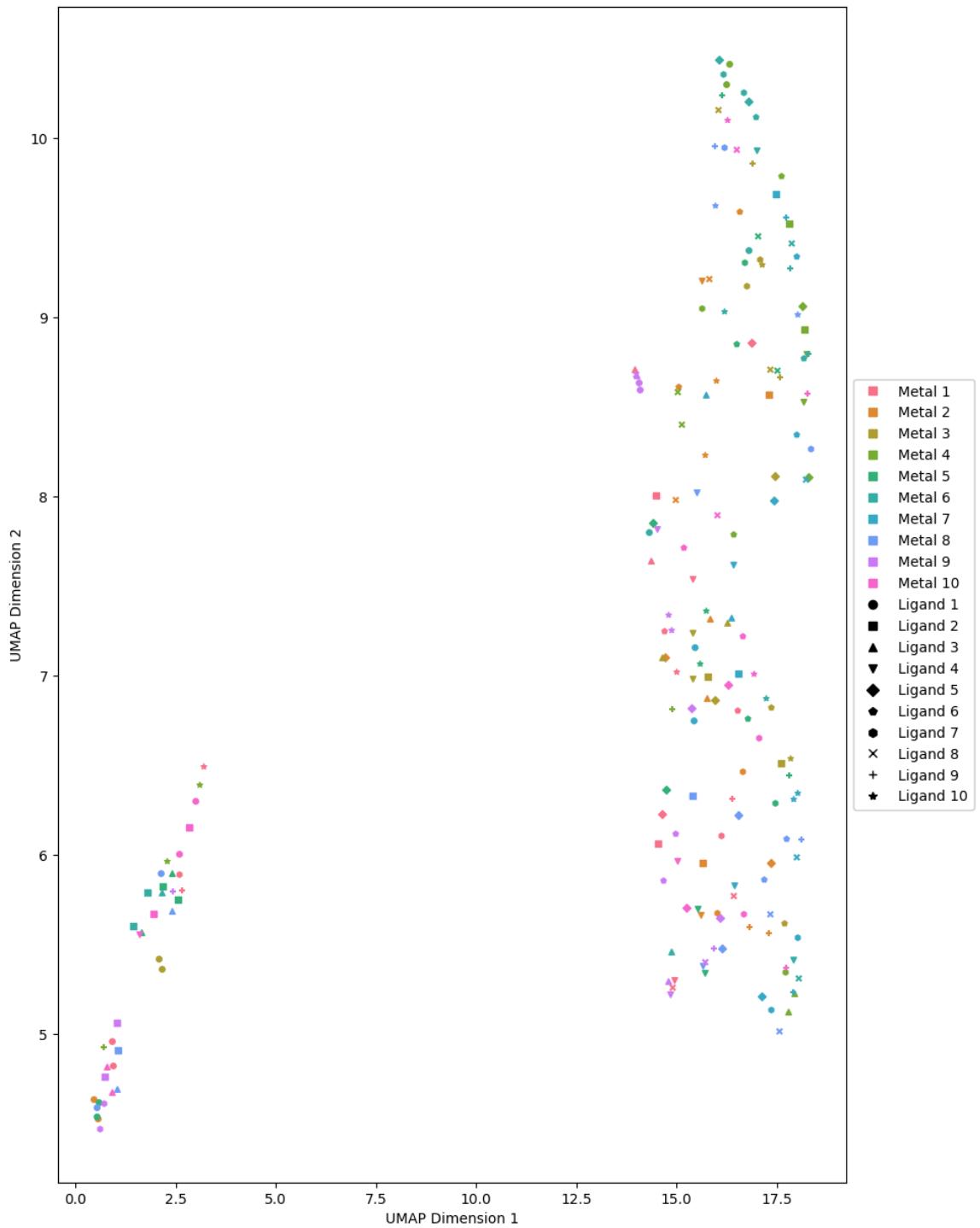


Figure 3.12: Differential Pulse Voltammetry UMAP Projection

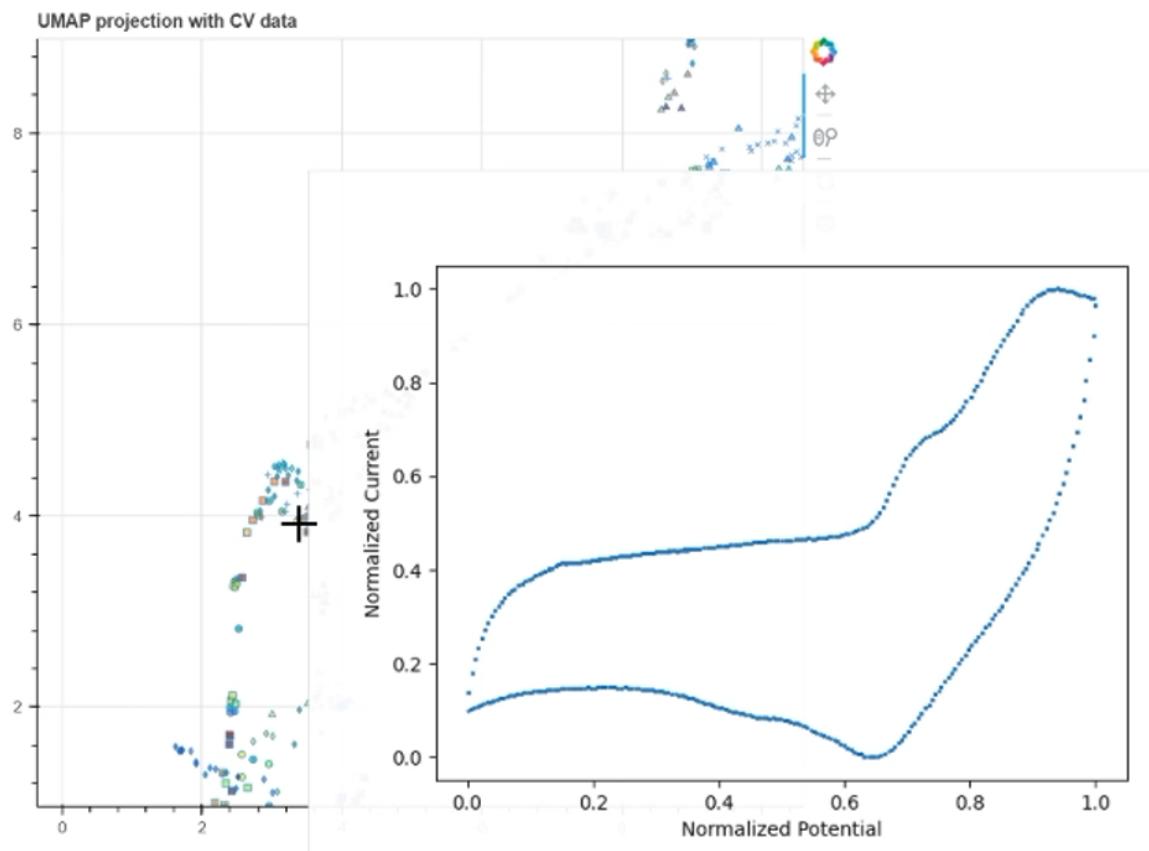


Figure 3.13: Bokeh Interactive Plot with CV Data

Chapter 4

Classification

4.1 Introduction

A classifier was trained to predict what metals and ligands were used to generate each voltammogram to demonstrate further the feasibility of using this encoding technique for various machine learning tasks. It is important to note that the dataset used is relatively small for a deep learning task. For training, the dataset was split with 80% for training, 10% for validation, and 10% for testing. 5-fold cross-validation, a technique for assessing the performance of a machine learning model by dividing the dataset into k subsets, training the model on $k-1$ subsets, and evaluating it on the remaining subset for each subset, is also used [26]. An important insight to consider is the similarity between voltammetry data and images. After all, each point has a potential and current value, similar to an image's RGB values. The main difference is that an image is 2-dimensional while voltammetry data is 1-dimensional. Many previous works have used convolutional neural networks (CNNs) for classification tasks [27]. Using this as inspiration, the proposed model architecture for voltammetry data classification uses 1-dimensional convolutional layers.

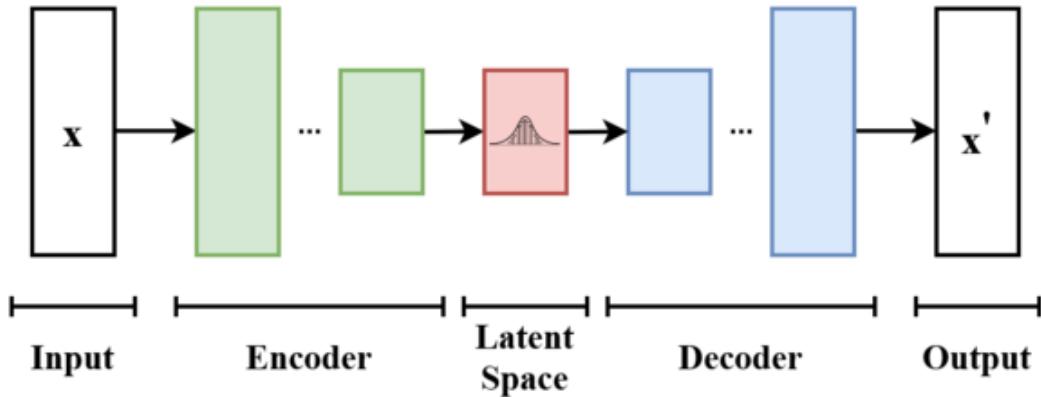


Figure 4.1: Autoencoder Diagram

Input is encoded into latent space and then recreated using decoder

4.2 Variational Autoencoders

Since one major challenge is the dataset size, one method to address this is to create synthetic data. A variational autoencoder (VAE) is similar to the autoencoder neural network architecture shown in Figure 4.1, with the main difference being that VAEs connects the encoder to its decoder through a probabilistic latent space that corresponds to the parameters of a variational distribution [28]. The encoder maps each point from the dataset into a distribution within the latent space rather than a single point in that space. The distribution is typically Gaussian with a mean and a variance. Once the VAE is trained, different points can be sampled from the learned latent space distribution. These samples represent different configurations of the input data in the latent space. The sampled points from the latent space are fed into the decoder network, which reconstructs the input data corresponding to those points

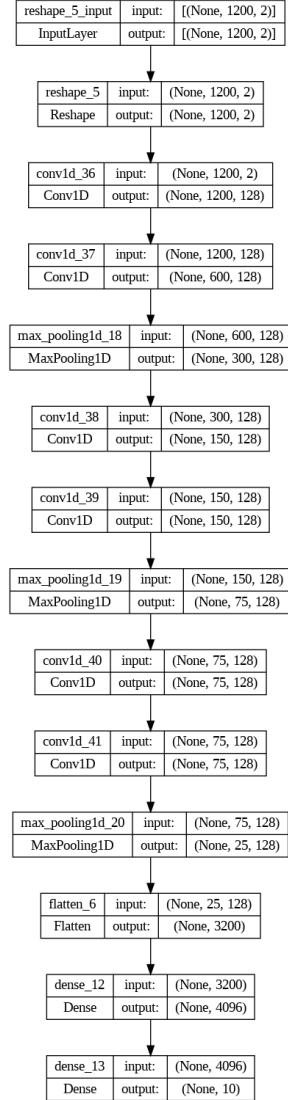
and generates a diverse set of synthetic data samples that resemble the original data distribution. The variability in the latent space allows for the generation of novel and diverse data samples that capture the underlying characteristics of the training data.

4.3 Conditional Variational Autoencoders

While traditional VAEs learn a latent space for the dataset, conditional variational autoencoders (CVAEs) expand this concept by introducing conditional dependencies between the input data and the latent variables. In the context of generating synthetic data, CVAEs offer a more controlled approach by allowing the generation process to be conditioned on additional information, such as class labels or other attributes associated with the data. By conditioning the generation process on known attributes or labels, CVAEs can generate synthetic data samples that capture the underlying data distribution and adhere to specific conditions or constraints defined by the conditioning variables. This enables the targeted generation of synthetic data for different classes or categories when labelled data is lacking. In this case, the metal and ligand are encoded using one-hot encoding and passed to the decoder to generate data from the same class.

4.4 Classifier Model Architecture

The classifier architecture can be seen in Figure 4.2. The model consists of several convolutional and max-pooling layers to encode the data and reduce dimensions. All


Figure 4.2: Classification Model Architecture

layers except for the output layer use the ReLU activation function. The output layer is a dense layer with ten units, one for each metal/ligand, and a softmax activation function. The Adam optimizer and categorical cross-entropy loss are used to train the model. Additionally, the model uses L2 regularization and early stopping to prevent overfitting and ensure smooth convergence. The Glorot uniform initializer is used for

weight initialization to facilitate better gradient flow and prevent exploding gradients.

4.5 Results and Discussion

Model	Fold 1 Acc	Fold 2 Acc	Fold 3 Acc	Fold 4 Acc	Fold 5 Acc	Avg Acc
CV Ligands	75.13%	80.29%	69.65%	76.40%	78.82%	76.06%
CV Metals	79.24%	82.50%	81.36%	80.11%	78.97%	80.44%
DPV Ligands	30.00%	35.00%	25.00%	25.00%	30.00%	29.00%
DPV Metals	25.00	30.00%	25.00%	20.00%	25.00%	25.00%

Table 4.1: Classification Accuracy

Separate classifiers were trained, each with a unique task of classifying CV ligands, CV metals, DPV ligands, and DPV metals. The accuracy of the classifiers can be seen in Table 4.1, and the results were much better for the CV data than the DPV data. This difference can likely be attributed to the size of the datasets. After

Model	Fold 1 Acc	Fold 2 Acc	Fold 3 Acc	Fold 4 Acc	Fold 5 Acc	Avg Acc
CV Ligands	77.86%	78.50%	80.94%	72.01%	73.02%	76.47%
CV Metals	85.00%	81.23%	84.19%	83.10%	78.75%	82.45%
DPV Ligands	25.00%	25.00%	15.00%	20.00%	20.00%	21.00%
DPV Metals	25.00%	20.00%	20.00%	25.00%	25.00%	23.00%

Table 4.2: Classification Accuracy with Synthetic Data

incorporating synthetic data generated with the CVAE into the training process, accuracy significantly improved for classifying CV data, as seen in Table 4.2. However,

the DPV ligands classifier saw a decrease in performance. Again, this is likely due to the size of the dataset. Several key considerations impact the quality of data generated by VAEs, especially when dealing with small datasets. Firstly, the quality and diversity of the original data influence the effectiveness of the synthetic data produced by VAEs. With limited variation or complexity in a small dataset, the VAE might struggle to accurately capture the proper underlying data distribution, potentially resulting in synthetic data that fails to represent the characteristics of the actual data fully. This mismatch can detrimentally affect classifier performance.

Additionally, the risk of overfitting is heightened in small datasets, where the classifier may excessively specialize in training data patterns that do not generalize well. Introducing synthetic data from a VAE can compound this issue if the VAE itself overfits the small dataset, producing synthetic data overly similar to the training data, which provides minimal additional information for the classifier and can lead to decreased performance on unseen data. VAEs implicitly learn the probability distribution of the input data. However, suppose the actual data distribution is significantly different from the distribution learned by the VAE due to the small dataset size. In that case, the synthetic data generated by the VAE may not accurately represent the true data distribution. This distribution mismatch can confuse the classifier, as it may encounter data points in the synthetic dataset that deviate from the real data distribution, leading to suboptimal performance. Table 4.3 provides insights into the precision, recall, and F1-score when classifying each metal type, along with the number of instances (support) for each metal type. Precision indicates the proportion of true positive predictions among all positive predictions, while recall measures the

	Precision	Recall	F1-Score	Support
Metal 1	0.88	0.88	0.88	8
Metal 2	0.80	1.00	0.89	8
Metal 3	1.00	1.00	1.00	4
Metal 4	1.00	0.83	0.91	12
Metal 5	1.00	0.71	0.83	7
Metal 6	0.88	0.78	0.82	9
Metal 7	0.82	0.90	0.86	10
Metal 8	0.50	0.40	0.44	5
Metal 9	0.78	1.00	0.88	7
Metal 10	0.82	0.90	0.86	10
Accuracy			0.85	80
Macro Avg	0.85	0.84	0.84	80
Weighted Avg	0.86	0.85	0.85	80

Table 4.3: CV Metals Classification Report

proportion of true positives that were correctly identified. F1-score, the harmonic mean of precision and recall, provides a balanced measure between the two. Overall, the classifier model achieved an accuracy of 85%, indicating its effectiveness in classifying different metal types. However, it is important to note variations in performance across metal types. For instance, Metal 3 achieved perfect precision, recall, and F1-score, suggesting the model's ability to classify this particular metal type accurately. On the other hand, Metal 8 exhibited lower precision and recall scores, indicating potential challenges in distinguishing this metal type from others. Both macro-average and weighted-average metrics hover around 0.85, indicating a reason-

ably balanced performance across all metal types. These metrics consider the average performance across all classes, with macro-average treating all classes equally, while weighted-average considers the contribution of each class based on its support. Ta-

	Precision	Recall	F1-Score	Support
Ligand 1	0.88	0.78	0.82	9
Ligand 2	0.88	0.88	0.88	8
Ligand 3	0.75	0.86	0.80	7
Ligand 4	0.45	0.71	0.56	7
Ligand 5	0.78	0.70	0.74	10
Ligand 6	1.00	0.86	0.92	7
Ligand 7	0.71	0.50	0.59	10
Ligand 8	0.67	0.89	0.76	9
Ligand 9	1.00	0.67	0.80	3
Ligand 10	1.00	0.90	0.95	10
Accuracy			0.78	80
MacroAvg	0.81	0.77	0.78	80
WeightedAvg	0.80	0.78	0.78	80

Table 4.4: CV Ligands Classification Report

ble 4.3 shows the classification report for classifying ligands. The classifier achieved an accuracy of 78% overall, indicating its capability to classify different metal types to some extent. However, upon closer examination, there are notable variations in performance across metal types. For instance, Metal 6 demonstrates excellent precision, recall, and F1-score, suggesting the model's proficiency in accurately classifying this metal type. Conversely, Metal 4 exhibits lower precision, recall, and F1-score,

indicating challenges in distinguishing this metal type from others. Table 4.5 and

	Precision	Recall	F1-Score	Support
Ligand 1	1.00	1.00	1.00	1
Ligand 2	0.00	0.00	0.00	2
Ligand 3	0.33	0.25	0.29	4
Ligand 4	0.33	0.33	0.33	3
Ligand 5	0.50	0.50	0.50	4
Ligand 6	0.00	0.00	0.00	1
Ligand 7	0.00	0.00	0.00	1
Ligand 8	0.00	0.00	0.00	0
Ligand 9	0.00	0.00	0.00	1
Ligand 10	1.00	0.33	0.50	3
Accuracy			0.30	20
MacroAvg	0.32	0.24	0.26	20
WeightedAvg	0.42	0.30	0.33	20

Table 4.5: DPV Ligands Classification Report

Table 4.6 show the classification reports for DPV ligands and metals. However, the small sample size makes it difficult to draw definitive conclusions from this data. To further assess the performance of these classification models, receiving operating characteristic (ROC) curves and area under the ROC curve (AUC) values can be used to gain valuable insights into the discrimination capabilities and robustness of the models when distinguishing between various metals and ligands. ROC curves and AUC values help assess the robustness of the classification model by showing how well it performs across different thresholds and levels of noise. A smooth ROC curve

	Precision	Recall	F1-Score	Support
Ligand 1	0.33	1.00	0.50	2
Ligand 2	0.00	0.00	0.00	1
Ligand 3	0.25	0.50	0.33	2
Ligand 4	0.00	0.00	0.00	3
Ligand 5	0.00	0.00	0.00	1
Ligand 6	0.50	0.25	0.33	4
Ligand 7	0.00	0.00	0.00	2
Ligand 8	1.00	0.50	0.67	2
Ligand 9	0.00	0.00	0.00	3
Ligand 10	0.00	0.00	0.00	0
Accuracy			0.30	20
MacroAvg	0.23	0.25	0.20	20
WeightedAvg	0.26	0.25	0.22	20

Table 4.6: DPV Metals Classification Report

with a high AUC suggests that the model can reliably discriminate between different metals and ligands even in the presence of noise or variability.

Furthermore, there may be a need to choose a classification threshold that balances sensitivity and specificity according to specific requirements or constraints. ROC curves provide a visual aid for selecting an appropriate threshold based on the desired trade-off between true positives and false positives. For example, when integrating with an SDL, minimizing false positives (misclassification of metals or ligands) might be prioritized over maximizing true positives. The ROC curves in Figure 4.3 and

Chapter 4. Classification

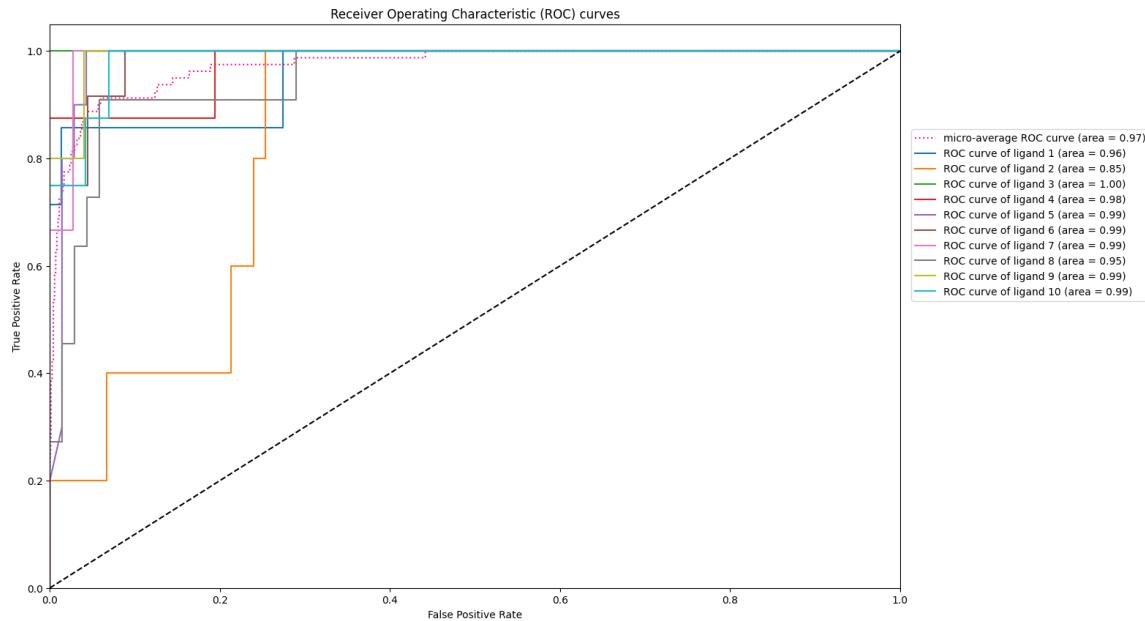


Figure 4.3: CV Ligand ROC Curves

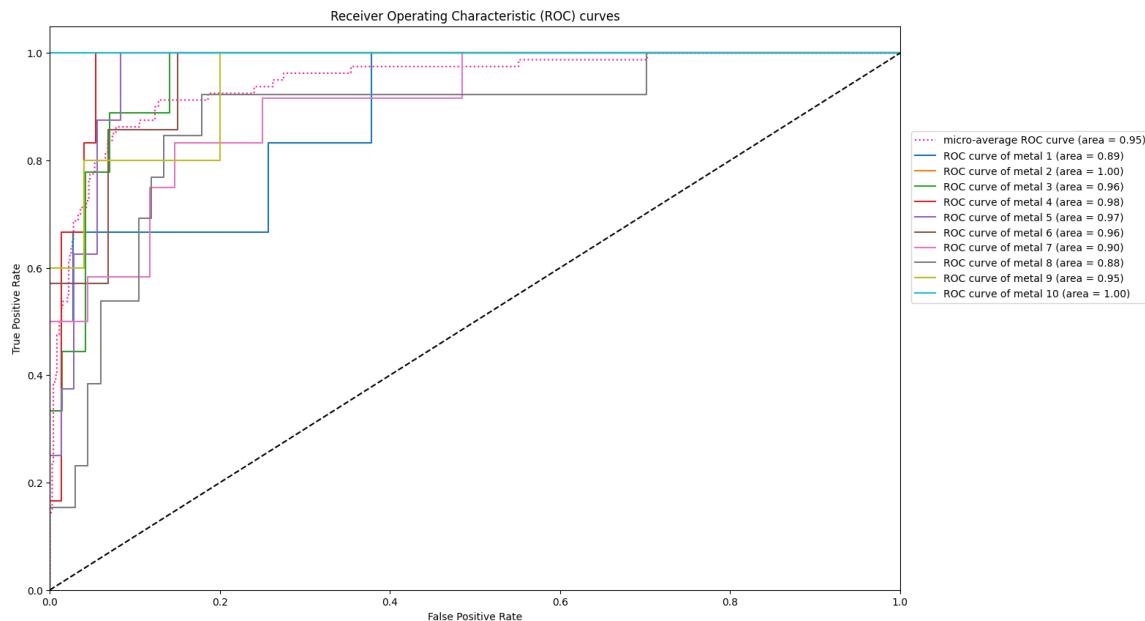


Figure 4.4: CV Metal ROC Curves

Chapter 4. Classification

Figure 4.4 show good results for both metals and ligands. The area under the ROC curve (AUC) calculation summarized the ROC curve analysis into a scalar value, which ranges between 0 and 1. The closer the AUC score to value 1, the better the application's overall performance. In Figure 4.5 and Figure 4.6, the ROC curves

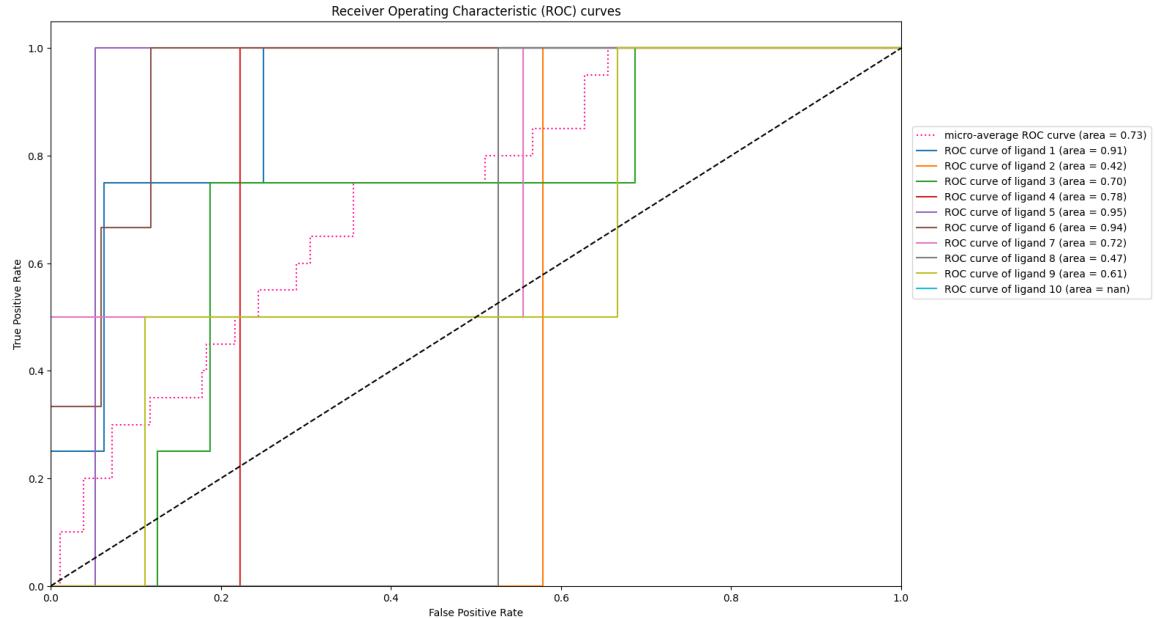


Figure 4.5: DPV Ligand ROC Curves

show that the classifier outperforms a random classifier by having an AUC value above 0.5. The data itself may cause issues with classification as some ligands and metals may be more difficult to distinguish than others. The confusion matrices are provided to investigate this. From the confusion matrix for ligands 4.7, ligand seven was often misclassified as ligand six. However, this misclassification is understandable. Figure 4.8 shows that the voltammograms for ligand six and ligand seven are quite similar. From the confusion matrix for metals 4.9, metal one was difficult to recognize, with many metals being misclassified as metal one. From the DPV confusion matrices

Chapter 4. Classification

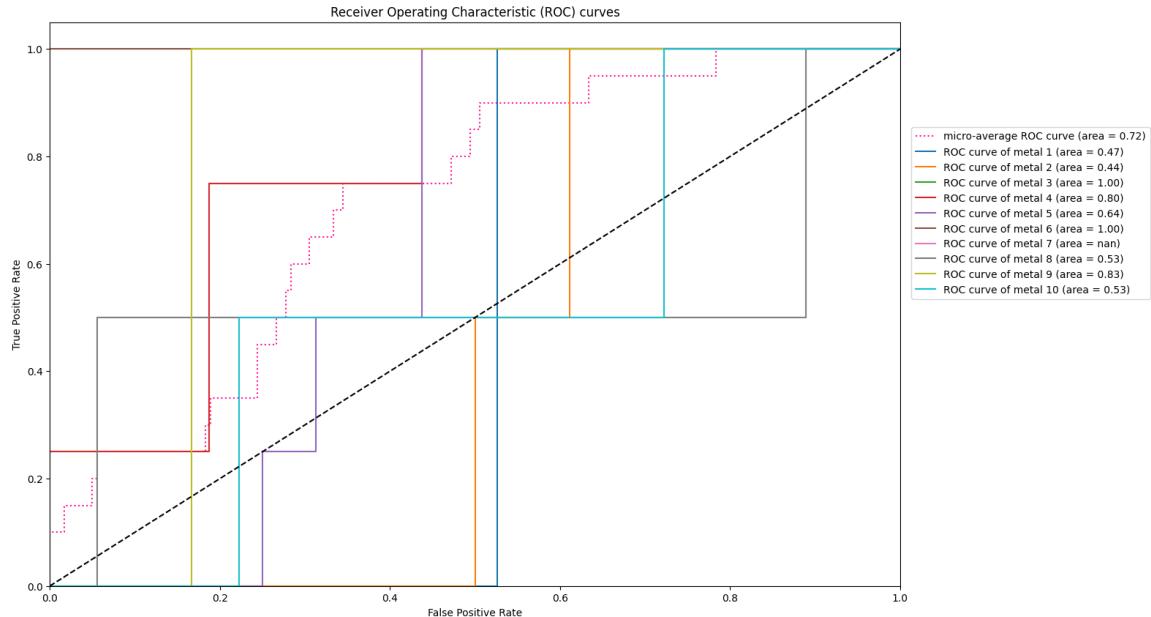


Figure 4.6: DPV Metal ROC Curves

seen in Figure 4.11 and Figure 4.10, it is hard to draw any definitive conclusions due to the dataset size. A major challenge in supervised learning is providing good examples during training. However, despite using a small dataset, these results are promising. This study establishes that crude CV data obtained from an economical potentiostat can be effectively encoded using CNNs. It was also shown that VAEs and CVAEs can generate high-quality, generalizable synthetic data. These findings align with recent research demonstrating that deep learning models can efficiently process CV data [29]. Although the ligand and metal labels are recorded. The results demonstrate that the encoding technique effectively captures the chemistry behind the measurements. Future research can incorporate group SELFIES within the decoder layer to predict or select from a pool of candidate redox groups identified through voltammetry or predict the compound used.

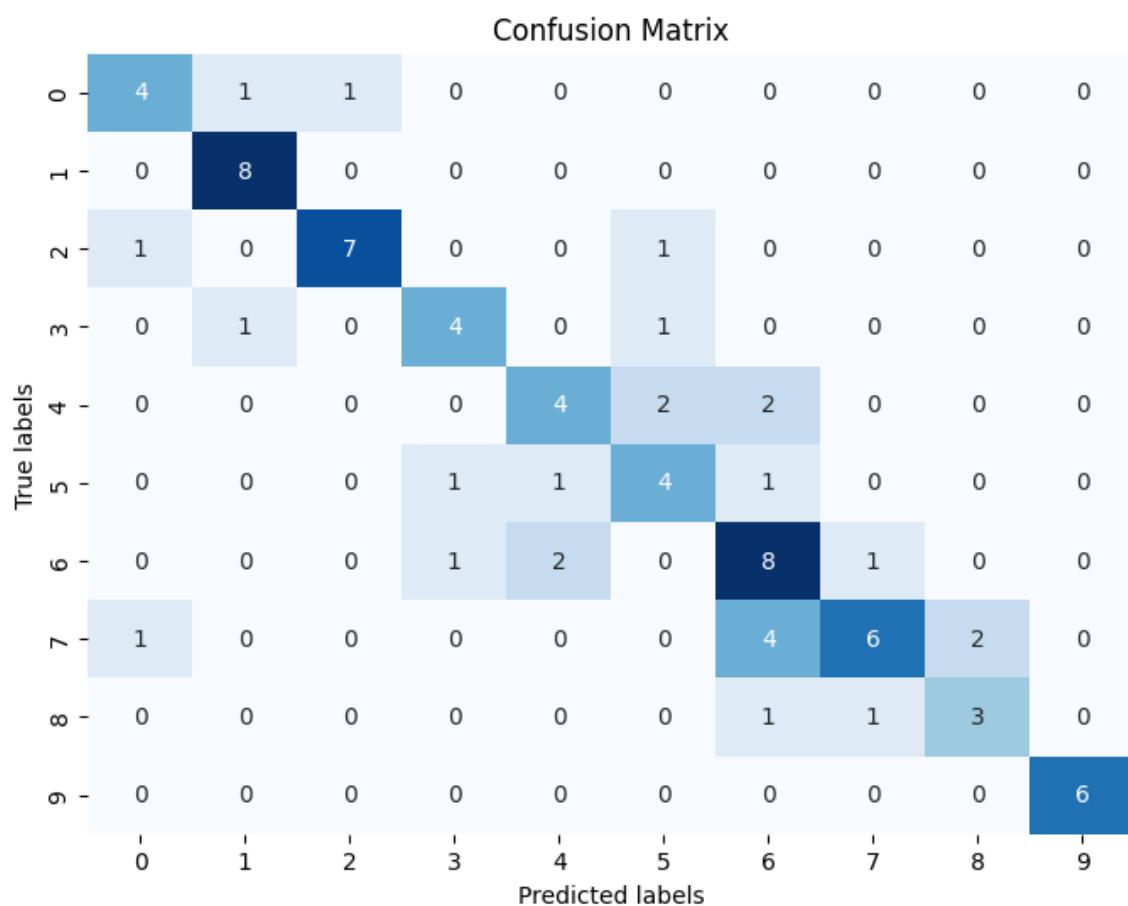


Figure 4.7: CV Ligand Confusion Matrix

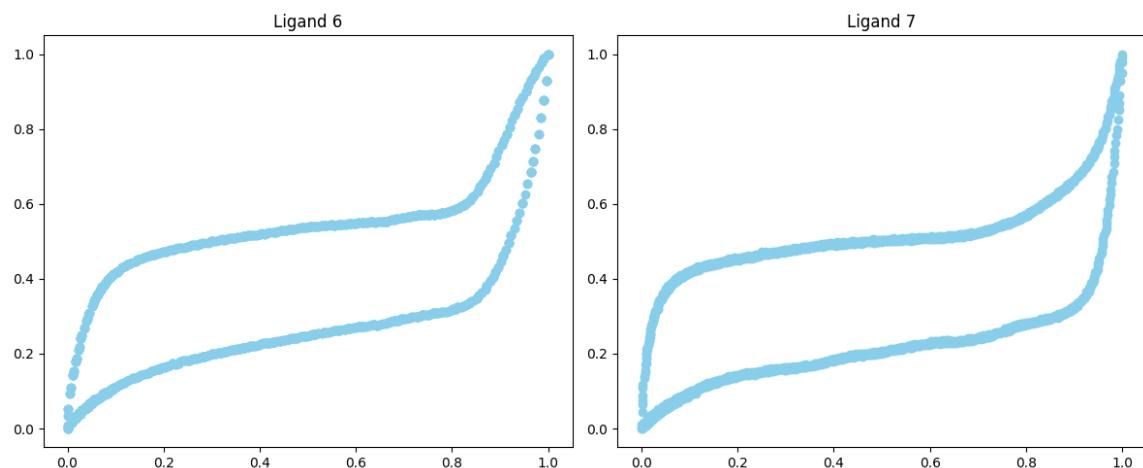


Figure 4.8: Ligand Six and Ligand Seven Voltammogram Comparison

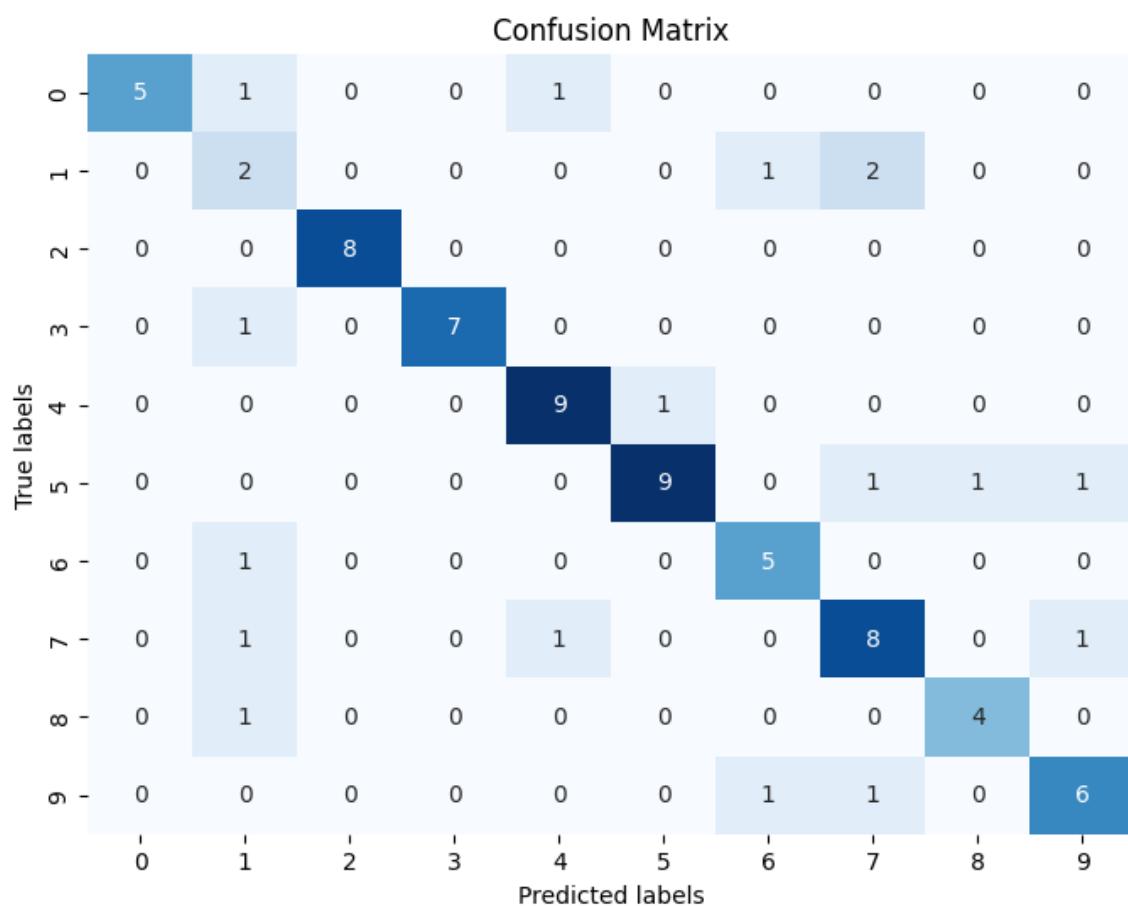


Figure 4.9: CV Metal Confusion Matrix

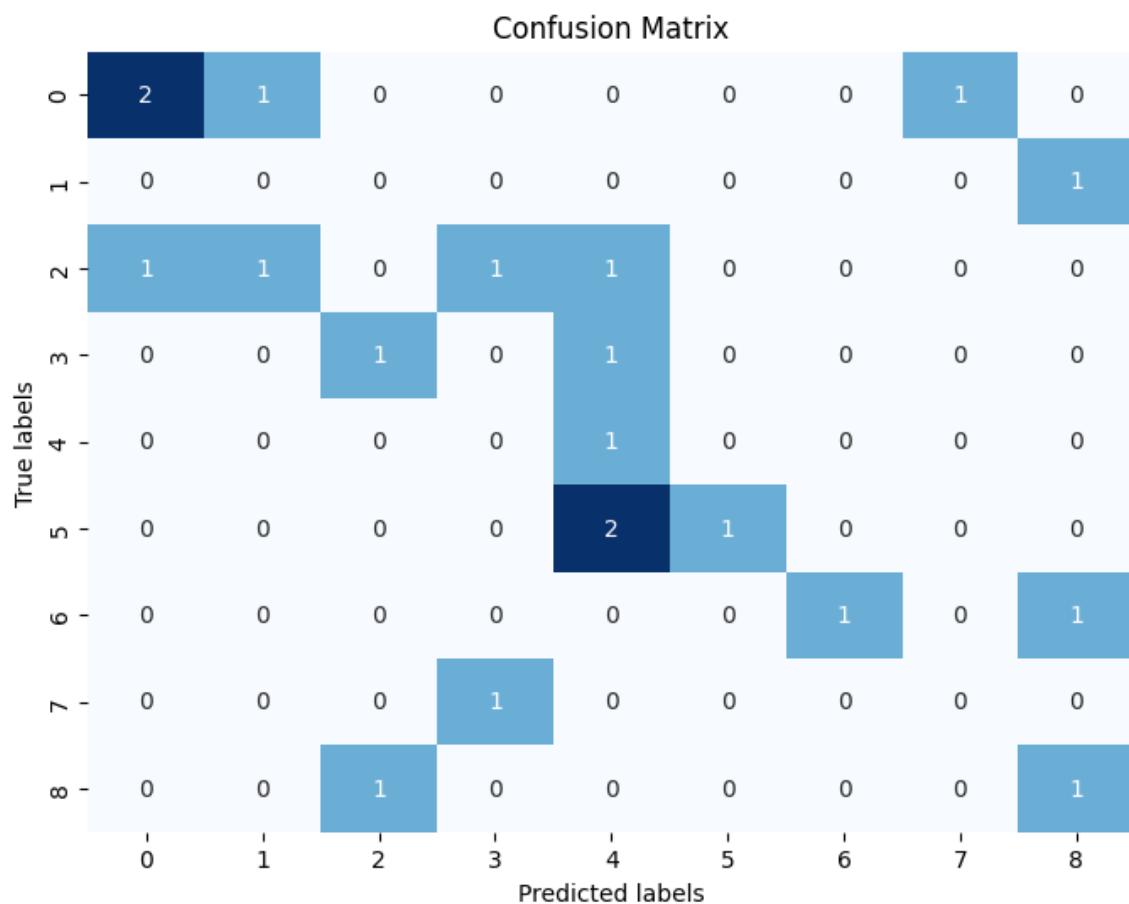


Figure 4.10: DPV Ligand Confusion Matrix

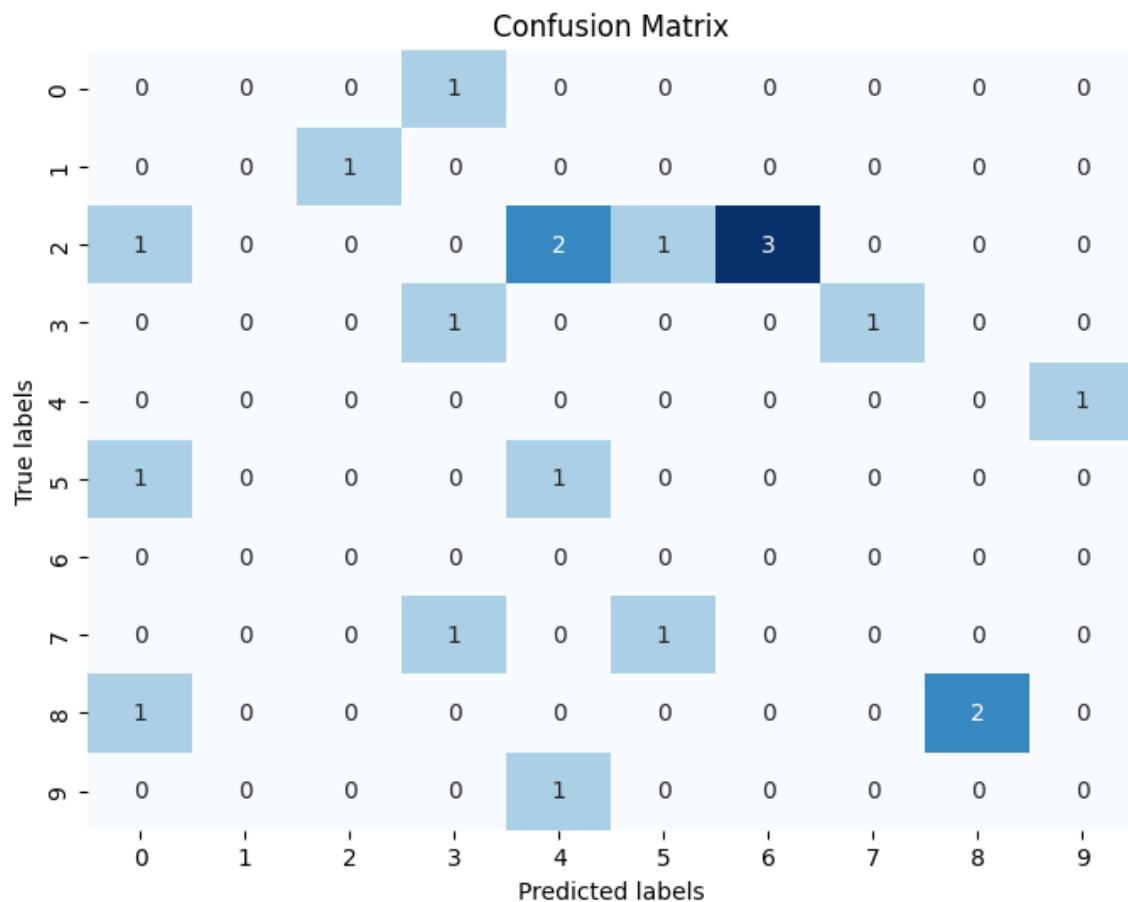


Figure 4.11: DPV Metal Confusion Matrix

Chapter 5

Denoising

5.1 Introduction

As the previously discussed dataset was generated using a low-cost potentiostat that lacks the accuracy of commercial options [11], we attempt to improve the quality of the data obtained by this potentiostat by applying ML to denoise the raw data with the commercial potentiostat data as a reference.

5.2 Autoencoder

As previously shown in Figure 4.1, an autoencoder is a neural network used to learn an efficient low-dimensional encoding of data. An autoencoder consists of an encoder and a decoder. The encoder transforms the input data into an encoded representation, and the decoder attempts to recreate the data from the encoded representation. Since the goal is to try and improve the data quality, the commercial potentiostat data is used for the decoder instead. This way, the low-cost potentiostat data creates an encoded representation, and an equivalent commercial potentiostat data is decoded. The main problem is how to pair results from the two potentiostats. While the metal and ligand used for each experiment are recorded, numerous other variables can influence the data. Therefore, the challenge revolves around accurately aligning

the outcomes generated by the low-cost potentiostat with their counterparts from the commercial one. This alignment is crucial for ensuring the reliability and validity of the encoded representations created by the autoencoder. Without precise pairing, the encoded representations may not accurately capture the underlying patterns in the data, leading to suboptimal performance of the autoencoder. The clustering technique previously described can be employed to address this challenge. Similar experimental results can be grouped by leveraging the recorded information on metals, ligands, and other relevant variables. The clustering process helps identify pairs of results with comparable characteristics, despite potential variations introduced by the different potentiostats.

5.3 Results and Discussion

In Figure 5.1, both the As seen in Figure 5.1, both the input and output are similar in overall shape. However, the output contains a much more defined duck-shaped voltammogram, which is typically expected. The results show promising outcomes and indicate that an autoencoder can effectively transform data from the low-cost potentiostat to resemble data from the commercial potentiostat. By leveraging the capacity of deep neural networks to learn complex patterns and relationships within the data, it becomes feasible to enhance the quality of measurements obtained from low-cost instruments, thereby expanding their utility in research and industrial applications. However, despite the promising results, several drawbacks and considerations must be acknowledged. Firstly, the effectiveness of the transformation heavily

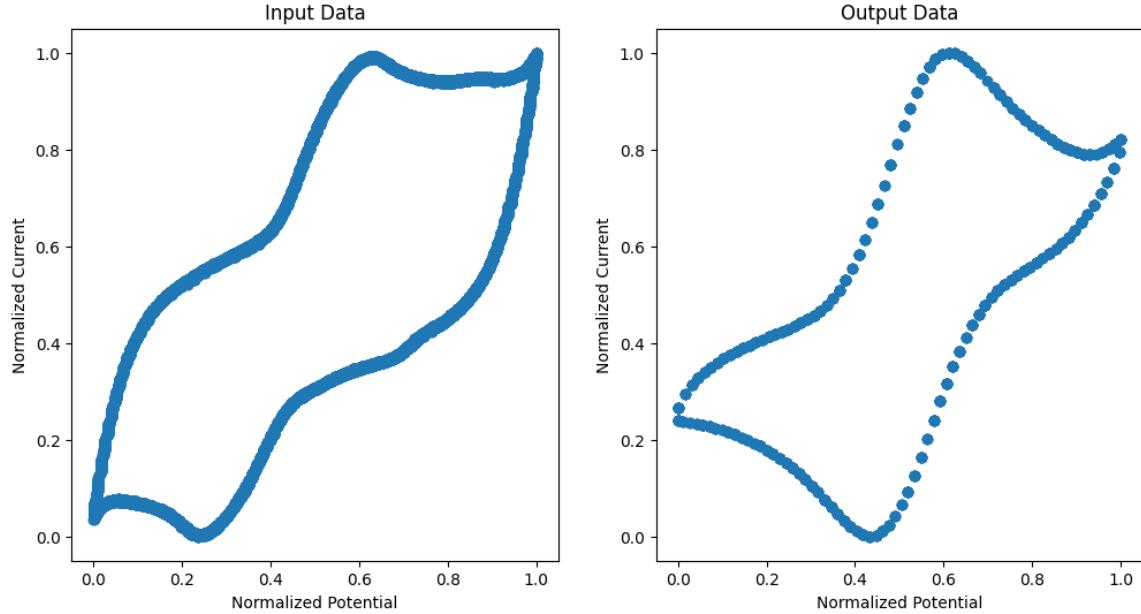


Figure 5.1: AutoEncoder Results

relies on the quality and diversity of the training data. Insufficient or biased training samples may lead to suboptimal performance and generalization issues, especially when dealing with complex electrochemical processes or diverse experimental conditions. While the autoencoder can effectively capture and replicate the dominant features present in the data, it may struggle with preserving subtle nuances or domain-specific characteristics inherent to the commercial potentiostat. Variations in hardware specifics, measurement protocols, or environmental factors could introduce discrepancies between the transformed and reference datasets. In conclusion, while autoencoders offer a promising avenue for enhancing the capabilities of low-cost potentiostats, their deployment must be accompanied by rigorous validation and consideration of the aforementioned limitations. Future research could focus on optimizing the autoencoder architecture, exploring alternative deep learning techniques,

Chapter 5. Denoising

and investigating strategies for addressing data heterogeneity to further improve the robustness and versatility of the proposed approach.

Chapter 6

Conclusion

In summary, the novel technique introduced for encoding CV and DPV data represents a pivotal advancement in the realm of SDLs. By effectively segmenting voltammograms according to their distinct characteristics and showcasing its effectiveness across a spectrum of machine learning applications, from clustering and classification to denoising and synthetic data generation, this technique signifies a significant step in improving the automation of custom low-cost devices in SDLs. Machine learning models able to precisely encode chemical data from characterization results may be used to enhance high-throughput operations by integrating multiple low-cost devices using the same trained model. This approach streamlines the adoption of SDL and HT setups and facilitates their integration into diverse research endeavours.

Looking ahead, there is a vast landscape for further exploration, particularly in investigating alternative curve simplification algorithms and seamlessly integrating the encoding technique into operational SDL frameworks. This approach not only promises to significantly enhance the efficiency and accuracy of SDL setups but also holds the potential to revolutionize access to such technologies. By significantly reducing the entry barriers for new research groups interested in embarking on SDL and high-throughput setups, this advancement opens the doors to a more inclusive and collaborative scientific landscape, sparking new possibilities and inspiring future research.

6.1 Data and Code Availability

All the relevant code can be found on [GitHub](#). The generated database containing the raw and processed CV and DPV measured results of all the measurements can be found in a Zenodo data repository (DOI:10.5281/zenodo.10633135) [11].

Bibliographic references

1. Tom, G. *et al.* Self-driving laboratories for chemistry and materials science. doi:[10.26434/chemrxiv-2024-rj946](https://doi.org/10.26434/chemrxiv-2024-rj946) (2024).
2. Hickman, R. *et al.* Atlas: a brain for self-driving laboratories. doi:[10.26434/chemrxiv-2023-8nrxx](https://doi.org/10.26434/chemrxiv-2023-8nrxx) (2023).
3. Strieth-Kalthoff, F. *et al.* Delocalized, asynchronous, closed-loop discovery of organic laser emitters. doi:[10.26434/chemrxiv-2023-wqp0d](https://doi.org/10.26434/chemrxiv-2023-wqp0d) (2023).
4. Lo, S. *et al.* Review of low-cost self-driving laboratories in chemistry and materials science: the “frugal twin” concept. *Digital discovery*. doi:[10.1039/d3dd00223c](https://doi.org/10.1039/d3dd00223c) (2024).
5. Bagotsky, V. S. *Fundamentals of electrochemistry* doi:[10.1002/047174199x](https://doi.org/10.1002/047174199x) (Wiley, 2005).
6. Electrode potential. doi:[10.1351/goldbook.E01956](https://doi.org/10.1351/goldbook.E01956) (2019).
7. Elgrishi, N., Rountree, K. J., McCarthy, B. D., Rountree, E. S., Eisenhart, T. T. & Dempsey, J. L. A practical beginner’s guide to cyclic voltammetry. *Journal of chemical education* **95**, 197–206. doi:[10.1021/acs.jchemed.7b00361](https://doi.org/10.1021/acs.jchemed.7b00361) (2018).
8. *Handbook of electrochemistry* (ed Zoski, C. G.) (Elsevier Science, London, England, 2006).
9. Yoshikawa, N., Akkoc, G. D., Pablo-García, S., Cao, Y., Hao, H. & Aspuru-Guzik, A. Does one need to polish electrodes in an eight pattern? automation provides the answer. doi:[10.26434/chemrxiv-2024-ttxnr](https://doi.org/10.26434/chemrxiv-2024-ttxnr) (2024).
10. Wain, A. J. & Dickinson, E. J. in *Nanoscale electrochemistry* (eds Wain, A. J. & Dickinson, E. J.) 1–48 (Elsevier, 2021). doi:<https://doi.org/10.1016/B978-0-12-820055-1.00008-3>.
11. Pablo-García, S. *et al.* An affordable platform for automated synthesis and electrochemical characterization. doi:[10.26434/chemrxiv-2024-cwnwc](https://doi.org/10.26434/chemrxiv-2024-cwnwc) (2024).

Bibliographic references

12. Nicholson, R. S. & Shain, I. Theory of stationary electrode polarography. single scan and cyclic methods applied to reversible, irreversible, and kinetic systems. *Analytical chemistry* **36**, 706–723. doi:[10.1021/ac60210a007](https://doi.org/10.1021/ac60210a007). eprint: <https://doi.org/10.1021/ac60210a007> (1964).
13. Heinze, J. Cyclic voltammetry—“electrochemical spectroscopy”. new analytical methods (25). *Angewandte chemie international edition in english* **23**, 831–847. doi:<https://doi.org/10.1002/anie.198408313>. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/anie.198408313> (1984).
14. Nicholson, R. S. & Shain, I. Theory of stationary electrode polarography. single scan and cyclic methods applied to reversible, irreversible, and kinetic systems. *Analytical chemistry* **36**, 706–723. doi:[10.1021/ac60210a007](https://doi.org/10.1021/ac60210a007) (1964).
15. Libretexts. *Cyclic voltammetry* 2023.
16. Grimshaw, J. in *Electrochemical reactions and mechanisms in organic chemistry* (ed Grimshaw, J.) 1–26 (Elsevier Science B.V., Amsterdam, 2000). doi:<https://doi.org/10.1016/B978-044472007-8/50001-X>.
17. Faulkner, L. R. Understanding electrochemistry: some distinctive concepts. *Journal of chemical education* **60**, 262. doi:[10.1021/ed060p262](https://doi.org/10.1021/ed060p262) (1983).
18. *Electroanalytical methods* 1st ed. en (ed Scholz, F.) (Springer, Berlin, Germany, 2005).
19. Laborda, E., González, J. & Molina, Á. Recent advances on the theory of pulse techniques: a mini review. *Electrochemistry communications* **43**, 25–30. doi:[10.1016/j.elecom.2014.03.004](https://doi.org/10.1016/j.elecom.2014.03.004) (2014).
20. Bellman, R. *Dynamic programming* (Dover Publications, Mineola, NY, 2003).
21. MacQueen, J. B. *Some methods for classification and analysis of multivariate observations* in *Proc. of the fifth berkeley symposium on mathematical statistics and probability* (eds Cam, L. M. L. & Neyman, J.) **1** (University of California Press, 1967), 281–297.
22. Thorndike, R. L. Who belongs in the family? *Psychometrika* **18**, 267–276. doi:[10.1007/bf02289263](https://doi.org/10.1007/bf02289263) (1953).

Bibliographic references

23. Ester, M., Kriegel, H.-P., Sander, J. & Xu, X. *A density-based algorithm for discovering clusters in large spatial databases with noise* in *Knowledge discovery and data mining* (1996).
24. Van der Maaten, L. & Hinton, G. Visualizing data using t-SNE. *Journal of machine learning research* **9**, 2579–2605 (2008).
25. McInnes, L., Healy, J. & Melville, J. *Umap: uniform manifold approximation and projection for dimension reduction* 2018. doi:[10.48550/ARXIV.1802.03426](https://doi.org/10.48550/ARXIV.1802.03426).
26. Hastie, T., Tibshirani, R. & Friedman, J. *The elements of statistical learning: data mining, inference and prediction* 2nd ed. (Springer, 2009).
27. Sharma, N., Jain, V. & Mishra, A. An analysis of convolutional neural networks for image classification. *Procedia computer science* **132**. International Conference on Computational Intelligence and Data Science, 377–384. doi:<https://doi.org/10.1016/j.procs.2018.05.198> (2018).
28. Pinheiro Cinelli, L., Araújo Marins, M., Barros da Silva, E. A. & Lima Netto, S. in *Variational methods for machine learning with applications to deep networks* 111–149 (Springer International Publishing, 2021). doi:[10.1007/978-3-030-70679-1_5](https://doi.org/10.1007/978-3-030-70679-1_5).
29. Hoar, B. B., Zhang, W., Xu, S., Deeba, R., Costentin, C., Gu, Q. & Liu, C. Electrochemical mechanistic analysis from cyclic voltammograms based on deep learning. *Acs measurement science au* **2**, 595–604. doi:[10.1021/acsmeasurescäu.2c00045](https://doi.org/10.1021/acsmeasurescäu.2c00045) (2022).
30. Faraday, M. S. (B. *On electro-chemical decomposition* 1970.
31. Dirac, P. A. M. *The principles of quantum mechanics* (Clarendon Press, 1981).
32. Knuth, D. *Knuth: computers and typesetting* <https://www-cs-faculty.stanford.edu/~knuth/abcde.html>.
33. Knuth, D. E. in. Chap. 1.2 (Addison-Wesley, 1973).
34. J., A. & Faulkner, L. R. *Student solutions manual to accompany electrochemical methods: fundamentals and applicaitons*, 2e en (John Wiley & Sons, Nashville, TN, 2002).

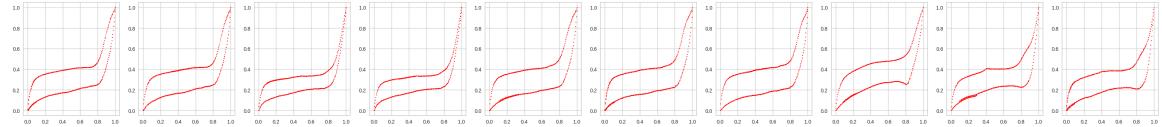
Bibliographic references

35. Rousseeuw, P. J. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics* **20**, 53–65. doi:[https://doi.org/10.1016/0377-0427\(87\)90125-7](https://doi.org/10.1016/0377-0427(87)90125-7) (1987).
36. Kramer, M. Autoassociative neural networks. *Computers & chemical engineering* **16**, 313–328. doi:[10.1016/0098-1354\(92\)80051-a](https://doi.org/10.1016/0098-1354(92)80051-a) (1992).

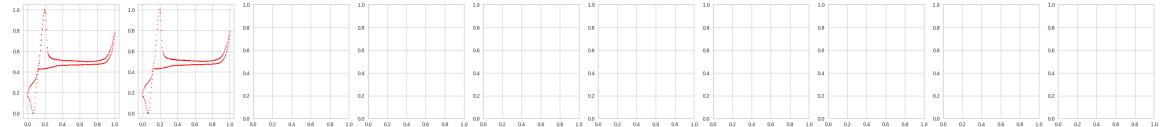
Appendix A

A.1 CV K-Means Cluster Results

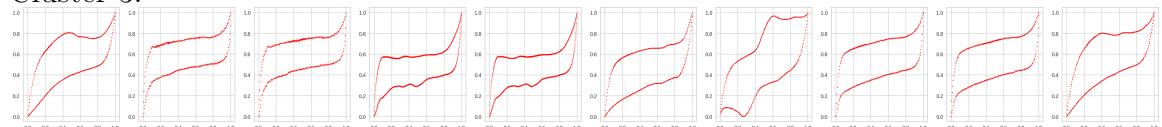
Cluster 1:



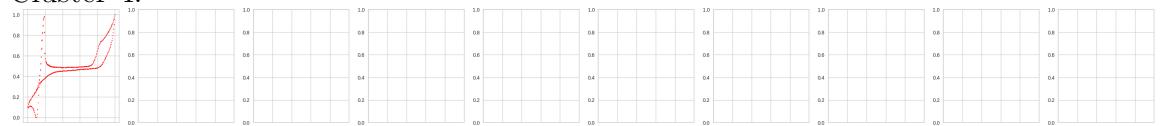
Cluster 2:



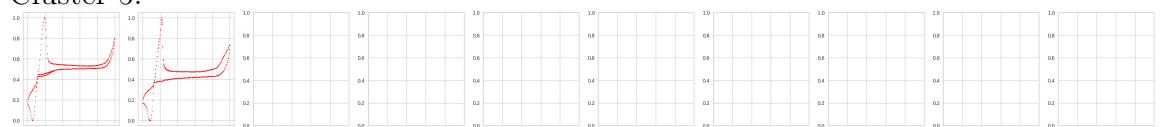
Cluster 3:



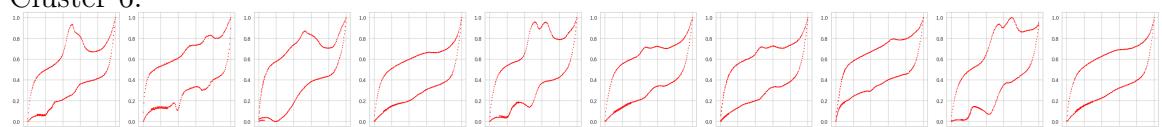
Cluster 4:



Cluster 5:

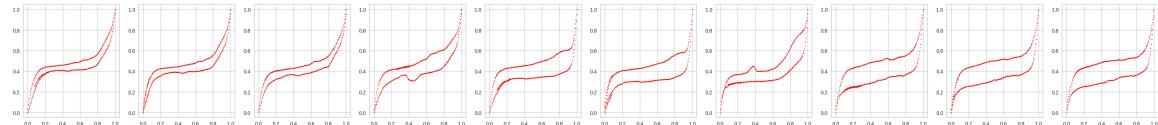


Cluster 6:

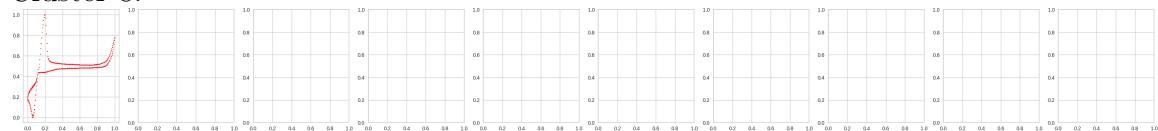


Cluster 7:

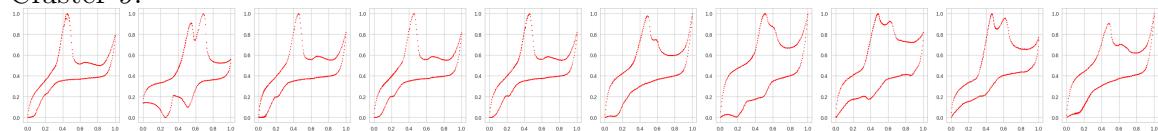
Appendix A.



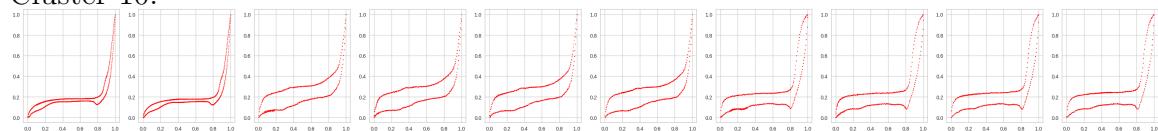
Cluster 8:



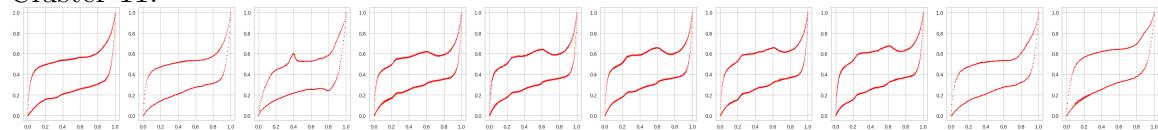
Cluster 9:



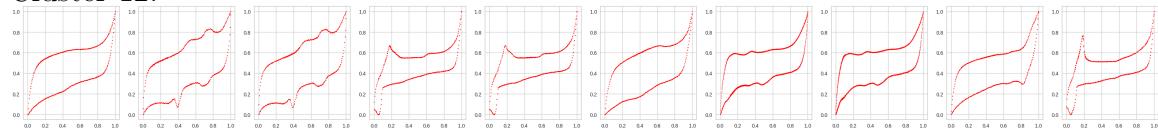
Cluster 10:



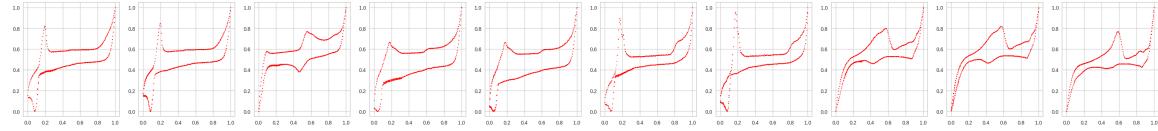
Cluster 11:



Cluster 12:



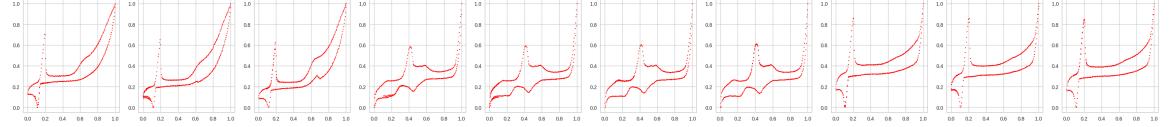
Cluster 13:



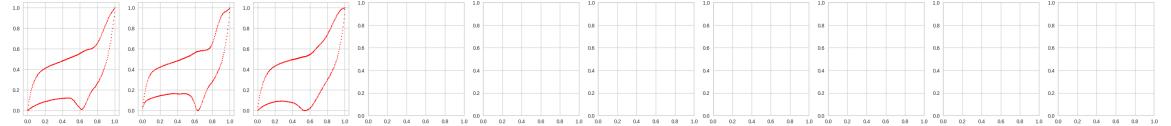
Cluster 14:

Appendix A.

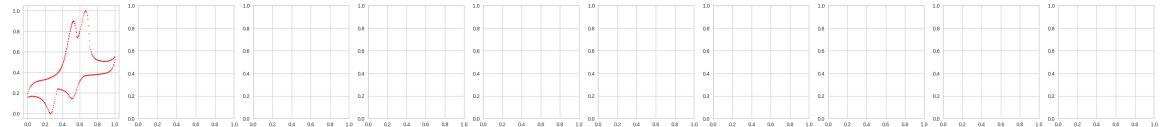
Cluster 15:



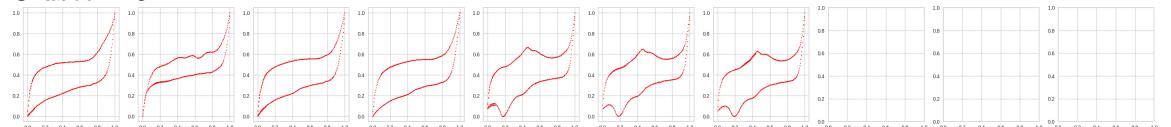
Cluster 16:



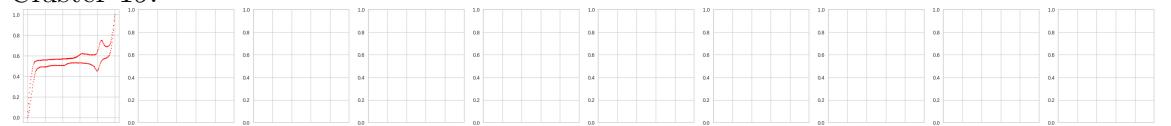
Cluster 17:



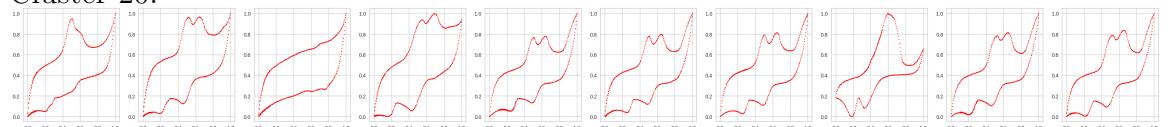
Cluster 18:



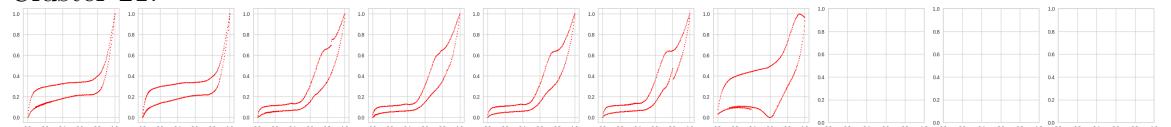
Cluster 19:



Cluster 20:

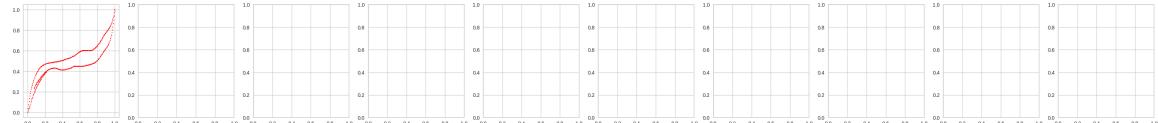


Cluster 21:

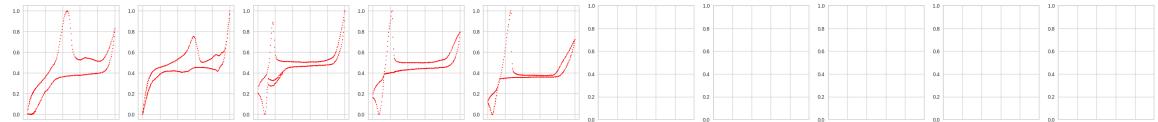


Cluster 22:

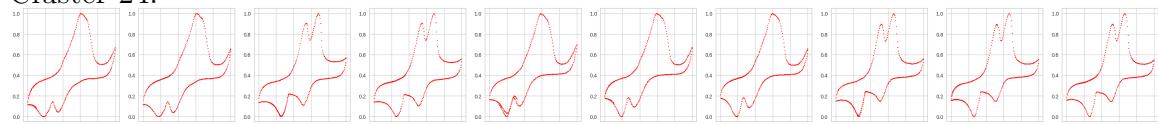
Appendix A.



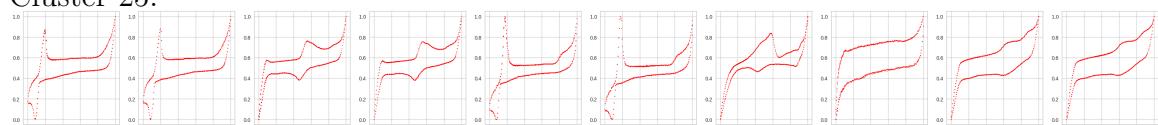
Cluster 23:



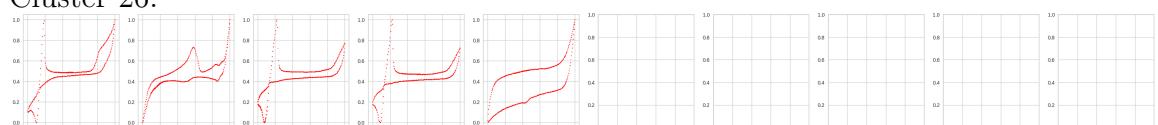
Cluster 24:



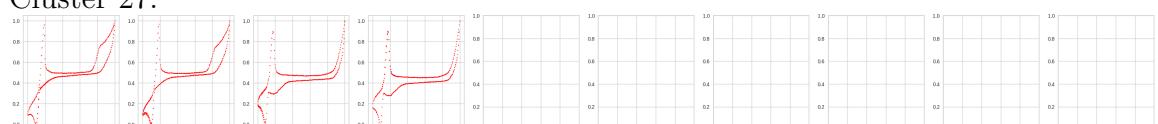
Cluster 25:



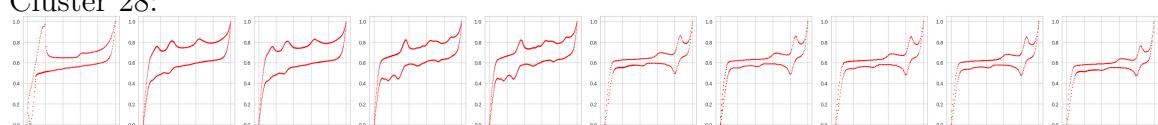
Cluster 26:



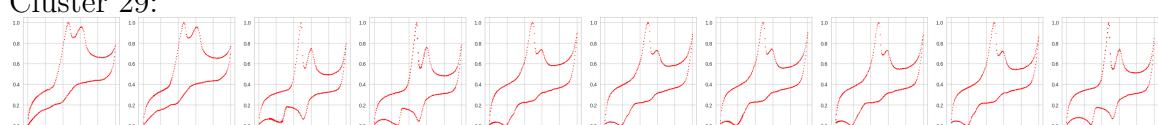
Cluster 27:



Cluster 28:

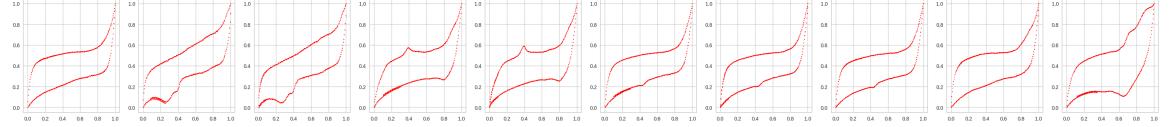


Cluster 29:

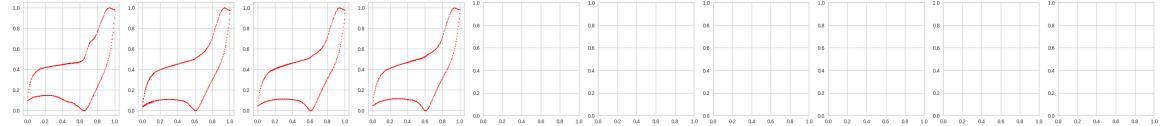


Appendix A.

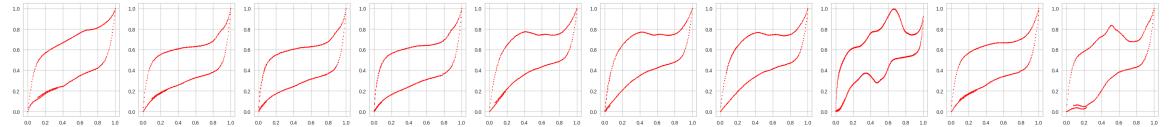
Cluster 30:



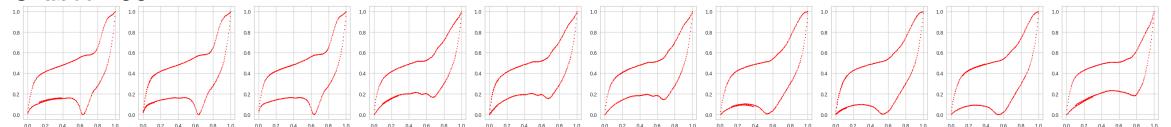
Cluster 31:



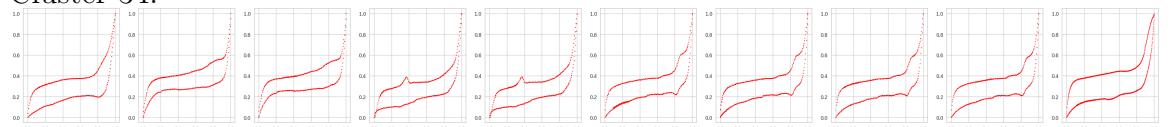
Cluster 32:



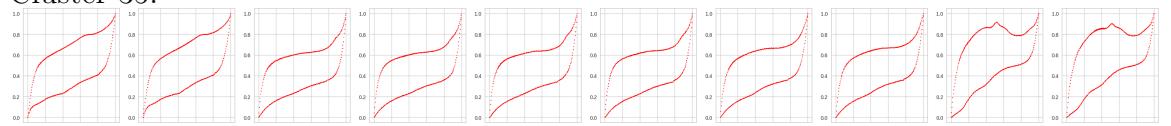
Cluster 33:



Cluster 34:



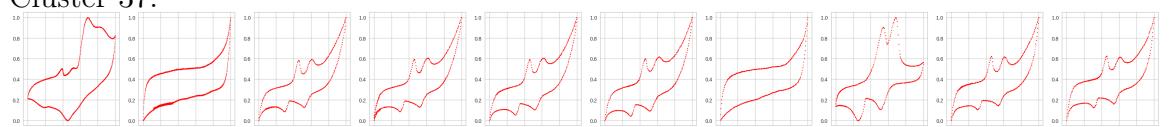
Cluster 35:



Cluster 36:

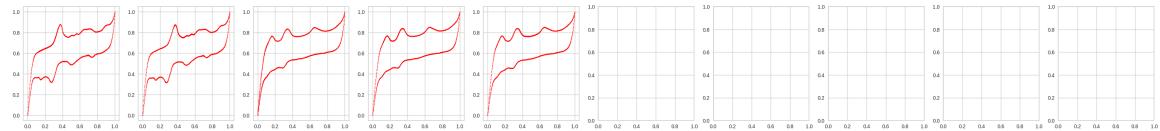


Cluster 37:



Appendix A.

Cluster 38:



A.2 Metals and Ligands

Entry	Metal	Form	CAS
1	V(IV)	VOSO ₄ xH ₂ O	123334-20-3
2	Cr(III)	CrK(SO ₄) ₂ 12H ₂ O	7788-99-0
3	Mn(II)	MnSO ₄ H ₂ O	10034-96-5
4	Fe(II)	FeSO ₄ 7H ₂ O	7782-63-0
5	Co(II)	CoSO ₄ 7H ₂ O	10026-24-1
6	Ni(II)	NiSO ₄ 6H ₂ O	10101-97-0
7	Cu(II)	CuSO ₄ 5H ₂ O	7758-99-8
8	Zn(II)	ZnSO ₄ 7H ₂ O	7446-20-0
9	Cd(II)	CdSO ₄ 8/3H ₂ O	7790-84-3
10	Pd(II)	Na ₂ PdCl ₄	13820-53-6

Table A.1: Table of Metals

Appendix A.

Entry	Ligand	SMILES	CAS
1	ammonia	N	1336-21-6
2	hydrazine	NN	7803-57-8
3	ethylenediamine	NCCN	107-15-3
4	ethanolamine	NCCO	141-43-5
5	diethanolamine	OCCNCCO	111-42-4
6	triethanolamine	OCCN(CCO)CCO	102-71-6
7	piperidine	N1CCCCC1	110-89-4
8	morpholine	N1CCOCC1	110-91-8
9	pyridine	n1ccccc1	110-86-1
10	2,2'-bipyridine (in HCl salt form)	c1ccc(nc1)c2cccn2	336-18-7

Table A.2: Table of Ligands

