# AN OVERVIEW OF ADVANCED TECHNOLOGIES APPLIED TO IDENTIFIED PRINTED AND HANDWRITTEN TEXT IN GURMUKHI SCRIPT: A REVIEW

**Rajinder Kaur[1], Anantdeep[2]**

*Department of Computer Science and Engineering, Punjabi University, Patiala[1,2]*
rkkhinda1999@gmail.com[1], *Boparai.dolly@gmail.com[2].*

## ABSTRACT

Handwriting recognition of Gurmukhi is an area of research which has gained traction in recent years. This paper reviews the state of the art of handwriting recognition for Gurmukhi. A brief overview of the various approaches used for Gurmukhi handwriting recognition, their benefits and drawbacks, the current challenges and the future research directions is presented. Specifically, the review will focus on the recent advances in deep learning models, feature extraction techniques, and preprocessing strategies. It further examines the challenges and potential areas of improvement in the field of Gurmukhi handwriting recognition. Finally, it provides insights into the potential applications of this technology and the opportunities for future research in the area. The paper concludes by highlighting the importance of handwriting recognition for Gurmukhi and the need for further research and development in the field.
**Keywords:** Feature extraction, CNN, Holistic approach, MLP, RDP algorithm

## INTRODUCTION

Handwriting is the individual and unique way a person forms letters and words when writing. It is often considered to be an expression of one's personality, as it can reveal details about the writer's style, character, and even emotional state. Artificial Intelligence (AI) is a powerful tool that can be used to analyse and interpret data in ways that would not be possible with manual methods. One area where AI has made significant strides is handwriting recognition, which involves using computer algorithms to identify patterns in handwriting and convert it into digital text. Handwriting recognition has become increasingly important in recent years due to advancements in technology that allow computers to read handwritten text more accurately than ever before. It is being used to help automate processes and make them more efficient. Handwriting recognition can be used

in a variety of ways, from helping to fill out forms faster, to reducing the amount of manual data entry required by businesses. Additionally, handwriting recognition can also be used for security purposes by verifying the authenticity of signatures or documents. With advances in artificial intelligence and machine learning algorithms, handwriting recognition accuracy has improved significantly over time and it continues to evolve with new technologies such as deep learning being developed every day.

The challenge faced in handwriting recognition is recognising the text from different scripts, as each script has its own unique characteristics. Gurmukhi script, which is used to write Punjabi language, differs significantly from English script, making it more difficult for computers to recognize it accurately. Gurmukhi consists of 35 characters written in two forms - vertical (as seen in Hindi) and horizontal (as seen in English). It also contains special characters such as punctuation marks and symbols. Unlike

in English, where each letter has a distinct shape and is usually written independently of other letters, Gurmukhi characters are often written together or adjacent to one another. This makes it difficult for traditional handwriting recognition algorithms to identify individual characters accurately. Additionally, there can be significant variation in how a given character is written by different people due to regional dialects and personal preferences. Furthermore, due to its complex writing system which combines multiple scripts (Devanagari & Sharada), recognising certain letters can be particularly challenging. To address these issues, researchers have developed various approaches that leverage deep learning techniques to improve accuracy in Gurmukhi handwriting recognition systems. These include using convolutional neural networks (CNNs) with custom-built feature extractors for extracting features from handwritten samples; recurrent neural networks (RNNs) for modelling temporal dynamics of text sequences; and reinforcement learning algorithms. By combining advanced machine learning methods with pre-processing steps such as character segmentation and noise reduction techniques, state-of-the-art results have been achieved on benchmark datasets. This demonstrates the effectiveness of these approaches in improving the accuracy of Gurmukhi handwritten text recognition systems when compared against traditional techniques used for English script recognition tasks.

## MOTIVATION

Handwritten Gurmukhi recognition is an important and challenging task that has gained a lot of attention in the field of pattern recognition. This research topic has been motivated by the need for efficient methods for automating document processing tasks such as handwriting recognition, optical character recognition (OCR), and other related tasks.

Gurmukhi script consists of 39 consonants and 11 vowels with various combinations, making it difficult to recognize due to its complex structure. Moreover, there are several variations in handwriting styles among different writers which further increases the complexity of the problem. As a result, automated techniques are required to accurately identify these characters from scanned documents or images. In recent years, numerous approaches have been proposed using machine learning algorithms such as Support Vector Machines (SVMs) and Deep Neural Networks (DNNs). However, existing approaches still face challenges due to lack of sufficient training data sets with large variety of fonts/styles available publicly or limited availability thereof.

Handwritten Gurmukhi recognition has the potential to provide a number of benefits and services to local language populations. It can help in preserving, digitizing, and sharing historical documents written in this script, providing access to information that would otherwise be lost or inaccessible. Additionally, deploying efficient and reliable systems for recognizing multilingual documents can allow organizations to process forms more quickly and accurately with minimal manual labor. Furthermore, such technologies can enable governments and other institutions to provide better services tailored specifically for their local language population by developing products like online databases or search engines for specific languages. Therefore, handwritten Gurmukhi recognition is an important technology that could have significant implications on how we interact with our past as well as how we serve our present needs.

## CHALLENGES

Pre-processing of Gurmukhi Handwriting: The pre-processing of Gurmukhi handwriting is one of the main challenges in recognition. This involves segmentation, normalization and feature extraction techniques that are specific to

this script. Segmentation is a process separate individual characters from each other, while normalisation helps reduce variability between different writers by making sure all characters have similar size, orientation and shape. Feature extraction is also important as it helps extract relevant features from the input data which can be used for classification purposes.

Feature Representation: Feature representation plays an important role in any handwriting recognition system since it determines how accurately a character can be recognised. In Gurmukhi handwriting recognition systems, various types of feature representations such as stroke-based features or structural features need to be explored for better performance results. Each type has its own advantages and disadvantages so finding the best combination for optimal performance requires careful consideration and experimentation with different approaches.

Ambiguity Between Similar Looking Letters: Another challenge faced during recognition process is ambiguity between similar looking letters especially if they appear close together such as ਗ vs ਦ or ਸ vs ਜ. Such cases require extra processing steps like context analysis or post-processing rules applied after initial detection stage so as accuracy can be improved significantly.

Lack of large dataset: Gurmukhi handwriting recognition is a challenging task due to the lack of available datasets. While there are some existing datasets, they tend to be small and limited in scope. This makes it difficult for researchers to develop robust algorithms that can accurately recognize Gurmukhi script from handwritten documents. To address this challenge, larger and more diverse datasets need to be created which include multiple writers, different writing styles, as well as different levels of complexity (e.g., cursive vs non-cursive). Additionally, these datasets should also contain various types of text (e.g.,

short stories or poems) so that the models can learn how to handle a variety of input data.

Training Data Collection: Gathering training data is another major challenge in developing a successful Gurmukhi handwritten text recognition system. This is due to limited resources available online or offline related to this script's writing style, compared to other languages like English or Hindi. To overcome this challenge, one possible solution is to create a dataset of handwritten Gurmukhi text using crowdsourcing tools such as Amazon Mechanical Turk. This would allow for the collection of large amounts of data from people with different writing styles and backgrounds, providing a comprehensive set of training data. Additionally, it may also be helpful to partner with local schools or universities in order to collect more samples from students who are native speakers and writers of Gurmukhi.

Limited Research Work: A major challenge faced by researchers working on developing robust handwriting recognition systems for Punjabi/Gurmukhi script is the intricate nature of these scripts compared to simpler languages like English. This makes it difficult to accurately recognize them and there has been a lack of sufficient research dedicated towards understanding the complexities involved. Furthermore, due to its complexity, obtaining enough data samples for training purposes can be challenging as well. Few efforts have been made to explore deep learning architectures and recognize Gurmukhi scripts precisely. However, much more needs to be done in this field. Active collaborations between academia and industry players would be beneficial for leveraging existing technologies being developed around the world into something useful.

Multi Script Environment: Last but not least, another issue faced when trying recognize Gurmukhi text is presence multiple coexisting

scripts alongside each other i.e. Hindi / Sanskrit / Urdu etc. In India, this type of scenario is quite common which leads to potential confusion among algorithms trained only to recognize one particular script. This results in errors during the decoding phase.

## TERMS USED

### 4.1 Convolutional Neural Network

A Convolution Neural Network (CNN) is a type of deep learning network used for image processing and recognition. It consists of several layers that make up the structure of the neural network, which can be used to process Gurmukhi handwriting recognition.

The first layer in a CNN is the convolutional layer which performs feature extraction on an input image. This layer applies a filter over the image, taking small portions from it and applying mathematical operations such as multiplication or addition to extract features from it. The output of this operation is then passed onto subsequent layers where further processing takes place until an output is produced by the network.

The next layer in a CNN is usually one or more pooling layers which reduce computational complexity by down sampling data while preserving important features extracted during convolution. Max-pooling and average-pooling are two common types of pooling layers used in CNNs. These layers are followed by fully connected (FC) layers, also called dense layers, which connect all neurons present in previous parts of the neural networks together forming what's known as a multilayer perceptron (MLP). The FC/dense layer typically contains multiple neurons that have weights associated with them and act as decision makers based on certain inputs given to them through training data sets fed into these neurons beforehand.

Finally, after passing through all these different steps, we reach our last step: classification using SoftMax function or sigmoid activation function depending upon our task requirement (binary classification vs multi class classification).s

SoftMax produces normalized probability scores for each class label whereas sigmoid outputs values between 0 & 1 indicating how likely an instance belongs to one particular class label among many others available at our disposal (e.g. - Gurmukhi letter "ਕ" vs other letters like "ਗ", "ਖ" etc.)

Once trained properly with training samples consisting handwritten Gurmukhi characters, our model should be able to classify new instances correctly with high accuracy even if they differ slightly from already seen examples!

### 4.2 Holistic approach v/s Analytical approach

Gurmukhi is a script used to write the Punjabi language. Handwriting recognition of Gurmukhi is an important task as it can help in digitizing documents and converting them into machine-readable formats. There are two approaches that can be used for this purpose - analytical approach and holistic approach.

In an analytical approach, the handwriting recognition system analyses each character separately, often using handcrafted features such as stroke count or aspect ratio of characters. This type of feature engineering process requires significant expertise from domain experts and may not work well when dealing with complex scripts like Gurmukhi which has many different writing styles across regions and writers. Furthermore, variations in handwriting styles due to age, gender or educational background cannot be addressed by this method since these factors are not considered during feature engineering phase.

On the other hand, a holistic approach uses deep learning techniques such as convolutional neural networks (CNNs) which take raw image data as input without any prior knowledge about the underlying script structure or writer's

style characteristics. These models learn patterns directly from images instead of relying on pre-defined features; thereby allowing them to capture subtle differences between various writing styles more accurately compared to analytics-based methods.

For example, consider two Gurmukhi words: ਸਤਿ and ਸੋਧੀ. An analytical approach would try to recognize them based on their separate components (i.e., ਸ + ਤ + ਿ and ਸ + ੋ + ਧ + ੀ). This approach may fail if the handwriting is not clear or if the characters are written in different sizes, shapes and styles. On the other hand, a holistic approach would take into account the entire context of how these words are written, such as their size and spacing relative to each other. This allows for more accurate recognition even when individual characters have variations in shape and style.

Holistic methods are faster than analytical approaches, as they do not require as much time to analyze individual elements before making decisions. This allows results to be obtained quickly and efficiently without sacrificing quality control measures such as double-checking data against known samples. Furthermore, holistic methods are more accurate since they take into account the whole context of a given piece of writing instead of relying on isolated elements in isolation. This makes them particularly well suited for applications where speed is important (e.g., automatic translation services).

## 4.3 Feature Extraction

Gurmukhi is an abugida script used to write Punjabi language. Handwriting recognition of Gurmukhi requires feature extraction from the handwriting samples. Feature extraction involves deriving useful information from the input data, which can be used as a basis for further processing. Feature extraction involves capturing the unique features of each character or symbol by analyzing its shape, form and

structure. The main aim of feature extraction is to reduce dimensionality while preserving the essential properties of an image that are necessary for accurate character identification or recognition.

### 4.3.1 Zoning

Zoning is a technique where each stroke or character is
divided into smaller regions called "zones". Zoning helps identify strokes by their size and shape so they can be classified more easily. This technique allows for better segmentation of characters from one another which makes them easier to recognize individually.

### 4.3.2 Transition

Transition is another technique used in handwriting recognition systems. It involves extracting information about how individual letters connect with each other while forming words or sentences. This is done by analysing their transitions patterns when using pen-based input devices such as tablets and styluses. Transition analysis helps to recognize characters or words written with a pen more accurately, compared to simply recognising the individual strokes of the letters separately.

### 4.3.3 Diagonal

Diagonal helps us identify diagonal lines within images; this allows us detect slanted writing which may otherwise not be detectable using zoning alone since some characters like "ਅ" do not always appear perfectly horizontal when written by hand but rather slant slightly downwards towards their right side depending on the writer's style/habit (this can also happen with other languages too!).

### 4.3.4 Peek extent

Peak extent is one such characteristic which helps recognize patterns within a letter or word. In simple terms, peak extent refers to the area around each character where its contours reach their highest point before descending again. This high point forms the 'peak' which is then followed by downward slopes on either side as it moves away from the center of the character. By measuring these peaks, feature extraction algorithms are able to detect subtle differences between characters even when they appear very similar at first glance.

### 4.3.5 Intersection and Open-End Point based feature extraction

This method works by detecting the intersection points between two different strokes in a character. For example, in Gurmukhi script, when two lines of a letter come together at an angle there is an intersection point which can be used for feature extraction. Similarly, open end points can also be detected where one line ends before another begins. Both these features are then used to classify the characters into their respective classes.

### 4.3.6 Curve-Based Feature Extraction

This technique relies on extracting curvature information from each Gurmukhi character by measuring its length and direction. Curves can be identified by examining the shape of the curves in each character; for instance, ਜ(ja) has an upward curve while ਫ(pha) has a downward curve. The curvature information can then be used to differentiate different characters with similar shapes such as ਝ(jha) and ਭ(bha).

### 4.4 Artificial neural network (ANN)

An Artificial Neural Network (ANN) is a type of computational model loosely based on the structure and functions of biological neural networks. It is composed of nodes arranged in layers, where each node contains an activation function that determines how it responds to input signals from its neighbouring nodes. This network can be trained to recognize patterns in data which are then used for various tasks such as handwriting recognition or image classification.

In Gurmukhi handwriting recognition, ANNs are used to identify characters written in the language by analysing their shape and orientation on a page. The ANN begins by extracting features from an image of Gurmukhi text, such as line width or curvature. These features are then fed into the network which uses them to classify different symbols according to their shapes and orientations. Once this process has been completed, the output is compared against known examples of Gurmukhi characters so that any unrecognised symbols can be identified correctly.

The advantage of using artificial neural networks for handwriting recognition lies in its ability to learn from mistakes; if character is misclassified during training, it will adjust itself accordingly until it achieves better accuracy over time with more data being provided for training purposes.

The main advantages of using ANNs for handwriting recognition of Gurmukhi include low complexity, low parameter requirements, and low learning and reprogramming requirements. In other words, they can quickly process data without needing a large amount of computing power or time spent training them before they can be deployed in production environments. This makes them ideal for applications where rapid results are needed but resources may be limited.

### 4.5 Multilayered perceptron model

A multilayered perceptron model is a type of artificial neural network which consists of multiple layers of interconnected nodes. Each layer contains an array of neurons that are connected to the input and output layers. The

first layer receives the inputs from the environment, while subsequent layers receive inputs from previous ones in order to generate outputs. These outputs can be used for classification or prediction tasks. The main advantage of using a multilayered perceptron model is its full connectivity, meaning each neuron in one layer is connected to all neurons in the next layer. This allows for more complex interactions between different levels, making it possible to learn nonlinear functions such as those found in handwriting recognition tasks like Gurmukhi characters recognition. Furthermore, by increasing the number of hidden layers, this type of network structure can learn more intricate patterns. This is because it increases its representational power exponentially with each additional layer added on top. As a result, these networks have better generalisation capabilities and improved fit when compared to linear models which are limited by their linear separability criteria. MLPs are well-suited for applications such as handwriting recognition, where there may not be a clear separation between character classes (such as Gurmukhi characters). This is because MLPs can handle high-dimensional data with ease and accuracy, even if these features cannot be linearly separated into distinct classes. Through training processes such as back propagation algorithms, weights associated with various connections between nodes can be adjusted. This adjustment of the weights allows for decision boundaries to become increasingly accurate over time. As a result, better generalisation performance across unseen test cases is achieved. In summary, Multilayer Perceptrons offer great potential for solving complex problems involving pattern detection and classification using only simple mathematical operations - making them ideal candidates when dealing with datasets containing nonlinearly separable information.

### 4.6 Random forest classifier

Random Forest Classifier is a supervised machine learning algorithm which uses ensemble learning to classify objects. It works by constructing a multitude of decision trees at training time and outputting the class that is most common or has the highest mean probability across all of them.

The way it works is that it randomly selects data points from the dataset, then builds multiple decision trees using those selected data points as root nodes. The randomness in selection helps reduce overfitting, since different subsets of data will lead to different tree structures. Each tree then makes its own prediction based on what it learned from the training set and this forms an "ensemble" of predictions which are combined together to form an overall prediction with better accuracy than any single tree could achieve on its own.

Random Forest Classifiers work best when dealing with large datasets where there are many features but only a few classes (i.e., binary classification). They also tend to perform well in cases may be non-linear relationships between predictors and outcomes, making them ideal for complex p problems like handwriting recognition of Gurmukhi characters or other languages containing many intricate shapes such as Chinese or Japanese characters. Random forest classifiers are able to handle high dimensional input spaces without sacrificing too much accuracy. This is due to their ability to select important features while discarding irrelevant ones through feature selection processes such as recursive feature elimination (RFE). These processes allow the algorithm to identify which features are most relevant and reduce the dimensionality of the data, thus improving accuracy while still maintaining a high-dimensional model.

### 4.7 Support Vector Machine

Support Vector Machine (SVM) is a supervised machine learning algorithm used for classification

and regression problems. It works by finding the best separating hyperplane between classes of data in a high-dimensional space. This concept can be extended to more than two dimensions, which makes it an effective tool for solving complex problems with multiple variables. SVM has been applied to many fields such as image recognition, text categorisation and bioinformatics.

In handwriting recognition of gurmukhi, SVM is used to identify patterns that exist in different fonts or styles of writing within the same language. The algorithm looks at features like line thickness, shape orientation and letter size in order to classify each character accurately. During training, these features are fed into the model so that it can learn how they relate to each other and build up a library of known characters that it can use when presented with unknown inputs during testing or real-world applications.

The main advantage of using SVM for handwriting recognition tasks is its ability to generalise well even when there are large amounts of variation present in datasets due to differences between writers' styles or font types being used. Additionally, SVMs have relatively low computational complexity compared with other models such as neural networks making them suitable for real time processing on small devices like smartphones where power consumption needs to be kept low.

Overall, Support Vector Machines offer an accurate solution for recognising Gurmukhi script from digital images while also providing good scalability across various platforms making them ideal candidates for tackling this problem efficiently in many scenarios.

## 4.8 k-Nearest Neighbor

The k-NN algorithm works by taking an input character and comparing it to a database of known characters. It then finds the k closest matches from the database and assigns them weights based on their similarity to the input character. Finally, it uses these weighted scores to decide which character best describes the input character.

The main advantage of using this technique for recognising handwritten Gurmukhi is that it does not require any prior training or knowledge about how each letter looks like as compared with other machine learning algorithms such as Support Vector Machines (SVM) or Artificial Neural Networks (ANN). Moreover, the k-NN algorithm is highly adaptable and can be used on different types of data such as images, text documents or audio files.

The main disadvantage of using this method for handwriting recognition is that it may not work well if there are too many similar characters in the dataset. Additionally, as more and more input characters are added to the database, it becomes increasingly difficult to find an accurate match since all characters must be compared against each other.

Overall, k-NN is a useful technique for recognising handwritten Gurmukhi but its limitations should be taken into consideration when deciding whether or not to use it for a particular application.

## 4.9 Ramer-Douglas-Peukar (RDP) algorithm

The Ramer-Douglas-Peucker (RDP) algorithm is self-controlled technique. It is used to simplify and optimise the process of recognising characters from Gurmukhi script. This algorithm works by reducing the number of points in a given polyline or curve, while preserving its overall shape. The main idea behind this algorithm is to identify and remove redundant points from the input data, thus reducing the number of points required for handwriting recognition.

In order to recognize Gurmukhi characters, each point must be identified as either being part of a character or not. To do so, RDP uses two parameters: epsilon ($\varepsilon$), which determines

how far away from a line segment any given point can be before it's considered irrelevant; and delta (δ), which defines how close together two consecutive lines should be before they are combined into one single line segment.

Once these parameters have been set, RDP begins by taking the first two points on the input data and calculating their Euclidean distance between them. If this distance is less than ε & then these two points are merged into one single line segment, otherwise they remain separate line segments. This process continues until all remaining pairs of adjacent lines have been evaluated against ε and δ respectively. In this way, every pair that does not meet both criteria will eventually get eliminated from consideration for further processing leaving only those that do meet both requirements intact for further analysis such as recognising patterns in strokes etc., thereby enhancing accuracy in reading handwritten text written in Gurmukhi script language.

## WORK DONE FOR RECOGNITION OF GURMUKHI SCRIPT

(Singh , Sharma and Chauhan, Online Handwritten Gurmukhi Word Recognition Using Fine-Tuned Deep Convolutional Neural Network on Offline Features 2021)In their work, the neural network architecture for handwriting recognition from both online and offline data are covered. The current study employed a network with three blocks of CNN layers for the recognition of online handwriting. The first block consists of two Conv1D layers with 64 filters and one MaxPooling1D layer, the second block of two Conv1D layers with 128 filters and one MaxPooling1D layer, and the third block of three Conv1D layers with 256 filters and one MaxPooling1D layer. To prevent overfitting of the network, the output of the CNN layers is flattened and sent through two fully connected Dense layers of 512 neurons, each followed by a dropout of 30%. Conv1Ds with a kernel size of 3, padding "same" and rectified linear unit (ReLU) as an activation

function were utilised in the current study. The network solves the optimisation problem by using the RMSprop optimiser using categorical entropy as the loss function. With a combined dataset of 6K words from Singh et al. (2016) and Singh and Sharma (2019), the validation of this work was conducted utilising both writer dependent and writer independent data in data dependent mode while splitting the dataset into 80:20 ratio for training and testing, respectively. Due to the limited dataset in comparison to the number of features, all algorithms (SVM, Logistic Regression, CNN-DNN) perform badly when learning from online data. Due to the online dataset's limited size, it is converted to an offline dataset so that transfer learning may be applied. On the offline dataset, handwriting recognition was performed using the two pre-trained neural networks VGG16-DNN and InceptionV3-DNN. Due to the use of dropout and transfer to deal with over-fitting, it exhibits very nice convergence in a few number of epochs. Due to the combining of the points on the online samples, offline data samples have greater continuity. The present results as 97.44%, 97.23% and 96.21% for 90:10, 80:20, 70:30 train and test data, respectively, are achieved for data dependent mode of online Gurmukhi handwriting. This study is a significant step in this approach because it offered benchmark results for online handwritten Gurmukhi words in data dependent writing mode, with the highest results achieving above 97% recognition accuracy, and it recognised handwritten text using a different DL architecture.

(Sharma, et al. 2022) In this study, recognition is carried out on 22 district names,22000 Gurmukhi handwritten dataset images are created, 1000 dataset image samples for every district name are created which will help in automatically reading the district name of Punjab. The Dataset is then divided into 80:20 ratio. A holistic approach is employed for the

purpose of implementing all operations on the entire word. A CNN model is developed to predict the accuracy, loss, recall, and precision of district names, and it consists of three layers: "convolution," "max-pooling," and "flattening". The relevant features of the given data images are preserved with the assistance of convolution and pooling. The flattening layer transforms the obtained features into a column so they can be easily fed to the very last layer of the model, which facilitates in output classification. With an average validation accuracy of 95.6%, the implemented CNN model has a validation accuracy of 99.0%.

(Dargan and Kumar, Writer Identification System Based on Offline Handwritten Text in Gurumukhi Script 2020) propose an effective and ideal system for writer identification. A dataset of 100 authors, or 10053x10=53000 Gurmukhi characters, was used. Zoning, Transition, Diagonal, and Peak Extent based feature extraction methods were employed. Artificial Neural Network (ANN) for its low complexity, low parameter requirements, and low learning and reprogramming requirements, Multi Layered Perceptron Model (MLP) for its full connectedness, presence of many layers, and best fit to the classification of data that cannot be linearly separable, and Random Forest (RF) classifier, which is an ensemble learning algorithm that works by making decision trees are used in their study. Three factors are utilised to evaluate the algorithm: accuracy, true positive rate, and false positive rate. Individual feature extraction approaches and hybridisation of these methods with classifiers are utilised, and the values of these three parameters are compared to determine the best identification accuracy rates. The experiment indicated a maximum identification accuracy of 93.6% when employing F1+F2+F3+F4 with Random Forest classifier after implementing feature extraction and classification algorithms.

(Dargan and Kumar, Gender Classification and Writer Identification System based on Handwriting in Gurumukhi Script 2021)The research offered a novel approach to the Gurmukhi(Punjabi) script by combining different feature extraction approaches and hybridising classification algorithms. A corpus of 200 writers, comprising 100 men and 100 women, was used for the experiment analysis. The objective of this paper is to develop a framework for writer identification system and gender classification system through handwritten text in Gurmukhi script. This paper considers two feature extraction methods: intersection and open-end point-based feature extraction and curve fitting-based feature extraction. Support Vector Machine (SVM), Multi-Layer Perceptron (MLP), K-Neural Network (NN), Random Forest, and hybrids of these classifiers are used for classification in order to identify the writer and categorise the gender, based on a sample of their handwriting. Using curve fitting-based features and classifier hybridisation, the maximum gender classification accuracy of 90.57% has been reported. As for writer identification, curve fitting-based features and classifier hybridisation are reported to have an accuracy of 87.76%.

(Kumar and Jindal 2019)This research presents a study of several features and classifiers for the recognition of multilingual pre-segmented characters made up of English, Hindi, and Punjabi. The objective of this paper is to create multi-lingual text recognition system. To create a dataset, 40 different writers are requested to write English, Hindi and Punjabi characters separately. For English, both upper case and lower case characters are considered so it consists of 52-class problem whereas Hindi and Punjabi contains 36-class and 35-class problems. As a result, 4920 samples of pre-segmented characters written in English, Hindi, and Punjabi have been collected. A combination of Linear-SVM, k-NN, and MLP

classifiers yielded a maximum accuracy of 92.89% for script identification. The relevant dataset is segregated in a 70:30 ratio to create training and testing sets for each of the scripts used in this work. For recognition of English, Hindi, Punjabi accuracy of 92.18%, 84.67%, 86.79% is achieved by using combination of zoning, diagonal and intersection and open end point based features and a combination of Linear-SVM, k-NN, and MLP classifiers.

(G. Singh 2021)This study addresses the issue of online handwriting recognition for bilingual scripts. The developed system receives input in the form of handwritten text that includes both English and Punjabi words. Almost 90% of the words in the input text are in Punjabi, and the remaining 10% are in English. However, in order to provide clarification, the meanings of a few challenging Punjabi words are also written within "(" and ")" using English. In this study, a dataset made up of 28 words written in Gurmukhi script is presented. All of the characters in the Gurmukhi script are covered by these particular words. This dataset was created with the help of 50 writers overall by providing more than 8400 handwriting, representing various age groups and professions. The input is given with the help of Tablet-PC (Acer-R7 model) which is used as Digitiser and active stylus which is used for writing on digitiser's surface. The sample collection procedure generates a number of strokes by utilising hardware technology such as a digitiser. A stroke is the collection of all the stylus-touchable spots on the digitiser's surface between the Pen Down() and Pen Up() events. Different symbols in the Gurmukhi script are represented by a single stroke class or by a particular arrangement of these classes. The database information includes the name of the word being considered, the sample number, the identity of the writer, the number of strokes used to write the word completely, and information on the coordinate points for each stroke. Pre-processing procedures are carried out to enhance the quality of the input supplied, such as the elimination of repetitive stroke points, the detection of missing points, and normalisation. When both English and Punjabi words are taken into account, the implemented system achieved a recognition accuracy of 93.07% for the bilingual system.

(Singh, Kumar and Bar, A self controlled RDP approach for feature extraction in online handwriting recognition using deep learning 2020)The current paper provides a feature extraction technique for online handwritten strokes based on a self-controlled Ramer-Douglas-Peucker (RDP) algorithm. The feature extraction challenge is tackled in the current work by changing the traditional chain code directional feature vector by detecting key points of the writing trajectory for OHWR. The suggested method is a self-controlled RDP based approach for recognising individual units in online handwriting trajectories. It uses key points to create a smaller length feature vector by computing the straight lines between them. This enables recognition without preprocessing, as the shorter feature vector requires less computational power and time than traditional methods. The aim of this work is to provide a method that employs deep learning to locate key points for online handwritten strokes and characters, which are then used as feature vectors in OHWR. The network that is used to recognise composed of three layers of Conv1Ds, two fully connected dense layers, one output layer, two 1-dimensional max pooling (MaxPooling1D), three dropout layers and one flatten layer. The algorithm was put to the test on several datasets to confirm the effectiveness of the suggested scheme. Gurmukhi stroke classes and Unipen for digits are two benchmarked online handwritten datasets that are used. The online handwritten Gurmukhi strokes dataset was created in an online handwriting environment with the help of 100 contributors. This research used a dataset of 34k samples of data for 62

classes of online handwritten Gurmukhi strokes. Datasets were randomly divided into ratios of 90:10, 80:20, 70:30, 60:40, and 50:50 for training and testing while experiments were being conducted on each dataset. Results are presented as the average of five random runs in each train-test data size. For the Gurmukhi and Unipen datasets, this method achieves 94.51% and 94.55% recognition accuracy, respectively.

( Aggarwal and Singh 2015) The purpose of this work is to discuss the offline recognition of handwritten Gurmukhi characters. It uses a handwritten Gurmukhi character database with 7000 sample character pictures. In order to extract features from handwritten offline Gurmukhi characters, a method utilising gradient and curvature hybrid features was used. For the purpose of recognition, an SVM classifier with an RBF kernel has been utilised. For computing the recognition rate, the experimental framework employs a 5-fold cross validation approach. This framework separates the image in the  experimental database into 5 equal-sized subsets, with each subset being tested while the other four subsets are used for training. The proposed composite feature extraction approach has been tested, and a recognition accuracy of 98.56% has been attained in this paper.

## FUTURE SCOPE OF THIS STUDY

Gurmukhi is a script used primarily to write the Punjabi language. Handwriting recognition of Gurmukhi has been a challenge due to its complex script and difficulties in distinguishing between similar letter shapes. Despite the challenges, advancements in technology have enabled better accuracy and recognition of Gurmukhi handwriting. In this, we will explore the current state of handwriting recognition of Gurmukhi and discuss 6 potential areas for future development.

**1. Development of a Comprehensive Dataset:** The development of a comprehensive dataset is essential for the effective implementation of Gurmukhi handwriting recognition. This dataset should contain a variety of handwritten samples from different sources, including books, magazines, newspapers, and online resources. Furthermore, the dataset should contain samples from different writers and different writing styles. This will help the computer learn more about the Gurmukhi handwriting style and provide better recognition results.

**2. Automated OCR:** Optical Character Recognition (OCR) is a method of automatically recognising text from images, such as scanned documents or handwritten notes. Automated OCR for Gurmukhi would allow users to quickly and easily convert their handwritten notes into digital text. This technology could be used to facilitate rapid search and retrieval of handwritten Gurmukhi documents, and could be especially useful for archival and research purposes.

**3. Advanced Pre-Processing Algorithms:** Pre-processing algorithms are essential for the successful implementation of Gurmukhi handwriting recognition. Currently, pre-processing algorithms are used to improve the accuracy of handwriting recognition by removing noise, detecting characters, and normalising the handwriting. However, more advanced pre-processing algorithms can be used to further improve the accuracy of the recognition. For example, algorithms can be developed to detect loops, curves, and other features in the handwriting.

**4. Incorporation of Contextual Information:** Contextual information, such as the type of document, the writer, and the writing style, can be used to improve the accuracy of Gurmukhi handwriting recognition. For example, the computer can be trained to recognize patterns in

the handwriting that are specific to a certain writer or writing style. This will enable the computer to
recognize Gurmukhi characters better.

**5. Support for Multilingual Writing**: Gurmukhi is often used to write other languages, such as Dogri, Braj Bhasha, and Sanskrit. In order to provide better support for multilingual writing, the handwriting recognition algorithms must be able to identify and classify characters in multiple languages. This will require the development of algorithms that can recognize and distinguish between letter shapes of different languages.

**6. Voice Recognition**: Voice recognition technology could be used to convert spoken Gurmukhi into digital text. This technology could be used to help users quickly and easily transcribe their spoken words into digital documents, or to convert speech into text for use in natural language processing applications. Overall, handwriting recognition of Gurmukhi is a promising field of research that has the potential to revolutionise the way we use and interact with the language. The six potential future scopes discussed above are just a few of the possibilities that this technology could open up. As further research and development is conducted in this field, it is likely that even more exciting applications will be discovered.

## CONCLUSION

Gurmukhi handwriting recognition has been an active research area in the past years. It has made significant progress in terms of accuracy and speed, but still more work needs to be done to make it a viable solution for real-world applications. In this review, we have discussed the state-of-the-art methodologies used for Gurmukhi handwriting recognition. We have also discussed the datasets used for training and testing models. We have presented an overview of the various approaches used for Gurmukhi

handwriting recognition, different types of feature extraction techniques and classifiers used for the same purpose. Overall, it can be concluded that Gurmukhi handwriting recognition has made significant progress in terms of accuracy and speed, but there is still room for improvement. Efforts towards research in this field should continue to be made, in order to make it a viable solution for real-world
applications.

## REFERENCES

Aggarwal, Ashutosh, and Karamjeet Singh. "Handwritten Gurmukhi character recognition." 2015 International Conference on Computer, Communication and Control (IC4). Indore,India: IEEE, 2015. 1-5.

Singh, Gurpreet. "A Bilingual (Gurmukhi-Roman) Online Handwriting Identification and Recognition System." International Journal of Recent Technology and Engineering (IJRTE) , 2021: 2093-1204.

Dargan, Shaveta , and Munish Kumar. "Gender Classification and Writer Identification System based on Handwriting in Gurumukhi Script." 2021 International Conference on Computing, Communication, and Intelligent Systems (ICCCIS). Greater Noida, India: IEEE, 2021. 388-393.

Dargan, Shaveta, and Munish Kumar. "Writer Identification System Based on Offline Handwritten Text in Gurumukhi Script." Sixth International Conference on Parallel, Distributed and Grid Computing (PDGC). Waknaghat, India: IEEE, 2020. 544-549.

Kumar, Munish , and Simpel Rani Jindal. "A Study on Recognition of Pre-segmented Handwritten Multi-lingual Characters." Archives of Computational Methods in Engineering (SPRINGER) 27 (2019): 577-589.

Sharma, Sandhya , et al. "Optimized CNN-Based Recognition of District Names of

Punjab State in Gurmukhi Script." Journal of Mathematics (Hindawi), 2022: 1-10.

Singh , Sukhdeep , Anuj Sharma, and Vinood Kumar Chauhan. "Online Handwritten Gurmukhi Word Recognition Using Fine-Tuned Deep Convolutional Neural Network on Offline Features." Machine Learning With Python (Elsevier) 5 (2021): 1-15.

Singh, Sukhdeep , Vinod Chauhan Kumar , and Elisa H. Bar. "A self controlled RDP approach for feature extraction in online handwriting recognition using deep learning." Applied Intelligence (SPRINGER ) 50 (2020): 2093-2104.