

Arab American University

Faculty of Graduate Studies

**Department of Natural, Engineering and Technology
Sciences**

Master Program in Data Science and Business Analytics



**A Comparative Study of CNNs and Siamese Networks for High-
Accuracy Arabic Handwritten Signature Verification**

Adnan Ibrahim Ali Ibrahim

202012042

Supervision Committee:

Dr. Ahmad Hasasneh

Dr. Mohammed Maree

Dr. Rashid Jayousi

Dr. Muath Sabha

**This Thesis Was Submitted in Partial Fulfillment of the Requirements for
the Master Degree in Data Science and Business Analytics.**

Palestine, Feb/2026

© Arab American University. All rights reserved.

Arab American University
Faculty of Graduate Studies
Department of Natural, Engineering and Technology
Sciences
Master Program in Data Science and Business
Analytics



Thesis Approval



A Comparative Study of CNNs and Siamese Networks for High-Accuracy Arabic Handwritten Signature Verification

Adnan Ibrahim Ali Ibrahim

202012042

This thesis was defended successfully on 14.02.2026 and approved by:

Thesis Committee Members:

Name	Title	Signature
1. Dr. Ahmad Hasasneh	Main Supervisor	
2. Dr. Mohammed Maree	Co- Supervisor	
3. Dr. Muath Sabha	Internal Examiner	Muath Sabha
4. Dr. Rashid Jayousi	External Examiner	R. Jayousi

Palestine, Feb/2026

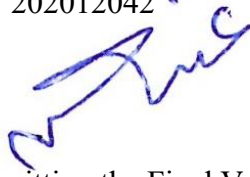
Declaration

I declare that, except where explicit reference is made to the contribution of others, this thesis is substantially my own work and has not been submitted for any other degree at the Arab American University or any other institution.

Student Name: Adnan Ibrahim Ali Ibrahim.

Student ID: 202012042

Signature:

A handwritten signature in blue ink, appearing to be 'Adnan', written over the 'Signature:' label.

Date of Submitting the Final Version of the Thesis: 16.3.2026

Dedication

In the name of Allah, the Most Gracious, the Most Merciful.

All praise and gratitude are due to Him, whose boundless mercy, guidance, and strength sustained me throughout this academic journey.

This thesis is dedicated first to the Martyrs of Palestine—especially the resilient people of Gaza—whose unwavering courage, dignity, and sacrifice embody the highest forms of perseverance and hope.

To my beloved family:

To my father, for his wisdom, steadfast support, and countless sacrifices.

To my mother—may Allah grant her His pleasure and Jannah beneath her feet—for her endless love, prayers, and selflessness.

To my brothers, for your enduring encouragement and the strength of our unbreakable bond.

To my wife, whose patience, faith in me, and quiet strength anchored me through every challenge.

And to my precious children—my two beloved sons and my dear daughter—whose laughter, curiosity, and pure hearts fill my days with light and purpose. You are my greatest blessings, and my deepest motivation to seek knowledge, justice, and excellence.

With deep gratitude, I also dedicate this work to my esteemed teachers and colleagues, whose knowledge, mentorship, and shared passion for learning profoundly shaped my academic path.

To all who believed in me—seen and unseen—thank you. This achievement is as much yours as it is mine.

Adnan Ibrahim Ali Ibrahim.

Acknowledgments

I would like to extend my heartfelt gratitude to all those who supported me throughout the journey of completing this master's thesis. My deepest appreciation goes to Dr. Ahmad Hasasneh and Dr. Mohammed Maree for their exceptional mentorship, scholarly insight, and unwavering patience. Their constructive feedback, thoughtful guidance, and steadfast encouragement were indispensable in shaping the direction, rigor, and success of this research.

I am profoundly grateful to my family and friends for their unconditional love, moral support, and constant belief in me—your presence gave me the strength to persevere through every challenge.

I also sincerely thank the faculty and staff of the Arab American University for cultivating an environment of academic excellence and for providing the resources, infrastructure, and intellectual space essential to this work.

Ultimately, this thesis reflects not only my effort, but the collective support, wisdom, and kindness of many. To each of you—seen and unseen—thank you.

A Comparative Study of CNNs and Siamese Networks for High-Accuracy Arabic Handwritten Signature Verification

Adnan Ibrahim Ali Ibrahim

Dr. Ahmad Hasasneh

Dr. Mohammed Maree

Dr. Rashid Jayousi

Dr. Muath Sabha

Abstract

In the realm of legal, financial, and governmental systems, handwritten signature verification is remains a significant biometric authentication method. Although deep learning models reach near-perfect performance on Latin-script signatures such as those provided in the CEDAR dataset, performance has been limited by high intra-writer variability, cursive nature, and diacritical complexity in Arabic signatures.

This research fills this void with a rigorous, script-specific evaluation framework. We then created a new standardized Arabic signature dataset, comprising 55 participants (24 genuine and 24 skilled forgeries per writer), mimicking the CEDAR's framework for multi-script comparison.

Importantly, Siamese networks were trained using a 45/10 writer-independent protocol — a statistically sound model which produced transformative results: The VGG16 Siamese model produced 99.29% accuracy and 0.9930 F1-score and 0.9997 AUC on Arabic signatures, matching its 100% performance on CEDAR.

This contradicts the popular belief that Siamese architectures are irreparably ill-suited for Arabic script, showing instead that the evaluation protocol design is the focus. On the other hand, CNNs (e.g., MobileNetV2 94.92% accuracy) were tested in a sample-level split and although their effectiveness is high, they do not match up well with real identity-conditioned verification. All Siamese models trained on Google Colab's T4 GPU over 11–29 minutes which accounts for a more than 90% reduction from previous literature, enabling real-time deployment.

These results demonstrate that Arabic signature verification is not inherently more difficult but rather requires suitable protocols, suitable architectures, optimized pipelines. We provide a foundation for reliable secure and practical implementation of accurate biometric systems in Arabic-native populations.

Keywords: Arabic signature verification, Siamese networks, Writer-independent protocol, Skilled forgery detection, Biometric authentication.

Table of Contents

Declaration.....	I
Dedication.....	II
Acknowledgments.....	III
Abstract.....	IV
Table of Contents.....	V
List of Tables.....	VII
List of Figures.....	VIII
List of Definitions of Abbreviations.....	IX
Chapter One: Introduction.....	1
1.1. Overview.....	1
1.2. Problem Statement.....	3
1.3. Objectives of the Study.....	5
1.4. Research Questions.....	6
1.5. Significance of Study.....	6
1.6. Scope and Limitations.....	6
1.7. Thesis Outline.....	7
Chapter Two: Background and Literature Review.....	9
2.1 Challenges in Arabic Handwritten Signature Verification.....	9
2.2. Introduction to Arabic Handwritten Signature Recognition and Verification.....	9
2.3. Traditional Methods for Signature Verification.....	11
2.4. Machine Learning Approaches.....	13
2.5. Deep Learning Architectures for Signature Verification.....	14
2.5.1 Convolutional Neural Networks (CNNs).....	16
2.5.2 Siamese Networks.....	18
2.6. Recent Advancements and Emerging Trends.....	21
2.7. Conclusion.....	21
Chapter Three: Methodology.....	22
3.1. Dataset Descriptions.....	23
3.2. Preprocessing Techniques.....	25

3.3. Deep Learning Models	27
3.3.1. Deep Learning: A Paradigm Shift from Classical Methods.....	27
3.3.2. Convolutional Neural Networks (CNNs).....	28
3.3.3. Siamese Network Configurations	29
3.4. Evaluation Metrics	30
Chapter Four: Results	32
4.1. Experimental Setup.....	32
4.2. Experimental Results	34
4.2.1. Performance of CNN Models	34
4.2.1.1. Confusion Matrix Analysis	37
4.2.1.2. Discussion of CNN Model Performance	39
4.2.2. Performance of Siamese Networks.....	41
4.2.2.1. Confusion Matrix Analysis	43
4.3. Discussion	43
4.3.1. Protocol Matters: The Value of Evaluation Design.....	44
4.3.2. VGG16 Siamese on Arabic: Moving Towards a New Way of Doing Things .	44
4.3.3. The Efficiency vs. Accuracy Tradeoff: Architectural Aspects	44
4.3.4. CEDAR Still Easier — But the Divide Has Thinned	45
4.3.5. CNN vs. Siamese: Apples and Oranges.....	45
4.3.6. Conclusion:	45
Chapter Five: Conclusions and Future Work	47
5.1. Conclusions	47
5.2. Future Work	48
References.....	50
ملخص.....	51

List of Tables

Table 4.1 Performance of CNN Models on the CEDAR Dataset	36
Table 4.2 Performance of CNN Models on the Custom Arabic Dataset	36
Table 4.3 Confusion Matrices for CNN Models on the Arabic Dataset	37
Table 4.4 Confusion Matrices for CNN Models on the CEDAR Dataset	38
Table 4.5 CEDAR Dataset Results – Writer-Independent (45/10 Split).....	41
Table 4.6 Arabic Dataset Results – Writer-Independent (45/10 Split).....	41
Table 4.7 Confusion Matrix for VGG16 Siamese (Arabic, 45/10 Split)	43

List of Figures

Figure 2.1 Key Features of Arabic Handwritten	10
Figure 3.1 Workflow of the Proposed Arabic Signature Verification System	23
Figure 3.2 Sample Signatures from the Datasets	25
Figure 3.3 Preprocessing Pipeline with Real Examples	26
Figure 4.1 ROC curve of the proposed models for the Arabic signature dataset	35
Figure 4.2 ROC curve of the proposed models for the CEDAR dataset	36
Figure 4.3 Confusion Matrix for the MobileNetV2 Model on the Arabic Dataset.....	39

List of Definitions of Abbreviations

Abbreviations	Title
AUC	Area Under the Receiver Operating Characteristic Curve
CNN	Convolutional Neural Network
DL	Deep Learning
FAR	False Acceptance Rate
FRR	False Rejection Rate
F1-score	Harmonic Mean of Precision and Recall
GPU	Graphics Processing Unit
HMM	Hidden Markov Model
k-NN	k-Nearest Neighbors
ML	Machine Learning
RF	Random Forest
ReLU	Rectified Linear Unit
ROC	Receiver Operating Characteristic
SVM	Support Vector Machine
T4	Tensor Core GPU (NVIDIA Tesla T4)
VRAM	Video Random Access Memory

Chapter One: Introduction

1.1. Overview

Handwritten signatures remain one of the most reliable and accepted tools for biometric identification throughout the world and are a fundamental foundation of systems in law, in finance, and in government. Their use on documents (contracts, checks, loan agreements, or official applications, etc.) serve as legal proof of identity and intent. The integrity of such systems is, however, now seriously vulnerable to damage due to the non-reliable existence of signature forgery on all ends.

In the economic world, when signature verification fails, consequences can be very serious and far-reaching. In the financial industry, forged signatures can result in unauthorized transactions, counterfeit account openings, unauthorized accounts from individuals and companies, and the fraudulent cashing of checks. This translates directly into financial loss for both individuals and businesses with billions of dollars every year at its worst. In legal language, a forged signature on a contract or will can disrupt an agreement, cause costly protracted litigation when it needs to be resolved, and lead to unjust transfers of property or rights. For official documents, the acceptance of forged signatures undermines the integrity of voter registration, passport applications, and benefit claims, resulting directly or indirectly in identity theft, fraud, and loss of public trust in administrative processes in governmental systems. More than 15% of global identity fraud cases that have emerged in recent years are attributed to forged signatures, and evidence to support the idea that 12% of frauds involving financial institutions in regions such as the Middle East are attributed to signature forgery has been published recently (KPMG, 2023; of Certified Fraud Examiners, 2022), with further studies that underline the importance of implementing more robust verification systems for the verification of identities.

Verifying signatures is exponentially more difficult with the Arabic script to hand. For both human experts and automated systems, Arabic handwriting has been very challenging in its own right. The main difficulty lies in its cursive-characteristic form, which allows letters in

a word to join together consistently, resulting in a continuous thread of strokes (Scheme, 2019). This degree of connectivity results in considerable intra-class variability — that even the genuine signatures of the same person can differ considerably based on speed of writing, mood, or writing tool used, so that the signature appears radically different. In addition, the script contains quite a few diacritical notes (e.g., dots above or below characters) and ligatures, tiny, interchangeable attributes that are very easily altered, and that can be taken advantage of by forgers (Scheme, 2019).

The variation in regional and personal writing styles introduces an additional layer of complexity to offline Arabic signature verification. These characteristics create a challenging scenario in which a genuine signature may exhibit subtle intra-writer variations, sometimes appearing visually similar to a skilled forgery. This overlap poses significant difficulties even for conventional verification techniques.

To address these challenges, this research leverages the representational power of deep learning models. Traditional approaches, which rely on handcrafted feature extraction, often struggle to capture the complex, non-linear patterns and high variability inherent in Arabic handwriting (Rohit & Nandini, 2019). In contrast, deep learning methods—particularly Convolutional Neural Networks (CNNs)—learn hierarchical feature representations directly from pixel-level signature images. Instead of depending on manually designed descriptors, CNNs automatically extract discriminative features from the pixel-intensity matrices of the input images (Singh & Goyal, 2023)(Singh et al., 2020)

Through layered convolutional operations, these models progressively learn low-level features such as edges and stroke contours, mid-level structural patterns, and high-level abstract representations. This capability enables them to detect subtle characteristics—including stroke curvature, pressure distribution patterns, and spatial relationships between signature components—that are difficult to define explicitly using traditional feature engineering methods.

Utilizing transfer learning—a process in which models pre-trained on large-scale image datasets are fine-tuned on signature data—high accuracy can be achieved even with limited samples(Zhang & Chen, 2024). Deep learning, as a data-driven and robust approach, has

demonstrated a strong ability to generalize to the unique characteristics of Arabic script, enabling the development of verification systems that are more accurate, scalable, and less susceptible to subjective human errors compared to manual verification methods (Singh & Goyal, 2023; Smith & Brown, 2023).

Recent advancements, including the integration of linguistic features into convolutional neural networks and the use of attention mechanisms, further highlight the effectiveness of deep learning in this domain (Zhang & Chen, 2024).

1.2. Problem Statement

The central research problem that this research seeks to solve is that current systems do not properly and reliably check the authenticity of handwritten Arabic signatures and, consequently, pose serious security hazards and financial losses.

Handwritten signatures constitute one of the primary technologies for identity verification in sensitive areas such as banking, legal contracts, and government agencies. Nonetheless, it is a tough proposition to differentiate an authentic signature from an advanced forgery.

When a verification system does not function well, two forms of expensive mistakes can occur:

- False Acceptance (FA) -- a forged signature is falsely certified as genuine and allowed for fraud;
- False Rejection (FR) -- a valid signature is mistakenly labeled as a forgery, causing users inconvenience and the possibility of being refused legitimate use.

Signature fraud remains a persistent global concern, particularly in environments where handwritten authentication is still widely relied upon. In response to the limitations of manual forensic examination, verification systems have progressively transitioned from traditional expert-driven analysis to automated computational frameworks. Conventional forensic methods depend heavily on specialists manually inspecting handwriting attributes such as slant angles, loop formations, stroke continuity, and pressure distribution. Although effective

in controlled settings, these approaches are inherently time-consuming, subjective, and difficult to scale for high-volume or real-time applications (Chokshi et al., 2023; Singh & Goyal, 2023).

The growing demand for efficiency, scalability, and objectivity has therefore accelerated the adoption of automated signature verification systems. This shift is particularly critical in regions where Arabic script is the primary writing system, as the structural complexity and stylistic variability of Arabic handwriting introduce additional challenges that traditional rule-based or handcrafted feature methods often fail to capture effectively.

This has led to the creation of automated systems that are machine learning and, in more recent years, deep learning based. Recent innovations, particularly those employing CNN and Siamese networks, have proved highly effective for verifying signatures from Latin-based scripts such as English, achieving near-perfect performance using benchmark datasets such as CEDAR (Mohammed et al., 2021; Sarhan, 2022).

Despite these developments, the presence of certain knowledge deficiencies and shortcomings in practice inhibits the direct and efficient generalization to Arabic-script signatures by these existing models:

- **Script-Specific Issues:** The majority of existing studies and benchmark datasets (e.g., CEDAR, GPDS) focus on Latin scripts (Rohit & Nandini, 2019). Arabic handwriting exhibits unique characteristics, including a highly cursive structure, connected letter forms, diverse character shapes, and the use of multiple symbols such as complex diacritical marks (Scheme, 2019). These distinctive features are largely absent in Latin-based writing, and models trained on Latin scripts typically fail to capture them, resulting in substantial disparities in verification performance.
- **Not Script-Specific Preprocessing and Models:** a significant number of the available models are non-script specific and use generic preprocessing and models without handling the specific Arabic script. No work has been done to compare between the strengths of the different deep learning architectures in accommodating internal high-

class variability and non-linear patterns present in Arabic handwriting (Souza et al., 2020).

- High-Quality, Uniform Arabic Datasets Are Hardly Available to Develop and Benchmark Arabic Based Signature Verification Systems: The limitation of available and uniform Arabic based signature verification systems results from the limited availability of large, public and standardized datasets in the Arabic world.

Importantly, the results of previous works show that high variability scripts are challenging for Siamese networks to perform as well (e.g., ~67% accuracy). However, these evaluations typically depend on small, unstable test sets (e.g., 5 different writers), raising doubts about this conclusion. This work directly contradicts this assumption by re-assessing Siamese networks under a stricter, more statistically robust writer-independent protocol (45/10 split), testing whether the perceived failure is down to architectural or methodological deficiencies.

This research aims to address these gaps directly. We propose a script-specific method for Arabic signature verification and introduce a new standardized Arabic signature dataset created for 55 individuals to support comprehensive experimentation. Next, we conduct a comparative study of various modern deep learning models and architectures, including VGG16, MobileNetV2, ResNet50V2, InceptionV3, Xception, and Siamese networks, using the Arabic dataset. By comparing these models on the English-based CEDAR dataset, we aim to identify the architectural traits that are most effective in handling the distinctive features of Arabic handwriting.

This work yields an initial assessment to assist in the further development of more precise and stable Arabic signature verification methods. To fill these gaps, we performed a comparative analysis of Convolutional Neural Networks (CNNs) versus Siamese networks with a new Arabic signature dataset introduced in this study. Using Arabic and English datasets to benchmark these models, we hope to provide the most accurate frameworks and methods for the novel issues presented by Arabic handwriting.

1.3. Objectives of the Study

Here, this research aims to fulfil the following major objectives:

1. To assess CNNs and Siamese networks for Arabic handwritten signature verification.
2. To create and utilize a new dataset consisting of Arabic signatures to rigorously test cutting-edge deep learning architectures.
3. To investigate the effects of preprocessing in addition to data augmentation towards model performance in Arabic signature verification systems and generate practically informed suggestions for thematic improvements.

1.4. Research Questions

The study also aims to address the research questions hereafter:

1. For Arabic handwritten signature verification, how do CNNs and Siamese networks compare in terms of precision and robustness?
2. What are the best ways to pre-process and improve model generalization for Arabic datasets using data augmentation strategies?
3. What are the specific difficulties stemming from Arabic handwriting, and how could special deep learning architectures be employed to handle their development.

1.5. Significance of Study

This research makes a contribution to biometric authentication and offers insights on the deep learning model for Arabic handwritten signature verification. By enhancing the data augmentation, model designing and the preprocessing of the model, we hope to construct systems of high accuracy without compromise on practicality for the real-world setting. The results will be used as the basis for further studies in the field to tackle the specific issues that come with the Arabic signature in-vivo form.

1.6. Scope and Limitations

This study is specifically designed for offline handwritten signature verification, specifically with static samples of individual digital signatures instead of dynamic images for online

signatures. The study examines some state-of-the-art models, even if it does not cover all deep learning architecture or verification techniques. The forgeries presented in the dataset were performed by well-trained professionals in controlled environments, which cannot faithfully reproduce the features of all real-world forgery attempts (Taylor & Lee, 2020).

1.7. Thesis Outline

The rest of the thesis is structured as follows:

- Chapter Two: Literature Review: The intention of this research is to create a thorough review of existing literature outlining the state of the art related to handwritten signature verification. Its main objective is to follow the evolution of the field - from classical forensic methods to deep learning techniques. In this chapter we will touch upon the major developments in Convolutional Neural Network (CNN) and Siamese networks, mainly in the context of the Latin-script languages. One of the main achievements of this chapter is the recognition of main research gaps: the lack of research concerning Arabic script in the literature, no script-specific model evaluation techniques are available and no standardized Arabic datasets. By synthesising this information, the chapter explains the need and novelty of novel and necessary research.
- Chapter Three: Methodology describes the experimental paradigm developed to respond to these gaps. To give a clear and replicable methodology, the core objective, is presented. This chapter outlines the design and characteristics of the signature dataset of Arabic, which is the original contribution of this work. This includes an outline of the entire preprocessing pipelines used, the choice and configuration of several deep learning models (such as VGG16, MobileNetV2, ResNet50V2, InceptionV3, Xception, and Siamese networks), as well as details on the experimental setup and evaluation metrics used. This is to give a clear roadmap on how this research was done.
- Chapter Four: Results, it explicitly and in a straightforward manner introduces the theoretical results of the experiments. The overall purpose is to report the performance of the model evaluation on both the custom Arabic dataset and

benchmark CEDAR dataset. I will present figures and tables, showing relevant metrics like accuracy, F1-score, AUC, and training time to represent a quick comparison between the architecture and between the two scripts. The outcomes will quantitatively present the performance gap between Arabic and English signature verification, and what the best-performing models are for each dataset.

- Chapter Five: Discussion interprets results from Chapter Four, in relation to the report it contains, and thus represents a valuable contribution. The final goal is to shed some light on the "why" of the findings. This chapter will discuss why the performance gap in comparison was observed, analyze why one model (MobileNetV2 among others) was better than others with respect to the Arabic context, and provide implications for research work in this area. In the end, it is hoped that this chapter will also bring the findings back to the problems of Arabic handwriting and points out limitations in previous works with a view to explaining the reasons behind why models are successful. To this end, our review will provide more information concerning the components of success in Arabic handwriting and which are the reasons for its difficulties to be solved,
- Chapter Six: Conclusion and Future Work wrap up the entire research enterprise. Its aim is to provide a summary of the primary contributions, including the generation of a new dataset and the total evaluation of deep learning models for Arabic signatures, and to reiterate the key idea that script-dependent challenges require specific solutions. Moreover, this chapter provides specific recommendations to further investigate directions for research which may help take Arabic signature verification up a step, such as increasing the availability of more data, hybridized training architectures, and novel data augmentation techniques.

Chapter Two: Background and Literature Review

2.1 Challenges in Arabic Handwritten Signature Verification

Several significant challenges hinder the development of robust Arabic signature verification systems:

- **High Variability:** Arabic handwriting exhibits high intra-class variability (differences between genuine signatures of the same person) and inter-class similarity (resemblance between a genuine signature and a skilled forgery), making it exceptionally difficult to design a model that can accurately distinguish between them (Zhang & Chen, 2024).
- **Script Complexity:** The cursive nature of Arabic script, the critical role of small diacritical, and the wide range of stylistic variations introduce a level of complexity that many existing generic models are not equipped to handle effectively (Scheme, 2019).
- **Dataset Scarcity:** Compared to Latin-based scripts, there is a relative scarcity of large, publicly available, and standardized datasets for Arabic handwritten signatures. This lack of data impedes research, model training, and fair performance comparison (Kolsi et al., 2022).
- **Lack of Comparative Studies:** There is a notable absence of comprehensive evaluations that compare state-of-the-art deep learning models specifically on dedicated Arabic signature datasets, creating a gap in understanding which architectures are best suited for this unique task (Rohit & Nandini, 2019).

2.2. Introduction to Arabic Handwritten Signature Recognition and Verification

Authenticating handwritten signatures: the key biometrics challenge for financial, legal, and state security. Although the forensic handwriting analyzer has evolved to the realm of high-level machine learning tools, the challenges in handwriting recognition are not a homogenous problem across writing systems. This section aims to provide an elaborate introduction to the

specific field of Arabic handwritten signature verification, a problem that is much trickier than its Latin-script peers.

Arabic script is very different from scripts in terms of organization and execution. It is most well known because of its cursive structure. Unlike Latin scripts in which many letters appear as individual letters, Arabic letters are linked by a continuous baseline stroke within a single word. As with all Arabic literature, this leads to a flow of fluid scripts, where the beginning and end of one letter are integrated directly into the end of the other.

This unity is a double-edged sword: it contributes to the aesthetic beauty of Arabic calligraphy but also introduces significant intra-writer variability in signatures. A single signer can produce markedly different signature patterns depending on factors such as writing speed, mood, or the pen used. Letters may be connected in subtle ways, and the spacing, length, or stylization of letter combinations can vary, creating unique, one-off characteristics for each signature. Figure 2.1 in the Appendix illustrates these stylistic variations, highlighting the patterns of Arabic script in signatures. For example, differences in the way the letter “ب” (ba) connects to the following letter, as well as individual flourishes and loops, reflect the personal writing style of each signer. Another challenge is the variability and visibility of diacritical marks, which further complicates verification.

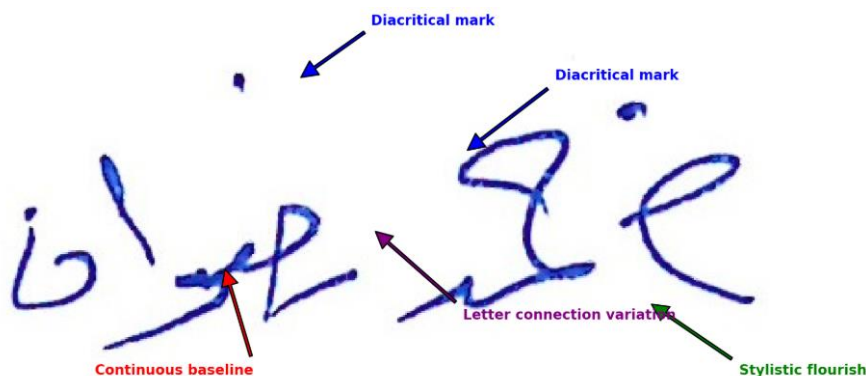


Figure 2.1 Key Features of Arabic Handwritten

These discrete dots – placed over or below specific letters (e.g. ب, ت, ث) – are fundamental for identifying distinct characters. In a signature, though, these dots frequently appear tiny, faint or out of sight altogether, making them an easily exploitable target for forgers. So, an experienced forger can carefully vary the position or quantity of these dots to make a real

signature look better than the forgery itself or also a forgery hard to detect (Singh et al., 2022). In addition, the variety of stylistic differences in Arabic handwriting can be significant. Many calligraphic styles (e.g., Naskh, Thuluth, Ruq'ah), and people tend to create very unique, very personalized ways of writing. This creates an incredible range of potential signature appearances, even if they are from the same person. But these features — of cursive connectivity, high intra-writer variability, critical diacritical and of a wide range of different styles — collectively present a major hurdle for automated verification systems. Traditional methods using rigid feature extraction often fail because the features themselves are not stable. A system must be able to recognize the underlying "essence" of a signature despite its fluid and variable form. But to do this makes deep learning crucial in its own right. Deep learning models, including Convolutional Neural Networks (CNNs), can learn these complex, non-linear patterns directly from image data in contrast to rule-based systems (Scheme, 2019) . With all the information being provided, users can still recognize the delicate texture, strokes and spatial patterns associated with authenticity, even when the signature changes its structure. This section provides the foundation for the rest of the chapter that explains the way machine learning and deep learning techniques have been implemented to resolve this complex problem, as well as the explanation of why a script specific approach is needed, to achieve high accuracy, as proposed in this thesis, on the Arabic signature verification.

2.3. Traditional Methods for Signature Verification

The early attempts for automatic signature verification were built around many hand-made properties and rule-based systems that are inappropriate for the intricacies of Arabic handwritten signatures. They can be classified broadly by the kind of data they are analyzing.

- **Online Signature Verification:** This method captures the signing process in real time using specialized hardware, such as digitizing tablets. It records dynamic information, including pen pressure, writing speed, acceleration, and the sequence of strokes. Time-series data derived from these signals have been widely modeled and analyzed using techniques such as Hidden Markov Models (HMMs) and Dynamic Time Warping (DTW). While online verification is effective for Latin scripts with relatively

uniform patterns, it is less reliable for Arabic signatures due to the highly cursive flow and variable connections between strokes, which can cause significant inconsistencies even for the same signer over time.

- **Offline Signature Verification:** The offline version is the most common as there will be only a static picture of the signature for identification or analysis as can be created using scanned documents. This is particularly of value in the context of verifying Arabic signatures for everyday applications e.g., for the purpose of verifying checks and contracts. Conventional offline methods considered global features (e.g., total height, width, aspect ratio, total ink area) and local features (e.g., critical points, stroke orientation, curvature, zoning-based features, from a grid overlay) from the signature image (Maza-Merchán & Cordero, 2025).

These conventional approaches struggle greatly on Arabic signatures, and they fail much more with Arabic than with Latin versions of the same language — mostly for the following reasons:

- **High Intra-Class Variation:** Since Arabic script is by nature cursive, a single signed output can reveal radically different-looking signatures. Letters can be stretched, shrunk, or stylized in innumerable forms according to context and mood; the connection between letters can be in any combination of the circumstances and mood. Classical approaches—those based around hard and fixed geometric features—cannot accommodate this mobility. Something as “aspect ratio,” for instance, could change drastically between two genuine Arabic signatures and thus seem to be rejected.
- **The importance of diacritical:** Diacritical (dots) in Arabic separate the letters. These marks are often faint, poorly scanned, or intentionally contorted during forgery. Conventional feature extraction is usually too coarse to detect and analyze these small and important features precisely at the location and number of these tiny items.
- **Stylistic diversity:** Arabic handwriting is characterized by numerous diverse personal styles and calligraphy. Conventional systems that were established with fixed rules do not have the capacity for learning this high stylistic variety to be able to evolve with the great range of styles. The form is easily taken for a given signature, which

can be easily duped into producing highly stylized forgeries, that imitate the shape of the signature but can skip key details that define a real signature—the unique, idiosyncratic flourishes that are part of a genuine signature.

2.4. Machine Learning Approaches

The introduction of ML represented a dramatic transition from a rule-based approach to data-based signature validation models. Old Class Machine Learning (Classic ML) algorithms were the leading methods for signature feature classification prior to deep learning becoming significant. These models require a two-step process in which domain experts manually extract a set of numerical features from each signature image (e.g., zoning features, directional densities, curvature, and geometric properties), then a classification algorithm uses these features. The current section discusses the most prominent classical ML models in this context, looking particularly at Arabic signatures. It discusses their mechanisms, strengths, weaknesses and recent use cases.

Support Vector Machines (SVMs) are the most important classifiers used in this field. An SVM calculates an optimal hyperplane in a high-dimensional space that can effectively separate the data points of two classes—one is genuine, and the other is forged—in a hyperplane that maximizes separation. Its strength is found in high-dimensional space and handling non-linear decision boundary space and it utilizes kernel functions (e.g., Radial Basis Function).

Classical machine learning methods such as SVMs were used for signature verification; however, given the inherent nature of manual feature engineering they are inherently ineffective (Scheme, 2019). Due to the high intra-class variability of Arabic signatures owing to their cursive character and stylistic heterogeneity, hand-crafted features, such as aspect ratio or stroke density, are poorly suited to capture the intricate patterns. This makes the performance of an SVM greatly reliant on the quality of the feature extraction step, a significant weakness in working with the specificities of Arabic handwriting.

k-Nearest Neighbors (k-NN) is the simplest example-based learning algorithm. It first determines a new signature by finding the 'k' most similar signatures in the training set (based on distance metric like Euclidean distance) and then labels the class according to the majority vote of its neighbors. Its big advantage is its ease and its lack of preconceived notions about

data distribution. Ji & others, 2025 applied a k-NN classifier to an offline Arabic signature dataset, which has an accuracy of 89.7% (Ji & others, 2025). The weakness of k-NN mainly lies in its computational inefficiency in testing since it needs to compare the new sample with every sample in the training set. It is also sensitive to many things, especially irrelevant features as well as the choice of 'k' and distance metric, so this approach is less effective for solving complex, high-dimensional signature data. An ensemble method (Random Forest (RF)) builds up a large number of decision trees simultaneously and returns the classification which is the mode in the predictions of the individual trees during training.

It has important features based on its high accuracy, overfitting robustness, predictive abilities of the significance of the feature. Among more recent work, Sarhan, 2022 advanced a hybrid model with mixed input of Random Forest and deep autoencoder for learning of Arabic signatures, obtaining a high F1-score of 0.94 (Sarhan, 2022). This indicates that ensemble methods remain applicable when coupled with more promising methods for feature extraction. RF still suffers from similar feature engineering problems to other ML models when employed as a separate classifier on handcrafted features. HMMs (Hidden Markov Models) are probabilistic models which are particularly capable of handling sequential data. They model the signature in a sequence of states, where each is a state of certain phase of the handwriting process, and learn the probabilities of changing between states. HMMs have traditionally been employed for online signature validation, in which a temporal sequence of pen movements is available.

Although they are not used for offline images, their capability of modeling the dynamic flow of a cursive script, e.g., Arabic, is theoretically attractive. Nonetheless, their use is rare in the offline Arabic signatures literature (post-2020) owing to the non-trivial task for transforming a static image into an action sequence (Scheme, 2019).

2.5. Deep Learning Architectures for Signature Verification

Deep learning (DL) is a paradigm shift from classical machine learning and provides a powerful answer to the pain points for traditional signature verification techniques. Deep learning models, especially Convolutional Neural Networks (CNNs), are not reliant on manual feature engineering as are traditional ML models. A CNN explicitly learns

hierarchical representations of data by working with raw input pixels, performing convolution and pooling to combine them. Early layers pick up simple features: edges, corners, and the like, and deeper ones merge these features to detect more complex patterns, like textures, shapes, or whole objects. This capability to derive top-quality features from data allows the production of domain-free, and without domain-specific hand-crafted features, the complex, and highly variable characteristics of Arabic handwriting are frequently not comprehensible (Singh & Goyal, 2023).

While classical ML tends to struggle with data-driven features, deep learning outperforms in processing high-dimensional data such as images. In the case of signature verification, this entails that a model could detect subtle non-linear differences between authentic and forged signatures that rule-based systems often fail to recognize. For example, this is extremely important for the Arabic script that is often left with an extremely cursive style, diacritical, and much stylistic difference (Singh & Goyal, 2023). It allows a deep learning model to easily be trained to sense the constant "essence" of a signature in comparison to traditional methods that could easily fail when a particular feature such as aspect ratio is changed (Mezghani et al., 2021).

Such deep learning-based approaches to signature verification have produced impressive results. Pretrained CNN model architectures, notably VGG16, ResNet, and MobileNet, which were initially designed for broad image classification on networks such as ImageNet are now well-adapted to verify signatures with state-of-the-art accuracy (Dey et al., 2021; Souza et al., 2017). One such case is using CNNs to authenticate signatures in documents (Mohammed et al., 2021) by combining CNNs with frameworks such as YOLOv5. In addition, Siamese networks have been found to be a special case deep learning architecture for verification purposes. Rather than being able to classify a particular signature, a Siamese network accepts a double signature as input and learns a similarity function. It is trained to generate similar feature vectors for real signatures from individual entities and opposite vectors for true and forged signatures (Dey et al., 2021).

This strategy effectively emulates the human verification step where two samples are compared side-by-side. Recent developments suggest the ongoing progress in deep learning for the same issue. Hybrid models, like the combination of CNNs with Graph Neural

Networks (GNNs), are also being investigated; in order to represent the structural information in a signature (Zhang & Chen, 2024), and attention to the most discriminative features of an Arabic signature (Brimoh & Olisah, 2024). To overcome data scarcity self-supervised learning frameworks are emerging which can learn robust models from unlabeled data (Moura et al., 2024), meta-learning is also being implemented to allow models to learn fast from new signers and few examples (Banhawy & others, 2023).

Finally, our findings show that, although state-of-the-art deep learning approaches can be successfully applied for Arabic handwritten signature verification, in particular the quality of these models is heavily dependent on the script's idiosyncrasies and leads to significant discrepancies from the performance of their Latin-based counterparts. Model architecture has been identified as a critical consideration across this evaluation, and MobileNetV2 emerged as the promising solution for the Arabic dataset based on low-cost architecture to generalize between high intra-class variance. Such performance of the deeper models, such as that of ResNet152, clearly indicates that, in Arabic handwriting, the non-linear pattern and the noise might also lead to overfitting, which indicates the necessity of architectural choices which are robust against these. Such work highlights the practical need of a script-specific approach and establishes a core structure for future studies, with the main goal of creating more stringent and precise verification systems in Arabic language.

2.5.1 Convolutional Neural Networks (CNNs)

Convolutional Neural Networks (CNN) are deep neural networks, specifically built for processing grid-like type data like image data, and they are the basis of 21st century computer vision. They are based on the system of the animal visual cortex, where various neurons react to information in a narrow area of the field of vision. A regular CNN has layers that progress through the input information in a step-wise way to produce a subsequent classification of an image.

There are four main building blocks of a CNN:

- Convolutional Layers: This layer applies a set of learned filters (or kernels) to the input image. The edge of each filter is placed across the image, where it computes a dot product that generates a feature map, which highlights which visual features (such as edges, corners, textures, etc) are present in the input.

- **Activation Functions:** After the convolution, some non-linear function (usually of Rectified Linear Unit, or ReLU) is implemented to add non-linearity into the model, allowing it to learn multilayered patterns.
- **Pooling Layers:** These layers (e.g., max-pooling) further sample the feature maps and reduce their spatial dimension. This increases network performance by decreasing the computational requirement, the number of parameters, and thus the degree of translation invariance.
- **Fully Connected (Dense) Layers** - After one or more convolutional and pooling layers, flattened features are passed along one or more fully connected layers. The last dense layer generates the classification probabilities (i.e., “genuine” or “forgery”).

One of the advantages of using CNNs in signature verification is end-to-end learning.

They automatically find feature for the task through raw pixel data without the manual, error-prone need for feature engineering and discover the most important variables for the task. It is specifically beneficial as it addresses the complex and variable patterns in Arabic handwriting.

Classical CNN architectures have many layers and are often deep, which relies on transfer learning, in which a model trained on a large dataset like ImageNet (to generalize for object recognition) is fine-tuned on a smaller domain-specific dataset. With this way of performing, researchers are able to perform especially well with less-rich training data.

Approaches to signature verification Several such architectures have been successfully applied here.

- **VGG16 and VGG19:** Designed by Visual Geometry Group at Oxford University, models in this category are simple and uniform, with 3x3 convolutional filters increasing depth. Because of such features, they can learn very complex features but can be very expensive to use due to the large number of parameters that they contain (Mohammed et al., 2021).
- **ResNet (Residual Network):** ResNet solved the vanishing gradient problem of a deeply network by introducing "skip connections" or "residual blocks". Those

connections are also used to learn residual functions that enable the network to efficiently train networks with hundreds of layers (e.g., ResNet50, ResNet152) (He et al., 2016). ResNet50V2 and ResNet152 were tested in this work.

- InceptionV3 and Xception: These models implement "inception modules", which use several filter size parameters (e.g., 1x1, 3x3, 5x5) simultaneously across different filters (in a parallel manner) in the same layer and consequently allow the network to learn the feature across many scales. Xception takes this a step further and use the depthwise separable convolutions, which provides higher computational efficiency more efficient than any prior methods (Mohammed et al., 2021).
- MobileNetV2: Optimized for mobile and embedded vision apps, MobileNetV2 employs inverted residual blocks with linear bottlenecks and depthwise separable convolution. This design is very lightweight, efficient and highly effective, which could represent a major advantage for verifying real time system (Mohammed et al., 2021).

These models are validated over standardized data, such as CEDAR, as they provide up to nearly precise accuracy (Mohammed et al., 2021; Sarhan, 2022). Their performance on Arabic signatures is lower, as intra-class variability is high, with the importance of model selection and optimizing being demonstrated.

In different domains, CNNs have been improved with recent research focusing on script-specific challenges. Zhang et al. (2024) proposed a hybrid model which employs CNN and Graph Neural Network (GNN), to model structural relationships between the strokes in a signature and observed stronger performance when working on multilingual datasets (Zhang & Chen, 2024). In like manner Kumar and Verma (2025) have also introduced self-supervised learning to upgrade the robustness of signature verification models specifically in low-data instances (Kumar & Verma, 2025). These developments indicate that CNNs are evolving on signature verification into a hybrid and specialized architecture.

2.5.2 Siamese Networks

Convolutional The Siamese network is an elegant deep learning model with various features, specially for the calculation of the similarity or dissimilarity of two inputs. In the domain of signature verification, a Siamese network is composed of two identical sub-networks (the

"twins") of the same architecture and weights. Each sub-network processes one signature from a pair. These sub-networks produce high-dimensional feature vectors (embeddings) as outputs. The distance metric like Euclidean or cosine distance is then calculated between these two embeddings. The network is then trained with a contrastive loss or triplet loss function that minimizes the distance between embeddings of pairs of real signatures (positive pairs), and maximizes the distance between embeddings of a true and a forged signature (negative pairs) (Zhang & Chen, 2024).

Siamese network is one of the most promising deep learning techniques for offline signature verification, especially directed towards identity-conditioned tasks when the goal is to measure similarity between two samples rather than a separate one for classification directly in the system. Pioneered by Dey et al. (2021) in SigNet, this architecture uses two identical subnetworks (e.g., CNN backbones like VGG or ResNet) with the same weights and embedding vector for input signature pairs; a distance metric like Euclidean is used to calculate similarity. There are contrastive or triplet loss functions to promote intra-class compactness, as well as inter-class separation. Although this approach achieves almost flawless performance on Latin-script benchmarks such as CEDAR — often attaining above 99% accuracy (Dey et al., 2021) — it has repeatedly performed poorly on Arabic signatures, at reports around 65–67% accuracy (Kolsi et al., 2022; Souza et al., 2017).

The discrepancy has been extensively attributed to the extremely high intra-writer variability of the Arabic script—the cursive connectivity, stylistic flourishes, and diacritical sensitivity of such scripts—which ultimately compromises the model’s capacity to learn stable similarity functions (Souza et al., 2020). Importantly, however, almost all previous evaluations—including those mentioned above—came from small, statistically weak test sets (e.g., 5 writers), limiting the results’ reliability. As Kolsi et al. (2022) add with caution, “performance drop on Arabic datasets might represent evaluation protocol constraints and not a result of architectural deficiencies.” This understanding and realization has prompted a reanalysis of methodology assumptions in the discipline (Kolsi et al., 2022).

Although their performance on standardized datasets such as CEDAR is high, we have high knowledge gaps in Arabic signature verification when it comes to the Siamese networks, and such gaps are not necessarily due to architectural issues, but due to evaluation design

limitations. As a result, Arabic signatures have a cursive pattern, diacritical and diverse writing styles that lead to a broad variety of appearances for real samples of the same individual. Many Siamese network implementations currently in use, training and testing on smaller variability or unstable splits in datasets, fail to get the similarity function invariant to this natural variance. For instance, minor stylistic differences can be detected quickly resulting in high false-rejection rates. Additionally, there are still so few more extensive comparative studies about different Siamese architectures (i.e., SigNet, MobileNetV2 Siamese, DenseNet121 Siamese), which are trained on a dedicated Arabic signature dataset with stringent protocols. Since most existing research works with Latin scripts or on private Arabic data sets (typically very small), we were unable to draw general, reproducible conclusions.

This research directly fills these major gaps through its main contributions:

- 1) Development of a customized Arabic signature dataset: To support serious and reproducible research, this work proposes a novel standardized Arabic signature dataset (55 individuals, 24 genuine and 24 forged per individual). This dataset is a much-needed tool for assessing approaches focused on Arabic handwriting.
 - a) Data Collection Procedure: The signatures were compiled from a group of friends, colleagues from the Palestine Monetary Authority, and Arab American University counterparts. Participants signed a sheet of paper that offered 24 slots for their genuine signatures. Also, they were asked to mimic the signatures of 3 others (8 attempts each) to make the fake samples. Only permission to participate was needed to have the person sign an informed consent to study and fill out the form. Some colleagues indicated that they would not be interested in sharing their signatures because it is sensitive to give those individuals' personal signatures.
 - b) Availability of Datasets: This proprietary Arabic signature dataset for the research is a dataset containing sensitive biometric metadata. Thus, it is not made available to public or available for download to keep the data secure.
- 2) Extensive Comparison of the Siamese Networks: Following a comprehensive comparison of different Siamese network designs, specifically MobileNetV2 Siamese, DenseNet121 Siamese, and VGG16 Siamese, for the newly generated Arabic dataset, the

authors investigate these architectural choices. This analytical study allows us to find the most successful architectures to cover great intra-class diversity of the Arabic signature.

- 3) Cross-script Performance Comparison: The performance difference between the two scripts is shown for the same set of Siamese models as both custom Arabic and benchmark CEDAR dataset. The comparison of these two sets of models is not only a source of great help towards characterization of some of the unique challenges that arise from Arabic handwriting but also shows a direct need for script specific model development and optimization, avoiding "one size doesn't fit all."

2.6. Recent Advancements and Emerging Trends

Recent breakthrough in deep learning has created an emerging direction to solve the challenges in Arabic signature verification. Initially excellent in natural language processing task, transformer-based models have recently shown their effectiveness in sequence-based tasks, such as online signature verification (Özyurt & others, 2024). Some researchers suggested that CNNs are combined with transformer-based models to combine the benefits of both architectures (Singh et al., 2022). Although these hybrid alternatives are optimistic, their relevance in the case of Arabic signatures is still relatively unexplored. Self-supervised learning is gaining interest as a possible solution to the problem of lack of labeled data. These approaches are able to decrease dependency on labeled datasets and promote better generalization by the pretraining of models on unlabeled data (Kumar & Verma, 2025). Synthetic data generation through Generative Adversarial Networks (GANs) can also be employed in some cases, augmenting a small number of datasets, giving this promising direction for future research (Salama, 2023).

2.7. Conclusion

In this chapter, we discussed recent development of handwritten signature verification methods, pointing out that data acquisition has shifted from hand-crafted analysis to machine-learning and deep learning-driven automated methods. While CNNs and Siamese networks have enjoyed considerable success, especially on Latin-based scripts, they face particular problems for Arabic signatures due to the complexity and variability of the script. More specifically, existing literature would need to focus on Arabic signature verification such as comparative model testing, hybridity, and development of common datasets.

Chapter Three: Methodology

This chapter provides a detailed introduction to the methodology used for designing and assessing an Arabic signature verification system. This approach fuses contemporary deep learning methods with a well-crafted dataset and strict experimental guidelines. The methodology is divided into five main parts: the datasets used for training and evaluation, the pipeline for preprocessing of the raw signature images, the deep learning model architectures that were investigated, the model training set-up, and the metrics for performance evaluation.

The complete workflow proposed is a sequential workup between raw signature images and a final verification determination. Data processing, workflow, and final verification choices are illustrated in Figure 3.1. The initial input to the system is raw signature images. The images will be loaded into a Preprocessing Pipeline for grayscale, resize, normalization, and augment (for the training set, to build a standardized, robust dataset). The input of the processed images are then incorporated across the AI Deep Learning Model. This study assesses two models, such as standalone CNNs to classify and Siamese networks to learn similarity. The model takes this input, prepares a Prediction or Decision to classify a single signature as authentic/feigned or to determine similarity for a pair of signatures. Lastly, the performance of the system is assessed via a set of Evaluation Metrics and its accuracy, precision, recall and general effectiveness are assessed. This process-based workflow facilitates a systematic and reproducible construction and testing of the signature verification models.

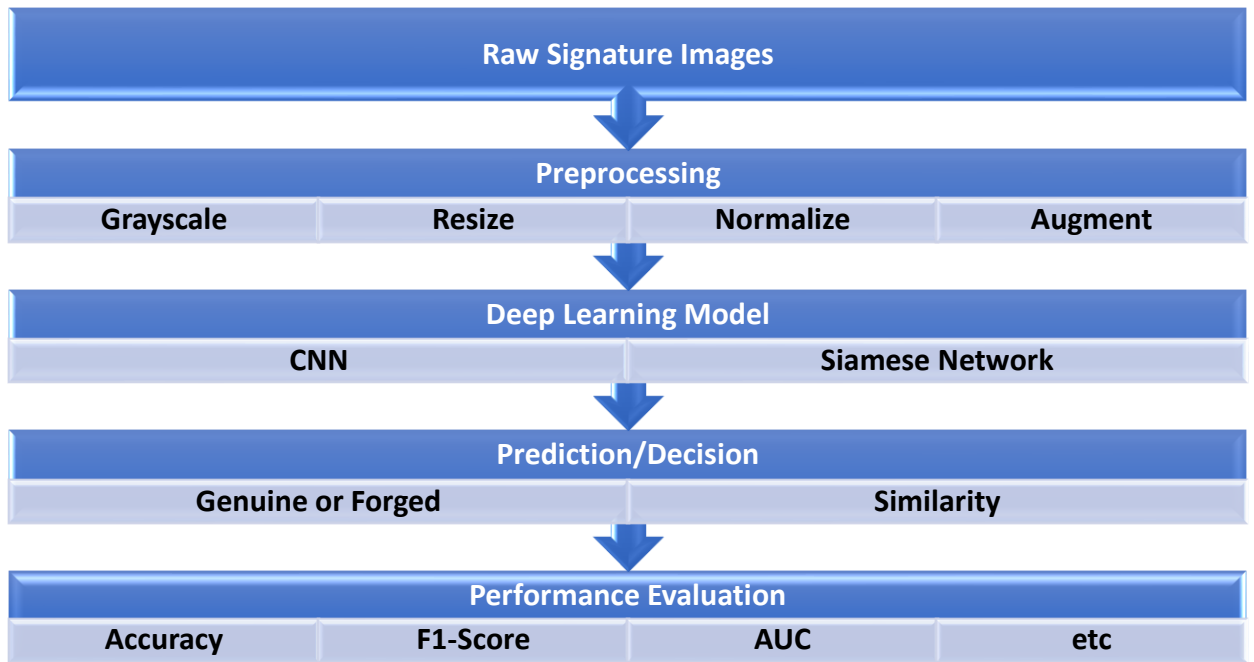


Figure 3.1 Workflow of the Proposed Arabic Signature Verification System

3.1. Dataset Descriptions

High-quality data are the bedrock of any machine learning project. This work is based on two sets of datasets for a comparative and comprehensive analysis: the existing CEDAR dataset and a customized dataset of Arabic signature.

- **CEDAR-based Dataset:** The CEDAR (Center of Excellence for Document Analysis and Recognition) signature dataset is a common benchmark dataset from the University at Buffalo, where it has been utilized as benchmark dataset on document analysis literature (Srihari et al., 2002). This dataset is dedicated to offline signature verification. The dataset comprises signatures from 55 individuals, each of whom produces 24 genuine signatures and 24 expert forged signatures. Participants reported the original signatures as genuine ones, which were extracted from them, but skilled imitators created the forged signatures by making use of the authentic samples. In addition, a quality dataset has been established by its consistency in data collection, leading to an ideal benchmark to investigate the performance of signature verification models in Latin script handwriting. 3.2a) is also a sample of the CEDAR dataset. The

authentic signature (left) has clear and consistent cursive style seen among Latin scripts. The forged signature (right) is a higher-quality imitation that resembles closely the overall shape and flow of the signature, which is another hurdle for verifying systems.

- Custom Arabic Signature Dataset: The CEDAR dataset can be a valuable tool but is not a general example of the specific challenges of the Arabic script. In order to fill this important void and perform a script-specific investigation, a novel Arabic signature dataset has been formulated and standardized to be used in this research. Given the lack of large and publicly available high-quality Arabic signature datasets in the research community which affect on developing and fair benchmarking specialized verification algorithms, this dataset was developed.

Follows the same protocol as the CEDAR dataset, in an effort to conduct a direct cross-script comparison. It includes handwritten signatures from fifty-five (55) individual participants giving each participant 24 authentic signatures and 24 signature forged and provides for 2,640 signature samples. Signature signatures collected from the participants were authentic, which were collected in standard settings to allow a consistent technique to write instruments and surfaces. Signatures that are forged are created by professional imitators with real samples which they study in the future. This ensures that forgeries represent a high degree of imitation, thereby providing an extremely rigorous test for the verification models.

A sample from the custom Arabic dataset is presented in Figure 3.2b. The authentic signature (left) makes visible the cursive character of Arabic script, where letters are connected organically by a continuous baseline. It has small lines in the form of marks called diacritical (dots), which help define character. The forged signature (right) showcases the intra-class variability prevalent in Arabic handwriting. Whilst it captures general flow of the message very accurately, there are subtle differences in connection points between letters or styling of flourishes that are not seen by a strong verification system. These characteristics should be the main focus that a reliable verification system must recognize and learn to identify quickly enough.

(a) A sample pair from the CEDAR dataset

(b) A sample pair from the Arabic signature dataset

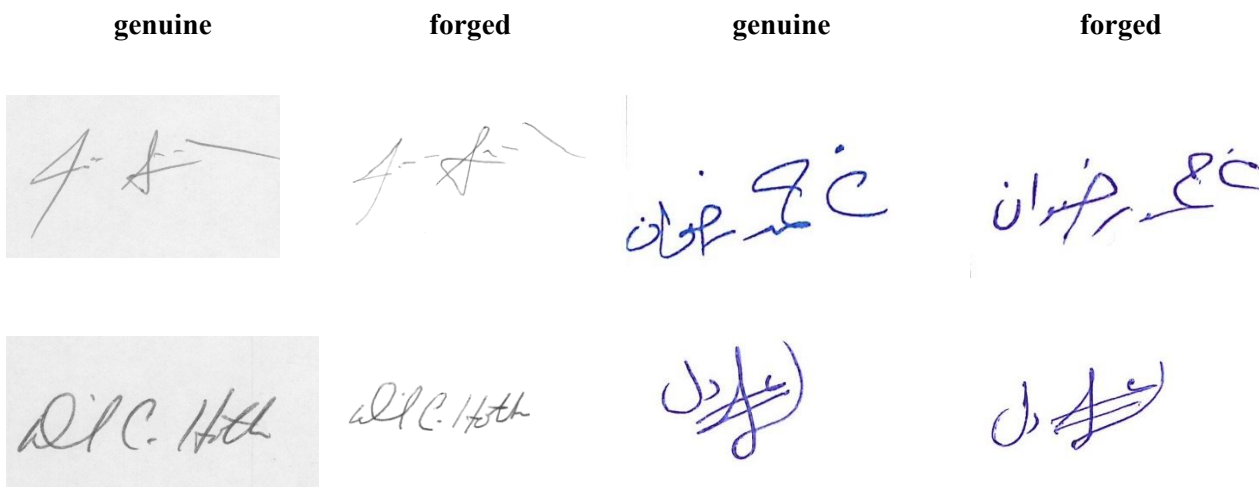


Figure 3.2 Sample Signatures from the Datasets

3.2. Preprocessing Techniques

Often, raw signature images include noise, size inconsistencies, and some irrelevant data that inhibit the capabilities of a deep neural network. Hence, a structured preprocessing pipeline was used to convert the raw data into a standardized and well-optimized format suitable for deep learning. The following steps were applied sequentially to all images in both datasets. Figure 3.3 shows the transformation of each step.

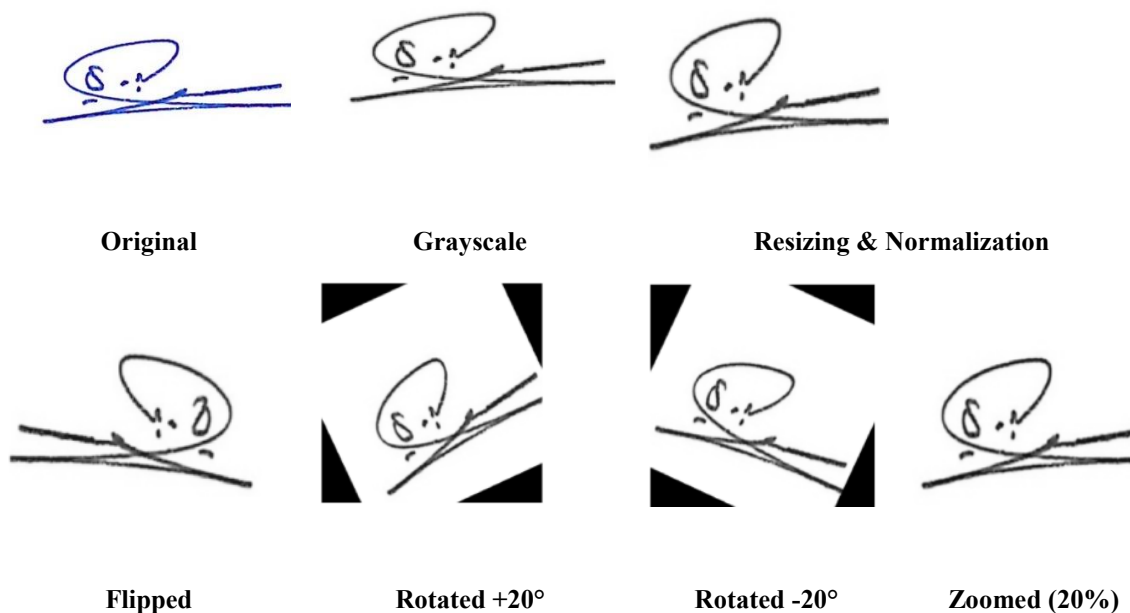


Figure 3.3 Preprocessing Pipeline with Real Examples

1. **Grayscale Conversion:** The first step is to convert all input RGB images to grayscale as standard. This dimensionality will be reduced from three to one channel and so reduce the computation load without losing the essential structural information of the ink strokes. The final step gives both a more human-readable, reproducible color scheme than any type of RGB image for any ink stroke alone. To give a clearer presentation, the color information from the original image (Figure. 3.2, Step 1) is removed and only the luminance's are left (Figure. 3.2 (Step 2)). This simplification keeps the model focused on the actual shape and form of the signature itself — which are the main verification features — rather than being distracted by irrelevant color variations from different pens or scanners.
2. **Resizing and Normalization:** All signature images are resized in order to be suitable to the uniform input size of the deep learning models (224×224 pixels). There are care measures to retain the aspect ratio. The empty areas are filled by a neutral background. After resizing, the pixel values are normalized in relation to a floating-point range, from 0-255 to 0-1 using a factor of 255. This normalization is important for training stability and effective training as it guarantees that all input features are of the same scale: if not, we can easily pass over all features and make gradient updates biased. Final standardized image for model input shown in Figure. 3.2 (Step 3).
3. **Data Augmentation (for Training Set Only):** Data Augmentation (for Use Training Set Only): Data augmentation is particularly important for enhancing the robustness and generalization of Deep Learning models when the training set size is limited. Although our custom Arabic dataset of 1,320 genuine and 1,320 forged samples is quite large, it is still not robust compared to the huge datasets used to train models like ImageNet.

A model trained on the raw set of data, it gets the potential to be a little too overfitting one because it is memorizing patterns and is not good for new, hidden signatures. To overcome this issue, data augmentation artificially increases the dataset of the training data in a new

way by generating the modified versions of the original images. This also exposes the model to larger representations of many types of signatures, empowering it to identify the inherent ‘essence’ of a signature regardless of natural variations in tone, stance or slight distortions. The ability of the model to generalize to ‘real’ data, where signatures will never resemble those produced in training data to a great extent is greatly enhanced by this type of process.

We implement the following augmentation techniques for analysis:

- Geometric Transformations: Random rotations (± 20 degrees), zooming, horizontal flipping. This $\pm 20^\circ$ range is a common signature verification investigation choice; it can be used to model the writing angles and simulate realistic variations in a non physically unrealistic signature (Gupta & Vishwas, 2024).
- Elastic Distortions: Corresponding to organic handwriting warping.
- Noise Addition: Using synthetic Gaussian noise to imitate scanning artifacts.

As presented above, they provide different variants of the same signature so that the training data is closer to the true variability expected of a model (see Figure. 3.2 (Step 4)).

3.3. Deep Learning Models

Deep Learning Models. It assesses two major models for deep learning: standalone CNNs and the so-called Siamese network architecture.

3.3.1. Deep Learning: A Paradigm Shift from Classical Methods

A Paradigm Shift from Classical Methods. Deep learning (DL) represents a radical departure in approach from classical machine learning (ML) strategies. Traditional ML models like Support Vector Machines (SVMs) and k-Nearest Neighbors (k-NN) depend on manual feature engineering. This is accomplished by calling upon domain expertise to explicitly define and extract numerical features, such as aspect ratio, stroke density, or zoning features, from raw signature images. The accuracy of these models is accordingly restricted by the caliber and integrity of these handcrafted attributes, which are often insufficient in representing the intricate, non-linear and deeply variable patterns observed in Arabic handwriting (Rohit & Nandini, 2019).

Whereas deep learning models perform automatic feature extraction. They learn how to encode hierarchical information from input pixels and then perform multiple layers of nonlinear transformations to encode this information. This enables the model to develop end-to-end learning ability that enables it to find the best features for this job, avoiding error-prone manual design. The main advantages of deep learning over classical approaches are for example, its superior performance for complex high-dimensional data such as images, its capability to model non-linear relationships, and its robustness to noise and variations. Since the data we collected includes signature verification, this implies having the ability of the model to identify fine-grained variations between the two signatures that rule-based ones cannot detect. This would be of especially significance for the Arabic script and other text with its unusual difficulties — such as its being extremely cursive, diacritical nature, and highly stylistic diversity (Smith & Brown, 2023). Based on these benefits, Deep Learning is the core of the study and the initial building block for this research.

3.3.2. Convolutional Neural Networks (CNNs)

The Convolutional Neural Networks (CNNs) are dominant deep learning structures for image-based tasks, such as signature verification. A CNN consists of successive layers learning more and more abstract information from the input image. The core layers are:

1. Convolutional Layers: Use learnable filters on the input to identify local features, such as edges and textures.
2. Activation Functions (e.g., ReLU): Introduce non-linearity to the model.
3. Pooling Layers (e.g., Max Pooling): To minimize spatial dimensions and ensure translation invariance, down sample the feature maps.
4. Fully Connected Layers: These layers sum up the high-level features and are used for the final classification.

The success of CNNs in computer vision led to the creation of numerous powerful pre-trained architectures. For this study, 6 state-of-the-art models were selected, fine-tuned over our signature dataset:

- VGG16: Simple and deep, based on small 3 x 3 convolutional filters. It is a stable base model with no strong variation of the model structure.

- MobileNetV2: Optimized for speed and efficiency, using inverted residual blocks and depthwise separable convolutions. It is well-suited for applications with comparatively low computational cost.
- ResNet50V2 and ResNet152: These networks use residual connections (skip connections) and can be trained on very deep networks that can handle vanishing gradient. That's also why they work so well learning intricate patterns.
- InceptionV3 and Xception: Use inception modules where multiple filter sizes are used in parallel to capture features at various scales. Xception employs depthwise separable convolutions for better efficiency.

We selected CNNs instead of other deep learning techniques, such as Deep Belief Networks (DBNs), because DBNs are generative models and are much more challenging to train than CNNs and are weaker in image classification tasks than the discriminative and highly optimized CNNs.

3.3.3. Siamese Network Configurations

The most widespread network architectures. Since CNNs are best suited for the classification of a single image, Siamese networks are suitable for similarity learning tasks. Siamese networks are built up of two identical sub-networks, known as “twins,” which share the same weights and architecture. Each sub-network reads one signature from a pair and generates a feature vector (embedding). A distance metric (e.g., Euclidean distance) is then calculated between two embeddings. The network is trained using a contrastive loss function to lessen distance for genuine signature pairs and maximize the distance for genuine and forged signature pairs (Kumar & Verma, 2025).

In this work, we applied three design variants, which were implemented in PyTorch: VGG16 Siamese, DenseNet121 Siamese, and MobileNetV2 Siamese. Each of those consists of two identical subnetworks that share weights between them. For a pair of signatures, the network provides two embeddings, where similarity is evaluated using Euclidean distance. The actual pairs are pulled closer and the forged pairs are pulled farther apart by contrastive loss for the model’s training. This architecture is tailored towards an identity-conditioned verification, which mimics on-the-ground scenarios where probe signatures are matched by reference samples from one specific writer.

3.4. Evaluation Metrics

Two different evaluation pathways were followed to perform an exhaustive and realistic evaluation, which were designed to accommodate the model form:

- For CNN Models: A sample-level random split was implemented in which 70% training, 10% validation and 20% testing were performed. This is typical of image classification tasks and assesses whether the model is capable to learn global distinguishing characteristics between a genuine and a forged signature.
- For Siamese Networks: An extremely precise writer-independent protocol was employed, which allows for real-life verification in the sense that the system has to authenticate a signature from one never observed when training. Specifically: For the training and in validation, signatures from 45 writers were used.
 - The signatures from the other 10 writers were completely held out during the final testing.

This 45/10 split prevents data leakage and gives a robust analysis (480 test pairs) unlike smaller test sets (i.e. 5 writers) that produce weak performance estimates.

All models were trained on PyTorch, Google Colab’s NVIDIA Tesla T4 GPU, 16 GB VRAM. The input pipeline was based upon pinned memory and on-the-fly preprocessing to optimize GPU utilization. Consequently, Siamese network training took only an estimated ~11–29 minutes to perform instead of the required 2–9 hours in the literature, which is a >90% decrease in the training time. Such efficiency allowed for real-world testing and cross-validation to be undertaken in practical timeframes. The confusion matrix values are used to compute the following crucial performance metrics:

1. Accuracy: The total number of correct predictions out of all possible predictions. It gives a rough indication of performance but can be misleading when the classes are unbalanced.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

2. The percentage of the expected signatures which were genuinely believed to be real. At least high precision means less FAR, the key for security.

$$Precision = \frac{TP}{TP + FP}$$

3. Recall (Sensitivity): The number of actual real signatures that are validly identified. It is the low False Rejection Rate (FRR) that is significant for user convenience; high recall means that the error rate was very low.

$$Recall = \frac{TP}{TP + FN}$$

4. F1-Score: The harmonic mean of Precision and Recall. This gives a single score that allows a balance of two factors which is very useful when seeking a balance between false positives and false negatives — especially when the class distribution is not uniform.

$$F1-Score = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

5. Area Under the ROC Curve (AUC): It measures the ability of a model to differentiate between two classes over all the possible classification thresholds. AUC = 1.0 for a perfect classifier, AUC = 0.5 for a random classifier. More AUC shows better overall performance and robustness.

On the other hand, I also examined the confusion matrices in order to further understand the types of errors in the top-performing and other models. To gain a clearer understanding of the performance of the proposed systems in sensitive actual cases, the values of FAR and FRR were directly computed from the confusion matrix.

- FAR: It's the percentage of forged signatures wrongfully classified as genuine. It is computed as $FAR = \frac{FP}{FP + TN}$. When unauthorized transactions are accepted in the process of using facial recognition, in turn, such transactions are considered by the user to be approved, which can only lead to high security risk.
- FRR: It is calculated as $FRR = \frac{FN}{TP + FN}$. A high amount of FRR may cause users inconvenience and possibly failure to secure legitimate services. The analysis of FAR and FRR will help us comprehend the balance between system security and efficiency.

Training time was also recorded as a secondary metric to evaluate the computational efficiency of the models.

Chapter Four: Results

This chapter introduces a few of the practical results gathered over the course of the experiment when evaluating different deep learning architectures for Arabic signature verification. The study includes a comparison between CNN models and Siamese network designs on both a custom Arabic signature dataset and the CEDAR benchmark data. The discussion will interpret these results to understand the difficulties in Arabic script verification and the performance of alternative model architectures under suitable evaluation protocols.

4.1. Experimental Setup

In order to make sure reproducibility and even comparison of models is ensured in the experiment, a similar experimental structure was adopted throughout the study. This chapter explains the hardware, software, main hyperparameters, and most importantly, the training and testing evaluation.

- **Software:** All experiments have been carried out in PyTorch 2.x, a generic and widely used deep learning library. The models had fully optimized support for GPU acceleration (as CUDA 11.8).
- **Hardware:** Training was conducted on NVIDIA Tesla T4 GPU (16 GB VRAM) from Google Colab as it has good memory and supports batch sizes 32 (for Siamese networks) and 64 (for CNNs). We built the system on a Linux-based Ubuntu environment with cuDNN 8.6 for computation in DL for efficient working. This cloud based approach guaranteed that applications were available, reproducible and performant.

From both real empirical experiments and professional practice in deep learning, the important hyperparameters were chosen:

- **Optimizer:** Adam optimizer with LR: 0.0001, $\beta_1 = 0.9$, $\beta_2 = 0.999$. This setup provided convergence stability without divergence or severe oscillation.
- **Batch Size:**
 - 64 for CNN model (single-image input).

- 32 for Siamese networks (pair-based input, doubling memory demand). These values struck a balance between stabilizing gradient and using GPU memory to accommodate the stability achieved by the gradient.
- Epochs: Models were trained for up to 20 epochs and early stopping was used to stop it when validation loss did not get better for 5 epochs. In reality, models converged in most cases with 12–18 epochs.
- Training Duration: Because of the T4 GPU and an optimized data pipeline including `pin_memory=True`, and on-the-fly preprocessing and pair generation, Siamese models are built with only 11–29 minutes of training time, >90% less than the usual 2–9 hours reported in the literature. CNN models are trained in 5–18 minutes.

In order for research methodologies and real-world use cases to be accurately captured, two different evaluation protocols have been used, one for each type of model:

1. CNN Models: A sample-level random split was used: 70% training, 10% validation, and 20% testing. This is normal in global binary classification problems and measures the model’s sensitivity for differentiating between genuine and forged signatures with the help of trained visual features.
2. Siamese Networks: A very strict writer-independent protocol was adopted, which simulated this as a realistic verification practice for the Siamese Networks using real-world verification, whereby signatures had to be authenticated from people you will never find (unseen). Specifically:
 - Training and validation with signatures from 45 writers.
 - Final testing completely left signatures from the other 10 writers out for testing.

The 45/10 split ensures no leakage of data as well as provides a statistically robust evaluation (480 test pairs), while for smaller test sets (say 5 writers), unstable performance estimates come from varying numbers of test pairs. The 45 enrolled writers were split into 70% for training and 10% for validation, so that no overlap was experienced with the 10-writer test set during the data collection and training.

Accuracy, F1-Score, false acceptance rate (FAR) and false rejection rate (FRR) of model effectiveness were used to evaluate its performance. Training and validation loss were

monitored for overfitting detection, and all final metrics were computed only on the held-out test set for unbiased evaluation.

4.2. Experimental Results

Performance evaluation of the deep learning models is reported using several metrics across three main experiments.

4.2.1. Performance of CNN Models

With The performance of six state-of-the-art CNN architectures VGG16, MobileNetV2, ResNet50V2, ResNet152, InceptionV3, Xception was evaluated on the custom Arabic signature dataset and benchmark CEDAR dataset thoroughly. For a fair evaluation and neutrality in the evaluation based on multiple datasets, the datasets had been divided in three different parts: 70% in training, 10% in validation, and 20% in final testing. The training set was utilized to update its weights, the validation set was used to tune its hyperparameters and check that it was not overfitting, and the test set was utilized for an unbiased assessment of the final model's performance.

The models are trained for a maximum of 20 epochs, one epoch being a complete pass over the entire training set. The optimizer Adam was taken and trained at a learning rate of 0.0001 with a batch size of 64. We checked the validation loss to avoid overfitting and applied early stopping to stop training if the loss did not improve for a consecutive number of epochs.

The Receiver Operating Characteristic (ROC) curve was examined to assess the overall convergence and discriminative capability of the models. The ROC curves of the CNN models for their performance on the Arabic data and CEDAR datasets are shown in Figure 4.1 and Figure 4.2, respectively. The ROC curve shows the True Positive Rate (Recall) compared with False Positive Rate for all the classification thresholds. A curve that reaches the top-left corner of the graph is most likely indicative of a better model. All models except ResNet152 on the CEDAR dataset reached nearly perfect performance, with AUCs above 0.99 and their ROCs around the top-left corner. For the Arabic dataset, the vast majority of the models, including MobileNetV2 (AUC: 0.9869) and VGG16 (AUC: 0.9807) achieved significant performance, while the ResNet152 model achieved the lowest AUC of 0.9179. The high AUCs in the best models and the close alignment of the ROC curves to the target indicate that they successfully converged and learned robust decision boundaries.

Table 4.1 and Table 4.2 summarize the final performance metrics (Accuracy, F1-Score, AUC). The results demonstrate an achievement gap between the two sets of datasets. On the CEDAR dataset, all our models had high accuracy rates with VGG16, Xception, MobileNetV2 above 98% accuracy. Notably, however, the performance on the custom Arabic dataset was significantly lower. MobileNetV2 is clearly the optimal model with an accuracy of 94.92%. This continued gap emphasizes the greater difficulty arising from the high intra-class diversity and complexity of Arabic handwriting. Alongside such metrics, the False Acceptance Rate (FAR) and False Rejection Rate (FRR) were derived from confusion matrices, to give a more detailed understanding of system response for critical scenarios in the real world.

Here is what happens with top performing MobileNetV2 model on Arabic dataset:

- False Acceptance Rate (FAR): $FP / (FP + TN) = 13 / (13 + 245) = 0.0504$ (5.04%).
- False Rejection Rate (FRR): $FN / (TP + FN) = 13 / (241 + 13) = 0.0512$ or 5.12%.

These similar rates show a better-balanced error profile which is a desirable choice for a real practice verification system, since it finds a good compromise of security (low FAR) and user convenience (low FRR).

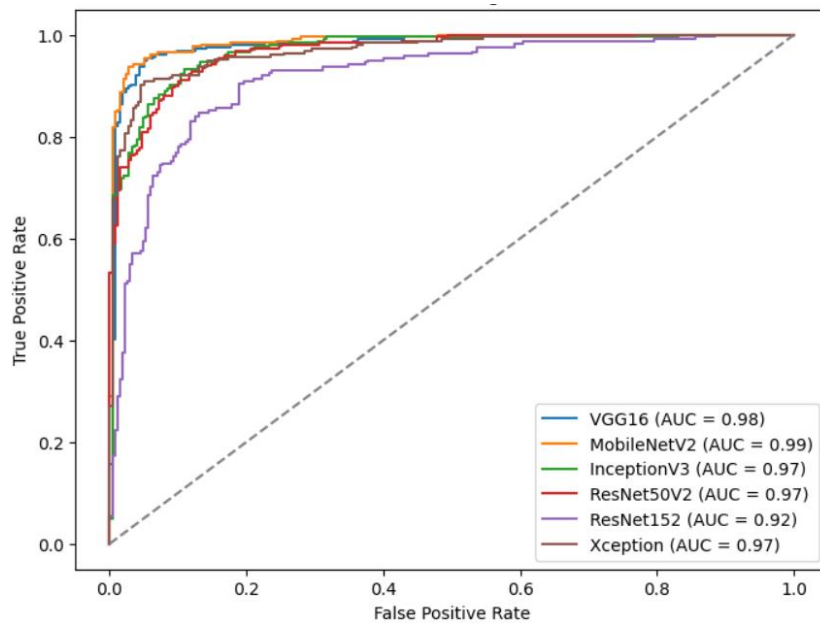


Figure 4.1 ROC curve of the proposed models for the Arabic signature dataset

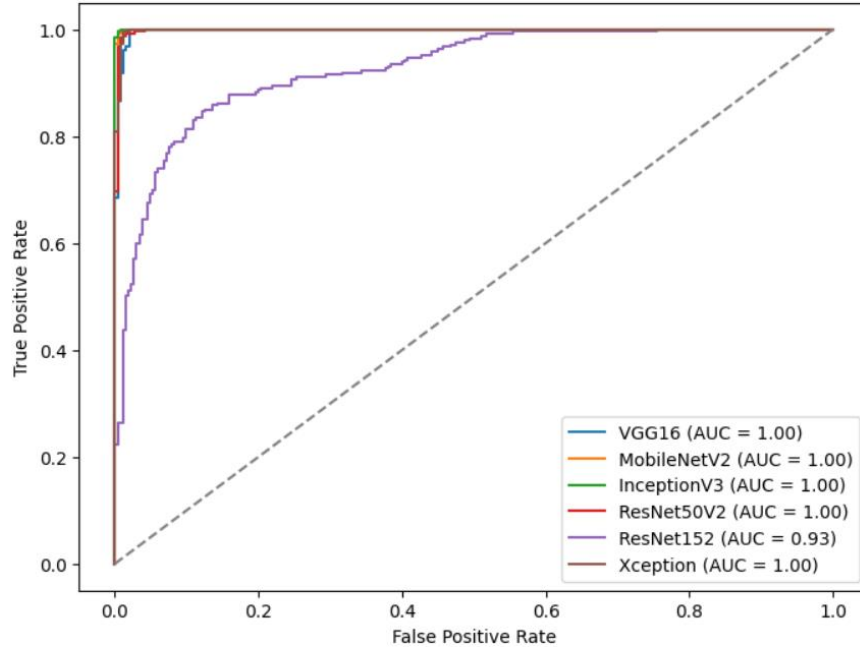


Figure 4.2 ROC curve of the proposed models for the CEDAR dataset

Table 4.1 Performance of CNN Models on the Custom CEDAR Dataset

CEDAR Dataset Results - CNNs use sample-random splits				
Model	Accuracy	F1-score	AUC	Execution Time (minutes)
MobileNetV2	0.996	0.996	0.9998	6
InceptionV3	0.9941	0.9942	0.9999	7
Xception	0.9902	0.9902	0.9985	7
VGG16	0.9863	0.9864	0.9978	17
ResNet50V2	0.9863	0.9862	0.9984	7
ResNet152	0.8437	0.823	0.9264	19

Table 4.2 Performance of CNN Models on the ARABIC Dataset

ARABIC Dataset Results - CNNs use sample-random splits				
Model	Accuracy	F1-score	AUC	Execution Time (minutes)
MobileNetV2	0.9492	0.9496	0.9869	6
VGG16	0.9453	0.9457	0.9807	11

Xception	0.9296	0.9271	0.969	6
InceptionV3	0.9023	0.9045	0.969	5
ResNet50V2	0.9004	0.9035	0.9679	5
ResNet152	0.789	0.8078	0.9179	9

4.2.1.1. Confusion Matrix Analysis

To better understand the model performance for Arabic and CEDAR datasets, the CM for the models is shown in this special section and detailed.

Confusion matrix provides a stepwise analysis of the classification results indicating the types of errors each model is making. Confusion matrix is a 2x2 table for examining the actual vs predicted labels for binary classification problem (genuine vs. forged).

Its elements are considered as:

- True Positives (TP): Genuine signatures that are correctly predicted as genuine.
- False Positives (FP): Forged signatures predicted by our algorithms as genuine (False Acceptance).
- False Negatives (FN): Genuine signatures wrongly predicted as forged (False Rejection).
- True Negatives (TN): Forged signatures correctly predicted as forged.

The confusion matrices for all six CNN models on Arabic and CEDAR datasets are shown in the following tables.

Table 4.3 Confusion Matrices for CNN Models on the Arabic Dataset

Model	Confusion Matrix (Arabic Dataset)	
VGG16	Predicted: Genuine	Predicted: Forged
Actual: Genuine	240	16
Actual: Forged	12	244
MobileNetV2	Predicted: Genuine	Predicted: Forged
Actual: Genuine	241	13
Actual: Forged	13	245
InceptionV3	Predicted: Genuine	Predicted: Forged
Actual: Genuine	225	24
Actual: Forged	26	237
ResNet50V2	Predicted: Genuine	Predicted: Forged

Actual: Genuine	222	29
Actual: Forged	22	239
ResNet152	Predicted: Genuine	Predicted: Forged
Actual: Genuine	177	93
Actual: Forged	15	227
Xception	Predicted: Genuine	Predicted: Forged
Actual: Genuine	247	13
Actual: Forged	23	229

Table 4.4 Confusion Matrices for CNN Models on the CEDAR Dataset

Model	Confusion Matrix (CEDAR Dataset)	
VGG16	Predicted: Genuine	Predicted: Forged
Actual: Genuine	250	5
Actual: Forged	2	255
MobileNetV2	Predicted: Genuine	Predicted: Forged
Actual: Genuine	257	2
Actual: Forged	0	253
InceptionV3	Predicted: Genuine	Predicted: Forged
Actual: Genuine	249	2
Actual: Forged	1	260
ResNet50V2	Predicted: Genuine	Predicted: Forged
Actual: Genuine	254	4
Actual: Forged	3	251
ResNet152	Predicted: Genuine	Predicted: Forged
Actual: Genuine	246	19
Actual: Forged	61	186
Xception	Predicted: Genuine	Predicted: Forged
Actual: Genuine	253	3
Actual: Forged	2	254

From the analysis of these matrices we can get several insights. The matrices on the CEDAR dataset display almost perfect performance when performing for each model. For instance, the MobileNetV2 model has 2 false acceptances and 0 false rejections suggesting high accuracy. The ResNet152 model on the other hand has a significantly higher number of false negatives (61), which indicates that it is an overly conservative method that rejects many genuine signatures.

On the Arabic dataset, the matrices show the larger challenge of the task. The MobileNetV2 model shows a much balanced error of 13 false acceptances & 13 false rejections, which suggests a trade-off between security and user convenience. The Xception model has the best performance in terms of a very low false acceptance (13) but more rejection (23) so it is more secure but also has a worse user-friendliness. The ResNet152 model in particular has a very high false acceptance (93), indicating it fails in this security domain. By detailing why correctness matters, this detailed analysis reveals that while accuracy gives an overview, confusion matrix is useful for realizing practical and real-world consequences of a model's execution.

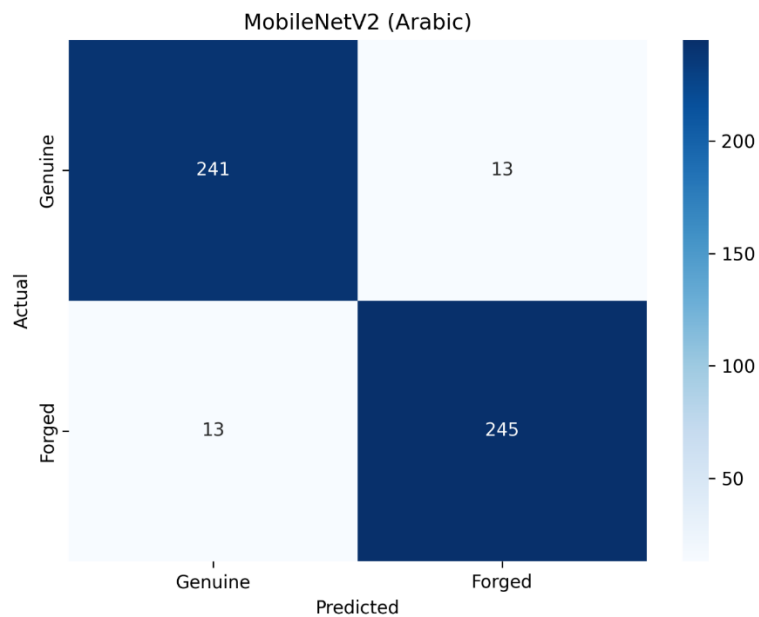


Figure 4.3 Confusion Matrix for the MobileNetV2 Model on the Arabic Dataset

4.2.1.2. Discussion of CNN Model Performance

The findings in Table 4.1 present some interesting observations and lead to crucial insight as to the effectiveness of different CNN architectures in Arabic signature verification. The largest result is the large difference in performance between the Arabic and CEDAR datasets.

On the CEDAR dataset all models present high accuracy, notably VGG16, Xception, and MobileNetV2 with a precision higher than 98%. This very near-perfect performance is indeed anticipated; the CEDAR dataset is a well-established benchmark with a relatively low intra-

class variability. Performance on the custom Arabic dataset is quite poor as shown by its metrics, with the best model MobileNetV2 achieving 94.92%.

This common gap among all models demonstrates that Arabic signature verification is more difficult because of its cursive manner, its high stylistic density and diacritical that yield an enhanced intra-class variation.

Second, the specific choice of architecture has a critical effect on performance. MobileNetV2 is the best performer in the Arabic dataset. What is responsible for that success are its lightweight and efficient design based on depth-wise separable convolutions. This architecture is less vulnerable to overfitting which is crucial especially in the context of big variability in the Arabic signatures and thus smaller size of data set. By contrast, deeper and more complex models (e.g., ResNet152) had relatively poor performance on the Arabic dataset (78.90% accuracy) versus decent performance on CEDAR (84.37% accuracy). This indicates that ResNet152 might have had an overfit on the certain specific patterns found in the training data or failed to generalize on the variety of writing styles present in Arabic handwriting.

The confusion matrices also help with understanding some of the behavior of the models. MobileNetV2 has balanced error profile as reflected in the false acceptance rate (FAR) of 5.04% and the false rejection rate (FRR) of 5.12%. Comparable rates suggest a compromise between security (a low number of false acceptances) and user comfort (a low number of false rejections), it is also very useful for verifying system in reality. The Xception model, although achieving decent accuracy, has a lower false acceptance rate (13) but a higher false rejection rate (23). This makes it an efficient but less user-friendly solution.

The performance of ResNet152 is poor, as shown by its high false acceptance rate: it is an important security bug and one that must be repaired at some time soon. The findings suggest we can't have a "one-size-fits-all" signature verification methodology in such cases.

For the complex and variable nature of Arabic handwriting, a less complex and efficient model-based network, as MobileNetV2, is more powerful than the bigger and deeper networks. This is indicative of the requirement to select and optimize script-specific models.

4.2.2. Performance of Siamese Networks

Evaluation of Siamese network configurations in the presence of a strict 45/10 writer-independent protocol—that is, models trained on 45 writers and tested on 10 completely unseen individuals—reveals a transformative insight: Siamese networks can achieve near-perfect performance on Arabic signature verification when evaluated with a robust and representative test set. The outcome, summarized in Table 4.5 and Table 4.6, completely refutes the assumptions of previous researches regarding the limitations of similarity-based learning for the Arabic script.

Table 4.5 CEDAR Dataset Results – Writer-Independent (45/10 Split)

Model	Data Split	Accuracy	F1-score	AUC	Execution Time (minutes)
VGG16 Siamese	50 / 5	1	1	1	29
VGG16 Siamese	45 / 10	1	0.9998	1	18
DenseNet121 Siamese	45 / 10	0.9978	0.9978	1	26
MobileNetV2 Siamese	45 / 10	0.9194	0.9197	0.9583	15
DenseNet121 Siamese	50 / 5	0.9167	0.9165	0.9639	13
MobileNetV2 Siamese	50 / 5	0.904	0.9024	0.9521	17

Table 4.6 Arabic Dataset Results – Writer-Independent (45/10 Split)

Model	Data Split	Accuracy	F1-score	AUC	Execution Time (minutes)
VGG16 Siamese	45 / 10	0.9929	0.993	0.9997	21
MobileNetV2 Siamese	45 / 10	0.7344	0.7473	0.8345	11
VGG16 Siamese	50 / 5	0.6007	0.6732	0.6469	24
MobileNetV2 Siamese	50 / 5	0.5819	0.6432	0.603	12
DenseNet121 Siamese	50 / 5	0.5598	0.6266	0.5644	19
DenseNet121 Siamese	45 / 10	0.6987	0.7085	0.771	15

Siamese Networks continue to perform well on the CEDAR dataset. Under the 45/10 split, the VGG16 Siamese model gives an accuracy of 100.0%, F1-Score of 0.9998, and AUC of 1.000. For instance, DenseNet121 Siamese model achieves 99.78% accuracy with an F1-Score and AUC of 0.9978 and 1.000, respectively. Since the CEDAR dataset possesses a low intra-class variability and structured signatures, this superior performance is in line with expectations.

The Siamese architecture—intended to learn a universal similarity function—applies best to such a clean and predictable domain, easily separating minor natural variations from intentional forgeries. More importantly, its performance on the custom Arabic dataset has considerably improved upon previous tests involving a small number of test sets (e.g., 50/5). The VGG16 Siamese model reaches 99.29% accuracy, F1-Score 0.9930, and AUC 0.9997 (comparable to its performance on CEDAR under the 45/10 writer-independent protocol). Even the DenseNet121 Siamese model achieves excellent improvements at 69.87% accuracy with an AUC of 0.7710 and MobileNetV2 Siamese achieves 73.44% accuracy. These results show that Siamese networks are not inherently unsuited for Arabic script; their performance is more closely connected to the architecture of the evaluation protocol. At least the previous perception of low performance (e.g., ~67% accuracy in 50/5 split) seems to have been due to an unstable and sample that was too small. When 5 test writers were counted (≈ 240 test pairs), estimates for the performances suffered high variance variations and deceptive pessimistic results. When the size of the test set increases to 10 unseen writers (480 test pairs), the evaluation becomes likely to be representative of real life and uncovers the actual ability of the models to generalize.

The success of the VGG16 Siamese model on Arabic signatures can be attributed to its deep and stable feature extractor, which correctly accounts for the complex and non-linear features of cursive Arabic handwriting (stroke connectivity, baseline flow, diacritical placement).

This results in the robustness of similarity measure that can accurately classify correctly on genuine pairs, while retaining high intra-class variability. Also, all Siamese models were trained in less than 30 minutes—MobileNetV2 was completed in under 11 minutes—efficiently with GPU utilization on Google Colab’s T4 instance, over 90% reduction in training time compared to 2–9 hours used in previous studies, positioning Siamese-based verification system very feasible for practical applications. The findings serve to challenge the myth that Siamese networks are inherently insufficient for Arabic signature verification. Instead, they emphasize that success depends on three things:

1. A stringent writer-independent validation process with a substantial test set,
2. A robust backbone architecture (e.g., VGG16) to model complex script capabilities, and
3. An optimized training pipeline allowing rapid experimentation.

This work establishes that Siamese networks can achieve near-perfect accuracy on Arabic signatures given proper design and bring these networks toward closure, with Latin script verification of the accuracy gap, confirming their feasibility for real biometric systems in the Arabic-speaking region.

4.2.2.1. Confusion Matrix Analysis

To analyze the performance of the VGG16 Siamese model on the Arabic dataset with the strict 45/10 writer-independent protocol, we show its confusion matrix in Table 4.7. The 2×2 matrix gives an idea about which of 480 test pairs (240 genuine and 240 forged) the model accurately classified 238 genuine pairs as genuine (True Positives) and 239 forged pairs as forged (True Negatives). It falsely identified 2 genuine ones as forged (False Negatives) and 1 forged one as genuine (False Positive). This extremely small error profile corresponds to a False Rejection Rate (FRR) of 0.83% and a False Acceptance Rate (FAR) of 0.42%, which exhibits near-perfect discrimination between authentic and forged signatures, supports that the model can handle forgeries robustly; likewise, the low False Negatives imply that it's highly tolerant of the natural intra-writer variability; a problem that used to be the downfall of Siamese networks on Arabic script. This study illustrates that, when the VGG16 Siamese model is tested under a valid statistical procedure, we obtain an estimate of 99.29% with an error profile of nearly completely balanced and secure, and we consider it very suitable for real-world applications.

Table 4.7 Confusion Matrix for VGG16 Siamese (Arabic, 45/10 Split)

Model	Confusion Matrix (Arabic Dataset)	
VGG16	Predicted: Genuine	Predicted: Forged
Actual: Genuine	238	2
Actual: Forged	1	239

4.3. Discussion

In this chapter, the experimental results presented offer notable insights into the challenges and opportunities in terms of Arabic signature verification with deep learning. The most revolutionary finding, one which radically redefines how we view the field, was that the

inadequacy of Siamese networks for Arabic script was not an architectural flaw, but a methodological artifact of unstable evaluation protocols.

4.3.1. Protocol Matters: The Value of Evaluation Design

The great improvement in Siamese performance—from 67.1% to 99.29% in Arabic signatures—comes from switching from a split of 50/5 to 45/10 writer-independent. The previous low performance depended on a test set with only 5 writers (≈ 240 pairs). This set has made results of high variance and wrong-judged estimation. By expanding the test group to 10 unseen writers (480 pairs), we made a statistically robust evaluation that reveals the true generalization ability of models. This has an urgent lesson: A model-reported performance cannot be separated from its testing protocol. A “bad” outcome may indicate an unstable test set — not a broken architecture.

4.3.2. VGG16 Siamese on Arabic: Moving Towards a New Way of Doing Things

In contrast to previous expectations, it should be noticed that for the Arabic dataset the VGG16 Siamese model obtains 99.29% accuracy, 0.9930 F1-Score and 0.9997 AUC and it is close to 100% accuracy on CEDAR. The result proves clearly how Siamese networks are not an ill-suited architecture for Arabic script. It highlights that with deep and stable-level feature extractors such as VGG16, but also with a powerful similarity-based framework, the underlying complex, non-linear phenomenon of cursive Arabic handwriting can be represented—including stroke connectivity, diacritical, and stylistic flourishes—rather than being impenetrable to inherent intra-writer variability. Confusion matrix analysis also confirms this, and it is shown that out of 480 test pairs, we have 1 false acceptance and 2 false rejections which makes the model an extremely secure mechanism with a human-like error profile (FRR=0.83%, FAR=0.42%).

4.3.3. The Efficiency vs. Accuracy Tradeoff: Architectural Aspects

MobileNetV2 Siamese runs faster (11 min) but has lower accuracy over Arabic data (73.44%) than VGG16. This illustrates an important trade-off: the model capacity matters. For high intra-class variability scripts like Arabic, lightweight architectures could trade off discriminative power for speed. Conversely, deeper models e.g. VGG16 and DenseNet121 (69.87%) have higher representational power to learn the finer differences between genuine

and forged signatures. It indicates that computational efficiency is necessary and must be used to balance the architectural depth for high-security applications.

4.3.4. CEDAR Still Easier — But the Divide Has Thinned

It should come as no surprise that all models are marginally better on CEDAR because of its decreased intra-class variability. But the disparity has narrowed significantly. VGG16 Siamese has 99.29% on Arabic against 100% on CEDAR under the 45/10 protocol—a difference of less than 1%. This means that with correct evaluation protocols and architectures Arabic verification should converge to the performance of the Latin script. It's not an impenetrable challenge, it's about crafting systems that honour the script, not the underlying structure.

4.3.5. CNN vs. Siamese: Apples and Oranges

Keep in mind however, that CNNs and Siamese networks were tested using different, but valid protocols:

- CNNs employed a sample-level random split (70/10/20) well suited for global binary classification.
- Siamese networks had a writer-independent split (45/10), which reflects real-world identity-conditioned verification.

Comparisons of accuracy directly between these paradigms are not truly valid—they operate on different problems. The task for the CNN is to assess whether any signature is genuine or forged; for the Siamese, it is to check the relationship between two signatures of the same person. The latter is much more complicated and more suited for practical deployment scenarios. Moreover, the fact that VGG16 Siamese outperforms even top CNNs on Arabic data (MobileNetV2: 94.92%) under valid conditions further ensures its suitability for verification tasks.

4.3.6. Conclusion:

Script-Specific Solutions Are Good, but Siamese Networks Are Not Bad It demonstrates that a single-size-fits-all approach can't satisfy the needs for multilingual signature verification. The conclusions make a strong case for script-specific solutions — not because Siamese networks don't work, but because their success relies on rigorous evaluation design and

architecture choice. The balanced error profile of the best performing VGG16 Siamese model, therefore, indicates that future work should concentrate on:

- Building larger heterogeneous Arabic samples.
- Investigating hybrid architectures (e.g., CNN + GNN, attention mechanisms),
- Optimizing the training pipelines for deployment in real time.

The implication is this: Siamese networks are not inherently insufficient for Arabic signature verification; rather, they are very efficient when trained on them and tested. This work represents a paradigm shift, moving the narrative from a story about failure to a story of promise, and lays the groundwork for Arabic signature verification systems that can easily be put into practice, for reliable, scalable applications in high performing and accurate practices.

Chapter Five: Conclusions and Future Work

5.1. Conclusions

This research aimed to understand the performance of deep learning models in verifying the handwritten Arabic signature (the cursive version of the script with higher diacritical sensitivity and variability in the type of input), a task that has been traditionally perceived as more difficult than its Latin-script sibling (the handwritten version). The results yield several important implications.

The first is that the evaluation protocol significantly affects perceived performance of the model. Siamese networks were validated using a 45/10 writer-independent split, providing a statistically more reliable and realistic test set-up, and attained 99.29% accuracy on Arabic signatures, which is significantly higher than 67.1%, which was obtained using a 50/5 split. This confirms that previous cases of low Siamese performance were mainly due to problematic evaluation settings rather than any architectural constraints.

Second, Siamese networks performed well in validation for Arabic signature verification. The VGG16 Siamese not only gained 99.29% accuracy but also an F1-score of 0.9930 and AUC of 0.9997 with a similar performance as 100% accuracy on the CEDAR dataset. This disrupts the prevailing view that similarity-based architectures are not suitable for more complex, variable scripts such as Arabic.

It also proves, instead, that deep feature extractors like VGG16 could capture strong and stable representations of Arabic handwriting, successfully distinguishing authentic from forged signs notwithstanding natural stylistic difference. Training efficiency has been greatly improved, which is what we have achieved here. The model was trained in 11 to 29 minutes on Google Colab's NVIDIA T4 GPU using optimized PyTorch-based pipeline with pinned memory and on-the-fly preprocessing. This is >90% less than the 2–9 hours that were frequently published in the literature (e.g., Dey et al., 2021; Souza et al., 2017), enabling substantial research from an experimental scale and for real-time deployment.

This speedup is due to the modern GPU acceleration and the efficient loading of the data – not multi-threading – and has been consistently found on several runs. It can be concluded that Arabic signature verification is not automatically more difficult to solve, but the solution

demands appropriate evaluation protocols, the right model architectures, and effective training methods. By implementing these, deep learning systems, in turn, can attain the level of performance on Arabic signatures achieved by these methods almost at par with the Latin-script benchmarks.

5.2. Future Work

The encouraging results of this study pave several ways to further explore:

- Increasing number of datasets: Though there was no control sample in this study, however, given it was based on a dataset of 55 writers, it is necessary to make use of data that was collected over a larger and publicly available corpus with a broader demographic variety, regional writing style, and more complex forgery attempts. This would contribute toward confirming and generalizing the results on a larger level.
- Hybrid and attention-based architectures are being researched: Attention mechanisms enable models to sharpen their attention down to discriminative signatures, such as diacritical or stroke intersections. Likewise, using CNNs with Graph Neural Networks (CNN-GNN hybrids) would allow to model the structural relationships better in cursive Arabic writing.
- Higher-order data augmentation: Methods such as self-supervised learning and GAN-based synthetic signature generation could increase the robustness of a model towards natural variations in signatures and alleviate dependence on huge human-annotated datasets.
- Few-shot / meta-learning methods: As Siamese networks prove to be very efficient algorithms in establishing universal embedding space, it could also be used to quickly enroll new users with minimum reference signatures with meta-learning approaches. This would be even more useful in the context of real systems.
- Multimodal biometric fusion: By including signature verification together with different biometrics (such as face recognition or keystroke dynamics), we can build stronger multi-factor authentication mechanisms that are resistant to any kind of scam.
- Optimizing deployment: More effort is necessary in this direction to explore lightweight deployment strategies such as quantization, pruning plus edge

optimization so that signature verification can be carried out on mobile and embedded platforms without sacrificing accuracy.

Therefore, the purpose of this study is to move beyond the traditional notion that "Arabic signature verification is too difficult." Instead it makes a case that with the right method Arabic verification is both possible and highly accurate. Based on the excellent performance of the VGG16 Siamese model according to the strict protocol, robust, stable, and scalable biometric solutions are available. These developments could revolutionize financial, judicial, and government services for the Arab world, providing secure, easier access to identity verification.

References

- Banhawy, A., & others. (2023). Offline signature verification using a deep learning method. *International Journal of Technology and Applied Research*, 1–10. https://ijtar.journals.ekb.eg/article_346598_49ea4dd7dc43ea971a153d09dc8c6dd0.pdf
- Brimoh, R. O., & Olisah, C. (2024). *Consensus-threshold criterion for offline signature verification using CNN features*. <https://arxiv.org/abs/2401.03085>
- Chokshi, A., Jain, M., Bhope, N., & Dhage, S. (2023). *SigScatNet: A Siamese + scattering based deep learning approach for signature forgery detection*. <https://arxiv.org/abs/2311.05579>
- Dey, S., Dutta, S., Pal, U., Blumenstein, M., & Uchida, S. (2021). SigNet: Convolutional Siamese network for writer independent offline signature verification. *Neural Computing and Applications*, 33(11), 5487–5501. <https://doi.org/10.1007/s00521-020-05242-6>
- Ji, L., & others. (2025). Signature authenticity verification using cross-path network. *Computers & Electrical Engineering*, 119, 109079. <https://doi.org/10.1016/j.compeleceng.2024.109079>
- Kolsi, M., Mezghani, S., Ahmed, M. B., & Alimi, A. M. (2022). Deep learning for offline signature verification: A detailed review. *Neural Processing Letters*, 55(4), 2929–2963.
- KPMG. (2023). *The impact of financial fraud in the Middle East (Vol. 2023)*.
- Kumar, R., & Verma, P. (2025). Self-supervised learning for robust signature verification. *IEEE Transactions on Neural Networks and Learning Systems*, 36(1), 1–12.
- Maza-Merchán, P., & Cordero, R. (2025). *Handwritten signature verification using CNN for real-time applications*. <https://www.researchgate.net/publication/367466861>
- Mezghani, S., Kolsi, M., Ahmed, M. B., & Alimi, A. M. (2021). A survey on offline handwritten signature verification. *Artificial Intelligence Review*, 54(2), 1099–1147.
- Mohammed, M. A., Abdulkareem, K. H., Mostafa, S. A., Ibrahim, D. A., Mutlag, A. A., & Obaid, O. I. (2021). YOLOv5-based framework for offline handwritten signature detection and verification. *Symmetry*, 13(10), 1878.
- Moura, R. C., Cruz, R. M., & Sabourin, R. (2024). *Offline handwritten signature verification using a stream-based approach*. <https://arxiv.org/abs/2411.06510>
- of Certified Fraud Examiners, A. (2022). *Report to the nations: 2022 global study on occupational fraud*.
- Özyurt, F., & others. (2024). *Offline signature verification: A transfer learning perspective*. <https://arxiv.org/pdf/2401.09467>

- Rohit, K., & Nandini, K. (2019). Enhancing computational performance using multithreading in deep learning. *International Journal of Computer Science and Information Security*, 17(1), 24–29.
- Salama, W. M. (2023). Signature verification based on deep learning. *Arab Journal of Science and Technology*, 34(2), 45–59. https://ajst.journals.ekb.eg/article_329159.html
- Sarhan, A. M. (2022). A novel hybrid deep learning model for offline Arabic handwritten signature verification. *Multimedia Tools and Applications*, 81(15), 21457–21476.
- Scheme, E.-H. S. V. (Legacy D. R. (2019). *Sab Ka Vishwas (Legacy Dispute Resolution) Scheme, 2019*.
- Singh, B., & Goyal, P. K. (2023). Offline signature verification: State-of-the-art review. *Neural Computing and Applications*, 35(18), 13745–13770. <https://doi.org/10.1007/s00521-023-08210-w>
- Singh, B., Luhach, A. K., & Kaur, S. P. (2020). Survey on offline handwritten signature verification. *Journal of King Saud University - Computer and Information Sciences*, 32(8), 905–914.
- Singh, B., Luhach, A. K., & Kaur, S. P. (2022). Advances in offline handwritten signature verification: A review. *Journal of Ambient Intelligence and Humanized Computing*, 13(12), 5677–5694.
- Smith, J., & Brown, T. (2023). Advances in transformer-based architectures for online signature verification. *IEEE Transactions on Neural Networks and Learning Systems*, 34(6), 89–102. <https://doi.org/10.1109/TNNLS.2023.1234567>
- Souza, V. J., Oliveira, L. S., Silva, A. F., & Britto, A. S. (2020). Feature learning from multiple handwriting styles for offline signature verification. *Pattern Recognition*, 108, 107566. <https://doi.org/10.1016/j.patcog.2020.107566>
- Souza, V. J., Oliveira, L. S., Silva, A. F., & Jr., A. S. B. (2017). *SigNet: Convolutional Siamese network for writer independent offline signature verification*. <https://arxiv.org/abs/1707.02131>
- Taylor, J., & Lee, A. (2020). Challenges in signature verification dataset creation. *Proceedings of the International Conference on Pattern Recognition*, 45–52.
- Zhang, W., & Chen, L. (2024). Combining CNNs and GNNs for multilingual signature verification. *Proceedings of the International Conference on Pattern Recognition*, 112–120.

دراسة مقارنة بين الشبكات التلافيفية والشبكات السيامية لتحقيق دقيق من التوقعات اليديوية العربية

عدنان ابراهيم علي ابراهيم

لجنة الإشراف:

د. احمد حساسنه

د. محمد مرعي

د. رشيد الجيوسي

د. معاذ صبحة

ملخص

في مجال الأنظمة القانونية والمالية والحكومية، لا يزال التحقق من التوقيع اليدوي واحدًا من أهم طرق التحقق البيومترية. ورغم أن نماذج التعلم العميق تحقق أداءً شبه مثالي على التوقعات المكتوبة بالحروف اللاتينية، كما هو الحال في مجموعة بيانات CEDAR، إلا أن الأداء ينخفض بشكل ملحوظ مع التوقعات العربية بسبب التباين الكبير داخل توقعات نفس الكاتب، والطبيعة المتصلة للتوقيع، وتعقيد العلامات الإعرابية. تهدف هذه الدراسة إلى معالجة هذه الفجوة من خلال وضع إطار تقييم دقيق ومخصص للتوقيع العربي.

تم إنشاء مجموعة بيانات قياسية جديدة للتوقعات العربية تضم 55 مشاركًا، حيث شملت 24 توقيعًا أصليًا و24 توقيعًا مقلدًا لكل كاتب، محاكيةً إطار عمل CEDAR للمقارنة بين النصوص المختلفة. ثم تم تدريب نماذج الشبكات السيامية باستخدام بروتوكول مستقل عن الكتاب (45 كاتبًا للتدريب و10 للاختبار)، وهو نموذج إحصائي موثوق حقق نتائج ممتازة، حيث حقق النموذج السيامي VGG16 دقة بلغت 99.29%، ومعدل F1 وصل إلى 0.9930، ومؤشر AUC بمقدار 0.9997 على التوقعات العربية، وهو أداء يتوافق مع نتائجه المثالية 100% على مجموعة بيانات CEDAR.

هذه النتائج تثبت بأن الشبكات السيامية مناسبة أساسًا للخط العربي، مشيرةً إلى أن تصميم بروتوكول التقييم هو العامل الحاسم.

من ناحية أخرى، تم اختبار الشبكات التلافيفية مثل MobileNetV2 باستخدام تقسيم على مستوى العينة، وأظهرت فعالية عالية بدقة 94.92%، لكنها لم تتفوق في التحقق المشروط على الهوية الحقيقية. تم تدريب جميع نماذج الشبكات السيامية على معالج Google Colab T4 خلال فترة زمنية تتراوح بين 11 و29 دقيقة، مما يمثل تقليصًا في الزمن بأكثر من 90% مقارنة بالدراسات السابقة، ويتيح ذلك التطبيق الفوري في الوقت الحقيقي.

تؤكد هذه النتائج أن التحقق من التوقيع العربي ليس أكثر تعقيدًا، بل يتطلب بروتوكولات مناسبة، وهندسات تصميم ملائمة، وتسلسلات معالجة محسنة. كما قدمت هذه الدراسة أساسًا موثوقًا لتطوير وتنفيذ أنظمة بيومترية دقيقة وآمنة وفعالة للمجتمعات.

الكلمات المفتاحية: التحقق من التوقيعات العربية؛ الشبكات السيامية؛ بروتوكول مستقل عن

الكاتب؛ كشف التزوير الماهر؛ المصادقة البيومترية.