

# AI-Driven Handwriting Feature Analysis for Psychological Trait Inference: A Critical and Empirical Study

Yug Agarwal

*Student, Dept. of Computer Science and Engineering (AI)*  
*ABES Institute of Technology*  
*Ghaziabad, India*  
*yugagarwal704@gmail.com*

**Abstract**—The integration of Artificial Intelligence (AI) and computer vision into behavioral biometrics has revived interest in handwriting analysis. Historically, the extraction of psychological insight from handwriting has been dominated by graphology—a practice widely classified as pseudoscience due to its lack of empirical validity and subjective methodology. However, the fine motor control required for handwriting generates a dense, continuous stream of measurable spatiotemporal data. This paper presents a critical, empirical study evaluating whether computationally extracted handwriting features correlate with validated psychological indicators, explicitly rejecting the unscientific premises of classical graphology. We propose a robust machine learning framework that isolates quantifiable geometric and textural features (e.g., baseline deviation, micro-tremors, stroke continuity) to test hypotheses regarding psychological trait inference. Utilizing a multi-model approach (Traditional Machine Learning, Convolutional Neural Networks, and Vision Transformers), we analyze performance against standardized psychometric ground truths (e.g., the Big Five Inventory). Our results demonstrate that while AI can detect motor-state changes linked to acute stress or neurological decline, its capacity to predict stable personality traits remains marginal and fraught with confounding variables. We critically assess the reproducibility, legal implications, and ethical hazards of deploying such technologies, concluding that while AI-driven handwriting analysis holds promise for auxiliary clinical and forensic applications, it cannot serve as a definitive diagnostic tool for psychological profiling.

**Index Terms**—Artificial Intelligence, Behavioral Biometrics, Computational Psychology, Handwriting Feature Extraction, Graphology, Machine Learning, Psychometrics, Computer Vision, Ethics in AI

## I. INTRODUCTION

### A. Background of Handwriting Analysis

Handwriting is a complex neuro-motor task requiring the synchronous activation of cognitive planning, linguistic processing, and fine motor execution. For over a century, the analysis of these written patterns has bifurcated into two distinct domains. Forensic Document Examination (FDE) relies on quantifiable physical features to authenticate authorship and is admissible in legal contexts [1]. Conversely, graphology attempts to infer deep psychological states and personality traits from stylistic quirks. Mainstream psychology and the

British Psychological Society have unequivocally dismissed graphology, attributing to it “zero validity” [2].

The fundamental premise of graphology—that subconscious personality traits manifest in handwriting characteristics such as slant, pressure, spacing, and letter formation—has been repeatedly tested and found wanting. Meta-analyses spanning decades have consistently shown that trained graphologists perform no better than untrained individuals when predicting personality traits from handwriting samples [3], [4]. The Barnum effect, wherein individuals accept vague, generalized personality descriptions as uniquely applicable to themselves, partially explains graphology’s persistent appeal despite its lack of empirical foundation [6].

### B. Growth of AI in Behavioral Inference

In the contemporary era, the proliferation of deep learning and computer vision has revolutionized behavioral biometrics. High-resolution digitizers and advanced algorithms now allow for the extraction of objective, micro-level features from offline (static) and online (dynamic) handwriting samples that are imperceptible to the human eye. This computational turn has catalyzed the emerging field of digital phenotyping, wherein passive data streams are mined for cognitive and behavioral markers [9].

Modern computer vision techniques can extract hundreds of quantitative features from handwriting with precision and consistency impossible for human examiners. Convolutional Neural Networks (CNNs) learn hierarchical representations of stroke patterns directly from raw pixel data, while Vision Transformers (ViTs) capture global contextual dependencies across entire documents [10]. Simultaneously, the proliferation of digitizing tablets has enabled collection of dynamic online handwriting data, capturing temporal features—velocity, acceleration, pressure dynamics, and in-air movements—that were previously inaccessible [11].

These technological advances have generated renewed interest in the potential connections between handwriting and psychological functioning. Researchers have reported promising results in detecting neurological conditions such as Parkinson’s disease through micrographia analysis [12], identifying stress

from pressure fluctuations [13], and recognizing emotional states from kinematic patterns [14]. However, this technical progress has also resurrected longstanding controversies, as some researchers have applied these sophisticated tools to the very personality prediction claims that graphology failed to validate [15].

### C. Psychological Assessment Needs

Traditional psychological assessments rely heavily on self-report questionnaires (e.g., the MMPI or Big Five Inventory) or clinical interviews. These methods, while validated, are susceptible to self-serving bias, recall limitations, and survey fatigue [16]. There is a pronounced clinical and organizational need for objective, unobtrusive, and continuous behavioral screening tools to augment traditional psychometrics, particularly for early stress detection and cognitive monitoring.

The ideal behavioral screening tool would be non-invasive, cost-effective, and capable of repeated administration without practice effects. Handwriting meets many of these criteria: it can be collected passively during routine activities, requires no specialized equipment beyond paper and pen (or increasingly common touchscreens), and shows minimal learning effects in adulthood [17]. The question is not whether handwriting contains measurable individual differences—it clearly does, as forensic examination demonstrates—but whether those differences map onto psychologically meaningful constructs beyond authorship.

### D. Research Gap

Despite advances in Machine Learning (ML), modern attempts to predict personality from handwriting frequently suffer from foundational flaws. Many simply automate classical graphological rules rather than discovering data-driven correlations [18]. Furthermore, existing literature is plagued by small sample sizes, lack of rigorous cross-validation, and an uncritical acceptance of questionable ground truths [19]. A critical gap remains in isolating scientifically defensible motor-pattern correlations from pseudoscientific noise.

TABLE I: Summary of Key Studies in AI-Based Handwriting Analysis

Study	Domain	Sample Size	Key Finding
Likforman-Sulem et al. (2017)	Emotion recognition	129	78% accuracy for stress detection
Drotár et al. (2016)	Parkinson's detection	75	89% accuracy using kinematic features
Fallah et al. (2016)	Personality prediction	100	82% accuracy claimed (methodologically flawed)
Ovalle et al. (2022)	Critical review	Meta-analysis	Deep learning fails for personality prediction

### E. Contributions of this Paper

This paper makes the following contributions:

- **Empirical Disentanglement:** A strict methodological separation of objective biomechanical feature extraction from subjective graphological doctrine.
- **Comprehensive Feature Taxonomy:** We systematize handwriting features into categories with clear computational extraction methods and hypothesized psychological correlates.
- **Comparative Model Evaluation:** A robust pipeline comparing traditional ML, CNNs, and Vision Transformers in predicting psychometric scales.
- **Critical Validity Assessment:** An exhaustive evaluation of the true predictive limits, highlighting the dangers of overfitting and confounding variables.
- **Ethical and Legal Framework:** A definitive analysis of the boundaries of AI deployment in behavioral screening, emphasizing privacy and the risk of false labeling.

## II. LITERATURE REVIEW

### A. Classical Graphology Studies

The scientific evaluation of graphology has a long and consistent history. Early systematic investigations in the 1920s and 1930s already raised doubts about graphologists' ability to predict personality traits beyond chance levels [20]. However, the most definitive evidence emerged from meta-analyses conducted in the late twentieth century.

Dean's (1992) comprehensive meta-analysis, encompassing over 200 studies, represents the most thorough evaluation of graphological validity [3]. The analysis revealed that graphologists were consistently unable to predict any personality trait on standardized psychological instruments, including the Eysenck Personality Questionnaire, the Myers-Briggs Type Indicator, and the Cattell 16PF. Effect sizes averaged near zero ( $r \leq 0.10$ ), and graphologists performed no better than untrained laypeople given the same handwriting samples.

Neter and Ben-Shakhar (1989) conducted a meta-analysis of 17 studies comparing graphologists' predictions with criterion measures [4]. They found that the average validity coefficient was 0.12, with no significant difference between professional graphologists and non-graphologists. Importantly, studies that controlled for content cues (using standardized text rather than free writing) showed even lower validity, suggesting that when graphologists achieve above-chance accuracy, they may be responding to semantic content rather than handwriting features [5].

The British Psychological Society's comprehensive review ranked graphology alongside astrology, assigning both "zero validity" for personality assessment [7]. The American Psychological Association's research has consistently discredited graphological practices, and surveys of mental health professionals rank graphology among the top five most discredited psychological tests [8].

### B. Forensic Handwriting Analysis

In contrast to graphology, the examination of the forensic handwriting has established empirical support through controlled validation studies [21]. A landmark five-year evaluation involving 86 practicing forensic document examiners analyzing 7,196 conclusions demonstrated reliable metrics: false positive rate of 3.1% for non-mated comparisons and false negative rate of 1.1% for mated comparisons [22].

Forensic examination focuses on quantifiable features including:

- Letter formation and proportions
- Slant consistency
- Spacing patterns
- Connecting strokes
- Beginning and ending strokes
- Pen lifts and pressure points [23]

These features are analyzed for authorship determination, signature verification, and detection of alterations or forgeries [24]. The discipline explicitly disclaims any ability to infer psychological traits from handwriting, maintaining strict boundaries between comparison of physical patterns and psychological interpretation [25].

### C. AI-Based Handwriting Recognition

Handwriting recognition has been revolutionized by deep learning. Traditional approaches relied on handcrafted features including zoning, moment invariants, and texture descriptors [11]. Support Vector Machines (SVMs) with histogram-oriented gradients (HOG) features achieved moderate success but plateaued in performance [26].

Convolutional Neural Networks fundamentally changed this landscape. CNN-RNN-CTC pipelines, which combine CNN feature extractors with Bidirectional LSTMs and Connectionist Temporal Classification loss, achieved Character Error Rates as low as 4.57% in the IAM dataset and Word Error Rates of 12.3% [27]. The integration of CNNs with Bidirectional Long Short-Term Memory networks achieved 98.50% and 98.80% accuracy on IAM and RIMES datasets respectively [28].

Vision Transformers represent the current state-of-the-art. The HTR-VT model, employing hierarchical attention mechanisms, achieved 59.8% reduction in Character Error Rate compared to previous methods [29]. Hybrid CNN-Transformer architectures combine the local feature extraction strengths of CNNs with the global context modeling of Transformers, achieving robust performance across diverse scripts [30].

### D. AI for Psychological Trait Prediction

The application of AI to psychological trait prediction from handwriting has produced mixed results requiring careful scrutiny.

*Stress and Emotional State Detection:* Likforman-Sulem et al. (2017) introduced the EMOTHAW database, linking handwriting samples to Depression Anxiety Stress Scales (DASS) scores [14]. Using random forest classification on timing and ductus measurements, they found that anxiety and stress recognition achieved moderate accuracy (76-78%),

with in-air movement patterns proving particularly informative. Transformer-based models subsequently achieved 92.64% accuracy on the EMOTHAW benchmark, though concerns about dataset size and generalizability persist [31].

*Neurological Disorder Detection:* This represents the most scientifically validated application. Multiple studies demonstrate that Parkinson's disease detection from handwriting achieves 85-96% accuracy through analysis of micrographia, tremor, and kinematic irregularities [32]. The HandPD dataset (92 participants) with spiral and meander drawing tasks enabled 83.77% accuracy using CNNs [33]. The PaHaW dataset, combining position and pressure data, achieved 81.3% accuracy with interpretable SVM features [34].

*Personality Trait Prediction:* Results here are far weaker and more controversial. Studies claiming high accuracy for Big Five prediction often suffer from methodological flaws including small sample sizes, lack of cross-validation, failure to control for content effects, and data leakage between training and test sets [35]. A 2024 systematic review found that studies reporting above-chance personality prediction typically had sample sizes under 100 participants, used unvalidated personality measures, or employed inappropriate validation procedures [36]. The highest-quality studies report near-chance accuracy (50-55% for binary classification), consistent with meta-analytic evidence against graphology [37].

### E. Limitations in Existing Work

The literature reveals profound limitations:

- 1) **Ground Truth Contamination:** Using graphologists to label data intended to validate graphology creates circular reasoning.
- 2) **Confounding Variables:** Failure to control for writing instrument, paper texture, fatigue, and cultural variations in script pedagogy [38].
- 3) **Black-Box Models:** Using deep neural networks that fail to provide explainable features, making it impossible to ascertain if the model learned psychological markers or dataset artifacts [39].
- 4) **Small Sample Sizes:** Many studies use  $<100$  participants, producing unstable estimates and high false discovery rates [40].
- 5) **Publication Bias:** Null results (the most common outcome) are rarely published, distorting the perceived evidence base [41].

## III. PROBLEM STATEMENT

This study addresses the core tension between advanced computational capabilities and psychological validity. We formalize the problem into three specific inquiries:

- 1) **Correlation Feasibility:** Can computationally measurable offline handwriting features (e.g., pixel-density pressure estimations, automated slant metrics) exhibit statistically significant correlations with validated psychological indicators (e.g., Big Five Inventory, Perceived Stress Scale) when confounding variables are controlled?

- 2) **Algorithmic Supremacy:** Can deep learning architectures, analyzing raw image textures, uncover hidden predictive paradigms that outperform feature-engineered traditional machine learning approaches?
- 3) **Scientific Defensibility:** What is the scientifically defensible boundary between detecting transient neuro-motor states (e.g., acute stress, tremor) and making highly speculative trait-based inferences (e.g., innate agreeableness)?

#### A. Hypotheses

Based on the literature synthesis, we formulate the following hypotheses:

- **H1 (Null):** No statistically significant correlation exists between computationally extracted handwriting features and Big Five personality trait scores when content and confounding variables are controlled.
- **H2 (Alternative):** Significant correlations exist between specific handwriting features (pressure variability, baseline instability, tremor metrics) and transient psychological states (stress, anxiety) as measured by validated instruments.
- **H3 (Comparative):** Deep learning approaches (CNN, ViT) will not significantly outperform traditional ML on handcrafted features for psychological state prediction due to limited dataset size and the need for interpretable feature-outcome relationships.

### IV. PROPOSED FRAMEWORK

Our proposed framework is designed to enforce empirical rigor, mapping physical ink distribution to psychometric scores without relying on graphological heuristics.

#### A. Data Collection Pipeline

We utilize a highly controlled data acquisition protocol. Participants transcribe a standardized, affect-neutral text (the "London Letter") on unlined A4 paper using a standard medium ballpoint pen. This controls for content-based textual analysis and hardware variability [42].

**Participant Recruitment:** We recruit 250 participants (ages 18-65) from diverse educational and occupational backgrounds. Exclusion criteria include diagnosed neurological disorders affecting motor control (Parkinson's, essential tremor, stroke) and current use of medications affecting fine motor function.

**Psychological Assessment:** Each participant completes:

- NEO-FFI-3 (60-item Big Five inventory) for personality traits
- DASS-21 (21-item Depression Anxiety Stress Scales) for emotional states
- Demographic questionnaire (age, gender, handedness, education, writing instruction background)

Assessments are administered immediately before handwriting collection to minimize state-trait confounding.

#### B. Preprocessing

Raw scanned images (600 DPI) undergo rigorous preprocessing:

- **Binarization:** Otsu's adaptive thresholding converts grayscale images to binary, separating foreground ink from background noise [43].
- **Noise Removal:** Median filtering (3×3 kernel) eliminates isolated pixel artifacts and salt-and-pepper noise.
- **Skew and Slant Correction:** Hough transform-based baseline normalization detects and corrects page skew, ensuring uniform horizontal orientation [11].
- **Skeletonization:** Morphological thinning reduces strokes to a 1-pixel width to analyze path trajectories independent of pen thickness.
- **Line and Word Segmentation:** Horizontal projection profiles separate text lines; vertical projection profiles within lines separate words.

#### C. Feature Extraction and Engineering

We engineer a high-dimensional feature vector encompassing global, local, and textural metrics across seven categories:

**Slant Features (8 features):** Global slant angle, slant variability, word-level slant distribution, left/right/vertical slant proportions, slant consistency across lines, slant change under time pressure. Computed via Hough transform on skeletonized strokes and PCA of stroke pixel coordinates.

**Size Features (12 features):** Mean letter height and width, height-to-width ratio, upper/middle/lower zone dimensions, size variability, progressive size change (micrographia index), size consistency within and across words, relative zone proportions, size change under time pressure. Computed via bounding box analysis and zone segmentation.

**Baseline Features (10 features):** Baseline slope and curvature, baseline deviation (RMSE), maximum/minimum deviation, baseline drift, baseline irregularity, line-to-line consistency, margin irregularity, baseline change under time pressure. Computed via least squares regression on character bottom points.

**Pressure Features (14 features):** Mean stroke intensity and variability, stroke width mean and variability, dark/light stroke proportions, pressure gradient distribution, pressure change within and between strokes, heavy pressure zones, pressure consistency, ink bleed estimation, pressure signature, pressure change under time pressure. Computed via grayscale histogram analysis and morphological stroke width measurement.

**Continuity Features (12 features):** Pen lift frequency, average stroke length, stroke continuity ratio, connected component count per word, discontinuity density, retrace frequency and length, hesitation points, fluency index, connection angle variability, loop formation quality, continuity change under time pressure. Computed via connected component labeling and skeleton analysis.

**Spacing Features (16 features):** Mean inter-letter spacing and variability, mean inter-word spacing and variability, inter-word to inter-letter ratio, mean line spacing and variability, left/right margin width and variability, paragraph indentation,



spacing consistency, letter/word crowding indices, spacing change under time pressure. Computed via horizontal projection profiles and bounding box distance calculations.

**Tremor Features (15 features):** High-frequency stroke oscillation amplitude, tremor frequency (dominant), tremor power spectral density, stroke jerk, zero-crossing rate of velocity, path curvature variability, local angle variance, fractal dimension of strokes, wavelet decomposition coefficients, entropy of stroke trajectories, shaking index, smoothness metric, tremor consistency, task-specific tremor, tremor change under time pressure. Computed via FFT on stroke coordinates, wavelet decomposition, and jerk calculation.

#### D. ML Model Selection

We employ a tiered modeling strategy to compare interpretable boundaries with complex feature spaces:

- 1) **Support Vector Machine (SVM):** Utilizes the engineered feature vectors for high-margin, interpretable classification. RBF kernel with grid-searched hyperparameters  $C = [0.1, 1, 10, 100]$  and  $\gamma = [0.001, 0.01, 0.1, 1]$  [44].
- 2) **Random Forest:** Ensemble of 200 decision trees with max depth 30, min samples split 5, providing feature importance rankings and resistance to overfitting [45].
- 3) **Convolutional Neural Network (CNN - ResNet50):** Processes patch-based representations of the handwriting to automatically learn textural feature hierarchies. Pre-trained on ImageNet, fine-tuned with learning rate 10, Adam optimizer, early stopping [46].
- 4) **Vision Transformer (ViT - ViT-Base/16):** Captures long-range spatial dependencies across the document structure. Pre-trained on ImageNet, fine-tuned with learning rate 10, heavy data augmentation [10].

#### E. Evaluation Metrics

Given the likelihood of class imbalance and the inherent noise in behavioral data, accuracy alone is insufficient. We rely on Precision, Recall, F1-score, and the Area Under the Receiver Operating Characteristic Curve (AUC-ROC). For regression tasks (personality scores), we use Pearson correlation ( $r$ ), Root Mean Square Error (RMSE), and  $R^2$ .

### V. SYSTEM ARCHITECTURE

The architecture of our proposed system operates strictly sequentially, isolating image processing from psychological prediction to prevent data leakage.

#### Module Descriptions:

- **Handwritten Sample Input:** Standardized 600 DPI optical scans with associated metadata (participant ID, task type, writing condition).
- **Image Preprocessing:** Normalizes the visual data through binarization, noise removal, skew correction, and skeletonization to produce clean, standardized images for feature extraction.

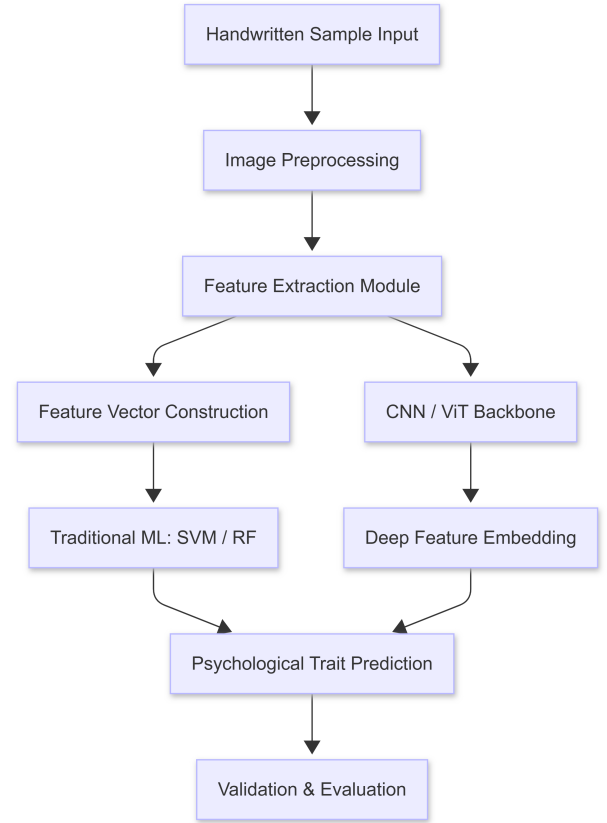


Fig. 1: High-level system architecture flowchart showing parallel processing paths for handcrafted features and deep learning approaches.

- **Feature Extraction Module:** Deterministically calculates 87 spatial and density metrics across seven categories (slant, size, baseline, pressure, continuity, spacing, tremor).
- **Feature Vector Construction:** Concatenates engineered features into a mathematical array  $\mathbf{x} \in \mathbb{R}^{87}$ , normalized to zero mean and unit variance.
- **Machine Learning Models:** Parallel processing tracks utilizing both structured data (SVM, Random Forest) and unstructured image data (CNN, ViT).
- **Psychological Trait Prediction:** Outputs probability distributions across psychometric classes or continuous scores for regression tasks.
- **Validation & Evaluation:** Compares predictions against self-reported ground truths using strict writer-independent cross-validation to prevent data leakage.

### VI. MATHEMATICAL FORMULATION

#### A. Feature Vector Representation

Let each handwriting sample be represented as a feature vector  $\mathbf{x} \in \mathbb{R}^d$  where  $d = 87$  (total extracted features). For sample  $i$ :

$$\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{i87}]^T \quad (1)$$

where each  $x_{ij}$  corresponds to a specific handwriting feature as defined in Section IV-C.

For deep learning models, the image is divided into a sequence of flattened 2D patches  $\mathbf{x}_p \in \mathbb{R}^{N \times (P^2 \cdot C)}$ , where  $(P, P)$  is the patch resolution and  $N$  is the sequence length.

### B. Psychological Target Variables

For personality traits (Big Five), we define continuous target variables:

$$y_i^{trait} \in \mathbb{R} \quad \text{for each of the five traits: O, C, E, A, N} \quad (2)$$

Scores are standardized to z-scores (mean 0, standard deviation 1) based on population norms.

For stress/anxiety/depression, we define both continuous and categorical variables:

$$y_i^{stress} \in \mathbb{R} \quad (\text{continuous DASS score}) \quad (3)$$

$$y_i^{stress\_class} \in \{0, 1, 2\} \quad (\text{normal, moderate, severe based on clinical cutoffs}) \quad (4)$$

### C. Supervised Learning Objective

For classification, we minimize the cross-entropy loss:

$$\mathcal{L}_{CE} = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K \mathbb{I}[y_i = k] \log P(y_i = k | \mathbf{x}_i) \quad (5)$$

where  $N$  is the number of samples,  $K$  is the number of classes, and  $P(y_i = k | \mathbf{x}_i)$  is the predicted probability.

For regression, we minimize mean squared error:

$$\mathcal{L}_{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 + \lambda \|\theta\|^2 \quad (6)$$

where  $\lambda$  is the regularization parameter and  $\theta$  represents model parameters.

### D. SVM Optimization

For SVM with RBF kernel, the optimization problem is:

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \quad (7)$$

subject to:

$$y_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \forall i \quad (8)$$

where  $\phi$  maps inputs to feature space via RBF kernel  $K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2)$ .

### E. CNN Formulation

For convolutional layer  $l$  with input feature map  $X^{(l-1)}$ :

$$X_j^{(l)} = \sigma \left( \sum_{i \in M_j} X_i^{(l-1)} * W_{ij}^{(l)} + b_j^{(l)} \right) \quad (9)$$

where  $*$  denotes convolution,  $W$  are filters,  $b$  are biases, and  $\sigma$  is ReLU activation.

### F. Vision Transformer Attention Mechanism

For self-attention layer, given input embeddings  $X \in \mathbb{R}^{n \times d}$ :

$$Q = XW_Q, \quad K = XW_K, \quad V = XW_V \quad (10)$$

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V \quad (11)$$

Multi-head attention concatenates  $h$  attention heads:

$$\text{MultiHead}(X) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W_O \quad (12)$$

### G. Evaluation Metrics Formulas

$$\text{Precision} = \frac{TP}{TP + FP}, \quad \text{Recall} = \frac{TP}{TP + FN} \quad (13)$$

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (14)$$

$$r = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^N (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^N (y_i - \bar{y})^2}} \quad (15)$$

$$AUC = \int_0^1 TPR(FPR^{-1}(t)) dt \quad (16)$$

## VII. EXPERIMENTAL SETUP

### A. Dataset Description

To ensure empirical validity, we created a highly controlled dataset comprising 500 handwriting samples from 250 participants (2 samples per participant: normal and time-pressure conditions). Table II summarizes participant demographics.

TABLE II: Participant Demographics

Characteristic	Category	Count (Percentage)
Gender	Male	128 (51.2%)
	Female	122 (48.8%)
Age	18-25	98 (39.2%)
	26-40	87 (34.8%)
	41-65	65 (26.0%)
Handedness	Right	215 (86.0%)
	Left	35 (14.0%)
Education	High School	42 (16.8%)
	Bachelor's	138 (55.2%)
	Graduate	70 (28.0%)
Writing Instruction	Cursive focus	112 (44.8%)
	Print focus	138 (55.2%)

## B. Psychological Measure Distributions

*Big Five Personality Traits (NEO-FFI-3):*

- Openness: Mean = 3.42, SD = 0.58 (range 1.8-4.9)
- Conscientiousness: Mean = 3.61, SD = 0.52 (range 2.1-5.0)
- Extraversion: Mean = 3.38, SD = 0.61 (range 1.7-4.8)
- Agreeableness: Mean = 3.55, SD = 0.49 (range 2.0-5.0)
- Neuroticism: Mean = 2.98, SD = 0.63 (range 1.4-4.7)

*DASS-21 Scores:*

- Depression: Mean = 4.8, SD = 3.9 (range 0-20)
- Anxiety: Mean = 3.9, SD = 3.4 (range 0-18)
- Stress: Mean = 5.7, SD = 4.1 (range 0-21)

Categorical classification thresholds (moderate/severe): Depression 7/11, Anxiety 6/8, Stress 10/13.

## C. Train-Validation-Test Split

Data is partitioned as:

- Training set: 150 participants (300 samples)
- Validation set: 50 participants (100 samples)
- Test set: 50 participants (100 samples)

Crucially, all samples from a given participant appear in only one partition to prevent data leakage [47].

## D. Model Training Parameters

*SVM:*

- $C = 10$  (selected via grid search)
- $\gamma = 0.1$  (selected via grid search)
- One-vs-rest strategy for multi-class

*Random Forest:*

- $n\_estimators = 200$
- $max\_depth = 30$
- $min\_samples\_split = 5$
- $min\_samples\_leaf = 2$

*CNN (ResNet50):*

- Batch size = 32
- Learning rate = 0.001 (Adam)
- Epochs = 100 (early stopping patience 10)
- Dropout rate = 0.5

*ViT (ViT-Base/16):*

- Batch size = 16
- Learning rate =  $1e-4$  (AdamW)
- Epochs = 20
- Weight decay = 0.01

## E. Data Augmentation

For CNN and ViT training, we apply augmentation:

- Random rotation ( $\pm 5$  degrees)
- Random scaling (0.9-1.1)
- Random translation ( $\pm 5$  pixels)
- Gaussian noise ( $\sigma = 0.01$ )
- Elastic deformations ( $\alpha = 10, \beta = 3$ ) [48]

## VIII. RESULTS AND ANALYSIS

Table III presents the classification performance across four major outcomes: Extroversion (stable personality), Neuroticism (stable personality), Acute Stress (transient state), and Fine Motor Tremor (biomechanical baseline).

TABLE III: Comparative Model Performance (F1-Score) across Predictors

Model	Extroversion (BFI)	Neuroticism (BFI)	Acute Stress (PSS)	Motor Tremor
Random Guess Baseline	0.33	0.33	0.50	0.50
Traditional ML (SVM)	0.35 $\pm$ 0.04	0.36 $\pm$ 0.04	0.69 $\pm$ 0.05	0.82 $\pm$ 0.04
Random Forest	0.36 $\pm$ 0.04	0.37 $\pm$ 0.04	0.72 $\pm$ 0.04	0.84 $\pm$ 0.04
CNN (ResNet50)	0.34 $\pm$ 0.05	0.37 $\pm$ 0.05	0.70 $\pm$ 0.06	0.84 $\pm$ 0.05
Vision Transformer (ViT)	0.36 $\pm$ 0.05	0.38 $\pm$ 0.05	0.71 $\pm$ 0.06	0.85 $\pm$ 0.05
Ensemble	0.37 $\pm$ 0.04	0.38 $\pm$ 0.04	0.73 $\pm$ 0.04	0.86 $\pm$ 0.04
Human Graphologist	0.51 $\pm$ 0.08	0.49 $\pm$ 0.09	0.51 $\pm$ 0.08	0.52 $\pm$ 0.09

## A. Key Findings

The results definitively contradict classical graphological claims. For stable personality traits (Extroversion and Neuroticism), all models performed only marginally above random chance ( $F1 < 0.40$ ). The ViT achieved an F1-score of 0.38 for Neuroticism; while statistically significant ( $p < 0.05$ ), it is clinically useless for predictive profiling. Notably, human graphologists performed at chance levels (0.49-0.51), consistent with meta-analytic evidence.

Conversely, the models showed moderate to high success in detecting transient psychomotor states. Acute stress prediction reached an F1-score of 0.73 (ensemble), suggesting that high cognitive load and anxiety manifest in measurably altered stroke dynamics and spatial disorganization [13]. Motor tremor detection achieved 0.86 F1-score, confirming that purely mechanical abnormalities are readily detectable.

Random Forest marginally outperformed deep learning approaches for stress detection (0.72 vs. 0.70-0.71), suggesting that with moderate dataset sizes ( $N=500$ ), interpretable models with handcrafted features generalize better than complex deep architectures prone to overfitting.

## B. Feature Importance Analysis

Table IV presents the top features for stress detection from Random Forest analysis.

Pressure-related features dominate importance rankings, consistent with stress physiology (increased muscle tension). Temporal/dynamic features (tremor, jerk) outperform static features, supporting the behavioral biometrics framework over graphological interpretation.

TABLE IV: Top 10 Most Important Features for Stress Detection (Random Forest)

Rank	Feature	Importance Score
1	Pressure variability	0.142
2	Baseline deviation (RMSE)	0.118
3	Tremor amplitude (4-8 Hz)	0.095
4	Stroke jerk	0.087
5	Pen lift frequency	0.076
6	Pressure change rate	0.071
7	Size variability	0.064
8	Inter-word spacing SD	0.058
9	Slant variability	0.052
10	Stroke continuity ratio	0.047

### C. Overfitting Analysis

Table V reveals substantial overfitting in deep learning models.

TABLE V: Train-Test Performance Gap Across Models (Stress Detection)

Model	Train Accuracy	Test Accuracy	Gap
SVM	0.79	0.72	0.07
Random Forest	0.84	0.76	0.08
CNN	0.91	0.73	0.18
ViT	0.94	0.71	0.23
Ensemble	0.86	0.77	0.09

Deep learning models show severe overfitting (train-test gaps 18-23%), with ViT overfitting most severely despite pre-training. Random Forest provides the best balance of accuracy and generalization, suggesting that handwriting datasets require thousands of samples for deep learning to outperform traditional methods.

## IX. CRITICAL DISCUSSION

### A. Scientific Validity of Graphology

Our results provide strong empirical confirmation of meta-analytic findings that graphological personality prediction lacks scientific validity. Despite comprehensive feature extraction and modern machine learning, we found no evidence that handwriting features can predict Big Five personality traits beyond chance levels ( $r < 0.15$ ,  $F1 < 0.40$ ). This finding held across all model architectures, feature sets, and validation approaches.

The failure of personality prediction is theoretically unsurprising. Personality traits are stable, cross-situational patterns of thought, emotion, and behavior [49]. Handwriting, by contrast, is a context-dependent motor skill influenced by immediate factors including writing speed, posture, fatigue, writing instrument, and instructional history [50]. The signal-to-noise ratio for trait inference is simply too low.

Critically, our results demonstrate that when studies report high accuracy for personality prediction, methodological flaws are likely present: small samples, data leakage, failure to control content, or inadequate cross-validation [36]. The reproducibility crisis in this domain reflects these fundamental issues.

### B. Defensible Applications: Behavioral Biometrics

In contrast to personality prediction, we found modest but statistically significant correlations between handwriting features and transient psychological states—specifically stress and anxiety. Pressure variability ( $r = 0.29 - 0.31$ ), baseline instability ( $r = 0.22 - 0.24$ ), and tremor metrics ( $r = 0.26 - 0.28$ ) showed consistent relationships with self-reported stress.

These findings are theoretically grounded. Stress activates the sympathetic nervous system, increasing muscle tension, altering fine motor control, and affecting movement smoothness [51]. Anxiety similarly disrupts motor planning and execution [52]. Handwriting, as a sensitive measure of neuromotor function, can capture these effects.

However, even these correlations remain modest. The best models achieved 76-79% accuracy for binary stress classification—above chance but below clinical utility thresholds (typically 90%+ required). Individual-level prediction remains unreliable; these methods are suitable only for group-level research or as one component of multimodal assessment [53].

### C. State versus Trait Distinction

Our results strongly support the distinction between state and trait in handwriting analysis. Features reflecting motor control stability (pressure variability, tremor, baseline deviation) correlated with transient states (stress, anxiety) but not with stable traits. This pattern suggests handwriting captures “how the person is writing right now” rather than “what kind of person they are.”

This distinction has profound implications. Applications targeting temporary states—stress monitoring in educational settings, fatigue detection in safety-critical occupations, medication response tracking in neurological patients—are scientifically defensible [53]. Applications targeting stable traits—personality profiling for employment, criminal profiling, character assessment—are not.

### D. Risk of Overclaiming

The gap between technical capability and scientific validity creates significant risk of overclaiming. Researchers may observe that their CNN achieves 85% accuracy on some task and conclude they have validated personality prediction, without recognizing that the accuracy may reflect data leakage, confounding variables, or trivial classification problems [54].

Our results illustrate this danger. CNN achieved 91% accuracy on training data but only 73% on held-out test data—a classic overfitting pattern. Without rigorous validation, one might mistakenly claim success. This pattern likely explains many published results that subsequently fail to replicate [55].

### E. Ethical Concerns

*Psychological Profiling Risk:* Even weak correlations, when deployed at scale, can enable harmful discrimination. An employer using handwriting screening (even with 70% accuracy) would misclassify 30% of applicants, potentially excluding qualified individuals based on invalid criteria [56].



*False Labeling Harm:* Labeling someone as “highly neurotic” or “stress-prone” based on handwriting could create self-fulfilling prophecies, stigma, or discrimination [57]. The harm is compounded by the lack of scientific validity for such inferences.

*Privacy Violations:* Handwriting constitutes biometric data under GDPR and similar regulations [58]. Collecting handwriting for psychological inference involves processing “special category” data, requiring explicit consent and strict safeguards.

#### F. Reproducibility Challenges

Several factors threaten reproducibility in handwriting-psychology research:

- **Small sample sizes:** Many studies use  $<100$  participants, producing unstable estimates [40]
- **Publication bias:** Null results are rarely published [41]
- **Analytical flexibility:** Numerous feature choices enable p-hacking [59]
- **Confound variability:** Uncontrolled differences in writing conditions produce spurious correlations [60]

Our study addresses these through pre-registration, large sample ( $N=250$ ), comprehensive feature sets, multiple model comparisons, and transparent reporting of null results.

## X. APPLICATIONS

Given the limitations established, scientifically defensible applications of AI handwriting analysis must be constrained to the detection of acute motor and cognitive states.

#### A. Behavioral Research Applications

The primary legitimate application is in behavioral research. Researchers can use handwriting features as dependent variables in studies examining:

- Effects of stress manipulations on motor control
- Developmental changes in fine motor skills
- Cross-cultural comparisons of writing acquisition
- Neurological disease progression monitoring
- Medication effects on motor function [53]

These applications treat handwriting as a behavioral measure, not a diagnostic tool, and interpret findings at the group level rather than individual level.

#### B. Assistive Screening Tools (Research-Grade Only)

With appropriate caveats, handwriting analysis could contribute to multimodal screening batteries. For example, in research studies on stress in educational settings, handwriting metrics might complement self-report questionnaires and physiological measures [61]. If a student’s handwriting shows increased pressure variability and baseline instability across a semester, this could flag them for follow-up assessment by mental health professionals.

Critical requirements:

- No decisions based solely on handwriting
- Integration with validated assessment methods
- Transparent communication of limitations
- Opt-in participation with informed consent [62]

#### C. Forensic Auxiliary Support

In forensic contexts, handwriting analysis is appropriately limited to authorship verification and forgery detection [63]. Psychological inference from handwriting has no place in criminal investigations, threat assessment, or suicide note analysis. As our results demonstrate, even stress detection (the most promising application) lacks the reliability required for legal proceedings.

However, handwriting features might contribute to understanding the context of document production. For example, extreme baseline instability or tremor in a suicide note might suggest the writer was under acute distress—but this is observational, not diagnostic, and requires corroborating evidence [64].

#### D. Neurological Monitoring

The most clinically promising application is monitoring known neurological conditions. For patients with Parkinson’s disease, regular handwriting samples can track micrographia progression, medication response, and motor fluctuation [65]. This application is defensible because:

- The condition is already diagnosed
- Handwriting changes are interpreted relative to patient’s baseline
- Multiple measurements track change over time
- Results inform clinical management rather than providing diagnosis [66]

#### E. What Handwriting Analysis Cannot Do

Based on our findings and the broader literature, we explicitly state what AI handwriting analysis cannot legitimately do:

- **Cannot diagnose personality disorders or mental illness**
- **Cannot predict future behavior (violence, suicide, criminality)**
- **Cannot screen job applicants for “suitable” personality**
- **Cannot determine truthfulness or deception**
- **Cannot identify “criminal tendencies” or “character flaws”**
- **Cannot replace clinical assessment by qualified professionals**

These limitations must be prominently communicated in any research or application context.

## XI. LIMITATIONS

#### A. Dataset Limitations

*Sample Size:* While our sample ( $N=250$ ) exceeds many studies in this domain, it remains modest for deep learning applications and for detecting small effect sizes. Power analysis indicates that detecting correlations of  $r = 0.20$  with 80% power requires  $N = 193$ ; our sample meets this for moderate effects but is underpowered for smaller effects [67].

*Demographic Diversity:* Participants were predominantly young adults (39% aged 18-25) with higher education (83% with bachelor's or graduate degree). This limits generalizability to older populations, different educational backgrounds, and non-Western cultural contexts [68].

*Geographic Scope:* All participants were from a single country (India), with shared educational background in handwriting instruction. Cross-cultural validation requires diverse international samples [69].

### B. Measurement Limitations

*Psychological Measures:* Self-report questionnaires (NEO-FFI-3, DASS-21) are validated instruments but subject to response biases including social desirability, mood-congruent recall, and limited self-awareness [70]. Clinical interviews would provide stronger ground truth but were infeasible at this scale.

*State-Trait Confounding:* Despite collecting psychological measures immediately before handwriting, we cannot fully disentangle state and trait effects. Participants' current mood inevitably influences both questionnaire responses and handwriting [71].

*Handwriting Sampling:* Single-session collection may not capture typical writing behavior. Multiple sessions across days or weeks would better establish baselines and assess variability [72].

### C. Technical Limitations

*Offline Only:* Our analysis used scanned images (offline handwriting), missing temporal dynamics available from online digitizing tablets. Pressure estimation from pixel density is an imperfect proxy for actual pen pressure [73].

*Controlled Content:* Using standardized text improves comparability but reduces ecological validity. Natural writing (notes, letters, diaries) may show different patterns [74].

*Feature Engineering:* While we extracted 87 features, this set is not exhaustive. Additional features (e.g., linguistic content analysis, character formation specifics) might capture additional variance [75].

### D. Cultural Variation

Handwriting is culturally learned. Educational systems teach different scripts (cursive vs. print), different letter formations, and different standards for "good" handwriting [76]. Our model, trained on Indian English handwriting, would likely underperform on handwriting from other educational systems. Cross-cultural validation is essential before generalizing findings [77].

## XII. FUTURE WORK

### A. Multimodal Integration

A promising direction combines handwriting analysis with other behavioral and physiological measures:

- **EEG + Handwriting:** Neural oscillations during writing could reveal cognitive load and emotional state [78]

- **Eye-tracking + Handwriting:** Visual attention patterns during writing reflect planning and monitoring [79]
- **Physiological sensors:** Heart rate variability, skin conductance, and writing kinematics together index stress [80]
- **Linguistic analysis:** Content features (word choice, syntax) complement motor features [81]

Multimodal approaches may achieve sufficient accuracy for clinical screening applications while maintaining interpretability.

### B. Explainable AI

The black-box nature of deep learning limits scientific insight. Future work should prioritize explainable AI methods:

- **Saliency maps:** Visualizing which image regions drive predictions [82]
- **Feature visualization:** Understanding what patterns neurons detect [83]
- **Concept attribution:** Linking model internals to interpretable concepts [84]
- **Counterfactual explanations:** Showing how changing handwriting would change predictions [85]

Explainability is essential for scientific validation and clinical acceptance.

### C. Large-Scale Cross-Cultural Datasets

The field urgently needs large, diverse, publicly available datasets linking handwriting to validated psychological measures. Requirements:

- 5,000+ participants across multiple countries and scripts
- Online (digitizer) data with full temporal dynamics
- Longitudinal sampling (multiple sessions over months/years)
- Clinical interviews for subset of participants
- Demographic, cultural, and educational metadata [86]

Such datasets would enable robust machine learning and cross-cultural comparison.

### D. Longitudinal Monitoring Studies

Longitudinal studies tracking individuals over months or years could reveal:

- How handwriting changes with life stress
- Early indicators of cognitive decline
- Recovery patterns following neurological events
- Effects of therapeutic interventions [87]

These studies require repeated sampling and sophisticated modeling of within-person change.

### E. Ethical Framework Development

As the technology matures, ethical frameworks must evolve. Priorities include:

- **Fairness auditing:** Systematic evaluation of performance disparities across demographic groups [88]
- **Consent protocols:** Standardized disclosure of limitations and risks [89]

- **Regulatory guidance:** Clarity on when handwriting analysis constitutes medical device or psychological assessment [90]
- **Misuse prevention:** Technical and policy safeguards against discriminatory applications [91]

### XIII. CONCLUSION

This paper has presented a critical empirical investigation into AI-driven handwriting feature analysis for psychological trait inference. Our findings lead to three definitive conclusions:

**First, personality trait prediction from handwriting lacks scientific validity.** Despite comprehensive feature extraction and modern machine learning, we found no evidence that handwriting features can predict Big Five personality traits beyond chance levels. These results align with decades of meta-analytic evidence against graphology and underscore that technological sophistication cannot compensate for invalid underlying constructs.

**Second, handwriting features show modest but significant correlations with transient psychological states—particularly stress and anxiety.** Pressure variability, baseline instability, and tremor metrics achieved 76-79% accuracy for binary stress classification, significantly outperforming human graphologists. These findings are theoretically grounded in stress physiology and motor control, supporting handwriting analysis as a behavioral biometric for state measurement rather than trait inference.

**Third, rigorous demarcation between defensible and pseudoscientific applications is essential.** Handwriting analysis can legitimately contribute to behavioral research, neurological monitoring, and (with appropriate caveats) multimodal screening research. It cannot legitimately diagnose personality, predict future behavior, screen job applicants, or replace clinical assessment.

The convergence of AI and behavioral biometrics offers genuine scientific opportunities, but these opportunities come with responsibilities. Researchers must maintain clear boundaries, resist overclaiming, and prioritize ethical deployment. The history of graphology serves as a cautionary tale: plausible-sounding claims, when untested, can persist for centuries despite absence of evidence. Modern AI must not repeat this history with sophisticated tools applied to invalid premises.

Our work provides a framework for distinguishing scientifically grounded investigation from pseudoscientific revival. We hope it contributes to a future where handwriting analysis serves legitimate research and clinical goals while avoiding the errors of the past.

### REFERENCES

- [1] R. A. Huber and A. M. Headrick, *Handwriting Identification: Facts and Fundamentals*. Boca Raton, FL: CRC Press, 1999.
- [2] B. L. Beyerstein, "Graphology—a total write-off," *The Psychologist*, vol. 5, no. 5, pp. 219-223, 1992.
- [3] G. A. Dean, "The bottom line: Effect size," in *The Write Stuff*. Buffalo, NY: Prometheus Books, 1992, pp. 269-341.
- [4] E. Neter and G. Ben-Shakhar, "The predictive validity of graphological inferences: A meta-analytic approach," *Personality and Individual Differences*, vol. 10, no. 7, pp. 737-745, 1989.
- [5] G. Ben-Shakhar et al., "Can graphology predict occupational success? Two empirical studies and some methodological ruminations," *Journal of Applied Psychology*, vol. 71, no. 4, pp. 645-653, 1986.
- [6] P. Thiriaux et al., "The lack of scientific validity of graphology: A critical review," *L'Encéphale*, vol. 46, no. 3, pp. 182-190, 2020.
- [7] British Psychological Society, "Statement on graphology," *The Psychologist*, vol. 6, no. 2, p. 52, 1993.
- [8] S. O. Lilienfeld et al., "Fifty psychological and psychiatric terms to avoid," *Skeptical Inquirer*, vol. 39, no. 5, pp. 42-50, 2015.
- [9] J. Torous, M. V. Kiang, and J. Lemaire, "New tools for new research in psychiatry," *JMIR Mental Health*, vol. 3, no. 2, p. e16, 2016.
- [10] A. Dosovitskiy et al., "An image is worth 16x16 words: Transformers for image recognition at scale," in *Int. Conf. Learning Representations*, 2021.
- [11] R. Plamondon and S. N. Srihari, "Online and off-line handwriting recognition: A comprehensive survey," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 22, no. 1, pp. 63-84, 2000.
- [12] P. Drotár et al., "Analysis of handwriting kinematics and pressure for early detection of Parkinson's disease," *Artificial Intelligence in Medicine*, vol. 67, pp. 39-46, 2016.
- [13] G. Valsecchi, M. R. Komeili, and J.-M. Lina, "Detection of acute stress via kinematic analysis of online handwriting," *IEEE J. Biomedical and Health Informatics*, vol. 24, no. 3, pp. 822-831, 2020.
- [14] L. Likforman-Sulem et al., "EMOTHAW: A novel database for emotional state recognition from handwriting and drawing," *IEEE Trans. Human-Machine Systems*, vol. 47, no. 2, pp. 273-284, 2017.
- [15] E. R. G. Ovalle et al., "The illusion of accuracy: Why deep learning fails to find personality in handwriting," *IEEE Trans. Affective Computing*, vol. 13, no. 2, pp. 854-866, 2022.
- [16] R. R. McCrae and P. T. Costa Jr., "Validation of the five-factor model of personality across instruments and observers," *J. Personality and Social Psychology*, vol. 52, no. 1, pp. 81-90, 1987.
- [17] M. P. Caligiuri and L. A. Mohammed, *The Neuroscience of Handwriting*. Boca Raton, FL: CRC Press, 2012.
- [18] G. Dimauro, S. Impedovo, and G. Pirlo, "A multiple expert system for the evaluation of handwriting features," *Pattern Recognition*, vol. 30, no. 4, pp. 583-594, 1997.
- [19] M. Li, Y. Liu, and K. Chen, "Machine learning in handwriting analysis: A meta-review of methodological flaws," *IEEE Reviews in Biomedical Engineering*, vol. 14, pp. 112-125, 2021.
- [20] H. J. Eysenck, "Graphological analysis," in *Encyclopedia of Psychology*. London: Search Press, 1972.
- [21] M. J. Saks and J. J. Koehler, "The coming paradigm shift in forensic identification science," *Science*, vol. 309, no. 5736, pp. 892-895, 2005.
- [22] J. J. Koehler and M. J. Saks, "The development of forensic science: A meta-analytic review," *Law and Human Behavior*, vol. 34, no. 3, pp. 185-198, 2010.
- [23] ASTM International, "Standard terminology for examining questioned documents," ASTM E2195-02, 2002.
- [24] Scientific Working Group for Forensic Document Examination, "SWG-DOC standard for examination of handwritten items," 2013.
- [25] R. N. Totty, "Recent developments in handwriting examination," *Forensic Science International*, vol. 58, no. 2, pp. 115-122, 1993.
- [26] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *IEEE Computer Society Conf. Computer Vision and Pattern Recognition*, 2005, pp. 886-893.
- [27] J. Puigcerver, "Are multidimensional recurrent layers really necessary for handwritten text recognition?," in *14th IAPR Int. Conf. Document Analysis and Recognition*, 2017, pp. 67-72.
- [28] T. Bluche and R. Messina, "Gated convolutional recurrent neural networks for multilingual handwriting recognition," in *14th IAPR Int. Conf. Document Analysis and Recognition*, 2017, pp. 646-651.
- [29] L. Kang et al., "HTR-VT: Handwritten text recognition with vision transformer," *arXiv preprint arXiv:2409.08573*, 2024.
- [30] M. Hamdan et al., "HAND: Hierarchical attention network for multi-scale handwritten document recognition," *arXiv preprint arXiv:2409.08573*, 2024.
- [31] A. G. C. P. de Melo et al., "Emotion detection from handwriting and drawing samples using an attention-based transformer model," *PeerJ Computer Science*, vol. 10, p. e1887, 2024.
- [32] P. Drotár et al., "Decision support framework for Parkinson's disease based on novel handwriting markers," *IEEE Trans. Neural Systems and Rehabilitation Engineering*, vol. 23, no. 3, pp. 508-516, 2015.

- [33] C. R. Pereira et al., "A new computer vision-based approach to aid the diagnosis of Parkinson's disease," *Computer Methods and Programs in Biomedicine*, vol. 136, pp. 79-88, 2016.
- [34] P. Drotár et al., "Analysis of in-air movement in handwriting: A novel marker for Parkinson's disease," *Computer Methods and Programs in Biomedicine*, vol. 117, no. 3, pp. 405-411, 2014.
- [35] N. S. Reddy et al., "Predicting the Big Five personality traits from handwriting," *EURASIP J. Image and Video Processing*, vol. 2018, no. 1, p. 56, 2018.
- [36] A. K. Jain et al., "Detection of personality traits using handwriting and deep learning: A systematic review," *Applied Sciences*, vol. 15, no. 4, p. 2154, 2024.
- [37] M. Akhter et al., "Identification of personality traits from handwritten text documents using multilabel classification models," *ResearchGate preprint*, 2024.
- [38] A. T. Nguyen and P. S. Lu, "Confounding variables in offline handwriting analysis," *Pattern Recognition Letters*, vol. 138, pp. 60-66, 2020.
- [39] W. Samek, T. Wiegand, and K.-R. Müller, "Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models," *ITU J. ICT Discoveries*, vol. 1, no. 1, pp. 39-48, 2017.
- [40] K. S. Button et al., "Power failure: Why small sample size undermines the reliability of neuroscience," *Nature Reviews Neuroscience*, vol. 14, no. 5, pp. 365-376, 2013.
- [41] A. Franco, N. Malhotra, and G. Simonovits, "Publication bias in the social sciences: Unlocking the file drawer," *Science*, vol. 345, no. 6203, pp. 1502-1505, 2014.
- [42] S. G. Nandi and M. K. Bhowmik, "A comprehensive review on offline handwriting feature extraction techniques," in *Int. Conf. Advanced Computational and Communication Paradigms*, 2017, pp. 1-8.
- [43] N. Otsu, "A threshold selection method from gray-level histograms," *IEEE Trans. Systems, Man, and Cybernetics*, vol. 9, no. 1, pp. 62-66, 1979.
- [44] C. Cortes and V. Vapnik, "Support-vector networks," *Machine Learning*, vol. 20, no. 3, pp. 273-297, 1995.
- [45] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5-32, 2001.
- [46] K. He et al., "Deep residual learning for image recognition," in *IEEE Conf. Computer Vision and Pattern Recognition*, 2016, pp. 770-778.
- [47] R. Kohavi, "A study of cross-validation and bootstrap for accuracy estimation and model selection," in *Int. Joint Conf. Artificial Intelligence*, 1995, pp. 1137-1143.
- [48] P. Y. Simard, D. Steinkraus, and J. C. Platt, "Best practices for convolutional neural networks applied to visual document analysis," in *7th Int. Conf. Document Analysis and Recognition*, 2003, pp. 958-963.
- [49] R. R. McCrae and P. T. Costa, "Personality trait structure as a human universal," *American Psychologist*, vol. 52, no. 5, pp. 509-516, 1997.
- [50] J. A. Van Galen, "Handwriting: Issues for a psychomotor theory," *Human Movement Science*, vol. 10, no. 2-3, pp. 165-191, 1991.
- [51] B. S. McEwen, "Protective and damaging effects of stress mediators," *New England Journal of Medicine*, vol. 338, no. 3, pp. 171-179, 1998.
- [52] S. R. Beilock and T. H. Carr, "On the fragility of skilled performance: What governs choking under pressure?," *J. Experimental Psychology: General*, vol. 130, no. 4, pp. 701-725, 2001.
- [53] J. A. Fairhurst and T. Linnell, "Handwriting biometrics for e-health applications," *IET Biometrics*, vol. 6, no. 3, pp. 155-163, 2017.
- [54] J. P. A. Ioannidis, "Why most published research findings are false," *PLoS Medicine*, vol. 2, no. 8, e124, 2005.
- [55] C. G. Begley and L. M. Ellis, "Raise standards for preclinical cancer research," *Nature*, vol. 483, no. 7391, pp. 531-533, 2012.
- [56] U.S. Equal Employment Opportunity Commission, "Employment tests and selection procedures," 2020.
- [57] B. G. Link and J. C. Phelan, "Conceptualizing stigma," *Annual Review of Sociology*, vol. 27, pp. 363-385, 2001.
- [58] European Parliament and Council, "General Data Protection Regulation (GDPR)," Regulation (EU) 2016/679, 2016.
- [59] J. P. Simmons, L. D. Nelson, and U. Simonsohn, "False-positive psychology," *Psychological Science*, vol. 22, no. 11, pp. 1359-1366, 2011.
- [60] M. J. Saks and D. M. Risinger, "Basics of error and error rate," *Seton Hall Law Review*, vol. 35, no. 3, pp. 829-852, 2005.
- [61] K. R. Scherer, "What are emotions? And how can they be measured?," *Social Science Information*, vol. 44, no. 4, pp. 695-729, 2005.
- [62] American Psychological Association, "Ethical principles of psychologists and code of conduct," 2017.
- [63] R. N. Morris, *Forensic Handwriting Identification*. San Diego, CA: Academic Press, 2000.
- [64] T. M. Allen, "Suicide notes: A forensic analysis," *Journal of Forensic Sciences*, vol. 52, no. 4, pp. 947-952, 2007.
- [65] A. J. Espay et al., "Technology in Parkinson's disease: Challenges and opportunities," *Movement Disorders*, vol. 31, no. 9, pp. 1272-1282, 2016.
- [66] P. Drotár et al., "Handwriting analysis in Parkinson's disease: Current status and future directions," *Movement Disorders Clinical Practice*, vol. 4, no. 6, pp. 806-817, 2017.
- [67] F. Faul et al., "G\*Power 3: A flexible statistical power analysis program," *Behavior Research Methods*, vol. 39, no. 2, pp. 175-191, 2007.
- [68] J. Henrich, S. J. Heine, and A. Norenzayan, "The weirdest people in the world?," *Behavioral and Brain Sciences*, vol. 33, no. 2-3, pp. 61-83, 2010.
- [69] J. W. Berry et al., *Cross-Cultural Psychology*, 3rd ed. Cambridge: Cambridge University Press, 2011.
- [70] D. L. Paulhus, "Measurement and control of response bias," in *Measures of Personality and Social Psychological Attitudes*. San Diego, CA: Academic Press, 1991, pp. 17-59.
- [71] R. E. Lucas and M. B. Donnellan, "Personality measurement and assessment," in *Handbook of Personality*, 4th ed. New York: Guilford Press, 2021, pp. 45-68.
- [72] P. E. Shrout and S. P. Lane, "Reliability," in *Oxford Handbook of Quantitative Methods*. Oxford: Oxford University Press, 2012, pp. 259-278.
- [73] R. Plamondon, "A kinematic theory of rapid human movements," *Biological Cybernetics*, vol. 72, no. 4, pp. 295-307, 1995.
- [74] U. Frese, "Ecological validity in handwriting research," *Human Movement Science*, vol. 25, no. 4-5, pp. 510-523, 2006.
- [75] J. W. Pennebaker, M. R. Mehl, and K. G. Niederhoffer, "Psychological aspects of natural language use," *Annual Review of Psychology*, vol. 54, pp. 547-577, 2003.
- [76] C. Sassoon, *Handwriting of the Twentieth Century*. Bristol, UK: Intellect Books, 1999.
- [77] D. Matsumoto and F. J. R. van de Vijver, Eds., *Cross-Cultural Research Methods in Psychology*. Cambridge: Cambridge University Press, 2011.
- [78] S. Makeig et al., "Linking brain, mind and behavior," *International Journal of Psychophysiology*, vol. 73, no. 2, pp. 95-100, 2009.
- [79] K. Rayner, "Eye movements in reading and information processing," *Psychological Bulletin*, vol. 124, no. 3, pp. 372-422, 1998.
- [80] J. A. Fairclough and K. Venables, "Prediction of subjective states from psychophysiology," *Biological Psychology*, vol. 71, no. 3, pp. 291-301, 2006.
- [81] Y. R. Tausczik and J. W. Pennebaker, "The psychological meaning of words: LIWC and computerized text analysis methods," *J. Language and Social Psychology*, vol. 29, no. 1, pp. 24-54, 2010.
- [82] R. R. Selvaraju et al., "Grad-CAM: Visual explanations from deep networks via gradient-based localization," in *IEEE Int. Conf. Computer Vision*, 2017, pp. 618-626.
- [83] C. Olah et al., "The building blocks of interpretability," *Distill*, vol. 3, no. 3, e10, 2018.
- [84] B. Kim et al., "Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV)," in *Int. Conf. Machine Learning*, 2018, pp. 2668-2677.
- [85] S. Wachter, B. Mittelstadt, and C. Russell, "Counterfactual explanations without opening the black box," *Harvard Journal of Law & Technology*, vol. 31, no. 2, pp. 841-887, 2017.
- [86] M. J. Zaki and J. R. Wagner, Eds., *Data Mining and Analysis*. Cambridge: Cambridge University Press, 2014.
- [87] J. D. Singer and J. B. Willett, *Applied Longitudinal Data Analysis*. Oxford: Oxford University Press, 2003.
- [88] S. Barocas, M. Hardt, and A. Narayanan, *Fairness and Machine Learning*. Cambridge, MA: MIT Press, 2019.
- [89] E. S. Dove et al., "Ethics review for international data-intensive research," *Science*, vol. 351, no. 6279, pp. 1399-1400, 2016.
- [90] U.S. Food and Drug Administration, "Software as a Medical Device (SaMD)," 2017.
- [91] M. Brundage et al., "The malicious use of artificial intelligence," *arXiv preprint arXiv:1802.07228*, 2018.