

Contrastive Learning for Continuous Touch-Based Authentication

Mengyu Qiao*
North China University of
Technology
Beijing, China
myuqiao@ncut.edu.cn

Yunpeng Zhai
North China University of
Technology
Beijing, China
2023322030149@mail.ncut.edu.cn

Yang Wang
Ultramain Systems, Inc.
Albuquerque, NM, USA
ywang@ultramain.com

Abstract

Smart mobile devices have become indispensable in modern daily life, where sensitive information is frequently processed, stored, and transmitted—posing critical demands for robust security controls. Given that touchscreens are the primary medium for human-device interaction, continuous user authentication based on touch behavior presents a natural and seamless security solution. While existing methods predominantly adopt binary classification under single-modal learning settings, we propose a unified contrastive learning framework for continuous authentication in a non-disruptive manner. Specifically, the proposed method leverages a Temporal Masked Autoencoder to extract temporal patterns from raw multi-sensor data streams, capturing continuous motion and gesture dynamics. The pre-trained TMAE is subsequently integrated into a Siamese Temporal-Attentive Convolutional Network within a contrastive learning paradigm to model both sequential and cross-modal patterns. To further enhance performance, we incorporate multi-head attention and channel attention mechanisms to capture long-range dependencies and optimize inter-channel feature integration. Extensive experiments on public benchmarks and a self-collected dataset demonstrate that our approach outperforms state-of-the-art methods, offering a reliable and effective solution for user authentication on mobile devices.

CCS Concepts

• Security and privacy → Biometrics; • Human-centered computing → Gestural input; • Mathematics of computing → Time series analysis.

Keywords

Contrastive Learning, Continuous Authentication, Self-Supervised Learning, Masked Autoencoder, Siamese Network, Temporal Convolutional Networks, Transformers, Multi-head Attention, Channel Attention, Touch Gesture

1 Introduction

The widespread adoption of smartphones has profoundly transformed human-device interaction in both personal and professional spheres. These devices now handle not only communication and entertainment but also sensitive tasks such as financial transactions and private data storage. As dependence on mobile technology continues to grow, safeguarding user data has become an increasingly critical priority.

Conventional authentication methods—such as PIN codes, passwords, and static biometric scans—are widely used for convenience. However, they remain susceptible to various attacks, including

credential theft and biometric spoofing [32]. To overcome these vulnerabilities, behavioral biometrics have emerged as a promising supplementary layer for continuous authentication [18]. Unlike traditional methods, behavioral biometrics continuously monitor users’ interactions with their devices, enabling persistent and dynamic identity verification.

Among the various behavioral biometric modalities, touch dynamics has garnered increasing scholarly attention owing to its unobtrusiveness and temporal consistency. Unlike gait or motion-based systems, which rely on sensors that are susceptible to changes in user posture or contextual conditions, touch-based authentication leverages interaction-specific features such as tap pressure, swipe velocity, and temporal rhythm [30, 31]. These characteristics exhibit relatively low sensitivity to environmental and physiological variability, thereby rendering touch dynamics a viable candidate for robust, real-world continuous authentication systems.

Despite progress in behavioral biometrics, many existing methods rely on handcrafted features and traditional classifiers like SVMs, without fully exploiting the temporal structure inherent in user interactions. This limits their ability to capture sequential dependencies critical for modeling user-specific behavior over time. Ignoring such temporal dynamics reduces the discriminative power and robustness of authentication systems [29].

To address these challenges, we propose TouchSeqNet—a Siamese time-series framework that integrates a pre-trained Temporal-Attentive Convolutional Network (TACN) with contrastive learning. The proposed architecture employs a self-supervised pretraining phase based on a Temporal Masked Autoencoder (TMAE), which reconstructs masked segments of time-series data to learn generalizable temporal representations. These representations are subsequently transferred to a Siamese network, where TACN refines the features through dilated causal convolutions and multi-head self-attention mechanisms. Additionally, a finger-channel attention module adaptively highlights the most discriminative features across input sequences. Finally, the representations of sample pairs are concatenated and passed through a classification head to determine identity similarity.

To evaluate the model, we introduce Ffinger, a new dataset comprising touch dynamics from 29 users. Ffinger captures diverse interaction patterns across users. Additionally, we benchmark performance using two widely adopted datasets—BioIdent and Touchalytics—ensuring fair comparison under standard protocols. It shows that our method is superior to existing gesture detectors and time series classification baselines and achieves state-of-the-art performance. In summary, the contributions of this paper are threefold as below:

*Corresponding author.

- We reformulate continuous authentication as a contrastive learning task by employing a Siamese network architecture, offering a novel perspective in this field.
- We propose a Temporal Masked Autoencoder for self-supervised learning, which effectively captures fine-grained temporal patterns associated with continuous motion and gesture dynamics, enabling the extraction of robust and generalizable representations.
- We propose a Temporal-Attentive Convolutional Network, which incorporates dilated causal convolutions, multi-head self-attention, and channel attention mechanisms to further enhance the network’s ability to capture long-range temporal dependencies and optimize feature integration.

2 Related Work

A wide range of behavioral biometric modalities have been explored for continuous user authentication on mobile devices, including touch dynamics [11, 12, 33, 38, 41], motion sensor signals [2, 6, 25, 27, 29], keystroke dynamics [1, 35–37], and gait patterns [26, 40]. These studies have demonstrated the effectiveness of mobile behavioral authentication in both constrained and unconstrained environments.

Among these modalities, touch-based biometrics have drawn significant attention due to their unobtrusiveness and stability. Frank et al. [12] pioneered large-scale evaluation of touch features, achieving 2% EER with multi-gesture fusion. Tolosana et al. [38] introduced MobileTouchDB and proposed a Siamese LSTM-DTW framework. More recent works [33, 41] leveraged CNNs and multi-modal representations to improve classification accuracy.

Motion and sensor-based methods also show promising performance, especially when combining CNNs or LSTMs with accelerometer and gyroscope data [2, 6]. DeepConvLSTM [25] and clock-variant RNNs [27] have been applied to capture temporal dependencies in such signals. Similarly, keystroke dynamics have been modeled with RNNs [1, 37] and Transformer hybrids [35], while gait recognition systems have adopted metric learning [40] and deep CNNs [26] for robust identification.

Despite encouraging progress, many existing methods rely on static or handcrafted features and shallow models, limiting their ability to capture the sequential and dynamic nature of user behavior. Moreover, few approaches incorporate self-supervised learning or pretraining strategies, which are crucial for generalization in real-world deployment.

3 Data Acquisition and Processing

3.1 Ffinger Dataset

Ffinger containing interaction data from 29 participants. Each participant performed both predefined gestures—seven structured multi-touch tasks denoted by $\{a, b, c, d, e, f, g\}$ —and free-form gestures, which allowed users to draw arbitrary patterns. For each sample, trajectories from all five fingers were simultaneously recorded.

Each finger’s trajectory is represented as a 7-channel time series, including the following features:

- x_i, y_i : spatial coordinates of the touch point at time t_i ,
- t_i : timestamp,

- p_i : applied pressure,
- s_i : touch area size,
- v_i : instantaneous velocity,
- d_i : movement direction,

This design captures both structured and natural usage scenarios, enabling fine-grained modeling of user behavior.

3.2 Data Processing

To ensure consistency across datasets, we selected five core features as model input: timestamp (T_i), horizontal and vertical positions (X_i, Y_i), applied pressure (P_i), and contact area (A_i). Each time step is represented as:

$$\mathbf{X}_i = [T_i, X_i, Y_i, P_i, A_i] \quad (1)$$

First-Order Differencing. To emphasize motion dynamics and reduce temporal redundancy, we compute the differences in time and position:

$$T'_i = T_i - T_{i-1}, \quad X'_i = X_i - X_{i-1}, \quad Y'_i = Y_i - Y_{i-1} \quad (2)$$

with $T_1 = X_1 = Y_1 = 0$ by default.

Z-Score Normalization. To mitigate scale disparities across feature dimensions, we apply Z-score normalization to pressure and contact area features within each gesture sample:

$$P'_i = \frac{P_i - \mu_P}{\sigma_P}, \quad A'_i = \frac{A_i - \mu_A}{\sigma_A} \quad (3)$$

where μ and σ are the mean and standard deviation computed within the current sample. This improves training stability while preserving relative intra-sample variation.

Figure 2 provides a comparison of touch dynamics for gestures performed by two users, showing the high similarity between gestures of the same user and low similarity between gestures of different users.

4 TMAE Pre-trained Model Architecture

In this section, we present the architecture of the Temporal Masked Autoencoder (TMAE), As illustrated in Figure 1, which serves as the self-supervised pretraining backbone for time-series representation learning. The model first processes the raw time-series input through a window-based convolutional projection to extract local feature sequences. We split the feature sequences into two branches: one branch generates discrete token embeddings through a learnable tokenizer, and these embeddings are divided into visible and masked parts; the other branch directly divides the feature sequences into visible and masked parts for representation prediction. Specifically, the visible features are encoded by a Transformer encoder to capture contextual dependencies, and the masked segments are processed by a momentum-updated encoder to produce target representations. A cross-attention regressor then predicts the masked representations using information from the visible tokens, while a decoder reconstructs the discrete codewords.

4.1 Touch Dynamic Feature Representation

We introduce the methods for representing touch dynamic sequences through window slicing and embedding via the Tokenizer [20].

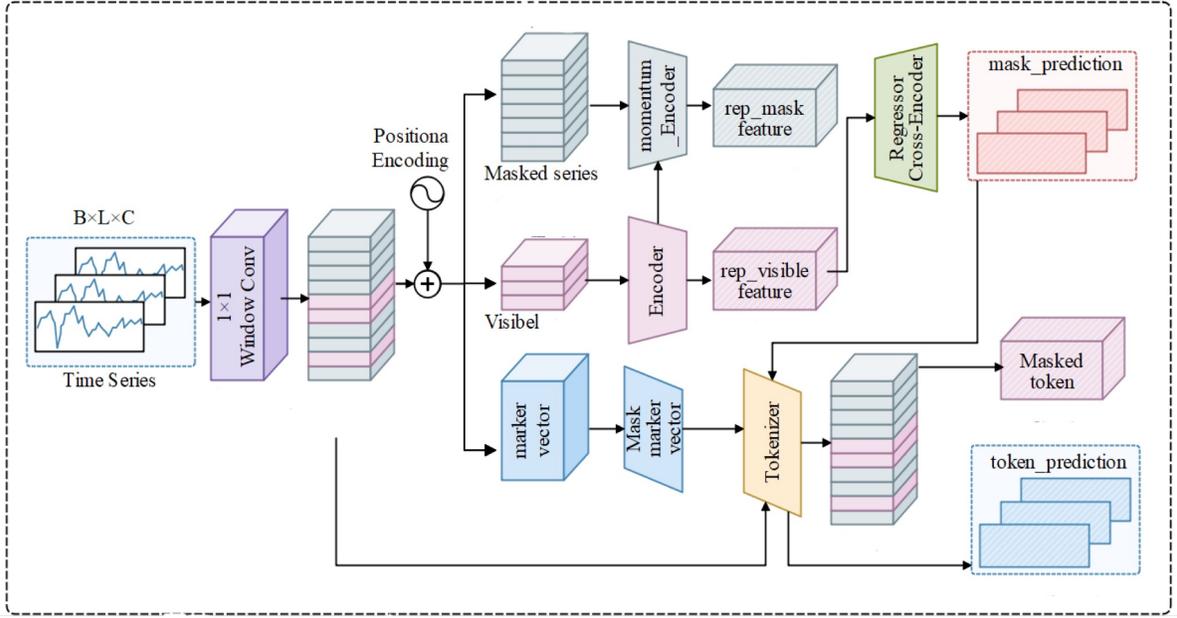


Figure 1: TMAE Model Architecture

In order to capture meaningful temporal dependencies across different time scales. We employ a window-slicing strategy to segment the continuous touch sequence into smaller, fixed-length sub-series [16].

Formally, let the touch sequence be represented as $X = \{x_1, x_2, \dots, x_T\} \in \mathbb{R}^{T \times C}$, where T is the length of the sequence and C is the number of channels (features such as x and y coordinates, pressure, etc.). We slice the sequence into non-overlapping windows of size σ , where each window $s_{i:j} = \{x_i, x_{i+1}, \dots, x_{i+\sigma}\}$ corresponds to a sub-series of length σ . The number of resulting sub-series is determined by $d = \lceil T/\sigma \rceil$. Then, the original sequence is encoded as $Z = \{z_1, z_2, \dots, z_d\} \in \mathbb{R}^{d \times m}$, where d is the new sequence length and m is the Model embedding dimension.

This strategy reduces temporal redundancy while ensuring that each sub-series contains enough semantic information, further enhancing the self-supervised learning process.

We convert each sub-series window Z into discrete embeddings using a Tokenizer [42], which maps raw segments into a compact latent space. Unlike handcrafted features, the Tokenizer supports end-to-end learning for effective representation of touch dynamics.

The Tokenizer module transforms each input sub-series $s_{i:j} \in \mathbb{R}^{\sigma \times C}$ into a continuous embedding representation $E \in \mathbb{R}^{\sigma \times m}$, where σ denotes the window length, C is the number of input channels, and m is the embedding dimension.

Each embedding $E_i \in \mathbb{R}^m$ is then projected into a vocabulary space of size K using a linear layer:

$$p_i = \text{softmax}(WE_i + b), \quad i \in [\sigma] \quad (4)$$

where $W \in \mathbb{R}^{m \times K}$ and $b \in \mathbb{R}^K$ are learnable parameters. This projection generates a probability distribution over the codebook entries.

To enable end-to-end differentiability, we adopt the Gumbel-Softmax trick to approximate discrete sampling during training. The final discrete token T_i for each position is selected via:

$$T_i = \arg \max_j p_i^{(j)}, \quad j \in [K] \quad (5)$$

Through this process, the windowed feature sequence Z is mapped into a discrete token sequence $T = \{T_1, T_2, \dots, T_d\}$, where each token encodes a local temporal pattern. These discrete representations provide a compact and informative abstraction of the raw input, facilitating downstream tasks such as classification and anomaly detection.

4.2 Masking Strategies

We describe the masking strategy used in the TMAE model as self-supervised learning [8]. The primary objective of this strategy is to reconstruct the hidden representations of masked windows and predict their corresponding discrete tokens, thereby enabling the learning of temporally structured semantic features from touch dynamic data.

To retain the positional information of each token after windowing, we first add positional encoding to the sequence embedding Z [39], resulting in:

$$Z^p = Z + \text{PositionEncoding}(Z), \quad Z^p \in \mathbb{R}^{d \times m} \quad (6)$$

Next, we split the window-convolved sequence Z into visible and masked representations. Let v_{index} denote the indices of the visible representations, and m_{index} denote the indices of the masked representations. The corresponding parts of the sequence can be denoted as:

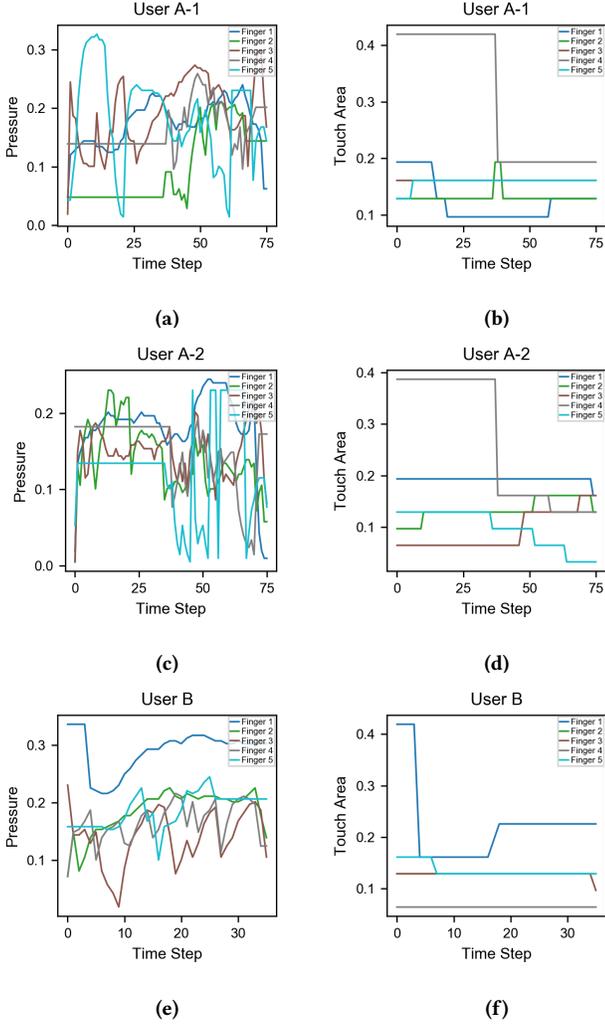


Figure 2: Comparison of touch dynamics for different gestures across users. (a), (b), (c), and (d) show the pressure and touch area of two instances of the same user. (e) and (f) represent the sequences of a different user.

$$Z_v = Z[v_{\text{index}}] \in \mathbb{R}^{d_v \times m}, \quad Z_m = Z[m_{\text{index}}] \in \mathbb{R}^{d_m \times m} \quad (7)$$

where $d_v + d_m = d$, indicating that the overall sequence length remains unchanged.

To ensure consistency, we similarly divide the discrete token sequence T into:

$$T_v = T[v_{\text{index}}] \in \mathbb{R}^{d_v}, \quad T_m = T[m_{\text{index}}] \in \mathbb{R}^{d_m} \quad (8)$$

During training, to represent the masked positions, we introduce a learnable mask token embedding vector $\mathbf{m} \in \mathbb{R}^m$ [13]. This vector is repeated for each masked position and combined with its positional encoding to form the masked input representation:

$$E_{\text{mask}} = \mathbf{m} \cdot \mathbf{1} + \text{PositionEncoding}(Z_m), \quad E_{\text{mask}} \in \mathbb{R}^{d_m \times m} \quad (9)$$

The mask token serves as a placeholder for the missing information and guides the model to learn to reconstruct the masked parts Z_m based on the contextual information from the visible part Z_v .

4.3 Self-supervised Regression

We perform two core pretext tasks: masked representation regression and discrete codeword prediction. In the following, we detail the implementation of these objectives and explain how our encoding strategy and momentum-based updates facilitate effective representation learning.

4.3.1 Multi-head Attention. The encoder in TMAE adopts the Multi-head Attention (MHA) mechanism [39, 42] to capture temporal dependencies across multiple subspaces.

Given the input sequence $X \in \mathbb{R}^{T \times m}$, each attention head h computes:

$$Q_h = XW_h^Q, \quad K_h = XW_h^K, \quad V_h = XW_h^V \quad (10)$$

$$\text{head}_h = \text{softmax}\left(\frac{Q_h K_h^\top}{\sqrt{d_{\text{head}}}}\right) V_h \quad (11)$$

The outputs of all heads are concatenated and linearly transformed:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_H)W^O \quad (12)$$

This structure allows the encoder to model complex temporal patterns in user interactions by leveraging multiple perspectives in parallel.

4.3.2 Masked Representation Regression. To reconstruct the representations of masked segments, we adopt a dual-encoder design consisting of a primary encoder and a momentum encoder. Both share the same Transformer architecture but differ in update strategies.

The primary encoder encodes the visible input:

$$R_v = \text{Encoder}(Z_v) \quad (13)$$

The momentum encoder, updated without gradients, encodes the masked tokens:

$$R_m = \text{Momentum_Encoder}(Z_m) \quad (14)$$

A cross-attention-based regressor predicts masked representations using visible context:

$$\hat{R}_m = \text{Regressor}(R_v, E_{\text{mask}}) \quad (15)$$

4.3.3 Discrete Codeword Prediction. To further enhance semantic learning, we introduce a discrete codeword prediction task. The tokenizer employs Gumbel-Softmax to discretize latent features. Given ground-truth discrete tokens $T_m \in \mathbb{R}^{d_m}$, the model predicts codeword distributions from reconstructed embeddings:

$$P(\hat{T}_m) = \text{Tokenizer.center}(\hat{R}_m) \quad (16)$$

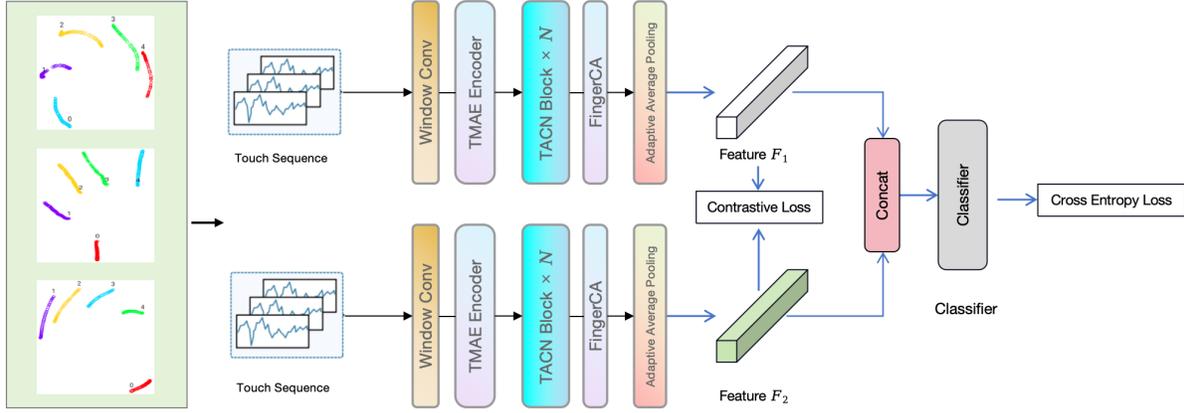


Figure 3: TouchSeqNet

4.3.4 *Momentum Encoder Updates.* The momentum encoder is updated via an exponential moving average of the primary encoder weights [7]:

$$\theta_m \leftarrow \mu\theta_m + (1 - \mu)\theta_e \quad (17)$$

Here, θ_m and θ_e denote the parameters of the momentum and primary encoders, respectively. A high smoothing factor μ (e.g., 0.99) ensures stable target representations for regression.

4.4 Self-supervised Loss

The self-supervised objective of our model consists of two components: an alignment loss and a discrete codeword prediction loss.

The alignment loss $\mathcal{L}_{\text{align}}$ minimizes the mean squared error (MSE) between the target representations from the momentum encoder and the predicted representations from the regressor:

$$\mathcal{L}_{\text{align}} = \text{MSE}(\text{rep_mask}, \text{rep_mask_prediction}) \quad (18)$$

The prediction loss $\mathcal{L}_{\text{pred}}$ uses cross-entropy to measure the discrepancy between the predicted token distributions and the ground-truth discrete tokens. We also monitor auxiliary metrics such as *Hits* and *NDCG@10* for evaluation.

The final loss function is a weighted sum of the two:

$$\mathcal{L} = \alpha\mathcal{L}_{\text{align}} + \beta\mathcal{L}_{\text{pred}} \quad (19)$$

where α and β control the contribution of each term.

5 TouchSeqNet Architecture

As illustrated in Figure 3, we propose TouchSeqNet, a contrastive learning framework designed for continuous user authentication based on dynamic touch data from mobile devices [21]. The architecture integrates pretrained temporal encoder via self-supervised learning on unlabeled behavioral sequences, a Temporal-Attentive Convolutional Network (TACN) module, and a hybrid loss that combines contrastive and cross-entropy objectives.

By leveraging a contrastive learning paradigm, TouchSeqNet extracts user-consistent yet discriminative representations that generalize well to real-world usage conditions.

5.1 Transfer Learning from TMAE

TouchSeqNet leverages the strengths of the pretrained TMAE model by transferring its temporal encoder into the downstream identity authentication framework. Specifically, we reuse two key modules from TMAE: (1) a window-based convolutional projection layer; (2) a multi-layer Transformer encoder pretrained via self-supervised masked representation regression.

These two modules act as the feature extractor in TouchSeqNet, transforming raw touch sequences into latent representations that encode both local and contextual temporal patterns. This transfer learning strategy enables TouchSeqNet to initialize from a representation space aligned with touch dynamics and temporal structure. As a result, the model starts with a strong inductive bias tailored for behavioral biometrics, allowing subsequent layers to focus on refining identity-specific discriminative features.

5.2 TACN block

To capture both local and global temporal patterns in dynamic touch sequences, the TACN module combines Temporal Convolutional Networks (TCN) [14] with Multi-head Attention to jointly model long-term dependencies and contextual correlations, while maintaining efficient and parallelizable computation. Compared with conventional CNNs, TCN supports longer effective memory via dilated convolutions without increasing parameter complexity [34].

5.2.1 *Temporal Convolutional Network.* The TCN is composed of stacked residual blocks, each containing two layers of dilated causal convolutions followed by weight normalization, ReLU activation, and dropout. A residual connection is added to stabilize training. The core computation in a residual block is given by:

$$y_i = \text{ReLU}(W_2 * (\text{ReLU}(W_1 * X_i + b_1)) + b_2), \quad (20)$$

where $W_1, W_2 \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in}} \times k}$ are convolutional kernels, and $*$ denotes the dilated causal convolution.

Dilated convolutions allow exponential expansion of the receptive field. Given a kernel f of size k and dilation factor d , the dilated convolution is defined as:

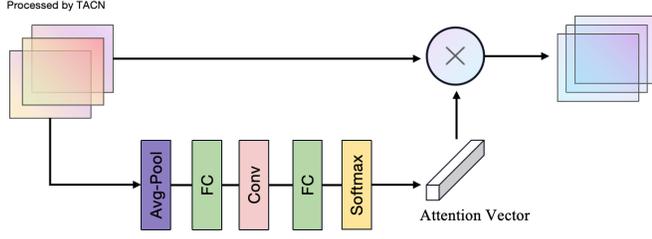


Figure 4: The structure of FingerCA module

$$F(s) = \sum_{i=0}^{k-1} f(i) \cdot X_{s-d \cdot i}. \quad (21)$$

By setting $d = 2^l$ for the l -th layer, the receptive field grows rapidly with depth while preserving the input length, which is essential for temporal alignment in authentication tasks.

Each residual block transforms the input sequence $X \in \mathbb{R}^{C_{in} \times L}$ as:

$$X^{(l)} = f(W^{(l)} *_d X^{(l-1)} + b^{(l)}), \quad X_{out}^{(l)} = X^{(l)} + X^{(l-1)}. \quad (22)$$

This structure allows efficient learning of long-range patterns while maintaining temporal consistency across layers.

5.2.2 Multi-head Attention and Hierarchical Temporal Fusion. While TCN are well-suited for extracting local patterns, they may struggle to fully capture long-range temporal dependencies that span across the entire sequence. To address this, we incorporate a Multi-head Attention mechanism after the TCN layers to enhance the model’s global temporal reasoning capabilities [28]. By leveraging multi-head attention, TACN can simultaneously attend to multiple temporal perspectives, enabling the model to better distinguish subtle variations in user touch dynamics, and laying a robust foundation for downstream identity authentication tasks.

5.3 FingerCA

As illustrated in Figure 1, To enhance discriminability after temporal modeling, we incorporate the FingerCA channel attention module [15, 22].

Let $X \in \mathbb{R}^{T \times C}$ be the output of the TACN blocks, where T is the number of time steps and C is the number of channels. A global average pooling is applied across time to produce a channel descriptor $z \in \mathbb{R}^C$:

$$z_c = \frac{1}{T} \sum_{t=1}^T X_{t,c} \quad (23)$$

This vector is passed through a two-layer MLP with non-linear activation to obtain attention weights $\alpha \in \mathbb{R}^C$, which are used to recalibrate the input features: $\tilde{X} = \alpha \circ X$.

5.4 Sample Pair Classification

For each input pair, we obtain two representations \tilde{F}_1 and \tilde{F}_2 , apply temporal average pooling, and concatenate the resulting vectors: $z = [z_1 || z_2] \in \mathbb{R}^{2C}$. A fully connected classification head then

predicts the probability that the pair belongs to the same user. This binary classification is trained end-to-end with a hybrid loss function.

5.5 Hybrid Loss Function

To jointly optimize representation learning and classification accuracy, we adopt a hybrid loss combining contrastive and cross-entropy terms [43]. Given embeddings (z_1, z_2) with label $y \in \{0, 1\}$:

$$\mathcal{L}_{\text{contrastive}} = y \cdot \|z_1 - z_2\|^2 + (1 - y) \cdot [\max(0, m - \|z_1 - z_2\|)]^2 \quad (24)$$

$$\mathcal{L}_{\text{CE}} = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y}) \quad (25)$$

The final training objective is a weighted sum of both terms:

$$\mathcal{L}_{\text{total}} = \lambda_1 \cdot \mathcal{L}_{\text{contrastive}} + \lambda_2 \cdot \mathcal{L}_{\text{CE}} \quad (26)$$

where λ_1 and λ_2 are hyperparameters balancing the two losses.

6 EXPERIMENTS

In this section, we describe our experimental protocol for evaluating TouchSeqNet. The evaluation consists of two stages: self-supervised pretraining of the TMAE encoder on unlabeled touch dynamics data, and fine-tuning the TouchSeqNet on labeled data. We conduct experiments on three datasets: our newly collected Ffinger dataset and two widely used public benchmarks, Touchalytics [12] and BioIdent [3].

6.1 Experimental Setup

6.1.1 Data Description. All datasets undergo the same preprocessing and normalization procedures to ensure consistency across experiments.

To handle variable-length sequences, we apply zero-padding along the temporal dimension. Specifically, for each sample $X \in \mathbb{R}^{T \times C}$, we pad it to length T_{pad} such that $T_{\text{pad}} \bmod \sigma = 0$, where σ is the convolution window size used in the pretraining stage. This ensures compatibility with window-based slicing.

We also construct a binary mask aligned with each sequence to mark valid positions. The mask is grouped into non-overlapping windows of size σ ; a window is marked as masked if more than half of its positions are padding. If masking is required, this window-level mask is then used during TMAE pretraining to guide the selection of visible and masked windows.

For all datasets, we adopt a contrastive pairing strategy. Each input to the model consists of a pair of samples, with the following labeling scheme:

- **Positive pair** ($y = 1$): both samples belong to the same user.
- **Negative pair** ($y = 0$): the two samples are drawn from different users.

This setup allows us to evaluate the model’s ability to learn identity-discriminative features under consistent experimental conditions.

6.1.2 Evaluation Methods. We assess the model using several standard classification metrics:

- **Accuracy:** the proportion of correctly predicted sample pairs over all pairs.

- F1 Score: the harmonic mean of precision and recall, computed as

$$F1 = \frac{2 \cdot P \cdot R}{P + R}. \quad (27)$$

- AUC: the Area Under the ROC Curve, indicating overall separability between positive and negative classes regardless of threshold.

6.1.3 Parameter settings. During the pretraining stage, we consistently set the embedding size to 64 for all models. The Adam optimizer is adopted as the default optimizer, with a fixed learning rate of 0.01 and no additional learning rate scheduling. The batch size is uniformly set to 128 across all experiments.

The Transformer encoder used in the TMAE model consists of 8 layers, each with 4 attention heads and a two-layer feed-forward network. A dropout rate of 0.2 is applied throughout the encoder [9]. In addition to the standard Transformer encoder, we further employ a 4-layer decoupled encoder to extract contextual representations from masked positions.

For each dataset, the slicing window size δ is selected from the candidate set {4, 8, 12}, and the default masking ratio is set to 40%. The vocabulary size of the discrete codebook in the tokenizer is fixed to 192 in our implementation.

In the fine-tuning stage, all hyperparameters shared with the pretraining stage are preserved. Additionally, the Temporal Convolutional Network (TCN) used in TouchSeqNet is configured with the input dimension `num_inputs` set to 64. The TCN consists of two residual blocks with output channel sizes defined as `num_channels` = [64, 128], and a dropout rate of 0.2. For each dataset, the convolutional kernel size is selected from {4, 5, 7} based on validation performance.

6.2 Experiment Results

6.2.1 Pre-training. Recent advances in self-supervised learning [5] have shown great promise in learning transferable representations from unlabeled data [24]. In many domains, it is common to pre-train a model on a large dataset and fine-tune it on target tasks. However, in the context of touch dynamics, datasets often differ significantly in terms of acquisition methods, gesture types, device specifications, and behavioral protocols.

To address this, we conduct self-supervised pretraining independently on the Touchalytics, BioIdent, and Ffinger datasets, and evaluate each model on its corresponding validation set. This setup ensures that the learned representations are adapted to the characteristics of each dataset and serve as a robust initialization for downstream fine-tuning.

6.2.2 Evaluations on TouchSeqNet. We evaluate the effectiveness of the proposed TouchSeqNet architecture on three representative touch dynamics datasets. Table 1 summarizes its performance on the held-out test sets. The results demonstrate that TouchSeqNet consistently achieves strong classification performance across different datasets, highlighting its robustness under diverse gesture and session conditions.

Across all evaluated datasets, TouchSeqNet consistently achieves high accuracy confirming its effectiveness in modeling fine-grained

Table 1: Performance of TouchSeqNet on All Datasets

Metric	Ffinger	BioIdent	Touchalytics
Accuracy	0.9769	0.9902	0.9908
F1 Score	0.9770	0.9908	0.9907
AUC	0.9769	0.9907	0.9908

touch dynamics for identity authentication. On the public benchmarks Touchalytics and BioIdent, it delivers near-perfect classification performance, demonstrating strong robustness and generalization in cross-user scenarios.

These results highlight the model’s ability to extract user-specific and gesture-aware representations, supporting its deployment as a unified solution for continuous authentication in both controlled and real-world environments.

6.2.3 Comparative Experiments. To further evaluate the effectiveness of the proposed TouchSeqNet model, we conduct comparative experiments against several strong baseline models under a consistent contrastive learning framework. Each model is used as a feature extractor for paired inputs, and a downstream classification head makes binary decisions. All models are trained and evaluated under the same experimental settings, and their performance is measured using classification accuracy across three datasets: Ffinger (our dataset), BioIdent, and Touchalytics.

The baseline models include:

- **TCN:** A dilated and causal convolutional network that captures long-range dependencies in time series data [19].
- **Gate-Transformer:** A lightweight attention-based model incorporating gating mechanisms to emphasize salient temporal features [23].
- **LSTM:** A recurrent neural network that models sequential dynamics through memory cells and gating structures [4].
- **InceptionTime:** A CNN-based model employing inception modules to extract multi-scale temporal features [17].
- **TSLANet:** A lightweight time series model featuring an Adaptive Spectral Block for Fourier-based denoising and an Interactive Convolution Block for efficient local feature extraction. [10].

As shown in Table 2, the proposed *TouchSeqNet* consistently outperforms all baseline methods across the three evaluated datasets. On the Ffinger dataset, it achieves the highest performance in all metrics, demonstrating strong robustness in modeling fine-grained and user-specific interaction patterns. In comparison, models such as TSLANet and LSTM show significant performance drops, indicating their limitations in capturing such behavioral variability.

On the public benchmarks BioIdent and Touchalytics, TouchSeqNet also reaches near-perfect results in terms of accuracy, F1 score, and AUC, outperforming or matching all competing approaches.

These results highlight the model’s superior generalization and discriminative capability, validating its effectiveness for continuous touch-based user authentication in practical scenarios.

6.2.4 Ablation Study. To assess the contributions of key components in TouchSeqNet, we conduct an ablation study across the Ffinger, BioIdent, and Touchalytics datasets. As summarized in

Table 2: Performance comparison of all models on Ffinger, BioIdent, and Touchalytics datasets (Accuracy / F1 / AUC).

Model	Ffinger			BioIdent			Touchalytics		
	Acc	F1	AUC	Acc	F1	AUC	Acc	F1	AUC
TCN[19]	0.9558	0.9582	0.9656	0.9052	0.9088	0.9050	0.9001	0.9033	0.9007
Gate-transformer[23]	0.9343	0.9384	0.9332	0.9794	0.9789	0.9792	0.9852	0.9852	0.9852
LSTM[4]	0.8851	0.8820	0.8812	0.9656	0.9638	0.9647	0.9713	0.9715	0.9713
InceptionTime[17]	0.8447	0.8644	0.8402	0.9886	0.9878	0.9885	0.9755	0.9754	0.9756
TSLANet[10]	0.7247	0.7453	0.7227	0.9831	0.9833	0.9831	0.9415	0.9426	0.9419
TouchSeqNet (ours)	0.9769	0.9770	0.9769	0.9902	0.9908	0.9907	0.9908	0.9907	0.9908

Table 3, we compare the full model against three variants: (1) removing the multi-head attention in TACN (*TACN w/o Attention*), (2) removing the Pretrained-module (*TACN w/o Pretrained-module*), and (3) using only the Pretrained-module (*Only Pretrained-module*).

Removing multi-head attention results in a clear drop in performance across all datasets (e.g., 97.69% to 94.05% on Ffinger), highlighting the importance of attention in capturing global temporal dependencies. Eliminating the Pretrained-module leads to an even more significant accuracy loss on BioIdent (from 99.02% to 83.56%) and Touchalytics, demonstrating the effectiveness of transfer learning from TMAE. Using only the Pretrained-module yields the lowest scores overall, confirming that pretraining alone is insufficient and must be complemented by downstream temporal modeling.

These results underscore the complementary benefits of the Pretrained-module and TACN components. While the encoder provides strong generalizable features, the attention-augmented temporal modeling in TACN is crucial for extracting task-specific discriminative representations.

Table 3: Ablation study results on classification accuracy

Model Variant	Ffinger	BioIdent	Touchalytics
TACN w/o Attention	0.9405	0.9806	0.9856
Only Pretrained-module	0.9325	0.8733	0.9366
TACN w/o Pretrained-module	0.9763	0.8356	0.8655
TouchSeqNet (Full)	0.9769	0.9902	0.9908

6.2.5 Summary of Experimental Findings. Experimental results across Ffinger, BioIdent, and Touchalytics confirm the effectiveness and robustness of the proposed TouchSeqNet framework for dynamic touch-based authentication.

In comparative experiments, TouchSeqNet consistently outperforms strong baselines such as TCN, Gate-Transformer, LSTM, and InceptionTime. Its performance gains are most evident on the challenging Ffinger dataset, demonstrating its strength in modeling fine-grained temporal and identity-specific patterns.

Ablation results highlight the complementary roles of the core components. The multi-head attention mechanism in TACN enhances temporal discriminability, while the pre-trained module

from TMAE provides transferable features that improve generalization, especially in low-data regimes. Removing either module leads to notable performance degradation.

Additionally, the model achieves near-perfect accuracy and F1 scores on public datasets, underscoring its strong generalization to different devices and behavioral contexts. The integration of contrastive learning and hierarchical temporal modeling enables robust discrimination between genuine and impostor pairs across a wide range of conditions.

7 Conclusion

We presented TouchSeqNet, a contrastive learning framework designed for continuous user authentication via touch dynamics. By integrating self-supervised pretraining, hierarchical temporal modeling, and attention mechanisms, the model learns rich, discriminative representations without requiring handcrafted features or domain-specific heuristics.

Our architecture demonstrates the synergistic value of transfer learning and structured temporal modeling. The pre-trained module offers a generalizable feature space, while the TACN block enhances temporal resolution and discriminability—together enabling robust identity verification across varying users and contexts.

Looking ahead, our framework provides a foundation for scalable behavioral biometrics. Its modular design makes it extensible to cross-device scenarios, federated authentication systems, and even multi-modal interaction signals (e.g., stylus, handwriting, or gesture input), pointing to broad applicability in real-world human-computer interaction systems.

References

- [1] Alejandro Acién, Aythami Morales, John V Monaco, Ruben Vera-Rodriguez, and Julian Fierrez. 2021. TypeNet: Deep learning keystroke biometrics. *IEEE Transactions on Biometrics, Behavior, and Identity Science* 4, 1 (2021), 57–70.
- [2] Sara Amini, Vahid Noroozi, Amit Pande, Satyajit Gupte, Philip S Yu, and Chris Kanich. 2018. Deepauth: A framework for continuous user re-authentication in mobile apps. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*. 2027–2035.
- [3] Margit Antal, Zsolt Bokor, and László Zsolt Szabó. 2015. Information revealed from scrolling interactions on mobile devices. *Pattern Recognition Letters* 56 (2015), 7–13.
- [4] Asrar Bajaber, Mai Fadel, and Lamiaa Elrefaei. 2022. Evaluation of Deep Learning Models for Person Authentication Based on Touch Gesture. *Computer Systems Science & Engineering* 42, 2 (2022).
- [5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural*

- information processing systems* 33 (2020), 1877–1901.
- [6] Mario Parreño Centeno, Yu Guan, and Aad van Moorsel. 2018. Mobile based continuous authentication using deep features. In *Proceedings of the 2nd international workshop on embedded and mobile deep learning*. 19–24.
 - [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*. PMLR, 1597–1607.
 - [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*. 4171–4186.
 - [9] Emadeldeen Eldele, Mohamed Ragab, Zhenghua Chen, Min Wu, Chee Keong Kwoh, Xiaoli Li, and Cuntai Guan. 2021. Time-series representation learning via temporal and contextual contrasting. *arXiv preprint arXiv:2106.14112* (2021).
 - [10] Emadeldeen Eldele, Mohamed Ragab, Zhenghua Chen, Min Wu, and Xiaoli Li. 2024. Tslanet: Rethinking transformers for time series representation learning. *arXiv preprint arXiv:2404.08472* (2024).
 - [11] Julian Fierrez, Ada Pozo, Marcos Martínez-Díaz, Javier Galbally, and Aythami Morales. 2018. Benchmarking touchscreen biometrics for mobile authentication. *IEEE transactions on information forensics and security* 13, 11 (2018), 2720–2733.
 - [12] Mario Frank, Ralf Biedert, Eugene Ma, Ivan Martinovic, and Dawn Song. 2012. Touchalytics: On the applicability of touchscreen input as a behavioral biometric for continuous authentication. *IEEE transactions on information forensics and security* 8, 1 (2012), 136–148.
 - [13] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. 2021. Mae: Masked autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377* (2021).
 - [14] Pradeep Hewage, Ardhendu Behera, Marcello Trovati, Ella Pereira, Morteza Ghahremani, Francesco Palmieri, and Yonghui Liu. 2020. Temporal convolutional neural (TCN) network for an effective weather forecasting using time-series data from the local weather station. *Soft Computing* 24 (2020), 16453–16482.
 - [15] Jie Hu, Li Shen, and Gang Sun. 2018. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7132–7141.
 - [16] Shima Imani and Eamonn Keogh. 2021. Multi-window-finder: domain agnostic window size for time series data. *Proceedings of the MileTS 21* (2021).
 - [17] Hassan Ismail Fawaz, Benjamin Lucas, Germain Forestier, Charlotte Pelletier, Daniel F Schmidt, Jonathan Weber, Geoffrey I Webb, Lhassane Idoumghar, Pierre-Alain Muller, and François Petitjean. 2020. Inceptiontime: Finding alexnet for time series classification. *Data Mining and Knowledge Discovery* 34, 6 (2020), 1936–1962.
 - [18] Parker Lamb, Alexander Millar, and Ramon Fuentes. 2020. Swipe dynamics as a means of authentication: Results from a bayesian unsupervised approach. In *2020 IEEE International Joint Conference on Biometrics (IJCB)*. IEEE, 1–9.
 - [19] Colin Lea, Michael D Flynn, Rene Vidal, Austin Reiter, and Gregory D Hager. 2017. Temporal convolutional networks for action segmentation and detection. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 156–165.
 - [20] Seunghan Lee, Taeyoung Park, and Kibok Lee. 2023. Learning to embed time series patches independently. *arXiv preprint arXiv:2312.16427* (2023).
 - [21] Seunghan Lee, Taeyoung Park, and Kibok Lee. 2023. Soft contrastive learning for time series. *arXiv preprint arXiv:2312.16424* (2023).
 - [22] Lin Lin, Jinlei Wu, Song Fu, Sihao Zhang, Changsheng Tong, and Lizheng Zu. 2024. Channel attention & temporal attention based temporal convolutional network: A dual attention framework for remaining useful life prediction of the aircraft engines. *Advanced Engineering Informatics* 60 (2024), 102372.
 - [23] Minghao Liu, Shengqi Ren, Siyuan Ma, Jiahui Jiao, Yizhou Chen, Zhiguo Wang, and Wei Song. 2021. Gated transformer networks for multivariate time series classification. *arXiv preprint arXiv:2103.14438* (2021).
 - [24] Mingsheng Long, Jianmin Wang, Yue Cao, Jianguang Sun, and Philip S Yu. 2016. Deep learning of transferable representation for scalable domain adaptation. *IEEE Transactions on Knowledge and Data Engineering* 28, 8 (2016), 2027–2040.
 - [25] Sakorn Mekruksavanich and Anuchit Jitpattanakul. 2021. Deep learning approaches for continuous authentication based on activity patterns using mobile sensing. *Sensors* 21, 22 (2021), 7519.
 - [26] Asif Iqbal Middy, Sarbani Roy, Saptarshi Mandal, and Rahul Talukdar. 2021. Privacy protected user identification using deep learning for smartphone-based participatory sensing applications. *Neural Computing and Applications* 33 (2021), 17303–17313.
 - [27] Natalia Neverova, Christian Wolf, Griffin Lacey, Lex Fridman, Deepak Chandra, Brandon Barbelo, and Graham Taylor. 2016. Learning human identity from motion patterns. *IEEE Access* 4 (2016), 1810–1820.
 - [28] Xuran Pan, Chunjiang Ge, Rui Lu, Shiji Song, Guanfu Chen, Zeyi Huang, and Gao Huang. 2022. On the integration of self-attention and convolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 815–825.
 - [29] Chao Shen, Yuanxun Li, Yufei Chen, Xiaohong Guan, and Roy A Maxion. 2017. Performance analysis of multi-motion sensor behavior for active smartphone authentication. *IEEE Transactions on Information Forensics and Security* 13, 1 (2017), 48–62.
 - [30] Chao Shen, Yong Zhang, Xiaohong Guan, and Roy A Maxion. 2015. Performance analysis of touch-interaction behavior for active smartphone authentication. *IEEE Transactions on Information Forensics and Security* 11, 3 (2015), 498–513.
 - [31] Zhihao Shen, Shun Li, Xi Zhao, and Jianhua Zou. 2022. MMAuth: A continuous authentication framework on smartphones using multiple modalities. *IEEE Transactions on Information Forensics and Security* 17 (2022), 1450–1465.
 - [32] Onsiri Silasai and Wachana Khowfa. 2020. The study on using biometric authentication on mobile device. *NU Int. J. Sci* 17 (2020), 90–110.
 - [33] Yunpeng Song and Zhongmin Cai. 2022. Integrating handcrafted features with deep representations for smartphone authentication. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 6, 1 (2022), 1–27.
 - [34] Yan Song, Shengyao Gao, Yibin Li, Lei Jia, Qiqiang Li, and Fuzhen Pang. 2020. Distributed attention-based temporal convolutional network for remaining useful life prediction. *IEEE Internet of Things Journal* 8, 12 (2020), 9594–9602.
 - [35] Giuseppe Stragapede, Paula Delgado-Santos, Ruben Tolosana, Ruben Vera-Rodríguez, Richard Guest, and Aythami Morales. 2023. Mobile keystroke biometrics using transformers. In *2023 IEEE 17th International Conference on Automatic Face and Gesture Recognition (FG)*. IEEE, 1–6.
 - [36] Giuseppe Stragapede, Ruben Vera-Rodríguez, Ruben Tolosana, Aythami Morales, Alejandro Acien, and Gaël Le Lan. 2022. Mobile behavioral biometrics for passive authentication. *Pattern Recognition Letters* 157 (2022), 35–41.
 - [37] Yu Sun, Qiyuan Gao, Xiaofan Du, and Zhao Gu. 2019. Smartphone User Authentication Based on Holding Position and Touch-Typing Biometrics. *Computers, Materials & Continua* 61, 3 (2019).
 - [38] Ruben Tolosana, Ruben Vera-Rodríguez, Julian Fierrez, and Javier Ortega-García. 2020. BioTouchPass2: Touchscreen password biometrics using time-aligned recurrent neural networks. *IEEE Transactions on Information Forensics and Security* 15 (2020), 2616–2628.
 - [39] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
 - [40] Cong Wang, Yanru Xiao, Xing Gao, Li Li, and Jun Wang. 2021. A framework for behavioral biometric authentication using deep metric learning on mobile devices. *IEEE Transactions on Mobile Computing* 22, 1 (2021), 19–36.
 - [41] Huanran Wang, Hui He, Chen Song, Hao Tang, Yanwei Sun, Yanchen Qiao, and Weizhe Zhang. 2022. Who Is Using the Phone? Representation-Learning-Based Continuous Authentication on Smartphones. *Security and Communication Networks* 2022, 1 (2022), 6339407.
 - [42] Jingyun Xiao, Ran Liu, and Eva L Dyer. 2024. Gafomer: Enhancing timeseries transformers through group-aware embeddings. In *The Twelfth International Conference on Learning Representations*.
 - [43] Iqra Zahid, Yue Chang, Tharindu Madusanka, Youcheng Sun, and Riza Theresa Batista-Navarro. 2024. Multi-Loss Fusion: Angular and Contrastive Integration for Machine-Generated Text Detection. In *Findings of the Association for Computational Linguistics: EMNLP 2024*. 7189–7202.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009