



# Vilniaus universitetas Duomenų mokslo ir skaitmeninių technologijų institutas

Informatikos krypties doktorantų atestacinė konferencija  
Veiklos ataskaita už 2023 m. rugsėjo 27 d. – 2024 m. kovo 28 d.

## ANOMALINIŲ ĮVYKIŲ IDENTIFIKAVIMAS IR JŲ UŽKARDYMAS KOMPIUTERIŲ TINKLUOSE TAIKANT MAŠININIO MOKYMOSI METODUS

**dokt. Arnoldas BUDŽYS** – Informatika N 009

**Studijų metai: IV**

Darbo vadovas: dr. Viktor Medvedev

Doktorantūros pradžios ir pabaigos metai: 2020–2024



## • STUDIJŲ PLANAS IR JO VYKDYMO SUVESTINĖ

Studijų metai	Egzaminai <sup>1</sup>	
	Planas	Įvykdyta
I (2020/2021)	1	1
II (2021/2022)	3	3
III (2022/2023)		
IV (2023/2024)		
<b>Iš viso:</b>	<b>4</b>	<b>4</b>

Studijų metai	Dalyvavimas konferencijose				Publikacijos					
	Tarptautinėse <sup>2</sup>		Nacionalinėse <sup>3</sup>		Su citav. rodikliu <sup>4</sup>			Be citav. rodiklio <sup>5</sup>		
	Planas	Įvykdyta	Planas	Įvykdyta	Planas	Įvykdyta <sup>6</sup>	Būklė <sup>7</sup>	Planas	Įvykdyta <sup>6</sup>	Būklė <sup>7</sup>
I (2020/2021)										
II (2021/2022)	1	1	1	1				1	0	
III (2022/2023)	1	1+1***	0	1	1	1	Publikuota	1	2*	Publikuota
IV (2023/2024)			<b>0</b>	<b>1</b>	<b>1</b>	<b>0</b>	<b>Įteikta**</b>			
<b>Iš viso:</b>	<b>2</b>	<b>2+1***</b>	<b>1</b>	<b>3</b>	<b>2</b>	<b>1</b>		<b>1</b>	<b>2</b>	

\*HCI2023 – publikuota (CA WoS, Springer), CISTI2023 – publikuota (CA WoS, IEEE )

\*\* Publikacija pataisyta atsižvelgus į recenzijas (major revision), pakartotinai įteikta 2024 m. vasario 03 d.

\*\*\* Prisidėta prie pranešimo bei straipsnio konferencijų medžiagoje parengimo.

# ATASKAITINIS STUDIJŲ PUSMETIS (IV: 2023/2024 – I pusmetis)

3

## Publikacijos 2023/2024 (I pusmetis)

Planas	Įvykdyta	Būklė	Publikacijos tipas
Informatics	Kurasova, O.; <b>Budžys, A.</b> ; Medvedev, V. <i>Exploring multidimensional embeddings for decision support using advanced visualization techniques</i> // Informatics. Basel : MDPI. eISSN 2227-9709. 2024, vol. 11, iss. 1, art. no. 11, p. [1-17]. DOI: 10.3390/informatics11010011. [Emerging Sources Citation Index (Web of Science)] [CiteScore: 4,80; SNIP: 1,023; SJR: 0,545; Q1 (2022, Scopus Sources)] [S1]	Publikuota	(IF – 3.1) CA WoS duomenų bazėje
Artificial Intelligence Review (Springer Nature)	<b>Budžys, A.</b> , Kurasova, O., Medvedev, V. <i>Deep Learning-Based Authentication for Insider Threat Detection in Critical Infrastructure</i> . Artificial Intelligence Review (2024). [Science Citation Index Expanded (SCIE) (Web of Science), IF – 12; Q1(Computer Science, Artificial Intelligence)]	Įteikta (gautos pirmos recenzijos): 2024 m. vasario 3 d.	(IF - 12) CA WoS duomenų bazėje



# ATASKAITINIS STUDIJŲ PUSMETIS (IV: 2023/2024 – I pusmetis)

4

## Publikacijos 2023/2024 (I pusmetis)

Planas	Įvykdyta	Būklė	Publikacijos tipas
COMPUTER STANDARDS & INTERFACES (ELSEVIER)	<b>Budžys, A.</b> , Kurasova, O., Medvedev, V. <i>Data Fusion Based Authentication Using Complex Keystroke Dynamics Analysis</i> . Computer Standards & Interfaces (2024). [Science Citation Index Expanded (SCIE) (Web of Science), IF – 5; Q1 (COMPUTER SCIENCE, HARDWARE & ARCHITECTURE) ; Q1(COMPUTER SCIENCE, SOFTWARE ENGINEERING)]	Bus įteikta 2024 m. balandžio 12 d.	(IF - 5) CA WoS duomenų bazėje



# DOKTORANTŪROS STUDIJŲ PASIEKIMAI

5

## Dalyvavimas tarptautinėse konferencijose

1. Budžys, A., Kurasova, O., and Medvedev, V., „Deep learning-based prevention of insider threats using user behavioral keystroke biometrics“, 32nd European Conference on Operational Research (EURO XXXII)], Espoo, Finland, July 3-6, 2022.
2. Budžys, A., Kurasova, O., and Medvedev, V., “Behavioral Biometrics Authentication Using Siamese Neural Networks”, in HCI for Cybersecurity, Privacy and Trust 5th International Conference, HCI-CPT 2023, Held as Part of the 25th HCI International Conference, HCI2023, 2023, pp. 1–14 (in press).

# DOKTORANTŪROS STUDIJŲ PASIEKIMAI

6

## Publikacijos (tik su citavimo rodikliu)

	Bibliografinis aprašas	Būklė
1.	Kurasova, O.; <b>Budžys, A.</b> ; Medvedev, V. Exploring multidimensional embeddings for decision support using advanced visualization techniques // Informatics. Basel : MDPI. eISSN 2227-9709. 2024, vol. 11, iss. 1, art. no. 11, p. [1-17]. DOI: 10.3390/informatics11010011. [Emerging Sources Citation Index (Web of Science)] [CiteScore: 4,80; SNIP: 1,023; SJR: 0,545; Q1 (2022, Scopus Sources)] [S1]	Publikuota
2.	<b>Budžys, A.</b> , Kurasova, O., Medvedev, V. Deep Learning-Based Authentication for Insider Threat Detection in Critical Infrastructure. Artificial Intelligence Review (2024). [Science Citation Index Expanded (SCIE) (Web of Science), IF – 12; Q1(Computer Science, Artificial Intelligence)]	Įteikta (gautos pirmos recenzijos): 2024 m. vasario 3 d.



# Doktorantūros mokslinių tyrimų ir disertacijos rengimo etapai

Darbo pavadinimas	Atlikimo terminai	Pastabos
<b>Mokslinių tyrimų disertacijos tema apžvalga ir analizė (Lietuvoje ir užsienyje):</b>	2020 m. spalio mėn. – 2021 m. rugsėjo mėn.	Atlikus literatūros analizę pavyko identifikuoti problemos sprendimo būdus panaudojant dirbtinius neuroninius tinklus.
<b>Mokslinio tyrimo vykdymas:</b> 2.1. Tyrimo metodikos sudarymas: 2.1.1. Tyrimo metodikos iškeltiems uždaviniams spręsti parinkimas; 2.1.2. Teorinio ir empirinio tyrimų suplanavimas pagal pasirinktą metodiką. 2.2. <b>Teorinis tyrimas:</b> 2.2.1. Mašininio mokymosi metodų, naudojamų kompiuterių tinkluose įsilaužimų prevencijai, tyrimas. 2.2.2. 2.3. <b>Empirinis tyrimas:</b> 2.3.1. Sudarytų metodų pritaikymas praktinių uždavinių sprendimui. 2.3.2. Gautų duomenų analizė, rezultatų apibendrinimas, išvadų parengimas.	2021 m. spalio mėn. – 2022 m. sausio mėn.  2022 m. vasario mėn. – 2022 m. rugsėjo mėn.  <b>2022 m. spalio mėn. – 2023 m. rugsėjo mėn.</b>  Pateiktas naujas autoriaus pasiūlytas metodas (GAFMAT) skaitiniams duomenims konvertuoti į vaizdinius. Siekiant įvertinti šio naujo metodo veiksmingumą, atlikta išsami lyginamoji analizė naudojant Siamo neuroninius tinklus, gauti rezultatai palyginti su esamomis metodikomis, aprašytomis literatūroje. Gauti eksperimentinio tyrimo rezultatai lyginami tarpusavyje. Tyrimai rodo, jog konvertavus skaitines reikšmes į vaizdus galima pagerinti vartotojų klasifikavimo rezultata, lyginant su mašininio mokymosi algoritmais, bei klasikinais dirbtiniais neuroniniais tinklais kuomet klasifikavimui naudojami skaitiniai duomenys.	

# Doktorantūros mokslinių tyrimų ir disertacijos rengimo etapai

Atskirų daktaro disertacijos dalių (tyrimo metodikos, rezultatų, ginamų teiginių, išvadų, ir kt.) parengimas:

- 3.1. Tikslų, uždavinių, tyrimo metodikos, ginamųjų teiginių patikslinimas;
- 3.2. Analitinės disertacijos dalies parengimas;
- 3.3. Teorinės disertacijos dalies parengimas;
- 3.4. Eksperimentinės disertacijos dalies parengimas;
- 3.5. Bendrųjų išvadų formulavimas.

**2023 m. spalio mėn. - 2024 m. gegužės mėn.**

Pasiūlyta dimensijos mažinimu pagrįsta vizualizavimo metodika, kurioje unikaliai integruojami dimensijos mažinimo metodai ir Siamo neuroniniai tinklai su trigubo nuostolio funkcija. Metodika sukurta sudėtingiems daugiamačiams giliųjų neuroninių tinklų įterpiniams interpretuoti ir vizualizuoti, siekiant pagerinti didelės apimties duomenų interpretavimą ir analizę. Siūlomos metodikos efektyvumas įvertintas panaudojant klavišų paspaudimų dinamikos duomenų rinkinį, skirtą naudotojo autentifikavimo problemai spręsti. Pasiūlytas naujas metodas GAFMAT (GAbor Filter MATrix Transformation), skirtas nevaizdiniais duomenimis transformuoti į vaizdus. Šis klavišų paspaudimų dinamikos transformavimas į vaizdus atskleidžia esminius elgesio bruožus, susijusius su slaptažodžio įvedimu. Pasiūlyto metodo efektyvumas įrodytas naudojant viešai prieinamus duomenų rinkinius. Pasiūlyta įvairių klavišų paspaudimų dinamikos ilgių suvienodinimo strategija skirtingiems duomenų rinkiniams. Tai leidžia sukurti giliuoju mokymusi pagrįstą modelį, kuris gali prisitaikyti prie bet kokio ilgio slaptažodžių, kad būtų galima sukurti tikslesnę ir bendresnę naudotojo autentifikavimo sistemą.

Daktaro disertacijos parengimas ir svarstymas padalinyje

2024 m. birželio mėn.

Daktaro disertacijos gynimas

2024 m. rugsėjo mėn.



# Disertacijos tema, tyrimo objektai ir tikslas

9

## Preliminari disertacijos tema:

- Anomalinių įvykių identifikavimas ir jų užkardymas kompiuterių tinkluose taikant mašininio mokymosi metodus.

## Tyrimo objektai:

- vartotojo sugeneruoti klaviatūros, pelės biometriniai duomenys, bei mašininio mokymosi metodų taikymas anomalinių įvykių identifikavimui ir neteisėtų veiksmų užkardymui.

## Tikslas:

- pasiūlyti metodiką sistemos vartotojui autentifikuoti pagal jo biometrinius elgsenos duomenis siekiant užkardyti insaiderio veiklą bei apsaugoti sistemą nuo jo neteisėtų veiksmų.



# Tyrimo uždaviniai

10

- Atlikti išsamią literatūros analitinę apžvalgą, siekiant identifikuoti tinkamus metodus anomalinių įvykių identifikavimui ir insaiderio užkardymui kompiuterių tinkluose;
- Atlikti skirtingų mašininio mokymosi metodų, skirtų anomalinių įvykių identifikavimui ir insaiderio užkardymui kompiuterių tinkluose, analizę ir tyrimą;
- Sukurti metodiką, apimančią mašininio mokymosi grįstus algoritmus, sistemos vartotojui autentifikuoti pagal jo biometrinius elgsenos duomenis;
- Įvertinti sukurtos metodikos efektyvumą realaus laiko duomenims atliekant eksperimentinius tyrimus;
- Atlikti gautų rezultatų analizę: rezultatų apibendrinimas, išvadų parengimas.



# Methodology

11

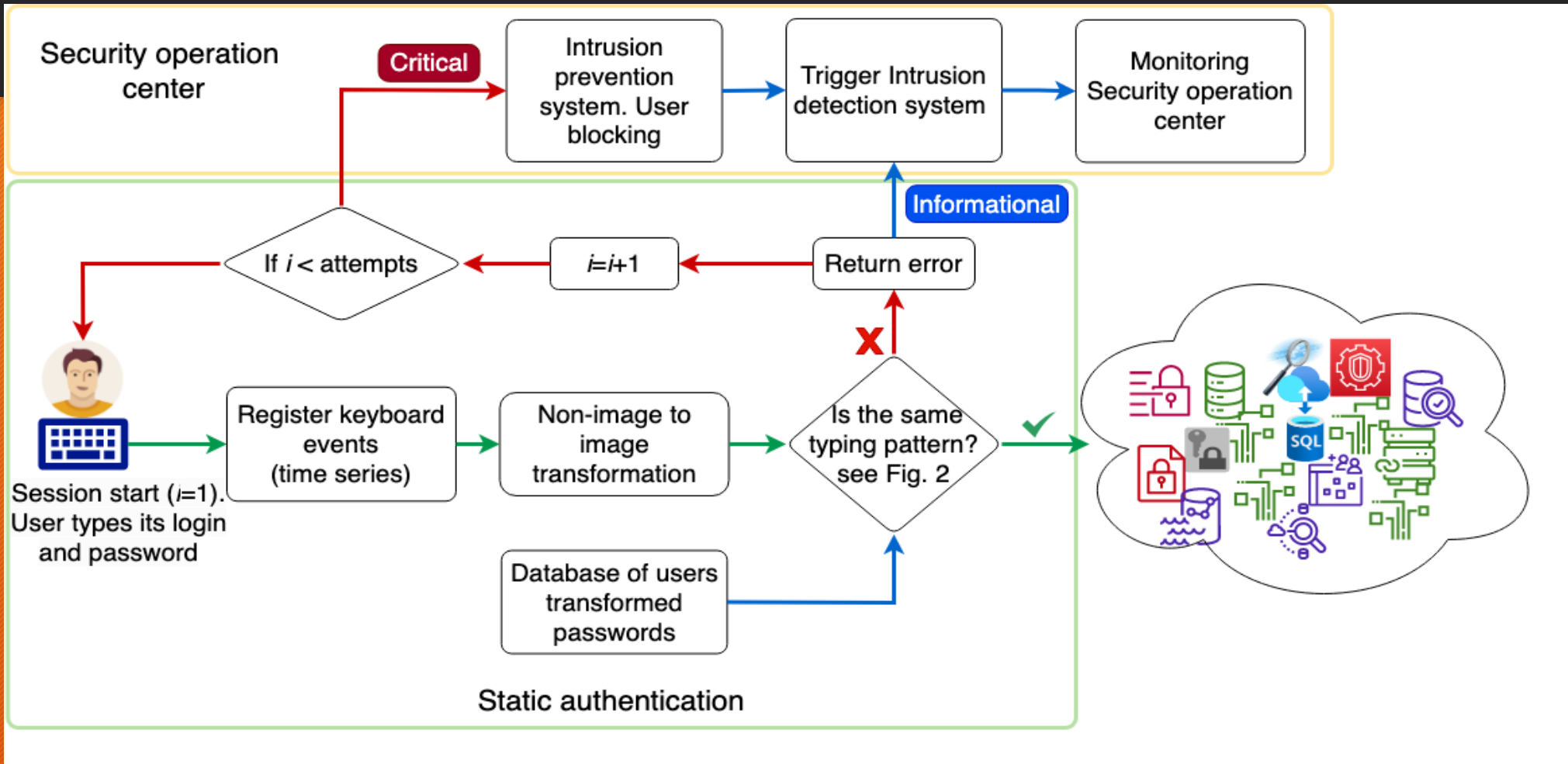


Fig. 1 Schematic representation of the user authentication process using an intrusion detection system and an intrusion prevention system based on user typing behavior.

# Methodology

12

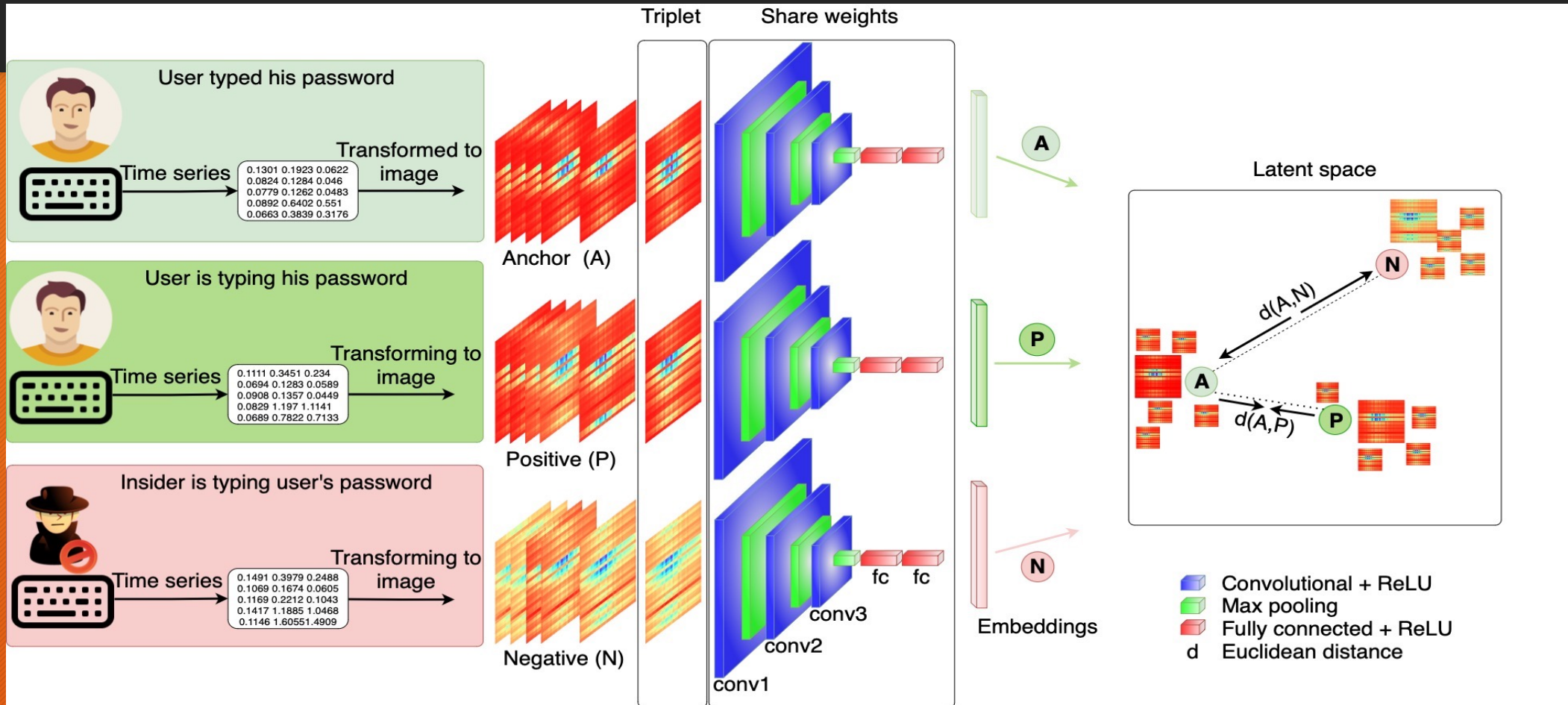
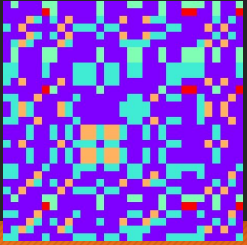


Fig. 2 Schematic representation of the proposed framework for time series transformation from keystroke biometric data features into images and training process of Siamese neural network with CNN branches

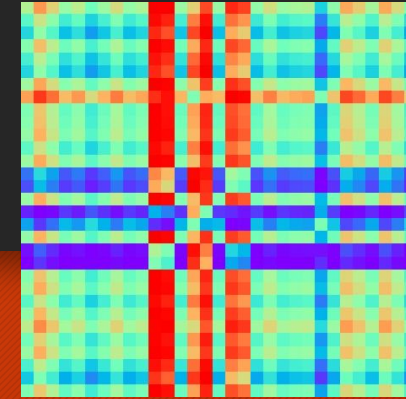


Markov Transition Field (MTF)

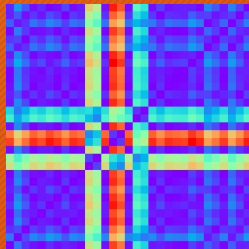


# Non-Image to Image

Gramian Angular Difference Field  
(GADF)



Recurrence Plots (RPs)



Gramian Angular  
Summation Field (GASF)

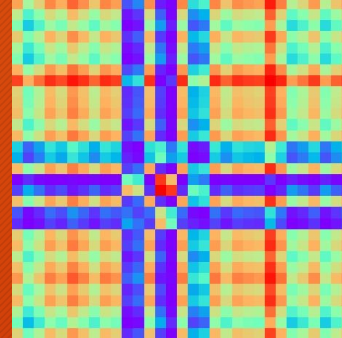
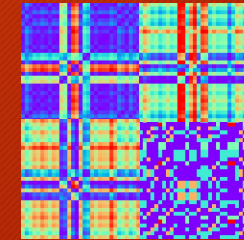


Image Fusion (combined)



Carnegie Mellon University dataset:  
Class: 51; Objects: 20400  
Objects: 20400

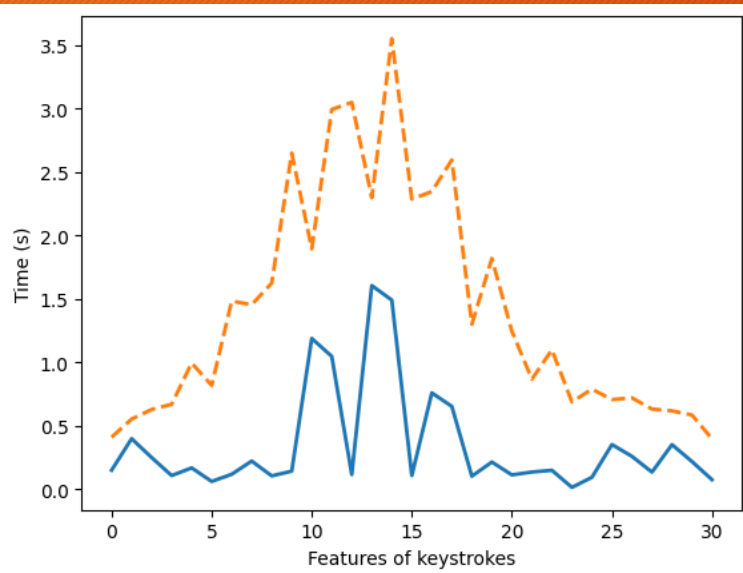


**Algorithm 1** Gabor Filter algorithm

```

1: function GaborFilter(discrete_signal,  $\sigma$ ,  $\theta$ ,  $\lambda$ ,  $\psi$ ,  $\gamma$ )
2:    $n \leftarrow$  length of discrete_signal
3:   Initialize  $x$  as array of size  $n$  generating evenly-
   spaced values in an interval  $(-3\sigma, 3\sigma)$ 
4:    $x \leftarrow x \cdot \cos(\theta)$ 
5:   Initialize gabor as an empty array of size  $n$ 
6:    $gabor \leftarrow \exp\left(-0.5 \cdot \left(\frac{x}{\sigma}\right)^2\right) \cdot \cos\left(2\pi \cdot \frac{x}{\lambda} + \psi\right)$ 
7:    $gabor \leftarrow \frac{gabor}{\sqrt{\sum_{i=0}^{n-1} gabor[i]^2}}$ 
8:    $gabor \leftarrow \text{Convolution}(\text{discrete\_signal}, gabor)$ 
9:   return gabor
10: end function

```



# Non-Image to Image

**Algorithm 2** GAFMAT algorithm

**Require:** *discrete\_signal*,  $\sigma\_list$ ,  $\theta\_list$ ,  $\lambda\_list$ ,  $\psi\_list$ ,  $\gamma\_list$

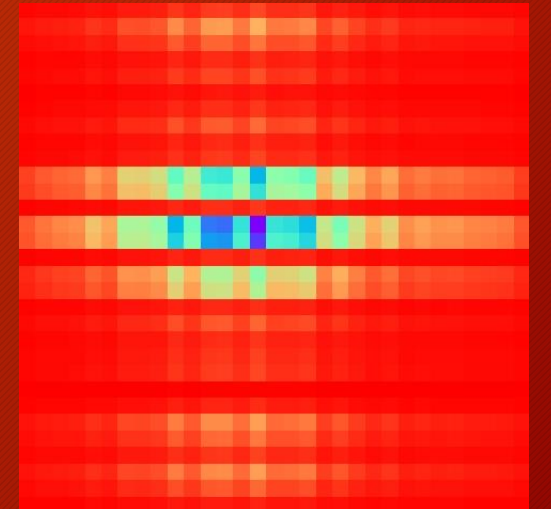
```

1:  $n \leftarrow$  length of discrete_signal
2: image  $\leftarrow$  create zero array of size  $n$ 
3: combinations  $\leftarrow$  CartesianProduct( $\sigma\_list$ ,  $\theta\_list$ ,
    $\lambda\_list$ ,  $\psi\_list$ ,  $\gamma\_list$ )
    $\triangleright$  The set of all possible pairs (see Table 2)
4: for each  $(\sigma, \theta, \lambda, \psi, \gamma)$  in combinations do
5:   gabortemp  $\leftarrow$  gabor(discrete_signal,  $\sigma$ ,  $\theta$ ,  $\lambda$ ,  $\psi$ ,  $\gamma$ )
    $\triangleright$  see Algorithm 1
6:   gabor  $\leftarrow$  transpose(gabortemp)
7:   image2D  $\leftarrow$  OuterProduct(image, gabor)
8: end for
9: return image2D

```

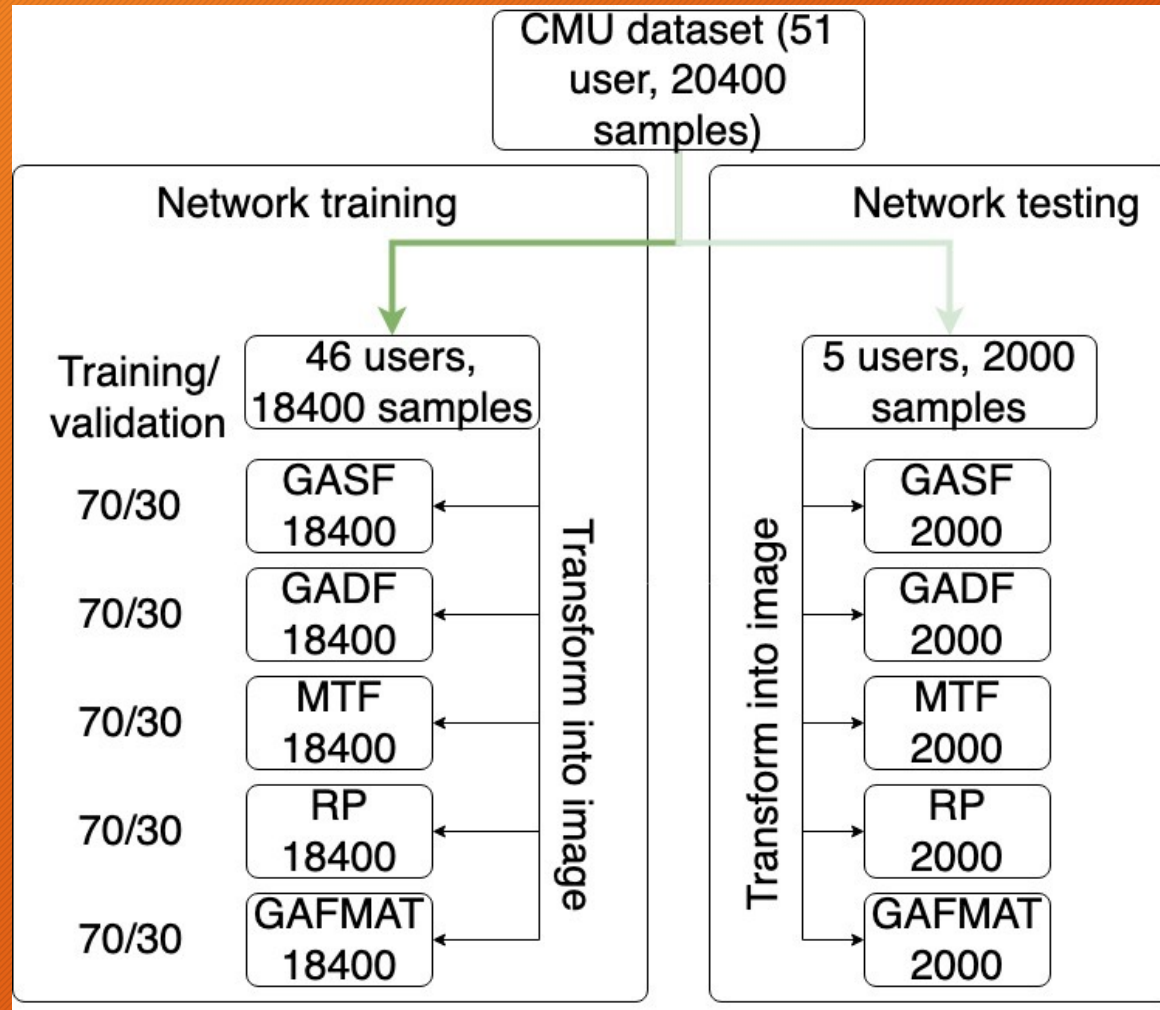
$$image2D = \begin{bmatrix} a_1b_1 & a_1b_2 & \cdots & a_1b_n \\ a_2b_1 & a_2b_2 & \cdots & a_2b_n \\ \vdots & \vdots & \ddots & \vdots \\ a_nb_1 & a_nb_2 & \cdots & a_nb_n \end{bmatrix}$$

## GABOR FILTER MATRIX TRANSFORMATION (GAFMAT)

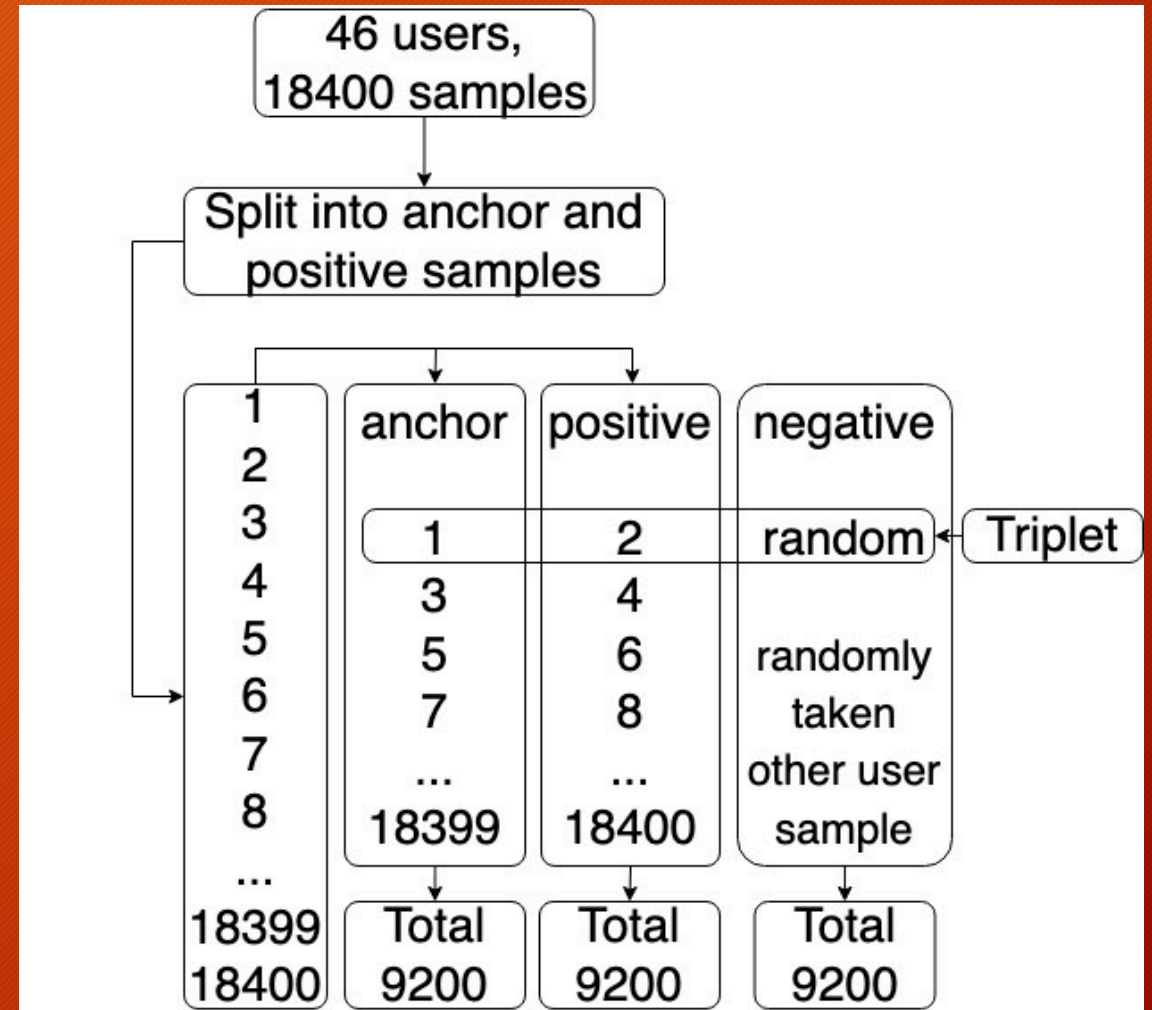




# Data Preparation



a)



b)

Fig. 3 a) The process of preparing data for model training/validation and testing; b) Splitting CMU data into an anchor and positive samples for each transformed dataset using GASF, GADF, MTF, RP, and GAFMAT methods for triplet preparation



# Results

Metrics	Non-Image to Image Transformation Methods				
	GADF	GASF	RP	MTF	GAFMAT
Accuracy↑	0.99077	0.98473	0.98331	0.94744	0.98935
EER↓	0.04794	0.05540	0.05327	0.12074	0.04545
AUC↑	0.98612	0.98290	0.98394	0.94862	0.98668
AP_ED↓	0.44127	0.47255	0.43633	0.56487	0.48600
AN_ED↑	1.72784	1.71689	1.68884	1.59469	1.76378
AP_STD↓	0.27487	0.29295	0.28245	0.36906	0.31383
AN_STD↓	0.32888	0.34455	0.34881	0.40005	0.31295
AN_CS↓	0.45772	0.45264	0.46871	0.46011	0.43755
AP_CS↑	0.77936	0.76373	0.78183	0.71756	0.75700

**Table 1** Results of image transformation methods on keystroke dynamics data from the CMU dataset using GADF, GASF, RP, MTF, and GAFMAT algorithms: Metrics-based evaluation on validation data.

Metrics	Non-Image to Image Transformation Methods				
	GADF	GASF	RP	MTF	GAFMAT
Accuracy↑	0.99077	0.98473	0.98331	0.94744	0.98935
EER↓	0.04794	0.05540	0.05327	0.12074	0.04545
AUC↑	0.98612	0.98290	0.98394	0.94862	0.98668
AP_ED↓	0.44127	0.47255	0.43633	0.56487	0.48600
AN_ED↑	1.72784	1.71689	1.68884	1.59469	1.76378
AP_STD↓	0.27487	0.29295	0.28245	0.36906	0.31383
AN_STD↓	0.32888	0.34455	0.34881	0.40005	0.31295
AN_CS↓	0.45772	0.45264	0.46871	0.46011	0.43755
AP_CS↑	0.77936	0.76373	0.78183	0.71756	0.75700

**Table 2** Results of image transformation methods on keystroke dynamics data from the CMU dataset using GADF, GASF, RP, MTF, and GAFMAT algorithms: Metrics-based evaluation on test data.



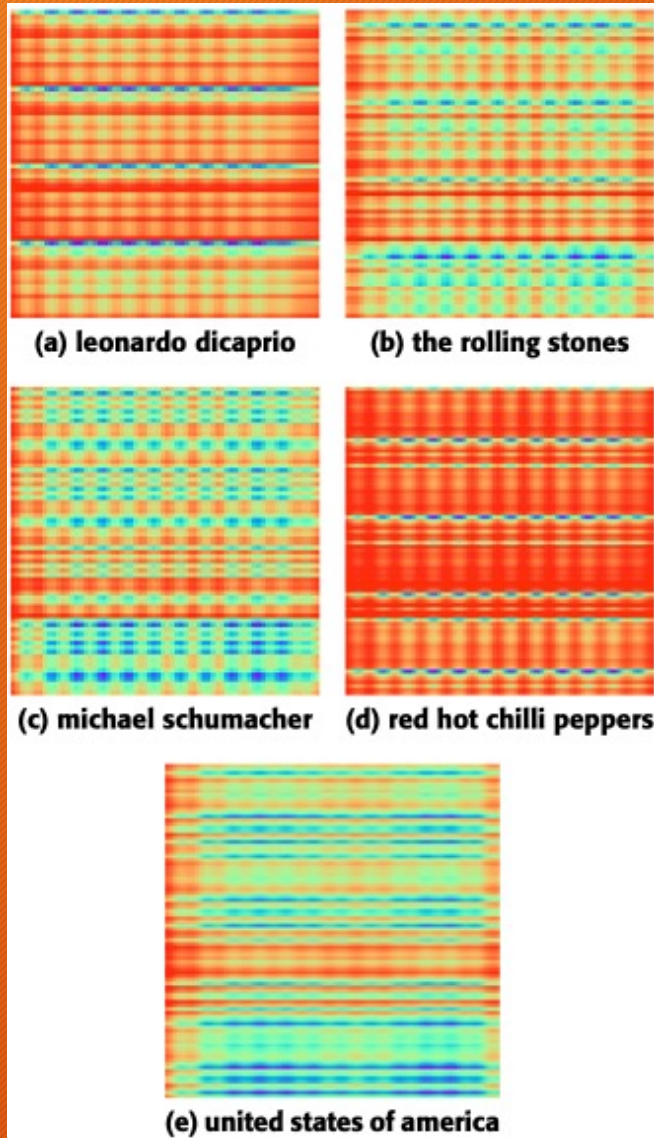
# Results

<b>Authors</b>	<b>Method</b>	<b>EER</b>
This Paper	GAFMAT	0.04545
	GASF	0.0554
	GADF	0.04794
	RP	0.05327
	MTF	0.12074
Killourhy (original) [11]	Manhattan (scaled)	0.096
Zhong et al. [8]	Nearest Neighbor (new distance metric) + outlier removal	0.084
Zhong et al. [8]	Nearest Neighbor (new distance metric)	0.087
Monaco et al. [36]	Inductive transfer encoder (Manhattan distance)	0.063
Hayreddin et al. [35]	Convolutional Neural Network	0.065
Ivannikova et al. [53]	Dependence Clustering with Manhattan	0.077
Sae-Bae et al. [54]	Manhattan (scaled with standard deviation)	0.0916

**Table 3** Performance evaluation for CMU dataset passwords on validation data: a comparison of results



# GAFMAT Method Validation Using the GREYC Dataset



- In the initial phase of the experiments section, using the <sup>18</sup>CMU dataset, we empirically established that our proposed method, called GAFMAT, provides the lowest EER value.
- The GREYC-NISLAB dataset was collected in 2013 and includes five passwords entered by 110 users. The passwords are as follows:
  - 1. "leonardo dicaprio"
  - 2. "the rolling stones"
  - 3. "michael schumacher"
  - 4. "red hot chilli peppers"
  - 5. "united states of america"
- The dataset of a single password consists of 2200 samples. In total, the dataset consists of 11000 data samples corresponding to 110 users, 20 samples per user.



# Results: GAFMAT method validation

Metrics	Passwords (GREYC-NISLAB)				
	leonardo dicaprio	the rolling stones	michaell schu-macher	red hot chilli peppers	united states of america
<b>Accuracy</b> ↑	0.97656	0.98698	0.99219	0.97778	0.99220
<b>EER</b> ↓	0.07552	0.04688	0.0651	0.04444	0.04688
<b>AUC</b> ↑	0.97824	0.98667	0.98771	0.98272	0.98847
<b>AP_ED</b> ↓	0.44736	0.43986	0.39958	0.45165	0.39566
<b>AN_ED</b> ↑	1.55644	1.61202	1.48864	1.63478	1.61275
<b>AP_STD</b> ↓	0.24318	0.21992	0.20467	0.21505	0.19676
<b>AN_STD</b> ↓	0.40601	0.37381	0.38351	0.38917	0.38013
<b>AN_CS</b> ↓	0.49905	0.48703	0.52795	0.47839	0.49790
<b>AP_CS</b> ↑	0.77632	0.78007	0.80021	0.77417	0.80217

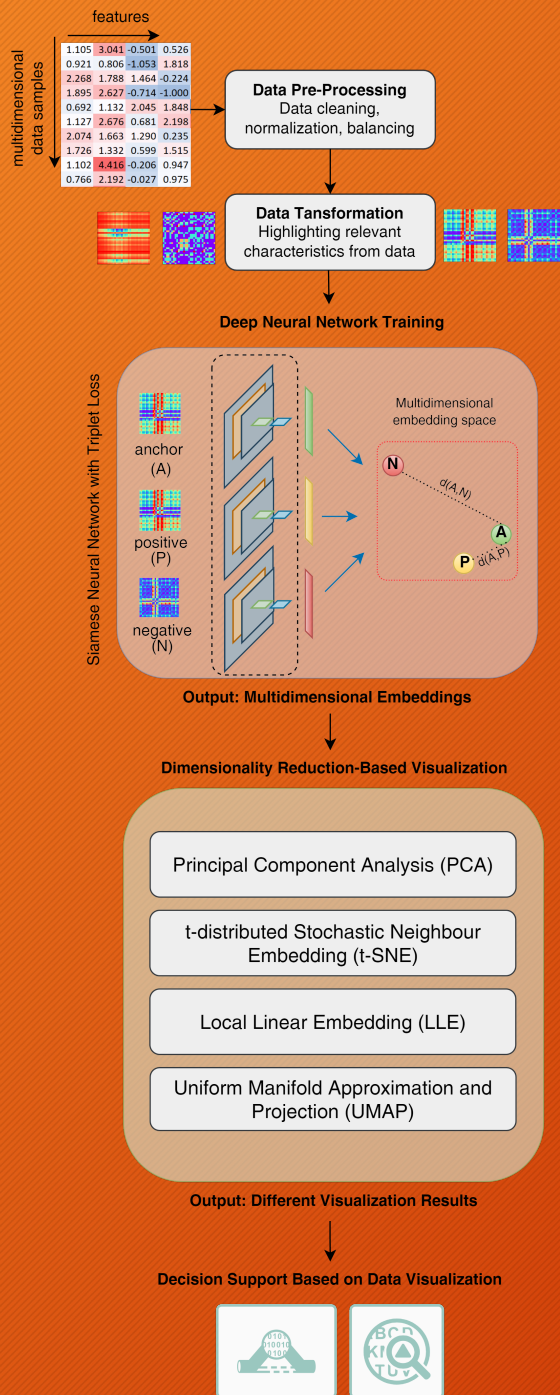
**Table 4** Results using different accuracy metrics for passwords from GREYC-NISLAB on validation dataset when transforming time series features of keystroke dynamics into an image using the GAFMAT algorithm.

Metrics	Passwords (GREYC-NISLAB)				
	leonardo dicaprio	the rolling stones	michaell schu-macher	red hot chilli peppers	united states of america
<b>Accuracy</b> ↑	0.84000	0.86000	0.86000	0.84000	0.92000
<b>EER</b> ↓	0.16000	0.20000	0.22000	0.22000	0.14000
<b>AUC</b> ↑	0.90320	0.85920	0.85400	0.86680	0.89240
<b>AP_ED</b> ↓	0.78894	0.86642	0.67407	0.87670	0.75085
<b>AN_ED</b> ↑	1.55808	1.49985	1.33055	1.55131	1.50073
<b>AP_STD</b> ↓	0.41371	0.40861	0.31141	0.44201	0.43587
<b>AN_STD</b> ↓	0.40956	0.41111	0.49554	0.40963	0.42794
<b>AN_CS</b> ↓	0.41324	0.40843	0.49884	0.39300	0.43711
<b>AP_CS</b> ↑	0.60553	0.56679	0.66297	0.56165	0.62458

**Table 5** Results using different accuracy metrics for passwords from GREYC-NISLAB on a test dataset when transforming time series features of keystroke dynamics into an image using the GAFMAT algorithm.



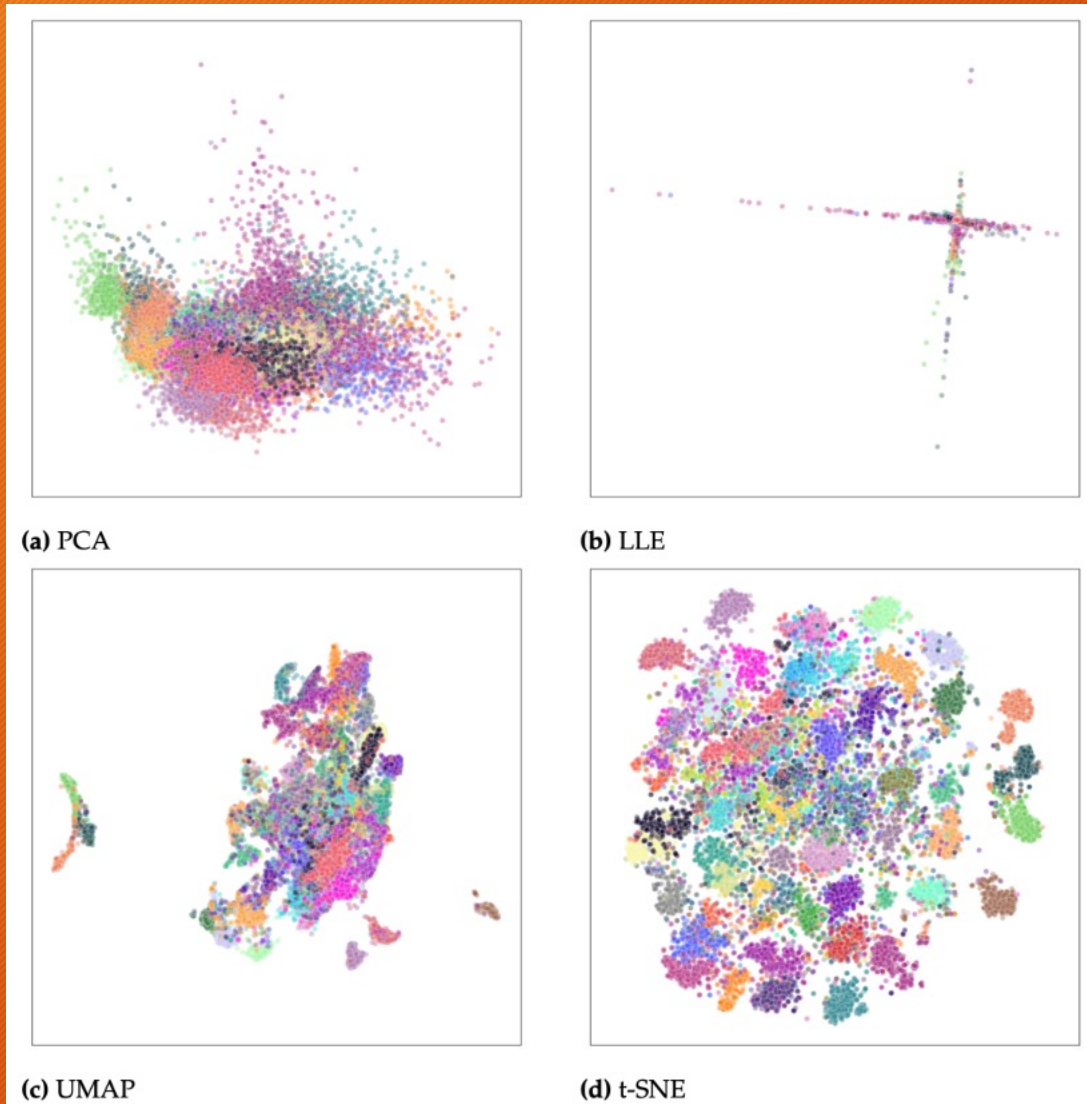
# Exploring Multidimensional Embeddings for Decision Support Using Advanced Visualization Techniques



**Fig 4.** The visualization framework based on dimensionality reduction for multidimensional embedding analysis in decision support



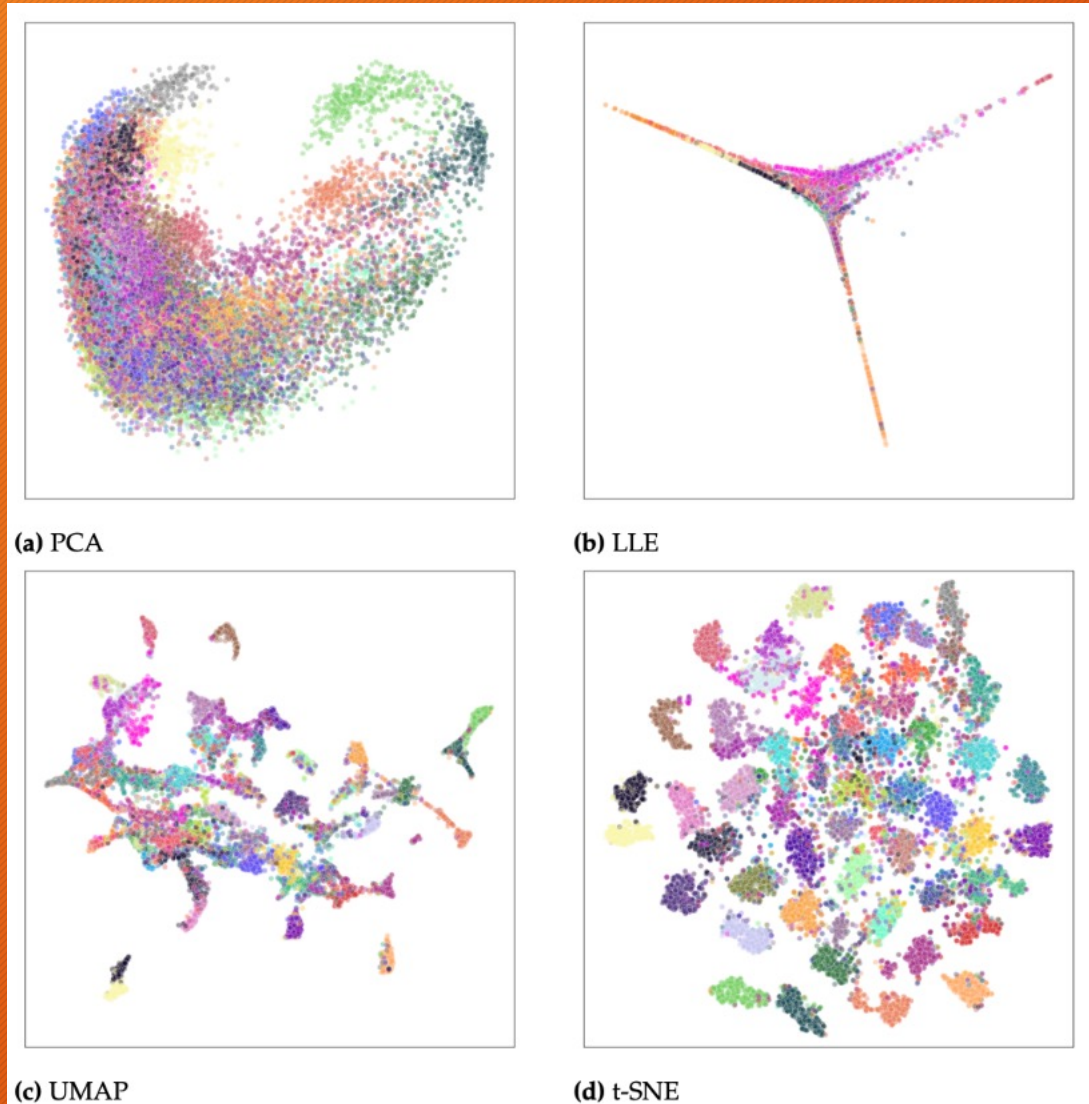
# Exploring Multidimensional Embeddings for Decision Support Using Advanced Visualization Techniques



**Fig 5.** Multidimensional data visualizations by using different dimensionality reduction techniques: (a) PCA, (b) LLE, (c) UMAP, (d) t-SNE. Each color corresponds to a different user in the CMU dataset



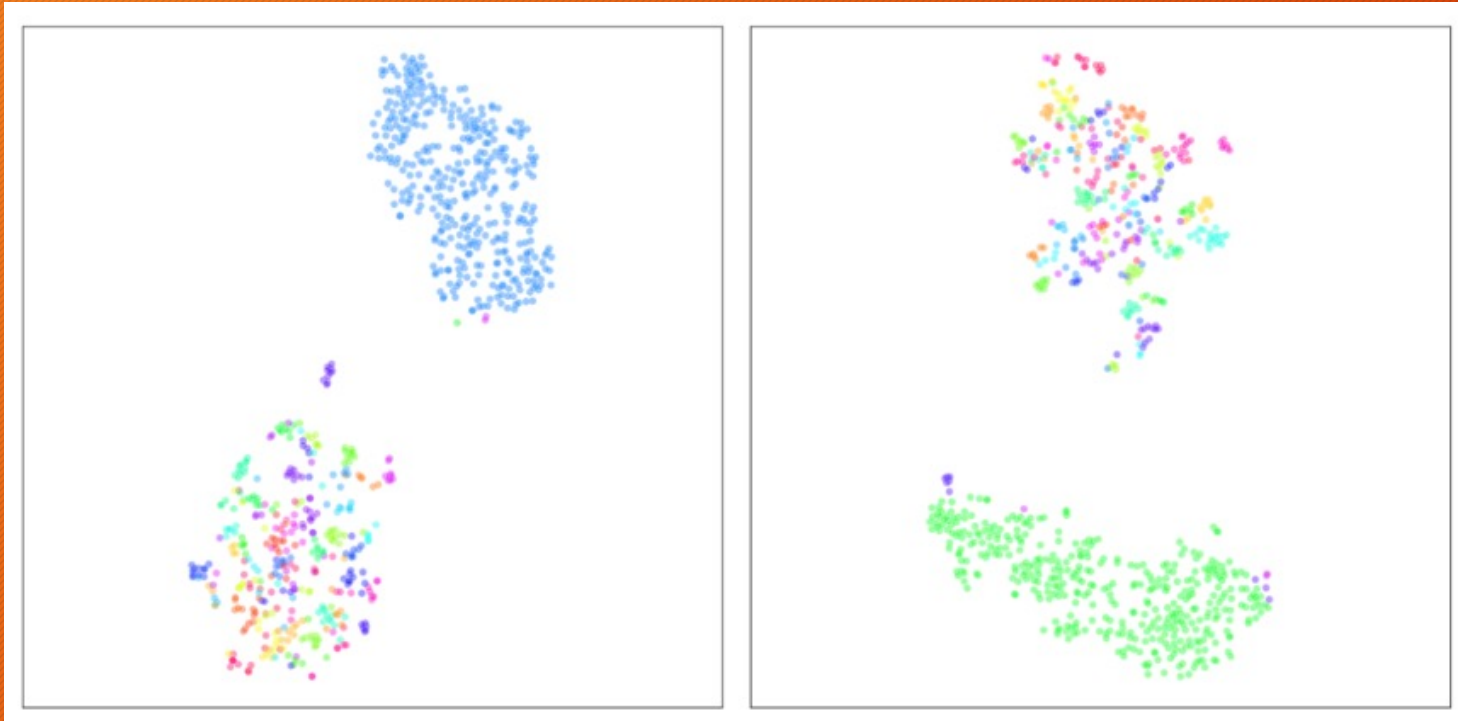
# Exploring Multidimensional Embeddings for Decision Support Using Advanced Visualization Techniques



**Fig 6.** Visualization of multidimensional embeddings obtained by Siamese neural network using different dimensionality reduction techniques ( $p = 256$ ): (a) PCA, (b) LLE, (c) UMAP, (d) t-SNE. Each color corresponds to a different user in the CMU dataset



# Exploring Multidimensional Embeddings for Decision Support Using Advanced Visualization Techniques



**Fig. 7** Examples of visualizations that show password typing patterns of the same user and the other randomly selected users



# Kito pusmečio darbo planas

24

- Publikacija mokslo leidinyje, turinčiame cituojamumo rodiklį Clarivate Analytics Web of Science duomenų bazėje (planuojama įteikti 2024 m. balandžio mėn., COMPUTER STANDARDS & INTERFACES);
- Tikslų, uždavinių, tyrimo metodikos, ginamųjų teiginių patikslinimas;
- Analitinės disertacijos dalies patikslinimas;
- Teorinės disertacijos dalies patikslinimas;
- Eksperimentinės disertacijos dalies patikslinimas;
- Bendrųjų išvadų formulavimas;
- Daktaro disertacijos parengimas.



Jeigu klausimas skamba kaip mįslė ir jūs neieškote sprendimo, nesakykite, kad atsakymas per sunkus. Ne mįslė kalta, kad nesuprantate, o jūsų neatidumas.

chatGPT

[arnoldas.budzys@mif.stud.vu.lt](mailto:arnoldas.budzys@mif.stud.vu.lt)