

A Key-Driven Framework for Identity-Preserving Face Anonymization

Miaomiao Wang, Guang Hua, Sheng Li*, Guorui Feng*

Shanghai University, Singapore Institute of Technology, Fudan University



#NDSSSymposium2025

Mask-Based

Background

DeepLearning-Based





Mosaic



High level of privacy protection •

Cannot be used for visual tasks

How to balance privacy and availability?









Original

DelD

DelDGAN (2021)

- High level of privacy preserving
- Identity is permanently lost •



Background



To balance privacy and recognizability

Identifiable virtual faces(2022)





Identity Hider(2024)

- Authentication in Virtual Domain
- The original face is not synchronized

- Facial feature synchronization
- Identity is easily leaked

> Background



To balance privacy and recognizability

Identifiable virtual faces(2022)



How to resolve the conflict between facial privacy and recognizability, while maintaining head posture and facial expressions?

> Research Objective



Overall Objective



> Research Objective



Overall Objective



Virtual faces anonymize human eyes and machines

Anonymity, synchronism, diversity, differentiation, interactivity, realism, and recognizability

Prevent unauthorized access

A key-driven face anonymization and authentication recognition (KFAAR) framework





Head posture-preserving virtual face generation (HPVFG) module

Generate virtual faces that maintain head posture and facial expressions

The key-controllable virtual face authentication (KVFA) module

Identity authentication using the correct key

Illustration of the proposed framework



Properties of virtual face

Anonymity: The virtual face has a different identity from the original face.

 $R\left(G\left(x_1, k_1\right)\right) \neq R\left(x_1\right)$

The virtual face The original face

R: Face recognizer G: Virtual face generator

Synchronism: Virtual faces generated from the same original identity belong to the same identity. $R(G(x_1,k_1)) = R(G(x_2,k_1))$ $x_1, x_2:$ Faces with the same identity collected at different times



Diversity: The virtual identities generated by the same original face but using different keys should be different.

 $R(G(x_1, k_1)) \neq R(G(x_1, k_2))$

Differentiation: For two original faces x and y which are from different identities, the generated virtual identities should be different even if they are derived from the same key.

 $R\left(G\left(x,k_{1}\right)\right)\neq R\left(G\left(y,k_{1}\right)\right)$

x, y: Faces are from different identities



HPVFG



Network Architecture

Training Strategy

Anonymity Loss:

 $L_{ano} = L_{cos} (R (G (x_1, k_1)), R (x_1), -1)$

Synchronism Loss:

 $L_{syn} = L_{cos} \left(R \left(G \left(x_1, k_1 \right) \right), R \left(G \left(x_2, k_1 \right) \right), 1 \right)$

Diversity Loss:

 $L_{div} = L_{\cos} (R (G (x_1, k_1)), R (G (x_1, k_2)), -1)$

Differentiation Loss:

 $L_{dif} = L_{\cos} (R(G(x, k_1)), R(G(y, k_1)), -1)$



Properties of KVFA



Authenticating virtual faces from different identities with the correct key

 $I(G(x, k_1), k_1) \neq I(G(y, k_1), k_1)$

x, y: Faces are from different identities





Properties of KVFA

Prevent False Rejection

When using the correct key for authentication

I(x) = I(G(x,k),k)

When using the correct key for authentication, two different virtual faces from the same original identity

 $I(G(x_1, k_1), k_1) = I(G(x_1, k_2), k_2)$ $I(G(x_1, k_1), k_1) = I(G(x_2, k_2), k_2)$



KVFA





Prevent Misidentification Loss:

 $L_{pmis1} = L_{cos} \left(I \left(G \left(x_1, k_1 \right) \right), I \left(x_1 \right), -1 \right)$

$$L_{pmis2} = L_{cos} \left(I \left(G \left(x_1, k_1 \right), k_2 \right), I \left(x_1 \right), -1 \right)$$

 $L_{pmis3} = L_{cos} \left(I \left(G \left(x, k_1 \right), k_1 \right), I \left(G \left(y, k_1 \right), k_1 \right), -1 \right)$

> Prevent False Rejection Loss:

$$\begin{split} L_{per1} &= L_{\cos} \left(I \left(G \left(x_1, k_1 \right), k_1 \right), I \left(x_1 \right), 1 \right) \\ L_{per2} &= L_{\cos} \left(I \left(G \left(x_1, k_1 \right), k_1 \right), I \left(G \left(x_2, k_1 \right), k_1 \right), 1 \right) \end{split}$$





Evaluation of Virtual Faces

Anonymity and Diversity

Methods	Anon	ymity ↑	Diversity↑		
	LFW	CelebA	LFW	CelebA	
CIAGAN [22]	0.674	0.0628	0.000	0.000	
IVFG [36]	0.988	0.889	0.750	0.787	
VFGM [33]	0.524	0.561	0.550	0.532	
Ours	0.962	0.922	0.783	0.728	

Head Posture and Facial Expression

Methods	Ya	w 🕹	Pite	ch ↓	Ro	oll ↓	Emo	otion ↑
	LFW	CelebA	LFW	CelebA	LFW	CelebA	LFW	CelebA
CIAGAN [22]	3.913	3.152	4.387	4.424	3.037	3.291	0.584	0.622
VFGM [33]	3.557	3.623	2.519	2.333	4.011	3.997	0.726	0.689
IVFG [36]	25.663	28.537	19.334	18.261	17.991	17.309	0.433	0.412
Ours	2.018	1.992	3.102	2.119	1.998	2.006	0.805	0.833

Synchronism and Detection Rate

Methods	AUC ↑		EER \downarrow		Detection Rate ↑	
	LFW	CelebA	LFW	CelebA	LFW	CelebA
CIAGAN [22]	-	<u>.</u>	-	-	0.986	0.998
IVFG [36]	0.929	0.933	0.103	0.122	1.0	1.0
VFGM [33]	0.889	0.920	0.182	0.139	1.0	1.0
Ours	0.973	0.992	0.092	0.089	1.0	1.0

Visual Quality

Methods	LFW	CelebA
CIAGAN [22]	7.64	6.90
VFGM [33]	6.99	6.91
IVFG [36]	6.19	6.78
Ours	7.29	6.82

Superior to existing SOTA methods in multiple metrics





Examples of the virtual faces.







Examples of the virtual faces.

Original Faces





Virtual Faces





æ









































Evaluation of KVFA

Recognition Accuracy

Dataset	CRR↑	FAR↓	AUC ↑
LFW	0.927	0.078	0.984
CelebA	0.956	0.064	0.989

Different Authentication Scenarios



"In-the-wild" experiment

Datasets		Perform of the virtu	Reco accuracy	gnition of KVFA		
	EER	Anonymity	Diversity	FID	CRR	FAR
LFW FFHQ	0.092 0.089	0.962 0.931	0.783 0.707	7.29 6.36	0.927 0.887	0.078 0.081

- Accurately authenticate the identity of a virtual face with the correct key
- Performed well on additional datasets





Ablation Studies

Head Posture Correction Module



Able to meet synchronization requirements

Each Loss of HPVFG

Methods	AUC	EER	Detection	Anonymity	Diversity	FID
w/o Lano	0.914	0.091	1.0	0.908	0.779	6.86
w/o Laun	0.897	0.104	1.0	0.974	0.787	7.73
w/o Ldiv	0.942	0.091	1.0	0.962	0	7.91
w/o Ldif	0.917	0.095	1.0	0.953	0.790	8.04
Ours	0.973	0.092	1.0	0.962	0.783	7.29

The performance is the best when all losses participate in training

Each Loss of KVFA

Methods	CRR	FAR	AUC
w/o Ltot1	0.756	0.298	0.833
w/o Ltot2	0.636	0.270	0.744
Ours	0.945	0.068	0.987

The performance is the best when all losses participate in training



0.97

0.96

≥ 0.95

0.94

0.93

0.92

0.91

8.0

7.8

7.6

₽ 7.4

0.2

0.4 0.6

Weight of the loss function

0.8

1.0

λano

0.8

Adiv

1.0

0.98

0.97 0.96

0.95

0.93

0.92

0.91

0.90

0.8

0.7

0.6

0.3

0.2

0.2

0.4

0.6

Weight of the loss function

0.8

All 0.5

, P.0 Dive

0.2

0.4 0.6

Weight of the loss function

0.94







Weight of the loss function

(b)

- The best performance is achieved when ٠ the threshold is set to 0.7
- The performance of KVFA loss function ٠ is best when the weight is all set to 1

> Experiment

Ablation Studies

1.0

0.8

0.6

0.4

0.2

0.0

Threshold

(a)

CRR

Threshold Setting for KVFA

Performance Analysis of Key

Key Generation

> Experiment

Use the secure random number generator KeyGen to generate the key on the client side, ensuring that the key is unique and not easily obtained by adversaries

Key Storage

The key is stored only on the client side, and FAS deletes the key after generating the virtual face

Fault tolerance of key

key length		1	Key error	bit	
	0 bit	1 bit	3 bits	5 bits	16 bits
8 bits	0.811	0.708	0.622	0.492	120
128 bits	0.803	0.691	0.640	0.523	0.522
256 bits	0.862	0.702	0.603	0.595	0.485

The impact of key length

key length	Anonymization	FID
8 bits	0.962	7.29
128 bits	0.956	6.95
256 bits	0.964	7.02



Conclusion



Performance Analysis of Key

Contributions

- Propose HPVFG to generate anonymous, synchronized, diverse, and high-quality virtual faces.
- Develop the KVFA module for authenticating virtual faces' original identities with the correct key.
- Successfully harmonize privacy with recognizability in virtual faces.

Future direction

- Enhance security by fortifying protection mechanisms, researching detection and prevention strategies against adversary model training, and bolstering system security.
- Streamline key strategies to fix random key issues and stabilize privacy.