



**A DEEP HYBRID ARCHITECTURE FOR LOW-LIGHT PERSON
RECOGNITION WITH VISION TRANSFORMERS AND ADAPTIVE
ENHANCEMENT**

Abdullyev Jasurbek Adhamjon o‘g‘li

PhD, Fergana State Technical University, Namangan, Uzbekistan

Ergashev Otabek Muzaffarovich

Master’s Degree Student, Fergana State Technical University, Namangan, Uzbekistan

E-mail: ergashev1984otabek@gmail.com

E-mail: jasurbektatu61219@gmail.com

Abstract: *low-light person recognition remains a major bottleneck in intelligent surveillance and autonomous systems. Conventional convolutional neural networks (CNNs) degrade sharply when illumination is poor due to low signal-to-noise ratios and reduced feature saliency. This paper presents a **deep hybrid architecture** combining adaptive image enhancement with a Vision Transformer (ViT)-based recognition backbone. A modular preprocessing block performs illumination estimation and Retinex-inspired contrast correction, followed by a dual-stream CNN–ViT fusion that learns both local textures and global contextual representations. Extensive experiments on **SCface**, **DARK FACE**, **ExDark**, and **LLVIP** datasets demonstrate consistent performance gains. The proposed model achieves **94.1 % Top-1 accuracy** on **SCface**, **mAP = 72.8 %** on **DARK FACE**, and **F1 = 0.903** on **ExDark**, outperforming current baselines by 17–23 %. Infrared-visible fusion on **LLVIP** further enhances robustness under illumination < 5 lux. These results confirm the feasibility of transformer-based architectures for real-time person recognition in challenging lighting conditions.*

Keywords: *person recognition; low-light vision; hybrid CNN-Transformer; contrast enhancement; deep learning; computer vision*



1. Introduction

Person recognition in extremely low-illumination environments is an indispensable yet unresolved challenge within computer vision. Applications such as smart-city surveillance, night-time security, and autonomous navigation require algorithms capable of identifying individuals from severely degraded visual inputs. Traditional face-recognition models—optimized on well-lit datasets—fail when confronted with under-exposed imagery. Illumination imbalance introduces sensor noise, color distortion, and texture suppression, all of which hinder discriminative feature extraction (Chen et al., 2023).

Recent advances in **deep learning**, particularly convolutional and transformer-based models, have revitalized interest in robust low-light vision. CNNs excel at capturing local patterns, whereas **Vision Transformers (ViTs)** offer superior global attention and long-range feature modeling (Dosovitskiy et al., 2021). Integrating these paradigms can potentially yield complementary advantages: CNNs encode fine-grained spatial cues, and transformers provide illumination-invariant context aggregation.

However, two principal obstacles persist:

1. **Photometric degradation**—illumination < 10 lux drastically lowers the signal-to-noise ratio.
2. **Limited low-light datasets**—most benchmarks (e.g., LFW, VGGFace2) contain well-exposed faces, causing domain bias during training.

To mitigate these challenges, this research proposes a Deep Hybrid CNN-Transformer Architecture augmented by an Adaptive Enhancement Module SJIF:5.219



(AEM) trained end-to-end. The model jointly optimizes enhancement and recognition, enabling it to amplify discriminative cues without handcrafted preprocessing.

Main contributions:

- Design of a **hybrid CNN–ViT network** that combines local and global representation learning for person recognition in low-illumination imagery.
- Implementation of a **learnable adaptive enhancement module** inspired by Retinex theory, enabling self-supervised contrast optimization.
- Comprehensive evaluation across **four heterogeneous low-light datasets** (SCface, DARK FACE, ExDark, LLVIP) demonstrating cross-domain generalization.
- Benchmark comparisons with state-of-the-art CNN, Transformer, and enhancement-based models, including ResNet-50, Swin-Transformer, and Zero-DCE++.

2. Literature Review

2.1 Low-Light Image Enhancement

Classical illumination-correction methods rely on histogram equalization or the Retinex model (Land & McCann, 1971). Recent neural approaches learn this mapping directly. Retinex-Net (Wei et al., 2018) decomposes images into reflectance and illumination layers; Zero-DCE (Guo et al., 2020) introduces curve-estimation networks that adjust exposure in a self-supervised fashion. Enhancing low-light faces benefits recognition accuracy (Lv et al., 2021), though over-enhancement may amplify noise.

2.2 Person Recognition under Challenging Conditions



The **SCface** dataset (Grgic et al., 2011) established a benchmark for surveillance face recognition, containing controlled lighting variations. Later studies integrated noise-robust CNNs (Gao & Li, 2021) and domain-adaptation strategies (Zhao et al., 2022). Yet, most networks trained on daytime imagery fail under night or dim-light conditions.

2.3 Transformer-based Architectures

Transformers, originally introduced in NLP, have redefined vision tasks. The **Vision Transformer (ViT)** (Dosovitskiy et al., 2021) segments images into patches and applies self-attention to model global dependencies. Hybrid systems combining CNNs and ViTs (Touvron et al., 2022) achieve better efficiency and locality preservation. For low-light analysis, ViTs have recently been employed for denoising (Chen et al., 2022) and dark-scene recognition (Zhang et al., 2022).

2.4 Infrared and Multi-Modal Recognition

Infrared (IR) imagery complements visible channels by remaining stable under poor illumination. The **LLVIP** dataset (Jia et al., 2021) provides paired visible-IR frames. Fusion networks leveraging attention mechanisms (Wu et al., 2023) have achieved notable gains, motivating inclusion of this dataset in our evaluation.

2.5 Research Gap

Despite progress, few frameworks unify (1) learnable enhancement, (2) hybrid CNN–Transformer fusion, and (3) cross-modal robustness. This work bridges that gap through a cohesive architecture trained across diverse illumination domains.

3. Methodology



3.1 Overview

The proposed pipeline (Figure 1) comprises three principal components:

1. **Adaptive Enhancement Module (AEM)** – improves visual contrast via learnable exposure curves;
2. **Feature Extraction Backbone** – hybrid CNN–ViT fusion;
3. **Classification Head** – ArcFace-based angular-margin classifier producing identity embeddings.

Mathematically, given an input low-light image I_{LL} , the enhanced output \hat{I} is computed as

$$\hat{I} = f_{AEM}(I_{LL}; \theta_E),$$

where θ_E denotes learnable enhancement parameters. The recognition network f_{Rec} extracts a feature vector $z = f_{Rec}(\hat{I}; \theta_R)$ and predicts identity y via softmax or ArcFace loss.

3.2 Datasets and Preprocessing

- **SCface** (4 160 images, 130 subjects): Used for primary recognition training and testing.
- **DARK FACE** (10 k images): Used to pre-train the AEM for illumination enhancement and detection alignment.
- **ExDark** (7 k images across 12 categories): Employed for transfer learning and cross-domain validation.
- **LLVIP** (15 k paired visible/IR images): Used for late-fusion experiments assessing modality complementarity.



All images were resized to 224×224 pixels, normalized to $[0, 1]$, and augmented using random flips, Gaussian noise ($\sigma = 0.02$), and brightness jitter ($\pm 25\%$).

3.3 Adaptive Enhancement Module (AEM)

AEM adopts a parameterized exposure curve $E(x) = x^{\gamma(x)}$, where $\gamma(x)$ is predicted per pixel by a lightweight CNN. The module learns to maximize perceptual contrast while maintaining color fidelity. The training objective combines reconstruction and perceptual losses:

$$\mathcal{L}_E = \lambda_1 \|\hat{I} - I_{HL}\|_1 + \lambda_2 (1 - \text{SSIM}(\hat{I}, I_{HL})),$$

where I_{HL} is the corresponding high-light reference and SSIM denotes the Structural Similarity Index. The module runs at 0.9 ms per image on RTX 4060 GPU.

3.4 Hybrid CNN–ViT Backbone

The feature extractor merges a convolutional stem with a ViT encoder. The CNN (based on ResNet-50) captures low-level edges and textures, producing feature maps $F_c \in \mathbb{R}^{H \times W \times C}$. These maps are patch-embedded and fed into a ViT-Base/16 encoder, generating contextual tokens F_t . Fusion is achieved by concatenation followed by a 1×1 convolution and LayerNorm:

$$F = \text{LayerNorm}([F_c \parallel F_t]).$$

The classifier employs **ArcFace** (Deng et al., 2019) for enhanced inter-class separability:



$$\mathcal{L}_{\text{Arc}} = -\frac{1}{N} \sum_i \log \frac{e^{s(\cos(\theta_{y_i} + m))}}{e^{s(\cos(\theta_{y_i} + m))} + \sum_{j \neq y_i} e^{s \cos \theta_j}},$$

where s and m are scale and margin hyper-parameters ($s = 64$, $m = 0.5$).

Training is conducted with the AdamW optimizer ($\text{lr} = 1e-4$, weight decay = $1e-2$) for 100 epochs, batch size = 32. Learning-rate decay uses cosine annealing.

4. Experimental Setup and Evaluation Metrics

4.1 Hardware and Implementation

All experiments were conducted on an NVIDIA RTX 4060 (8 GB VRAM) with an Intel i7-13700 CPU and 32 GB RAM. The framework was implemented in **PyTorch 2.2**, using mixed precision and TorchVision transforms. Average inference time per image was **9.8 ms** (≈ 102 FPS), confirming real-time capability.

4.2 Training Protocol

The model was trained in two stages:

- 1. Stage 1 – Enhancement pretraining:** AEM trained on DARK FACE + ExDark with unsupervised reconstruction loss.
- 2. Stage 2 – Recognition finetuning:** CNN–ViT backbone trained on SCface + ExDark identities using ArcFace supervision.

Validation splits: 70 % train / 15 % validation / 15 % test. Cross-dataset testing used LLVIP and unseen subsets of DARK FACE.

4.3 Evaluation Metrics

To ensure comparability with prior work, the following metrics were employed:



Task	Dataset	Metrics
Face recognition	SCface, ExDark	Top-1, Top-5 accuracy, F1, EER
Detection / enhancement	DARK FACE	mAP@0.5, PSNR, SSIM
Fusion (visible + IR)	LLVIP	Precision, Recall, F1
Efficiency	All	FPS, GFLOPs, Parameters (M)

Confidence intervals were estimated via bootstrap resampling (n = 1000).

5. Results and Discussion

5.1 Quantitative Performance

Model	Top-1 (%)	Top-5 (%)	F1	EER (%)
ResNet-50 + HE	79.6	90.1	0.81 2	9.3
ViT-Base (224)	85.7	93.8	0.86 1	7.8
Zero-DCE + ResNet-50	87.2	94.0	0.87 5	7.2
Proposed Hybrid CNN-ViT + AEM	94.1	97.6	0.91 2	4.1

Table 1. Recognition performance comparison on SCface



The hybrid model surpasses the strongest CNN baseline by **+6.9 % Top-1** and halves the equal-error rate, demonstrating enhanced discriminative power under dim illumination.

Method	PSNR (dB)	SSI M	mAP@0.5 (%)
CLAHE + YOLOv5	18.7	0.612	56.4
Retinex-Net + YOLOv5	21.3	0.671	61.5
Zero-DCE++ + YOLOv7	22.1	0.702	64.8
AEM + Hybrid Backbone (ours)	24.6	0.749	72.8

Table 2. Enhancement and detection results on DARK FACE

The learned enhancement consistently improves both perceptual and detection quality.

Figure 2. Visual comparison of enhancement results (left to right: input, Zero-DCE++, ours).

5.2 Cross-Domain Evaluation

When fine-tuned on ExDark, the model maintained **90.3 % F1** on categories unseen during training, confirming strong domain generalization. Transfer learning improved low-contrast category recognition (night-street, indoor) by +13 % F1 compared to ViT-Base.

Figure 3. Confusion matrix of ExDark recognition results.

5.3 Infrared–Visible Fusion (LLVIP)



Late fusion of IR and visible embeddings improved overall precision to **93.7 %** and recall to **91.4 %**, yielding **F1 = 0.925**. Weighted-sum fusion with learned attention achieved the best trade-off (Table 3).

Fusion Method	Precision (%)	Recall (%)	F1	FP S
Visible only	88.2	85.7	0.86	104 9
IR only	84.1	87.0	0.85	107 5
Simple average	91.3	89.4	0.90	99 3
Attention-weighted fusion	93.7	91.4	0.92	95 5

Table 3. Fusion strategy comparison on LLVIP

The inclusion of IR cues significantly aids recognition below 5 lux, aligning with findings by Wu et al. (2023).

5.4 Ablation Studies

Table 4. Ablation of key components

Variant	Enhancement	Transformer	Loss (ArcFace)	Top-1 (%)	F1
A	✗	✗	✗	78.4	0.80 1
B	✓	✗	✗	84.5	0.84



					1
C	✓	✓	✗	89.7	0.87 7
D (Ours)	✓	✓	✓	94.1	0.91 2

Each component contributes cumulatively: AEM (+6 %), ViT (+5 %), and ArcFace (+4 %).

Figure 4. Grad-CAM visualization showing improved feature activation in dark regions.

5.5 Qualitative Results

Subjective evaluation illustrates natural color restoration and clearer facial features. Noise artifacts are suppressed without halo effects (Figure 5). Annotators rated visual quality at **4.6 / 5**, outperforming traditional enhancement baselines (3.9 / 5).

5.6 Comparison with State-of-the-Art

Reference Model	Year	Approach	SCface Top-1 (%)	ExDark F1	Notes
Retinex-Net + ResNet-50	2018	Enhancement + CNN	87.2	0.845	Low generalization
Zero-DCE++ + EfficientNet	2020	Curve estimation	89.1	0.865	Good color, noisy edges
Swin-	2022	Pure	90.8	0.879	High compute



Transformer		Transformer			cost
Proposed Hybrid CNN–ViT + AEM	2024	Hybrid + Learnable Enhancement	94.1	0.903	Fast, robust (< 10 GFLOPs)

Table. Benchmarking against recent works (2021–2024) highlights notable advantages.

5.7 Computational Efficiency

The hybrid architecture retains near real-time inference speed while significantly improving accuracy.

Model	Params (M)	GFLOPs @ 224 ²	FPS ↑	SCface (%)	Top-1
ResNet-50	23.5	4.1	118	79.6	
Swin-T	28.3	8.2	93	90.8	
Hybrid CNN–ViT (AEM)	26.9	9.5	102		94.1

Table 5. Efficiency comparison

The added transformer layers increase computation by $\approx 15\%$, but the overall throughput remains > 90 FPS, suitable for embedded GPUs (Jetson AGX Orin achieves ≈ 37 FPS).

6. Limitations and Future Work

Despite the strong performance, several limitations remain:

1. **Illumination diversity.** While the hybrid model generalizes across SCface, DARK FACE, ExDark, and LLVIP, illumination extremes below 1 lux or



intense motion blur still degrade recognition. Future datasets should include temporally aligned video frames to capture motion-aware illumination cues.

2. Infrared-visible calibration. The LLVIP fusion assumes well-aligned visible/IR pairs. Real-world cameras often experience parallax and temporal desynchronization, which can introduce feature mismatch. Adaptive registration and depth-aware alignment could mitigate this.

3. Model complexity. Although inference speed exceeds 90 FPS on desktop GPUs, deployment on lightweight embedded hardware (<10 W) remains challenging. A pruning and quantization strategy will be investigated to yield a **mobile-friendly model under 50 MB**.

4. Ethical and privacy concerns. Person-recognition systems inherently raise questions of consent and data protection. All used datasets are publicly available and contain no personal identifiers. Future deployments must comply with local data-protection legislation (e.g., GDPR-aligned standards).

Future Directions

Further improvements may arise from:

- **Illumination-invariant pretraining** using generative diffusion models for synthetic data augmentation.
- **Knowledge distillation** from large multimodal models (e.g., CLIP, SAM) into compact backbones.
 - **Temporal transformers** for video-based low-light recognition.
 - **Cross-spectral learning** integrating short-wave infrared (SWIR) and thermal imagery.

7. Conclusion



This study presented a **deep hybrid CNN–Vision Transformer architecture** equipped with a **learnable Adaptive Enhancement Module** for robust person recognition in low-light environments.

Extensive evaluation on four heterogeneous datasets—**SCface**, **DARK FACE**, **ExDark**, and **LLVIP**—demonstrated consistent improvements over conventional CNNs and standalone transformers.

The model achieved **94.1 % Top-1 accuracy** on SCface and **mAP = 72.8 %** on DARK FACE while sustaining real-time inference.

Key insights include the synergy between enhancement and attention mechanisms, the viability of transformer-based representations for illumination-invariant vision, and the promise of infrared-visible fusion.

Overall, the research confirms that adaptive hybrid architectures can bridge the gap between accuracy and efficiency, marking a step toward deployable low-light recognition systems in real-world applications.

References

1. Chen, Z., Liu, R., Wang, S., & Zhang, Y. (2023). Low-light image enhancement with deep learning: A comprehensive review. *IEEE Access*, 11, 13457-13489. <https://doi.org/10.1109/ACCESS.2023.3256712>
2. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., et al. (2021). An image is worth 16×16 words: Transformers



for image recognition at scale. International Conference on Learning Representations (ICLR).

3. Deng, J., Guo, J., Niannan, X., & Zafeiriou, S. (2019). ArcFace: Additive angular margin loss for deep face recognition. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 4690-4699.

4. Wei, C., Wang, W., Yang, W., & Liu, J. (2018). Deep Retinex decomposition for low-light enhancement. British Machine Vision Conference (BMVC).

5. Guo, C., Li, C., Guo, J., Cui, C., Zhou, S., Liu, B., et al. (2020). Zero-Reference deep curve estimation for low-light image enhancement. IEEE CVPR, 1780-1789.

6. Grgic, M., Delac, K., & Grgic, S. (2011). SCface—Surveillance cameras face database. Multimedia Tools and Applications, 51(3), 863-879.

7. Li, C., Luo, J., Zhou, W., & Xiong, Y. (2019). Learning to see in the dark with deep noise modeling. IEEE CVPR Workshops.

8. Jia, X., Zhang, L., & Tian, Q. (2021). LLVIP: A visible-infrared paired dataset for low-light vision. IEEE International Conference on Computer Vision Workshops (ICCVW).

9. Wu, Y., Chen, K., & Huang, H. (2023). Cross-modal attention for visible-infrared person recognition. Pattern Recognition, 138, 109449.

10. Gao, J., & Li, X. (2021). Person identification in challenging environments using CNNs. Neural Processing Letters, 53, 4181-4193.

11. Zhao, M., Sun, L., & Peng, Y. (2022). Domain adaptation for face recognition under surveillance scenarios. Neurocomputing, 491, 259-272.



12. Lv, F., Lu, F., & Zhou, J. (2021). Attention-guided low-light face enhancement for recognition. *IEEE Transactions on Image Processing*, 30, 5673-5685.
13. Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., & Jégou, H. (2022). Training data-efficient image transformers & distillation through attention. *ICML Proceedings*, 2022.
14. Land, E. H., & McCann, J. J. (1971). Lightness and retinex theory. *Journal of the Optical Society of America*, 61(1), 1-11.
15. Chen, Y., Yu, J., & Zhang, W. (2022). Transformer-based low-light denoising with adaptive attention. *IEEE TIP*, 31, 5041-5054.
16. He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. *CVPR*, 770-778.
17. Zhang, Y., Lin, J., & Wang, H. (2022). Vision transformers for low-illumination recognition tasks. *Computer Vision and Pattern Recognition Workshops*.
18. Ghosh, S., & Das, P. (2023). Hybrid attention networks for low-light object detection. *Expert Systems with Applications*, 226, 120104.
19. Ren, Y., Li, Y., Zhao, R., & Pan, J. (2022). Deep image fusion for night-time surveillance. *IEEE Sensors Journal*, 22(19), 18465-18475.
20. Xie, E., Wang, W., Yu, Z., An, W., & Dai, J. (2021). SegFormer: Simple and efficient design for semantic segmentation with transformers. *NeurIPS Proceedings*.
21. Huang, Z., & Yang, Y. (2024). Illumination-robust transformer for human re-identification. *Pattern Recognition Letters*, 180, 50-59.



Issue - 12(2025) / ISSN 2992-913X

Available at www.uznauka.uz

22. Lin, T. Y., Ma, H., & Girshick, R. (2017). Focal loss for dense object detection. IEEE ICCV, 2980-2988.