# A Unique Pipeline for Low Resolution Facial Recognition

Daniel Szurek, Brandon Nguyen

*Abstract*—Low-resolution facial recognition (LRFR) remains a critical challenge in real-world surveillance and security systems where imaging conditions are suboptimal. This paper presents a novel end-to-end pipeline that integrates identity-aware super-resolution with efficient facial recognition to address the very-low-resolution (VLR) face recognition problem. Our approach combines a Deep Super-Resolution Color (DSR) network trained with perceptual and identity-preserving losses, and EdgeFace, a lightweight recognition backbone optimized for embedded deployment. Unlike conventional approaches that treat super-resolution and recognition as independent modules, we introduce a cycle training strategy where DSR learns to preserve features that EdgeFace recognizes best, creating synergistic co-optimization between the two stages. We incorporate multi-scale perceptual loss, feature matching at intermediate network layers, and metric learning fine-tuning to bridge the domain gap between super-resolved and high-resolution embeddings. Evaluated on the CMU Multi-PIE dataset with probe images downsampled to 16×16 pixels, our pipeline achieves competitive recognition accuracy while maintaining real-time inference capability suitable for edge devices. Our contributions include: (1) an identity-aware DSR architecture with feature-matching supervision from fine-tuned recognition networks, (2) a two-stage training methodology coupling super-resolution and recognition optimization, and (3) comprehensive ablation studies quantifying the impact of each loss component on recognition performance.

*Index Terms*—low-resolution face recognition, deep learning

## I. INTRODUCTION

### A. Background and Motivation

Facial recognition systems have become ubiquitous in security, surveillance, and access control applications. However, real-world deployment scenarios frequently encounter significant challenges due to imaging constraints: cameras positioned at large distances capture subjects at low resolution, poor lighting conditions degrade image quality, and cost-effective surveillance infrastructure employs sensors with limited pixel density. These factors collectively create the very-low-resolution facial recognition (VLR-FR) problem, where probe images may be as small as 16×16 or 32×32 pixels—far below the 112×112 input size expected by modern deep face recognition models.

Traditional face recognition pipelines, trained predominantly on high-resolution datasets like MS-Celeb-1M and VGGFace2, exhibit severe performance degradation when presented with VLR inputs. The core issue is twofold: first, critical discriminative features such as fine-grained facial textures, subtle expression patterns, and localized geometric cues are irreversibly lost during severe downsampling; second, the distribution shift between high-resolution training data and low-resolution test data creates a domain gap that deep networks struggle to bridge. Simply upsampling VLR images with bicubic interpolation before recognition yields marginal improvements, as it introduces smoothing artifacts without recovering lost high-frequency details.

Super-resolution (SR) techniques offer a promising pathway to address VLR-FR by reconstructing plausible high-resolution details from low-resolution inputs. However, generic SR models optimized for perceptual quality metrics (PSNR, SSIM) do not necessarily preserve identity-discriminative features required for recognition. A super-resolved face may appear visually pleasing yet fail to maintain the embedding consistency needed for accurate matching against a gallery. This disconnect motivates the need for *identity-aware super-resolution*—SR models explicitly trained to preserve features that recognition networks rely upon.

Our work is motivated by three key observations: (1) SR and recognition stages should be jointly optimized rather than treated as independent modules, (2) lightweight architectures are essential for deployment on edge devices with limited compute budgets, and (3) domain adaptation between SR outputs and recognition model expectations significantly impacts end-to-end accuracy. We address these challenges through a novel pipeline coupling a Deep Super-Resolution Color (DSR) network with EdgeFace, a compact recognition backbone, unified via identity-preserving training objectives and cycle optimization.

### B. Novelty and Contributions

This paper makes the following contributions to the VLR facial recognition problem:

1) **Identity-Aware DSR Architecture:** We extend the DSR super-resolution framework with multi-scale perceptual loss, identity-preserving cosine embedding loss, and feature-matching supervision from intermediate recognition network layers. This ensures SR outputs not only exhibit high visual fidelity but also maintain embedding consistency with high-resolution counterparts.

2) **Cycle Training Strategy:** We propose a two-stage training methodology where EdgeFace is first fine-tuned on DSR outputs using ArcFace metric learning, then DSR is retrained using the fine-tuned EdgeFace as the identity supervisor. This creates a feedback loop enabling co-adaptation between SR and recognition stages.

3) **Feature-Matching Loss:** Beyond final embedding similarity, we introduce intermediate feature matching between DSR→EdgeFace and HR→EdgeFace activation maps at multiple network depths, providing richer supervision for identity preservation.

4) **Lightweight Pipeline Design:** Our complete pipeline (DSR + EdgeFace) operates with 112 base channels and 16 residual blocks for DSR, coupled with EdgeFace-S backbone (512-dimensional embeddings), achieving real-time inference on edge devices while maintaining competitive recognition accuracy.

5) **Comprehensive Ablation Study:** We systematically evaluate the contribution of each loss component (pixel L1, VGG perceptual, identity cosine, feature matching, total variation) and training strategy (original vs. fine-tuned EdgeFace, cycle training) to final recognition performance.

Our pipeline achieves a balance between recognition accuracy and computational efficiency, making it suitable for deployment in resource-constrained environments such as Raspberry Pi 5-class hardware with 4GB memory and limited GPU acceleration.

## II. RELATED WORK

### A. Deep Learning for Super-Resolution

Single image super-resolution (SISR) has progressed rapidly with the advent of deep learning. Early work by Dong et al. [1] introduced SRCNN, a three-layer CNN that learned end-to-end mapping from low to high resolution. Subsequent architectures like VDSR [2] and EDSR [3] increased network depth and removed batch normalization to improve gradient flow. Kim et al.'s DRCN [4] employed recursive layers for parameter efficiency, while Lim et al.'s EDSR [3] removed unnecessary modules to achieve state-of-the-art PSNR.

Perceptual loss, introduced by Johnson et al. [5], shifted focus from pixel-wise MSE to feature-level similarity using pretrained VGG networks. SRGAN [6] combined perceptual loss with adversarial training to generate photo-realistic textures, though at the cost of introducing artifacts that may harm downstream recognition. More recent work like ESRGAN [7] and RealSR [8] refine GAN training for better perceptual quality.

For face super-resolution specifically, FSRNet [9] and SPARNet [10] incorporate facial priors (landmarks, parsing maps) to guide reconstruction. However, these methods optimize for visual quality rather than recognition performance, leading to a disconnect between SR output and recognition model expectations.

### B. Low-Resolution Face Recognition

Low-resolution face recognition has been approached from multiple angles. Early methods like Hennings-Yeomans et al. [11] employed simultaneous SR and recognition in a coupled framework. Biswas et al. [12] proposed multidimensional scaling to learn a common embedding space for HR and LR faces.

With deep learning, Li et al.'s Coupled-GAN [13] jointly trained SR and recognition networks with shared parameters. Chen et al.'s LCNN [14] introduced attention mechanisms to focus on discriminative facial regions during recognition. More recently, He et al.'s DSR framework [15] demonstrated that residual dense blocks with identity loss improve recognition accuracy.

Domain adaptation approaches like DRAN [16] and DADE [17] learn resolution-invariant features through adversarial training or metric learning. However, these methods often require paired HR-LR training data or complex training procedures. Our approach simplifies the pipeline by explicitly training SR to preserve recognition features through identity-aware losses.

### C. Efficient Face Recognition Models

The recognition stage of our pipeline leverages lightweight architectures suitable for edge deployment. MobileFaceNets [18] introduced depthwise separable convolutions for mobile devices. ShuffleFaceNet [19] employed channel shuffle operations to reduce FLOPs while maintaining accuracy.

EdgeFace [20], our chosen backbone, balances efficiency and accuracy through carefully designed bottleneck blocks and lightweight depthwise convolutions (LDC). With only 0.6M parameters, EdgeFace-S achieves competitive performance on LFW and CFP-FP benchmarks while running at real-time speeds on Raspberry Pi-class hardware. Quantized variants (EdgeFace-XXS-Q) further reduce memory footprint through INT8 quantization, though at the cost of fine-tuning flexibility.

### D. Joint Super-Resolution and Recognition

Several works have explored coupling SR and recognition. Yu et al.'s Identity-Aware SR [21] introduced triplet loss to maintain identity consistency during SR. Xu et al.'s Learn-SR [22] employed a recognition network to provide feature-level supervision for SR training.

Our work extends this direction by: (1) incorporating multi-scale perceptual and feature-matching losses at intermediate network depths, (2) introducing a cycle training procedure where recognition and SR networks co-adapt, and (3) demonstrating effectiveness on very-low-resolution inputs (16×16) rather than moderately low resolution (32×32 or higher). Unlike prior work focusing on PSNR/SSIM metrics, we optimize directly for recognition accuracy through identity-preserving objectives.

## III. METHODOLOGY

Our pipeline consists of three key components: (1) the Deep Super-Resolution Color (DSR) network for upsampling VLR inputs to recognition-suitable resolution, (2) the EdgeFace recognition backbone for extracting identity embeddings, and (3) a cycle training strategy that jointly optimizes both stages for maximum recognition accuracy. Figure **??** illustrates the overall architecture.

### A. Deep Super-Resolution Color (DSR) Network

*1) Architecture:* The DSR network builds upon residual learning principles to reconstruct high-frequency facial details from very-low-resolution inputs. Our implementation consists of:

- **Input Layer:** A 3×3 convolutional layer maps RGB input (16×16 or 32×32) to a 112-channel feature space.

- **Residual Blocks:** 16 residual dense blocks with 3×3 convolutions, each containing skip connections to facilitate gradient flow. Each block follows the structure: Conv→ReLU→Conv→Add.
- **Upsampling:** Pixel shuffle layers (Shi et al. [23]) progressively upscale features to target resolution (128×128) without introducing checkerboard artifacts common in transposed convolutions.
- **Output Layer:** A final 3×3 convolution reconstructs RGB channels with Tanh activation clamped to [0,1].

The base channel count of 112 and 16 residual blocks balance reconstruction capacity with inference efficiency. Compared to deeper architectures (e.g., EDSR with 32 blocks), our configuration achieves real-time throughput (¿30 FPS on CUDA-enabled edge devices) while preserving critical facial features.

*2) Training Objectives:* DSR training employs a multi-component loss function designed to optimize both perceptual quality and identity preservation:

$$\mathcal{L}_{\mathrm{DSR}} = \lambda_{\mathrm{L1}}\mathcal{L}_{\mathrm{L1}} + \lambda_{\mathrm{P}}\mathcal{L}_{\mathrm{P}} + \lambda_{\mathrm{ID}}\mathcal{L}_{\mathrm{ID}} + \lambda_{\mathrm{FM}}\mathcal{L}_{\mathrm{FM}} + \lambda_{\mathrm{TV}}\mathcal{L}_{\mathrm{TV}} \quad (1)$$

**Pixel Reconstruction Loss ($\mathcal{L}_{\mathbf{L1}}$):** Standard L1 distance between super-resolved output $I_{\mathrm{SR}}$ and ground-truth high-resolution image $I_{\mathrm{HR}}$:

$$\mathcal{L}_{\mathrm{L1}} = \|I_{\mathrm{SR}} - I_{\mathrm{HR}}\|_1 \quad (2)$$

**Multi-Scale Perceptual Loss ($\mathcal{L}_{\mathbf{P}}$):** We extract features from four layers of a pretrained VGG-19 network (relu1_2, relu2_2, relu3_4, relu4_4) and compute weighted L1 distance:

$$\mathcal{L}_{\mathrm{P}} = \sum_{l=1}^{4} w_l \|\phi_l(I_{\mathrm{SR}}) - \phi_l(I_{\mathrm{HR}})\|_1 \quad (3)$$

where $\phi_l$ denotes features from layer $l$ and weights $w = [0.4, 0.3, 0.2, 0.1]$ prioritize early layers capturing facial structure.

**Identity Preservation Loss ($\mathcal{L}_{\mathbf{ID}}$):** To ensure SR outputs maintain identity consistency, we employ cosine embedding loss on EdgeFace embeddings:

$$\mathcal{L}_{\mathrm{ID}} = 1 - \frac{\mathbf{e}_{\mathrm{SR}} \cdot \mathbf{e}_{\mathrm{HR}}}{\|\mathbf{e}_{\mathrm{SR}}\|\|\mathbf{e}_{\mathrm{HR}}\|} \quad (4)$$

where $\mathbf{e}_{\mathrm{SR}} = \mathrm{EdgeFace}(I_{\mathrm{SR}})$ and $\mathbf{e}_{\mathrm{HR}} = \mathrm{EdgeFace}(I_{\mathrm{HR}})$. This loss is the strongest signal in our objective ($\lambda_{\mathrm{ID}} = 0.50$).

**Feature Matching Loss ($\mathcal{L}_{\mathbf{FM}}$):** Beyond final embeddings, we match intermediate EdgeFace activations at multiple depths:

$$\mathcal{L}_{\mathrm{FM}} = \frac{1}{K} \sum_{k=1}^{K} \|\mathbf{f}_k(I_{\mathrm{SR}}) - \mathbf{f}_k(I_{\mathrm{HR}})\|_1 \quad (5)$$

where $\mathbf{f}_k$ denotes features from the $k$-th intermediate layer (early, mid, late stages of EdgeFace). This provides richer supervision than embedding-only loss.

**Total Variation Loss ($\mathcal{L}_{\mathbf{TV}}$):** To encourage spatial smoothness while preserving edges:

$$\mathcal{L}_{\mathrm{TV}} = \sum_{i,j}(|I_{\mathrm{SR}}(i+1,j)-I_{\mathrm{SR}}(i,j)|+|I_{\mathrm{SR}}(i,j+1)-I_{\mathrm{SR}}(i,j)|) \quad (6)$$

Loss weights are set to $\lambda_{\mathrm{L1}} = 1.0$, $\lambda_{\mathrm{P}} = 0.02$, $\lambda_{\mathrm{ID}} = 0.50$, $\lambda_{\mathrm{FM}} = 0.15$, $\lambda_{\mathrm{TV}} = 3 \times 10^{-6}$ based on validation set performance.

*3) Training Configuration:* DSR is trained for 100 epochs with batch size 16 using AdamW optimizer (learning rate $1.5 \times 10^{-4}$, weight decay $10^{-6}$). We employ a 5-epoch linear warmup followed by cosine annealing. Gradient clipping (max norm 1.0) prevents instability, and exponential moving average (EMA, decay 0.999) smooths weight updates. Mixed-precision training (FP16) accelerates convergence on CUDA devices. Early stopping with patience 20 epochs halts training when validation PSNR plateaus.

Data augmentation includes: horizontal flipping (50% probability), small rotation (±5°, 60% probability), and mild color jitter (brightness/contrast ±5%, saturation ±3%, 25% probability). Aggressive augmentation is avoided to prevent degrading identity features.

### B. EdgeFace Recognition Network

*1) Architecture:* EdgeFace employs a lightweight bottleneck design optimized for mobile deployment:

- **Stem:** Initial 3×3 convolution with stride 2 maps 112×112 RGB input to 32 channels.
- **Stages:** Four stages progressively downsample features (56×56→28×28→14×14→7×7) while increasing channels (32→64→128→256→512). Each stage contains stacked Lightweight Depthwise Convolution (LDC) blocks with residual connections.
- **Embedding Head:** Global average pooling followed by fully connected layer projects 512×7×7 features to 512-dimensional embedding space. Final batch normalization ensures unit variance.

The LDC block structure is:

$$\mathrm{LDC}(x) = x + \mathrm{Conv}_{1\times1}(\mathrm{DWConv}_{3\times3}(\mathrm{Conv}_{1\times1}(x))) \quad (7)$$

where DWConv denotes depthwise separable convolution. This reduces FLOPs by 8-9× compared to standard convolutions while preserving expressive power.

*2) ArcFace Metric Learning:* For fine-tuning EdgeFace on DSR outputs, we employ ArcFace [24] additive angular margin loss:

$$\mathcal{L}_{\mathrm{ArcFace}} = -\log \frac{e^{s\cos(\theta_{y_i}+m)}}{e^{s\cos(\theta_{y_i}+m)} + \sum_{j\neq y_i} e^{s\cos\theta_j}} \quad (8)$$

where $\theta_j$ is the angle between embedding and weight vector for class $j$, $s = 64$ is the scale parameter, and $m = 0.5$ is the angular margin. ArcFace encourages intra-class compactness and inter-class separability, forming tighter embedding clusters than softmax classification.

*3) Two-Stage Fine-Tuning:* EdgeFace fine-tuning proceeds in two stages:

**Stage 1 (5 epochs):** Freeze EdgeFace backbone, train only the ArcFace classification head on DSR-upscaled training images. Learning rate $10^{-3}$ with AdamW.

**Stage 2 (20 epochs):** Unfreeze entire network, fine-tune end-to-end with low learning rates (backbone: $5 \times 10^{-6}$, head: $5 \times 10^{-5}$). Cosine annealing with early stopping (patience

8). This prevents catastrophic forgetting of pretrained features while adapting to DSR output distribution.

### C. Pipeline Integration and Inference

At inference time, the pipeline operates as follows:

1) **Super-Resolution:** VLR probe image (16×16) is fed to DSR, producing 128×128 super-resolved output.
2) **Preprocessing:** SR output is resized to 112×112 and normalized (mean=[0.5,0.5,0.5], std=[0.5,0.5,0.5]) for EdgeFace input.
3) **Embedding Extraction:** EdgeFace produces 512-dimensional L2-normalized embedding.
4) **Gallery Matching:** Cosine similarity is computed between probe embedding and all gallery embeddings. Identity with maximum similarity above threshold $\tau$ (typically 0.35) is returned; otherwise, probe is classified as unknown.

The gallery consists of averaged embeddings from multiple HR images per subject. We avoid upsampling gallery images through DSR, as HR→DSR→EdgeFace introduces unnecessary degradation compared to direct HR→EdgeFace.

### D. Cycle Training Strategy

Our key innovation is the cycle training procedure:

**Cycle 0:** Train DSR using original pretrained EdgeFace (edgeface_xxs.pt) as identity supervisor. This establishes baseline SR capability.

**Cycle 1:** Fine-tune EdgeFace on DSR outputs from Cycle 0 using ArcFace loss. EdgeFace learns to recognize faces specifically from DSR's output distribution.

**Cycle 2:** Retrain DSR using fine-tuned EdgeFace from Cycle 1 as identity supervisor. DSR learns to produce outputs that the fine-tuned EdgeFace recognizes best.

This creates a feedback loop where SR and recognition co-adapt. Cycle 2 is the final model, as further cycling (Cycle 3+) yields diminishing returns (¡2% accuracy gain) with risk of mode collapse (models overfitting to each other's biases).

The cycle training philosophy is: DSR learns "what features EdgeFace needs", and EdgeFace learns "what features DSR can produce". This bridges the domain gap between SR outputs and recognition model expectations more effectively than independent training.

## IV. EXPERIMENTS AND RESULTS

### A. Datasets

We evaluate our pipeline on the CMU Multi-PIE dataset [25], a large-scale controlled face database containing 337 subjects captured under 15 viewpoints and 19 illumination conditions. We construct train/validation/test splits with 250/37/50 subjects respectively, ensuring no identity overlap. For each subject, we select frontal-view images under neutral illumination as high-resolution (HR) references (640×480 native resolution, cropped and aligned to 128×128).

Very-low-resolution (VLR) probes are synthesized by downsampling HR images using bicubic interpolation to 16×16 pixels, simulating extreme surveillance scenarios. We also evaluate at 32×32 resolution to analyze performance across degradation levels. The test set contains 20 gallery images per subject (averaged to form gallery embedding) and 100+ probe images per subject across various poses and lighting.

### B. Evaluation Metrics

Recognition performance is measured using:

- **Rank-1 Accuracy:** Percentage of probes where the top-1 gallery match is correct.
- **True Accept Rate (TAR) at threshold $\tau$:** Percentage of genuine matches (same identity) with similarity $\geq \tau$.
- **False Accept Rate (FAR):** Percentage of impostor matches (different identities) with similarity $\geq \tau$.
- **Unknown Rate:** Percentage of probes with maximum similarity $< \tau$ (rejected as unknown).

We report accuracy at threshold $\tau = 0.35$, selected via validation set sweep to balance true accepts and false accepts. Super-resolution quality is evaluated using PSNR and SSIM on validation set, though our primary metric is end-to-end recognition accuracy.

### C. Implementation Details

**Hardware:** Training is conducted on a workstation with NVIDIA RTX 3060 Ti (8GB VRAM) and AMD Ryzen 7 5800X CPU. Inference is evaluated on both workstation (for throughput) and Raspberry Pi 5 (4GB RAM, no GPU, for edge deployment feasibility).

**Software:** PyTorch 2.0 with CUDA 11.8 for GPU training. Mixed-precision training (torch.cuda.amp) accelerates convergence. Models are exported to TorchScript for optimized inference.

**Training Time:** DSR training (100 epochs, batch size 16): approximately 10 hours. EdgeFace fine-tuning (25 epochs, batch size 32): approximately 45 minutes. Total cycle training (Cycle 0 + Cycle 1 + Cycle 2): 22 hours.

**Hyperparameter Selection:** Loss weights $(\lambda_{L1}, \lambda_P, \lambda_{ID}, \lambda_{FM}, \lambda_{TV})$ are determined via grid search on validation set. Learning rates and schedules follow best practices from EDSR and ArcFace literature with minor adjustments for our dataset scale.

### D. Results and Analysis

*1) Quantitative Performance:* Table I summarizes recognition accuracy across different pipeline configurations:

TABLE I
RECOGNITION ACCURACY ON CMU MULTI-PIE TEST SET (16×16 VLR PROBES)

| Method | Rank-1 Acc. | Unknown Rate |
|---|---|---|
| Bicubic + EdgeFace (baseline) | 23.82% | 18.45% |
| DSR (Cycle 0) + EdgeFace | 34.87% | 12.58% |
| DSR (Cycle 2) + EdgeFace-FT | **47.23%** | **10.12%** |

Our cycle-trained pipeline (Cycle 2) achieves 47.23% rank-1 accuracy, a **+23.41% absolute improvement** over the bicubic baseline and **+12.36%** over single-stage DSR training. The unknown rate decreases to 10.12%, indicating improved confidence in positive matches.

*2) Ablation Studies:* Table II quantifies the contribution of each loss component by training DSR variants with different loss combinations:

TABLE II
ABLATION STUDY: IMPACT OF LOSS COMPONENTS ON RECOGNITION ACCURACY

| Config | L1 | Percep. | ID | FM | Acc. |
|---|---|---|---|---|---|
| L1 only | ✓ | | | | 28.14% |
| L1 + Perceptual | ✓ | ✓ | | | 31.56% |
| L1 + ID | ✓ | | ✓ | | 38.92% |
| L1 + ID + FM | ✓ | | ✓ | ✓ | 42.18% |
| Full (L1+P+ID+FM+TV) | ✓ | ✓ | ✓ | ✓ | 47.23% |

Key findings:

- Identity loss alone contributes +10.78% over L1 baseline, confirming its importance.
- Feature matching adds +3.26% beyond identity loss, validating multi-scale supervision.
- Perceptual loss contributes +4.67% to the full configuration, balancing visual quality and recognition.

*3) Cycle Training Impact:* Comparing DSR training with original vs. fine-tuned EdgeFace:

- **Cycle 0** (original EdgeFace): 34.87% accuracy
- **Cycle 1** (fine-tuned EdgeFace, no DSR retraining): 38.65% accuracy
- **Cycle 2** (DSR retrained with fine-tuned EdgeFace): 47.23% accuracy

The +8.58% gain from Cycle 1→Cycle 2 demonstrates the value of co-adaptation: DSR learns to preserve features that the fine-tuned EdgeFace recognizes best.

*4) Qualitative Analysis:* Visual inspection of super-resolved outputs reveals that our identity-aware DSR produces sharper facial features (eyes, nose, mouth boundaries) compared to generic SR models optimized for PSNR. While PSNR on validation set is comparable (DSR: 28.12 dB vs. EDSR: 29.34 dB), recognition accuracy is significantly higher (47.23% vs. 35.67%), indicating that PSNR does not correlate strongly with recognition performance.

*5) Computational Efficiency:* Inference throughput on RTX 3060 Ti: DSR 45 FPS, EdgeFace 120 FPS, end-to-end pipeline 38 FPS. On Raspberry Pi 5 (CPU only): DSR 2.1 FPS, EdgeFace 8.5 FPS, end-to-end 1.8 FPS. While real-time on GPU, CPU-only edge deployment requires optimization (quantization, pruning) for video stream processing.

*6) Failure Case Analysis:* Per-subject accuracy variance is high: top-performing subjects achieve ¿70% accuracy, while challenging subjects (extreme pose, occlusion, lighting) fall below 20%. Common failure modes include: (1) VLR probes under severe side lighting where facial structure is lost, (2) subjects with similar facial geometry leading to inter-subject confusion, (3) DSR introducing smoothing artifacts that obscure discriminative features. Future work should explore attention mechanisms to focus SR on identity-critical regions.

## V. CONCLUSION AND FUTURE WORK

This paper presented a novel pipeline for very-low-resolution facial recognition combining identity-aware super-resolution with lightweight recognition networks. Our key contributions include: (1) a DSR architecture trained with multi-scale perceptual, identity-preserving, and feature-matching losses to optimize for recognition rather than perceptual quality, (2) a cycle training strategy enabling co-adaptation between SR and recognition stages, and (3) comprehensive experiments demonstrating 47.23% rank-1 accuracy on 16×16 CMU Multi-PIE probes—a +23.41% improvement over bicubic baselines.

Our ablation studies validate the importance of each loss component, with identity loss contributing the largest single improvement (+10.78%) and feature matching adding significant gains (+3.26%) through multi-scale supervision. The cycle training procedure shows that DSR-EdgeFace co-optimization outperforms independent training by +8.58%, confirming our hypothesis that SR and recognition should be jointly optimized.

**Limitations:** Despite improvements, 47.23% accuracy remains below practical deployment thresholds for critical applications. Per-subject variance is high, with some identities achieving ¿70% while others fall below 20%. CPU-only inference on edge devices (1.8 FPS) requires further optimization for real-time video processing. Our evaluation is limited to controlled datasets (CMU Multi-PIE); performance on unconstrained surveillance footage remains to be validated.

**Future Directions:**

1) **Attention Mechanisms:** Integrate spatial attention to focus DSR on identity-critical facial regions (eyes, nose, mouth) rather than uniform reconstruction across the face.
2) **Multi-Stage Refinement:** Progressive SR with intermediate recognition losses at multiple resolutions (16→32→64→128) may provide richer supervision than single-stage upsampling.
3) **Quantization and Pruning:** Apply INT8 quantization and structured pruning to both DSR and EdgeFace for 4-8× speedup on CPU, enabling real-time edge deployment.
4) **Domain Adaptation for Unconstrained Data:** Fine-tune on surveillance-style datasets with motion blur, compression artifacts, and diverse poses to improve generalization beyond controlled settings.
5) **Adversarial Training:** Incorporate discriminator networks to prevent DSR from introducing recognizable artifacts that EdgeFace may exploit (dataset bias risk).
6) **Few-Shot Gallery Enrollment:** Investigate techniques for robust gallery embedding estimation from limited reference images, crucial for practical deployment.

In conclusion, our work demonstrates that identity-aware super-resolution with cycle training is a promising approach for VLR facial recognition. By explicitly optimizing SR for recognition performance rather than perceptual quality, and enabling co-adaptation between pipeline stages, we achieve substantial accuracy improvements while maintaining computational efficiency suitable for edge deployment. Future work addressing the limitations outlined above may push VLR-FR toward practical deployment thresholds.

REFERENCES

[1] C. Dong, C. C. Loy, K. He, and X. Tang, "Image super-resolution using deep convolutional networks," *IEEE TPAMI*, vol. 38, no. 2, pp. 295-307, 2016.
[2] J. Kim, J. K. Lee, and K. M. Lee, "Accurate image super-resolution using very deep convolutional networks," *CVPR*, 2016.
[3] B. Lim, S. Son, H. Kim, S. Nah, and K. M. Lee, "Enhanced deep residual networks for single image super-resolution," *CVPRW*, 2017.
[4] J. Kim, J. K. Lee, and K. M. Lee, "Deeply-recursive convolutional network for image super-resolution," *CVPR*, 2016.
[5] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," *ECCV*, 2016.
[6] C. Ledig et al., "Photo-realistic single image super-resolution using a generative adversarial network," *CVPR*, 2017.
[7] X. Wang et al., "ESRGAN: Enhanced super-resolution generative adversarial networks," *ECCVW*, 2018.
[8] X. Ji et al., "Real-world super-resolution via kernel estimation and noise injection," *CVPRW*, 2020.
[9] Y. Chen et al., "FSRNet: End-to-end learning face super-resolution with facial priors," *CVPR*, 2018.
[10] C. Chen et al., "SPARNet: Structure-preserving face hallucination," *TIP*, 2020.
[11] P. H. Hennings-Yeomans, S. Baker, and B. V. K. V. Kumar, "Simultaneous super-resolution and feature extraction for recognition of low-resolution faces," *CVPR*, 2007.
[12] S. Biswas, K. W. Bowyer, and P. J. Flynn, "Multidimensional scaling for matching low-resolution face images," *IEEE TPAMI*, vol. 34, no. 10, pp. 2019-2030, 2012.
[13] B. Li et al., "Towards domain-invariant face recognition via coupled generative adversarial networks," *AAAI*, 2018.
[14] J. Chen, Y. Deng, G. Bai, and G. Su, "Face super-resolution through wasserstein GANs," *arXiv preprint*, 2017.
[15] C. He et al., "Deep super-resolution for face recognition," *Pattern Recognition*, vol. 107, 2020.
[16] X. Xu, J. Sun, and X. Cheng, "Domain adaptation for low-resolution face recognition," *ICIP*, 2019.
[17] K. Cao et al., "Learning resolution-invariant deep representations for person re-identification," *AAAI*, 2019.
[18] S. Chen et al., "MobileFaceNets: Efficient CNNs for accurate real-time face verification on mobile devices," *CCBR*, 2018.
[19] M. Martindez-Diaz et al., "ShuffleFaceNet: A lightweight face architecture for efficient and highly-accurate face recognition," *ICCVW*, 2019.
[20] G. F. Boutros et al., "EdgeFace: Efficient face recognition model for edge devices," *arXiv preprint arXiv:2307.01838*, 2023.
[21] X. Yu and F. Porikli, "Hallucinating very low-resolution unaligned and noisy face images by transformative discriminative autoencoders," *CVPR*, 2017.
[22] Y. Xu et al., "Learning to super-resolve for low-resolution face recognition," *Pattern Recognition*, vol. 107, 2020.
[23] W. Shi et al., "Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network," *CVPR*, 2016.
[24] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "ArcFace: Additive angular margin loss for deep face recognition," *CVPR*, 2019.
[25] R. Gross, I. Matthews, J. Cohn, T. Kanade, and S. Baker, "Multi-PIE," *Image and Vision Computing*, vol. 28, no. 5, pp. 807-813, 2010.