

Deep Learning-based Face Super-resolution: A Survey

Junjun Jiang · Chenyang Wang · Xianming Liu · Jiayi Ma

Received: date / Accepted: date

Abstract Face super-resolution, also known as face hallucination, which is aimed at enhancing the resolution of low-resolution (LR) one or a sequence of face images to generate the corresponding high-resolution (HR) face images, is a domain-specific image super-resolution problem. Recently, face super-resolution has received considerable attention, and witnessed dazzling advances with deep learning techniques. To date, few summaries of the studies on the deep learning-based face super-resolution are available. In this survey, we present a comprehensive review of deep learning techniques in face super-resolution in a systematic manner. First, we summarize the problem formulation of face super-resolution. Second, we compare the differences between generic image super-resolution and face super-resolution. Third, datasets and performance metrics commonly used in facial hallucination are presented. Fourth, we roughly categorize existing methods according to the utilization of face-specific information. In

each category, we start with a general description of design principles, present an overview of representative approaches, and compare the similarities and differences among various methods. Finally, we envision prospects for further technical advancement in this field.

Keywords Face super-resolution · Face hallucination · Deep learning · Convolution neural network

1 Introduction

Face super-resolution, a domain-specific image super-resolution problem, refers to the method of recovering high-resolution (HR) face images from given low-resolution (LR) face images. This method enhances the resolution of LR face images of low quality, and recovers the details of face images. In many real-world scenarios, limited by physical imaging systems and imaging conditions, the face images are always of low quality. Thus, face super-resolution has a wide range of applications and notable advantages, and it has always been a hot topic since its birth in the field of image processing and computer vision society.

The concept of face super-resolution was first proposed by Baker and Kanade [6] in 2000, which is the pioneer of face super-resolution technique. They develop a multi-level learning and prediction model based on Gaussian image pyramid to improve the resolution of LR face images. Liu *et al.* [100] proposed to integrate a global parametric principal component analysis (PCA) model with a local nonparametric Markov random field (MRF) model for face super-resolution. Since then, a number of innovative methods have been proposed and face super-resolution has been the subject of active research efforts. Researchers super-resolve the LR face images by using global face statistical models [50, 154, 20, 125, 99, 169], local

J. Jiang
School of Computer Science and Technology, Harbin Institute of Technology
Harbin 150001, China
E-mail: jiangjunjun@hit.edu.cn

C. Wang
School of Computer Science and Technology, Harbin Institute of Technology
Harbin 150001, China
E-mail: wangchy02@hit.edu.cn

X. Liu
School of Computer Science and Technology, Harbin Institute of Technology
Harbin 150001, China
E-mail: csxm@hit.edu.cn

J. Ma
Electronic Information School, Wuhan University
Wuhan 430072, China
E-mail: jyima2010@gmail.com

Table 1 Summary of face super-resolution survey in the past decade.

No.	Survey title	Year	Venue	Content
1	A survey of face hallucination [98]	2012	CCBR	
2	A comprehensive survey to face hallucination [150]	2014	IJCV	
3	A review of various approaches to face hallucination [5]	2015	ICACTA	
4	Face super resolution: a survey [73]	2017	IJIGSP	Traditional methods.
5	Super-resolution for biometrics: a comprehensive survey [122]	2018	PR	
6	Face hallucination techniques: a survey [128]	2018	CICT	
7	Survey on GAN-based face hallucination with its model development [103]	2019	IET	GAN-based methods.

patch-based representation methods [22, 115, 71, 67, 41], or hybrid ones [202, 58]. With the rapid development of deep learning technique, attractive advantages over previous attempts have been obtained and have been applied into image or video super-resolution, and many comprehensive surveys review the recent achievements in these fields, *i.e.*, general image super-resolution survey [156, 2, 174], and video super-resolution survey [102]. Towards face super-resolution, a domain-specific image super-resolution, a few surveys are listed in Table 1. In the early stage of research, [98, 150, 5, 73, 122, 128] provide a comprehensive review of traditional face super-resolution methods (mainly including patch-based super-resolution, PCA-based methods, *etc.*), while Liu *et al.* [103] offers a generative adversarial network (GAN) based face super-resolution model. However, so far no review literature is available on deep learning super-resolution specifically for human faces. Thus, we present a comparative study of different deep learning-based face super-resolution methods.

In this exposition, our focus is on deep learning based face super-resolution. The main contributions of this survey are threefold:

- The survey provides a comprehensive review of recent techniques for face super-resolution, including the characteristics of face super-resolution, benchmark datasets, deep learning-based face super-resolution methods, commonly used loss function, related applications, and others.
- The survey summarizes how existing deep learning-based face super-resolution methods explore the potential of network architecture and take advantage of characteristics of face images, and compare the similarities and differences between these methods.
- The survey discusses the challenges and envisions the prospects for future research in the face super-resolution field.

The overall organization is presented in Fig. 1. Section 2 introduces the problem definition of face super-resolution. Section 3 compares the differences between general image super-resolution and face super-resolution, and presents the characteristics of face images and reviews mainstream face datasets and evaluation metrics. In Section 4, we discuss face super-resolution methods. To avoid exhaustive enumeration and take facial characteristic into consideration, face super-resolution methods are categorized according to face-specific information used. In Fig. 1, six major categories are presented: general face super-resolution methods, prior-guided face super-resolution methods, attribute-constrained face super-resolution methods, identity-preserving face super-resolution methods, reference face super-resolution methods, and audio-guided super-resolution method. Depending on network architecture or the utilization of face-specific information, every category is further divided into several subcategories. Besides, Section 4 also reviews the commonly used loss function in face super-resolution, and some methods dealing with joint tasks and face super-resolution related applications. Section 5 concludes the face super-resolution and further discusses the limitations and envisions prospects for further technical advancement.

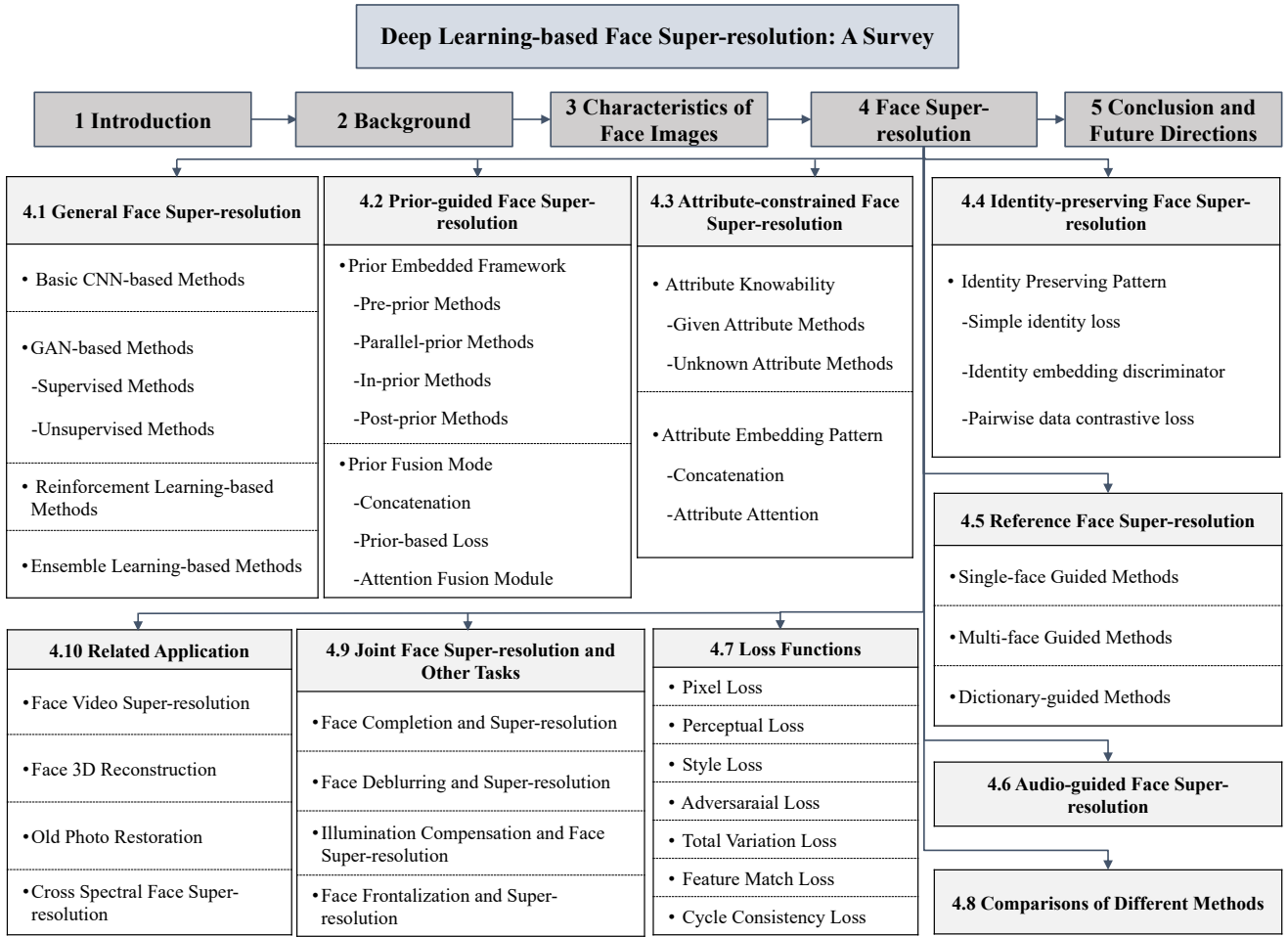


Fig. 1 Structure of this survey.

2 Background

2.1 Problem Definition

Face super-resolution focuses on recovering the corresponding HR face images from LR face images. Obtaining LR can be modeled as follows:

$$I_{LR} = \Phi(I_{HR}, \theta), \quad (1)$$

where I_{LR} is the LR face image, I_{HR} is its corresponding HR face image, Φ and θ are the degradation model and degradation parameters (*e.g.*, downsampling operation, noise, and others.) respectively. Nevertheless, face super-resolution is devoted to simulating the inverse process of the degradation model and recovers the I_{SR} from I_{LR} such that I_{SR} is as close as possible to I_{HR} , which can be expressed as:

$$I_{SR} = \Phi^{-1}(I_{LR}, \delta) = F(I_{LR}, \delta), \quad (2)$$

where F is the super-resolution model (inverse degradation model) and δ represents the parameters of F , I_{SR} represents

the super-resolved (SR) result. Toward the training of the parameters δ , it can be defined as

$$\hat{\delta} = \underset{\delta}{\operatorname{argmin}} \mathcal{L}(I_{SR}, I_{HR}), \quad (3)$$

where \mathcal{L} represents the loss between I_{SR} and I_{HR} and $\hat{\delta}$ is the parameter of the trained model. (In face super-resolution, MSE loss is the most popular loss in face super-resolution, and some models tend to use a combination of multiple loss functions, which are reviewed in section 4.7.)

In a real-world environment, the degradation model and parameters are all unavailable, and I_{LR} is only given information. At this point, face super-resolution model is difficult to train. Thus, researchers tend to use mathematical models to simulate the degradation model and generate the LR and HR pair to train the model. The simplest mathematical model is

$$I_{LR} = (I_{HR}) \downarrow_s, \quad (4)$$

where \downarrow denotes downsampling operation (always imresize function in MATLAB by default setting, sometimes maybe

other functions), and s is the scaling factor. However, this pattern is too simple to match the real-world degradation process. To simulate the real degradation process better, researchers design a degradation process with the combination of many operations (*e.g.*, downsampling, blur, noise and compression) as follows:

$$I_{LR} = ((I_{HR} \otimes k) \downarrow_s + n)_{\text{JPEG}}, \quad (5)$$

where k is the blurring kernel, \otimes presents the convolutional operation, n denotes the noise, and JPEG is the JPEG compression. Various combinations of different operations are also used in face super-resolution. However, they are not introduced in detail in this study.

2.2 Image Quality Assessment

In face super-resolution, two main methods of quality evaluation are qualitative and quantitative evaluation. Qualitative evaluation relies on the judgement of humans, and qualitative evaluation invites readers or interviewers to see and assess the quality of the generated images, leading to the results always consistent with human perception. Quantitative evaluation mainly utilizes statistical data to reflect the quality of the generated images. In general, the quantitative evaluation methods produce different results from qualitative evaluation methods, because the starting point of quantitative evaluation methods is mathematics instead of human visual perception, which leaves assessment in dispute image quality. In most cases, these two evaluation methods are both used to assess the quality of images. Here, we introduce some commonly used evaluation metrics.

2.2.1 Peak Signal-to-Noise Ratio (PSNR)

PSNR is commonly used in face super-resolution. Given reference face image (ground truth) I_{HR} and test image (super-resolved result) I_{SR} , the first step is to calculate the mean squared error (MSE) between them. Based on the MSE and the maximum possible pixel value M , PSNR is obtained, which is defined as:

$$\text{MSE} = \frac{1}{hwc} \sum_{i,j,k} (I_{SR}^{i,j,k} - I_{HR}^{i,j,k})^2, \quad (6)$$

where h , w , and c denote the height, width, and channel of the image, then

$$\text{PSNR} = 10 \log_{10} \left(\frac{M^2}{\text{MSE}} \right). \quad (7)$$

Generally, M is fixed and is usually equal to 255, then PSNR is only related with MSE between two images. The lower MSE is, the higher PSNR is. The higher PSNR is

better. In this pattern, PSNR focuses on distance between every pair pixel in two images, which is not consistent with human perception, resulting in poor performance in cases where human perception is more important. However, PSNR remains the most popular and commonly used metric in face super-resolution.

2.2.2 Structural Similarity

Different from PSNR which only measures distance between pixels, structural similarity index (SSIM) [199] is proposed to measure structural similarity between structural information. To be specific, SSIM measures similarity from three aspects: luminance, contrast, and structure. Given face image I_{HR} (I_{SR}) with height h , width w and channel c , luminance and contrast are estimated as the mean and standard of the image intensity:

$$\mu_{I_{HR}} = \frac{1}{hwc} \sum_{i,j,k} I_{HR}^{i,j,k}, \quad (8)$$

$$\sigma_{I_{HR}} = \left(\frac{1}{hwc-1} \sum_{i,j,k} (I_{HR}^{i,j,k} - \mu_{I_{HR}})^2 \right)^{\frac{1}{2}}, \quad (9)$$

where $\mu_{I_{HR}}$ is mean and $\sigma_{I_{HR}}$ is standard of I_{HR} . The similarity between luminance and contrast is defined as

$$\mathcal{L}(I_{HR}, I_{SR}) = \frac{2\mu_{I_{HR}}\mu_{I_{SR}} + C_1}{\mu_{I_{HR}}^2 + \mu_{I_{SR}}^2 + C_1}, \quad (10)$$

$$\mathcal{C}(I_{HR}, I_{SR}) = \frac{2\sigma_{I_{HR}}\sigma_{I_{SR}} + C_2}{\sigma_{I_{HR}}^2 + \sigma_{I_{SR}}^2 + C_2}, \quad (11)$$

where C_1, C_2 are constants to avoid division by zero. The image structural similarity between two images is estimated as the correlations between normalized pixel value, which is expressed as

$$\sigma_{I_{HR}, I_{SR}} = \frac{1}{hwc-1} \sum_{i,j,k} \left(I_{HR}^{i,j,k} - \mu_{I_{HR}} \right) \left(I_{SR}^{i,j,k} - \mu_{I_{SR}} \right), \quad (12)$$

$$S(I_{HR}, I_{SR}) = \frac{\sigma_{I_{HR}, I_{SR}} + C_2/2}{\sigma_{I_{HR}}\sigma_{I_{SR}} + C_2/2}. \quad (13)$$

SSIM is obtained by

$$\text{SSIM}(I_{HR}, I_{SR}) = \mathcal{L}(I_{SR}, I_{SR}) * \mathcal{C}(I_{HR}, I_{SR}) * S(I_{HR}, I_{SR}). \quad (14)$$

SSIM varies from 0 to 1, where a higher value is better. Considering the uneven distribution of the image, SSIM is less reliable than assessing images locally. Thus,

multi-scale structural similarity index measure (MS-SSIM) [157] is proposed to enhance reliability, which divides the whole images into multiple windows, and first assesses SSIM for every window separately, and then converges them to obtain the final MS-SSIM.

2.2.3 FID

PSNR measures the distance between pixels, SSIM assesses similarity between luminance, contrast and structure, and frechet inception distance score (FID) [52] compares the difference between feature maps. Different from PSNR and SSIM, FID is an evaluation metric that is used to assess the generative ability of generator adversarial network [46] (which is introduced in section 4.7), and is always applied to assess the visual quality of face images. FID is defined as

$$\text{FID} = \left\| \mu_{f_{I_{HR}}} - \mu_{f_{I_{SR}}} \right\|^2 + \text{Tr} \left(\Sigma_{f_{I_{HR}}} + \Sigma_{f_{I_{SR}}} - 2 \left(\Sigma_{f_{I_{HR}}} \Sigma_{f_{I_{SR}}} \right)^{1/2} \right), \quad (15)$$

where $f_{I_{HR}}, f_{I_{SR}}$ are corresponding features extracted from pre-trained Inception V3 [132], and μ is mean, Tr is trace of the matrix, and Σ is the covariance matrix. The smaller FID is, the better the visual quality is.

Certainly, except PSNR, SSIM, FID, there are still many evaluation metric, such as MS-SSIM, NIQE [119], mean opinion score, and others. However, PSNR and SSIM are the mostly widely used in face super-resolution,

3 Characteristics of Face Images

General image super-resolution aims to enhance the resolution of general images which can be various and diverse. The image can be that of a tree, a cat, a person, or any other object. However, the aim of face super-resolution is to improve the resolution of face images and all the images have faces, which enables face super-resolution to become a special form of general image super-resolution. The face image is a highly structured object that has its own unique information, which can be explored and utilized. In this section, we simply introduce the face-specific information.

3.1 Prior Information

First, obvious structural prior are found in face images, such as facial parsing maps, facial landmarks, and facial heatmaps.

- facial parsing maps: semantic segmentation of facial images separating the facial components from face images, including eyebrows, eyes, nose, mouth, skin, ears, hair, and others.
- facial landmarks: these locate the key points of facial components. The number of landmarks vary in different datasets, such as CelebA [107], which provides 5 landmarks while Helen [85] offers 194 landmarks.
- facial heatmaps: these are generated from facial landmarks. Facial landmarks give accurate points of the facial components, while heatmaps give the probability of the point being a facial landmark. To generate the heatmaps, every landmark is represented by a Gaussian kernel centered on the location of the landmark.

These prior information can provide the location of facial components and facial structure information, which helps to recover the face shape. To understand them clearly, we provide concise illustrations of these prior information in Fig. 2.

3.2 Attribute Information

Second, attribute information is also the unique information in face images, which is semantic-level information. Different from general images, the principal part of face images is the human face, which has natural properties such as gender, hair color, and others. In face super-resolution, because of one-to-many maps from LR to HR, recovered face images may contain artifacts and even generate wrong attributes, for example, the face in the recovered result does not wear eyeglasses but the ground truth wears eyeglasses. At this time, attribute information can remind the network which attribute should be covered in the results. From a different perspective, attribute information also contains facial details. If we take the eyeglasses as an example, the attribute of wearing eyeglasses provides the details of the facial eyes. To provide a more comprehensive overview of all existing attribute labels, we not only provide a concise example of attribute information in Fig. 2, but also list commonly used 40 attributes in Table 2, including 11 facial component-related attributes, 14 hair-related attributes, 10 style-related attributes and 5 accessory-related attributes. Moreover, these attributes are always binary in face datasets, 1 denotes that the face image match the attribute information while 0 means mismatch. Notably, some datasets use -1 to denote mismatch, *e.g.* CelebA [107].

3.3 Identity Information

Third, every face image has a face that can identify a person, which is enabled by identity information. This type

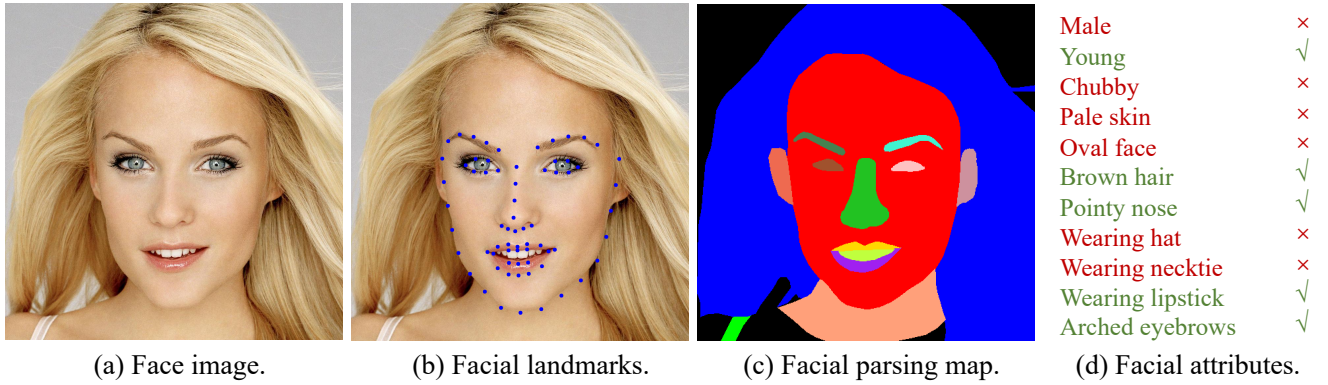


Fig. 2 Face-specific information.

Table 2 An overview of commonly used attributes.

Commonly used 40 attributes					
Facial component-related	Oval face	Arched eyebrows	Bushy eyebrows	Narrow eyes	Bags under eyes
	High cheekbones	Big nose	Pointy nose	Mouth slightly open	Big lips
	Double chin				
Hair-related	Receding hairline	Black hair	Blond hair	Brown hair	Gray hair
	Straight hair	Wavy hair	Bald	Bangs	No beard
	5 o'clock shadow	Mustache	Goatee	Sideburns	
Style-related	Heavy makeup	Attractive	Blurry	Chubby	Young
	Male	Smiling	Pale skin	Rosy cheeks	Wearing lipstick
Accessory-related	Wearing hat	Wearing earrings	Wearing necklace	Wearing necktie	Eyeglasses

of information is always used for keeping super-resolved results and ground truth identity consistency. On the one hand, visually the person should not be changed after super-resolution. On the other hand, face super-resolution should facilitate the performance of face recognition. Similar to attribute information, identity also offers identity-preserving facial details, which are beneficial to face restoration.

3.4 Datasets for Face Super-resolution

Many face image datasets are used for face super-resolution, which differ in many aspects (*e.g.*, image amounts, face-specific information contained, and others.). In Table 3, we list a number of commonly used face image datasets and simply indicate their amount, and the face-specific information they offer. For parsing maps and

identity, we only present whether they are provided, while for attributes and landmarks, we offer the specific amount. These datasets always provide HR face images, which means that we have to generate the corresponding LR by ourselves using degradation process. However, we show a few face images of the most commonly used two datasets: CelebA [107] and Helen [85] in Fig. 3.

Aside from the preceding datasets, many other face datasets are used in face super-resolution, including CACD200 [23], VGGFace2 [19], UMDFaces [7], and others. Some of them are facial recognition datasets that provide identity information and high-quality face images used in identity-preserving face super-resolution methods.

Table 3 List of public face image datasets for face super-resolution benchmark. Parsing maps and identity represent whether the dataset has facial parsing maps and identity information of face images.

Dataset	Amount of face images	Number of attributes	Number of landmarks	Parsing maps	Identity
CelebA [107]	202,599	40	5	×	✓
CelebAMask-HQ [74]	30,000	×	×	✓	×
Helen [85]	2,330	×	194	✓	×
FFHQ [75]	70,000	×	68	×	×
AFLW [81]	25,993	×	21	×	×
300W [131]	3,837	×	68	×	×
LS3D-W [11]	230,000	×	68	×	×
Menpo [187]	9,000	×	68	×	×
LFW [57]	13,233	73	×	×	✓
LFWA [158]	13,233	40	×	×	✓
CASIA-WebFace [176]	494,141	×	×	×	✓
VGGFace [126]	3,310,000	×	×	×	✓
WiderFace [173]	32,203	×	×	×	✓

4 Face Super-resolution Methods

At present, various deep learning face super-resolution methods are proposed. On the one hand, these methods explore the potential of efficient network for face super-resolution but set aside face particularity, *i.e.*, develop basic convolution neural network (CNN) or generative adversarial network (GAN) for face reconstruction. On the other hand, they focus on utilization of face-specific information, *e.g.*, using prior information to facilitate face repair, and so on. Furthermore, additional high-quality face image or audio sequences in face imaging also can be used to assist the restoration. Here, according to the type of face image special information, we divide face super-resolution methods into six categories: general face super-resolution, prior information guided face super-resolution, attribute-constrained face super-resolution, identity-preserving face super-resolution, reference face super-resolution, and audio-guided face super-resolution. In this section, we concentrate on every kind of face super-resolution methods and introduce each one in detail.

4.1 General Face Super-resolution Methods

General face super-resolution methods always design an efficient network for face images without any face-specific information. These methods exploit the potential of efficient network structure to super-resolve face images. Since they don't take prior information into consideration, we call them general face super-resolution methods. Depending on the design of their model, we divide general face super-resolution methods into four categories: basic CNN-based methods, GAN-based methods, reinforcement learning-based methods and ensemble learning-based methods. We show their frameworks in Fig. 4. Moreover, some of these categories can be further divided. Aiming to present a clear and concise overview, we create and summarize the general face super-resolution methods in Fig. 5.

4.1.1 Basic CNN-based Methods

Basic CNN-based methods develop CNN-based face super-resolution network for face reconstruction. Depending on whether they consider the global information and local differences, we can further divide basic CNN-based methods into three categories: global methods



(a) Examples of CelebA [107].



(b) Examples of Helen [85].

Fig. 3 Examples of commonly used datasets: CelebA [107] and Helen [85].

that feed the entire face into the network and recover face images globally, local methods that crop face images into regions or patches then recover face images locally or consider local differences, global and local methods that combine global and local methods and recover face images locally and globally.

Global Methods: In early years, researchers view a face image as a whole and recover it globally. Inspired by powerful representation of CNN, Zhou *et al.* [198] firstly adopt CNN and propose bi-channel convolutional neural network (BCCNN) to learn a mapping from LR face images to HR face images. Specifically, LR is fed into feature extraction layer (comprised of three convolution layer) followed by two branches, generating a parameter α and coarse SR results I_{Rec} respectively for adaptively integrating I_{LR} and I_{Rec} as follows:

$$I_{\text{SR}} = \alpha(I_{\text{LR}}) \uparrow_s + (1 - \alpha)I_{\text{Rec}}, \quad (16)$$

where I_{LR} is the LR, and \uparrow_s means upsampling $\times s$ by bicubic interpolation.

Similar to image super-resolution using deep convolutional networks (SRCNN) [35], the work of Huang

et al. [61] also consists of three layers for face reconstruction. Then, inspired by the performance gain of iterative back projection (IBP) in general super-resolution, IBP is introduced for face super-resolution as an extra post-processing step based on SRCNN, which is called SRCNN-IBP [56]. In depth, after repairing a SR result I_{CSR} by SRCNN [35], IBP projects I_{CSR} into the low-resolution space by downsampling operation and produces $I_{\text{CSR}}^{\downarrow s}$, and then calculates the residual map I_{Re} between $I_{\text{CSR}}^{\downarrow s}$ and I_{LR} , which is the error that should be complemented in I_{CSR} . To complement I_{CSR} and generate refinement results, I_{Re} is upsampled and then added to I_{CSR} . In this way, I_{CSR} is self-corrected and refined.

Cascaded model face hallucination (CDFH) [101] is a cascaded model that first denoises and recovers low-frequency information, and then compensates the high-frequency information using the second subnetwork. Similar to CDFH, face hallucination via convolution neural network (FHCNN) [124] is also a cascaded face super-resolution network that applies pixel-wise loss function at every step. Self-attention residual network

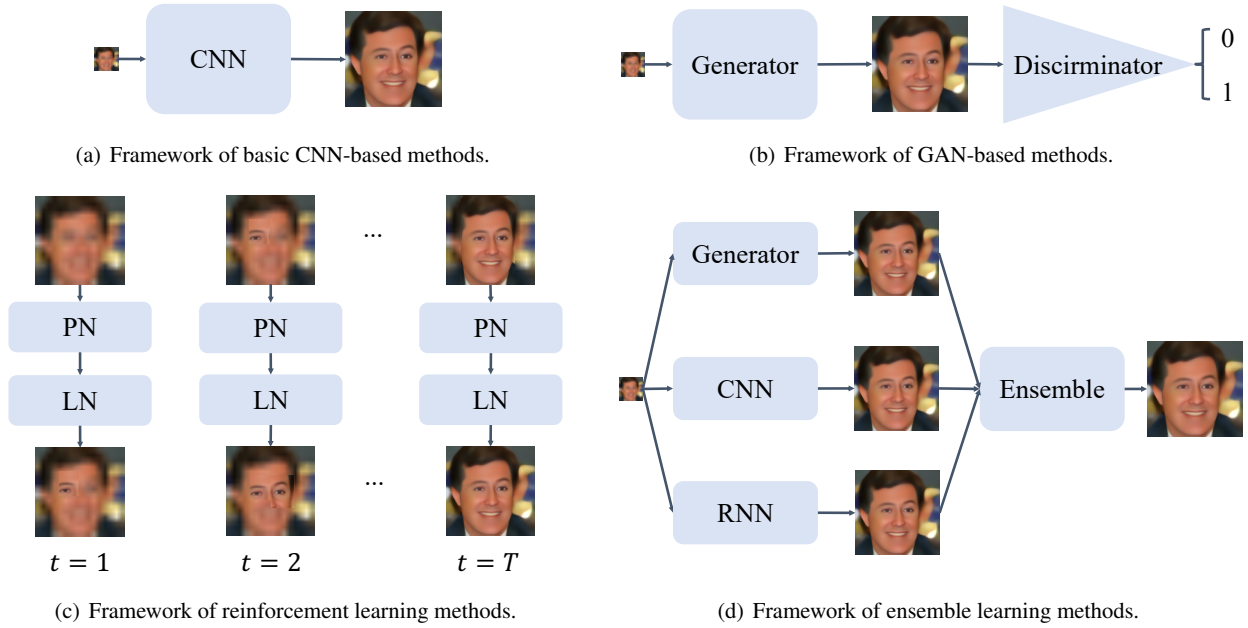


Fig. 4 Four frameworks of prior information guided face super-resolution methods. The wax yellow block denotes prior information. PN is policy network, while LN means local network in [18]. RNN is Recurrent neural network.

(SARN) [105] applies self-attention mechanism for face restoration. Residual back-projection network (RBPNet) [153] designs a residual back-projection two-branch network, including a base model, to extract features for super-resolution and edge map prediction branches. Split-attention in split-attention network (SISN) [155] designs internal-feature split attention block to generate internal-feature attention and to enhance the facial features correlation in channel dimension.

Considering that face images in the real world are always in multi-scale but the previous methods lack multiple scale representation ability, sequential gating ensemble network (SGEN) [28] combines multiple encoders and decoders in bottom-top and top-bottom pattern to extract multi-scale features and captures multi-scale details, and designs a sequential gating unit to improve the tasks of selecting and fusing information. It is observed that super-resolution in image domain produces smooth results without high-frequency detail. Considering that wavelet transform can depict the textural and contextual information of the images, WaveletSRNet [59] aims to super-resolve face images in wave coefficient, which first transforms face images into wavelet coefficients and super-resolves the face images in wavelet coefficient domain to address the problem of over-smooth results. WaveletSRNet is comprised of embedding, wavelet prediction, and reconstruction networks. The LR is fed into the embedding network to extract features that are imported into the wavelet prediction network predicting wavelet

coefficient, and then the reconstruction network transforms the wavelet into SR images. Then, the wavelet-based loop architecture face hallucination (WaveLAFH) [45] also predicts the wavelet coefficient for face restoration.

Local Methods: Although global methods can capture global information, they ignore the differences among different facial regions, thereby to difficulty of learning and a decrease in the performance of the model. To settle this problem, adaptive aggregation network (AAN) [51], a method to deal with noise face super-resolution, is proposed, including two generator branches and an aggregation branch. To deal with each case on its merits, the aggregation branch is designed to generate a mask M with input LR, and then M and $1 - M$ are used as the mask to be applied on the output of two generator branches respectively. Thus, two generator branches can pay attention to different regions. Finally, masked outputs of two generator branches is added to produce the final SR results. The process can be expressed as

$$I_{SR} = G_1(I_{LR}) \otimes M + G_2(I_{LR}) \otimes (1 - M), \quad (17)$$

where G_1, G_2 are the two generator branches, and \otimes denotes pixel-wise product. That is, M adaptively divides the whole image into two categories, and the two generator branches super-resolve different regions individually.

Super-resolution technique based on definition-scalable inference (SRDSI) [54] also divides face images and recovers different components. In contrast to AAN, SRDSI decomposes the face into a basic face with low-frequency

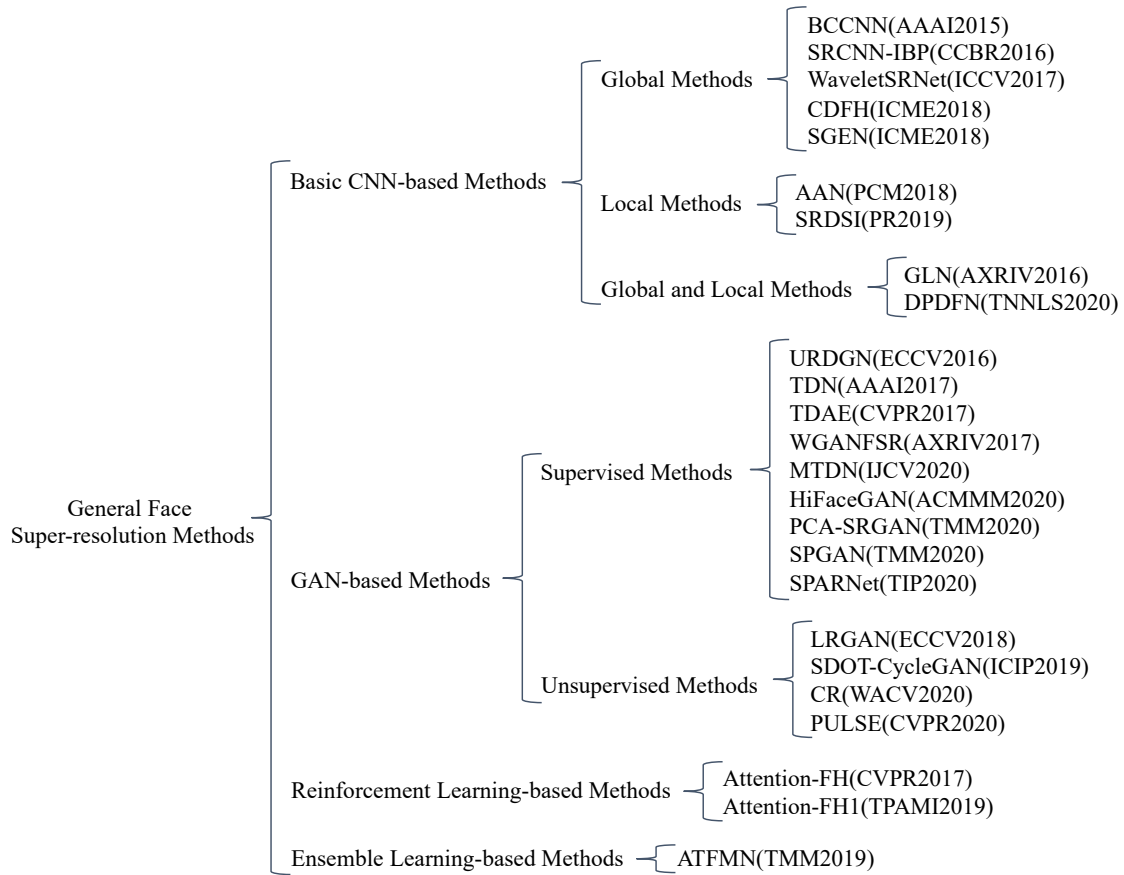


Fig. 5 Overview of general face super-resolution methods. Considering the space limitation, we only list some representative works.

and enhance face with high-frequency through PCA, and recovers the basic face with very deep convolutional networks (VDSR) [79], repairs enhance face with sparse representation, and finally fuses the recovered enhance face and the basic face. Thereafter, many patch-based methods occurs, including [80, 193, 111, 110, 123], all of which crop face images into several patches and train models for recovering corresponding patches.

Global and Local Methods: Considering that global methods can capture global structure but ignore local details, although local methods focus on local details but lose global structure, we can naturally combine global and local methods for capturing global structure and recovering local details. At first, global-local network (GLN) [144] proposes a global-local face super-resolution network, including two sub-networks: global upsampling network (GN) which models global constraints and local enhancement network (LN) that learns face-specific details. GN is a two-stream parallel network. The first stream uses deconvolutional network to simply upsample the LR, generating a smooth image with the same resolution as HR. As the output of the first stream shows limited details, the second stream, which is made up of fully connected network, aims to produce high-frequency detail. Then, the

concatenation of output of two streams is fed into the LN to enhance the local details. To simultaneously capture global clues and repair local details, dual-path deep fusion network (DPDFN) [69] constructs two individual branches for learning global facial contours and local facial component details, and then fuses the results of the two branches generating final SR results.

4.1.2 GAN-based Methods

GAN, first proposed by Goodfellow *et al.* [46], can be trained with paired or unpaired data, and generate real images with details, which inspire researchers to utilize GAN to recover face images with additional details. Thus, GAN-based methods are proposed. According to paired or unpaired training data, GAN-based methods can be divided into supervised methods with paired data and unsupervised methods with unpaired data.

Supervised Methods: In the early stage, Yu *et al.* [182] developed an original GAN-based face super-resolution network: ultra-resolving face images by discriminative generative networks (URDGN), which consists of two subnetworks: a discriminative model to distinguish a real HR face image or an artificially super-resolved output, and

a generative model to generate SR face images to fool the discriminative model and match HR and distribution of HR. The competition between two subnetworks enables the generative model to generate HR face images with better perceptual quality. However, URDGN acquires LR to be aligned, which may be unsatisfactory. Thus, Yu *et al.* [183] proposed an improved version by inserting multiple spatial transformer networks (STN) [65] in a generator to achieve face alignment, transformative discriminative networks (TDN). Nevertheless, the input of TDN is fixed-sized, which is inflexible. Thus, multiscale transformative discriminative networks (MTDN) [184] is developed, which downsamples LR into 16×16 as low-frequency and upsamples it into 32×32 as high-frequency, and feeds them into two branches to extract features that are concatenated for the following layers. As LR can be noisy and unaligned, Yu *et al.* built a decoder-encoder-decoder network named transformative discriminative autoencoders (TDAE) [164]. The first decoder upsamples and coarsely aligns LR by STN producing I_{CSR} , then to reduce the influence of noise and remaining misalignment, the encoder downsamples I_{CSR} and obtains I_{CLR} , which is imported into the second decoder for final reconstruction. Otherwise, super-resolving tiny faces with face feature vectors (FFVFSR) [113] designs a face feature vector-based method that embeds face feature vector (generated from SphereFace [106]) into super-resolution model and discriminator to boost face reconstruction. Inception residual network (IRN) [63] builds an inception-residual block for face super-resolution and the work of Chen *et al.* [26] designs a generator based on attention mechanism. Recently, Chan *et al.* [21] took advantage of pre-trained StyleGAN [76] to do large factor super-resolution, which is called GLEAN. GLEAN builds an encoder to extract latent vector and multi-scale features of LR, then feed the features and vector into pre-trained subnetwork of StyleGAN to produce another multi-scale features for the decoder that is designed to generate final output.

Although good performance has been achieved, the original GAN involves several serious issues. First, the training of the original GAN is unstable. Second, generated images lack diversity. Third, the loss function fails to facilitate model training. To resolve these problems, wasserstein generative adversarial networks (WGAN) [3], an improved version of the original GAN, is proposed. Based on WGAN, face super-resolution through wasserstein GANs (WGANFSR) [29], has been used on numerous experiments on different adversarial loss and network architecture to verify the effectiveness of WGAN and WGAN-GP in face super-resolution. Except WGAN, boundary equilibrium generative adversarial networks (BEGAN) [9] is also an improved conditional GAN, based on which face super-resolution using conditional generative

adversarial networks called FCGAN [10] is developed. In FCGAN, face super-resolution is formulated as an HR face image generation problem conditioning on LR face images. Subsequently, improved versions of GAN were introduced, including enhanced discriminative generative adversarial network (EDGAN) [175], DBGAN [167], multi scale gradient GAN with capsule network as discriminator (MSG-CapsGAN) [116], improved super-resolution generative adversarial network (ISRGAN) [149], feature-preserving face super-resolution (FPFSR) [196] and others. Recently, Chen *et al.* [24] develops spatial attention residual network (SPARNet) which introduces spatial attention into the generator and takes advantage of multi-scale discriminator to improve image quality. In the wavelet coefficient domain, Huang *et al.* [60] also developed WaveletSRGAN for super-resolution based on WaveSRNet.

Without facial prior and degradation, HiFaceGAN is built [172, 171], which consists of a suppression module to select and integrate useful information, and a replenishment module that makes the best use of selected information for recovery details. To further explore GAN for face super-resolution, PCA-SRGAN [36] takes advantage of PCA to obtain the orthogonal space and split it into two orthogonal subspaces: W (contains the first n vectors), V (contains the remaining vectors), and then project normalized face (subtracted average face) into w space by projection matrix P_w . The projection will have increasing facial information by progressively increasing n . Based on this view, PCA-SPGAN imports progressively increasing projections of the face into the discriminator rather than the entire face image, which can reduce the learning difficulty of the discriminator and help the generator learn progressively.

The commonality of these types of GAN is that the discriminator outputs a single number to represent an image. However, Zhang *et al.* [189] assume that a single number is too fragile to represent a whole image, so they design supervised pixel-wise GAN (SPGAN) whose discriminator outputs a discriminative matrix with the same resolution as input images. SPGAN designs a supervised pixel-wise adversarial loss as follows:

$$\mathcal{L}_{\text{SPGAN(D)}} = - \sum_{i,j} [\min(0, \mathcal{D}_{\text{sp}}(I_{\text{LR}}, I_{\text{HR}})_{i,j} - |I_{\text{HR}} - I_{\text{SR}}|)], \quad (18)$$

$$\mathcal{L}_{\text{SPGAN(G)}} = \sum_{i,j} [\alpha_{i,j} * (\mathcal{D}_{\text{sp}}(I_{\text{LR}}, I_{\text{HR}})_{i,j})], \quad (19)$$

$$\alpha_{i,j} = \begin{cases} 1.0 & \mathcal{D}_{\text{sp}}(I_{\text{LR}}, I_{\text{HR}})_{i,j} \geq 0 \\ s & \mathcal{D}_{\text{sp}}(I_{\text{LR}}, I_{\text{HR}})_{i,j} < 0 \end{cases}, \quad (20)$$

$$\mathcal{D}_{\text{sp}}(I_{\text{LR}}, I_{\text{HR}}) = \mathcal{D}(I_{\text{LR}}, I_{\text{HR}}) - \mathcal{D}(I_{\text{LR}}, I_{\text{SR}}), \quad (21)$$

where $\mathcal{L}_{\text{SPGAN(G)}}$ presents loss function of the generator, $\mathcal{L}_{\text{SPGAN(D)}}$ denotes loss function of discriminator, \mathcal{D} indicates the function of the discriminator, $\mathcal{D}_{\text{sp}}(I_{\text{LR}}, I_{\text{HR}})$ is the final output of the pixel-wise discriminator, and s is used to control SPGAN's pattern. In general, the generator aims to minimize the \mathcal{D}_{sp} to force I_{SR} to be as similar as possible to I_{HR} , while the discriminator aims at maximizing \mathcal{D}_{sp} to distinguish I_{HR} and I_{SR} . However, they assume that when I_{SR} is very close to I_{HR} , it has no sense to maximize \mathcal{D}_{sp} for the discriminator. Therefore, they add $|I_{\text{HR}} - G(I_{\text{LR}})|$ to $\mathcal{L}_{\text{SPGAN(D)}}$.

Towards s , when $s = 1$, $\mathcal{L}_{\text{SPGAN(G)}}$ is equal to original \mathcal{L}_G (which is introduced in section 4.7), leading to an unsupervised pattern. When $s = 0$, if $\mathcal{D}(I_{\text{LR}}, I_{\text{HR}}) < \mathcal{D}(I_{\text{LR}}, I_{\text{SR}})$, then the gradient of $\mathcal{L}_{\text{SPGAN(G)}}$ is 0, a weakly supervised pattern. When $s = -1$, if $\mathcal{D}(I_{\text{LR}}, I_{\text{HR}}) < \mathcal{D}(I_{\text{LR}}, I_{\text{SR}})$, the generator forces $\mathcal{D}(I_{\text{LR}}, I_{\text{SR}})$ to be close to $\mathcal{D}(I_{\text{LR}}, I_{\text{HR}})$. That is, when $s = 1$, the generator forces I_{SR} to be close to I_{HR} , a strong supervised pattern.

Unsupervised Methods: The aforementioned methods are always trained on artificial LR-HR pair generated by simple downsampling method such as bicubic or bilinear downsampling (a few of them may add noise or blur). However, these models fail to recover good results from real-world LR. The quality of the real-world LR is affected by a wide range of factors such as the weather, leading to the unknown complicated degradation of real LR. The gap between real LR and artificial LR is wide and results in the decrease of face super-resolution performance in [47], and needs to be bridged. It is necessary to bridge this gap. To achieve this, many scholars have proposed unsupervised and self-supervised face super-resolution methods.

LRGAN [13] proposes learning the degradation before super-resolution from unpaired data. LRGAN designs a high-to-low GAN to learn real degradation process from unpaired LR-HR and create paired LR-HR for training Low-to-High GAN. Specifically, with HR as input, the high-to-low GAN generates LR (GLR) that should belong to real LR distribution and close to corresponding downsampled HR by average pooling. Then, for low-to-high GAN, GLR are fed into the generator to repair the SR results which have to be close to HR and match the real HR distribution. Goswami *et al.* [133] further developed a robust face super-resolution method.

Discrepancies between GLR in training phase and real LR in testing phase still exist, then characteristic regularisation (CR) [31] is formulated. Different from LRGAN that learns the degradation process to create real LR from HR and super-resolves in real LR space, CR transforms real LR into artificial LR and then super-resolves in artificial LR space. Based on CycleGAN, CR learns the mapping between real LR and artificial LR. Then, it uses the generated artificial LR from real LR (RALR) to fine-tune the super-resolution model, which is pre-trained by artificial pair, producing SR results. Thereafter, it feeds SR and HR into pre-trained face recognition network to extract features of SR and HR. To constrain the model, features of SR need to match those of HR and a downsampled version of SR should be also close to RALR. Then, Zheng *et al.* [195] used semi-dual optimal transport to guide model learning and develop semi-dual optimal transport CycleGAN (SDOT-CycleGAN).

The aforementioned methods take unsupervised learning to learn the map from real LR to real HR, while PULSE [118] is a self-supervised method. Instead of learning map between LR and HR, PULSE formulates face super-resolution as a generation problem without HR in which a given generative model generates high-quality face images (we also call them SR to maintain consistency) so that the downsampled version of SR is close to LR. To achieve the destination, the problem can be expressed as:

$$\min_G \|G(z) \downarrow_s - I_{\text{LR}}\|, \quad (22)$$

where z is the latent vector and input of pre-trained StyleGAN [76], s is the downsampling factor, G denotes the function of generator, and the downsampling is employed by bicubic interpolation. PULSE solves face super-resolution from a new perspective and inspires others. However, there are still many problems to be addressed, such as bicubic interpolation, which may not be the best degradation, and so on.

4.1.3 Reinforcement Learning-based Methods

Inspired by human perception process in which humans start from whole images and then explore a sequence of regions with an attention shifting mechanism, attention-aware face hallucination via deep reinforcement learning (Attention-FH) is proposed [18, 137]. Specifically, Attention-FH has two subnetworks: policy network (PN) that locates the region needed to be enhanced in current step t , and local enhancement network (LEN) that enhances the selected region. After T steps, the final SR results are

recovered, and the reward r_t can be described as:

$$r_t = \begin{cases} 0 & t < T \\ -\mathcal{L}_2 & t = T, \end{cases} \quad (23)$$

where \mathcal{L}_2 is MSE loss, which is introduced in section 4.7.

4.1.4 Ensemble Learning-based Methods

Different from above methods that only use a single model (e.g., CNN, GAN, and others), ensemble learning is used in adaptive-threshold-based multi-model fusion network (ATFMN) [70] for compressed face super-resolution. Considering that different models have various special advantages, ATFMN takes multiple deep learning networks (CNN, GAN, and RNN) to integrate their advantages together. Specifically, ATFMN uses three models (CNN-based, GAN-based, and RNN-based) combined with their designed dense block to generate candidate SR faces that are fed into corresponding attention subnetworks to learn the attention matrices, then attention matrices and candidate SR faces are fused to reconstruct SR results. In contrast to the preceding methods, ATFMN exploits the potential of ensemble learning for face super-resolution instead of focusing on a single model.

4.2 Prior-guided Face Super-resolution

General face super-resolution methods aim to design efficient networks for face super-resolution. Nevertheless, different from the general images, face images have face-specific information, for example, prior information which includes facial landmarks, facial parsing maps, and facial heatmaps. General face super-resolution methods ignore the facial prior information, generating face images with fuzzy facial structure. Therefore, to recover facial images with much clearer facial structure, researchers begins to develop prior-guided face super-resolution network. (Note that some methods uses a depth map or 3D facial prior to boost face restoration, because the depth map and 3D prior also provide facial structure information. Thus, we categorize them as prior-guided face super-resolution methods.)

4.2.1 Prior Embedded Framework

Prior information guided face super-resolution methods refer to methods that extract facial prior information (from LR, intermediate results or features, or final SR results) and utilize prior information to facilitate face reconstruction. Considering the order of prior information extraction and

face super-resolution, we further divide prior information guided face super-resolution methods into four parts: i) pre-prior methods that extract prior information followed by face super-resolution; ii) parallel-prior methods that extract prior information and face super-resolution simultaneously; iii) in-prior methods extract prior information from the intermediate results or features at middle of face super-resolution, and iv) post-prior methods that extract prior information from face super-resolution results. We illustrate the main frameworks of four categories in Fig. 6 and outline the development of prior-guided face super-resolution methods in Fig.7.

Pre-prior Methods: These methods first extract prior information and then face super-resolution. That is, they always extract prior information from LR by an extraction network which can be a pre-trained network or a subnetwork in the model, then take advantage of the prior information to facilitate face super-resolution. In the early stage, cascaded bi-network (CBN) [201] first super-resolves face images and combines a dense correspondence field progressively. At every scale, the network predicts the dense correspondence field which is then fed into the SR network for recovering intermediate super-resolved results. Based on intermediate results of former steps, the latter steps proceed. Extracting prior information directly from LR by pre-trained network is a straightforward solution to utilize facial prior information. Learning to hallucinate face images via component generation and enhancement (LCGE) [141] simply uses a pre-trained network to extract facial landmarks to divide facial components and feed the five components into different branches to recover corresponding components, thereby generating enhanced results of LR. Then, LCGE takes advantage of HR training images to transfer high-frequency detail to enhanced results of LR. The drawback of LCGE is that it only super-resolves LR in a local way but disregards global information. Jiang *et al.* [68] develop a two-step method deep CNN denoiser and multi-layer neighbor component embedding for Face Hallucination (MNCEF) that first recovers the global face images and then compensates missing details for every component (the partition of facial components is the same as that of LCGE). However, LCGE and MNCEF divide the face into components by accurate landmarks, which is usually unavailable when LR is tiny (*i.e.*, 16×16). Thus, researchers turn to facial parsing maps.

Parsing map guided multi-scale attention network (PMGMSAN) [147], consists of two subnetworks: ParsingNet and FishSRNet. In detail, PMGMSAN first pre-trains ParsingNet to extract the parsing map from LR, and then feeds the concatenation of the parsing map and LR into FishSRNet, and finally trains FishSRNet and finetunes ParsingNet to produce SR results. Chen *et al.* [25] propose a multi-scale progressively face restoration method named

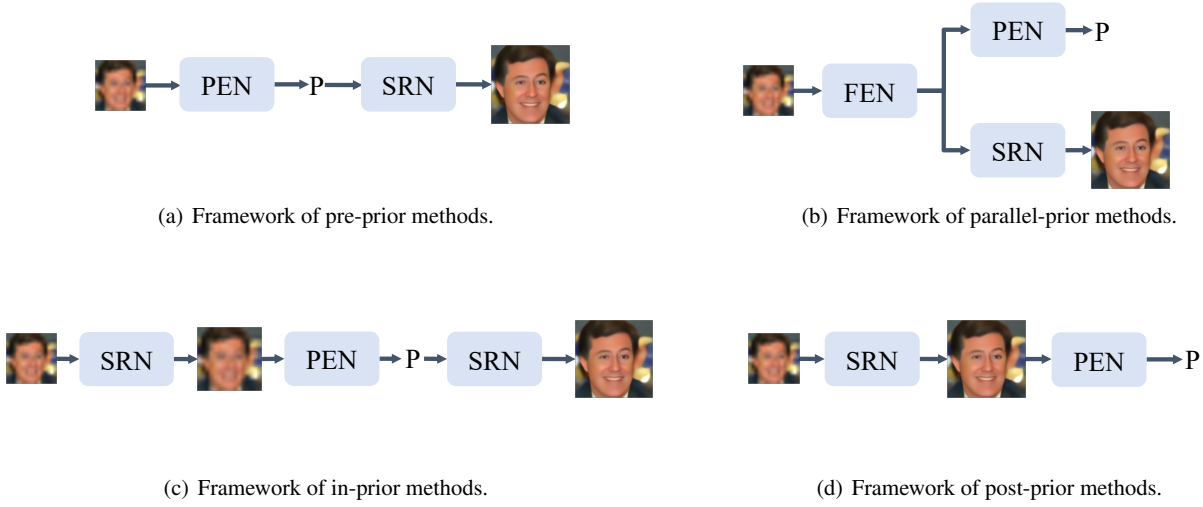


Fig. 6 Four frameworks of prior-guided face super-resolution methods. PEN is prior estimation network, SRN is super-resolution network, FEN is a feature extraction network, and P is prior information.

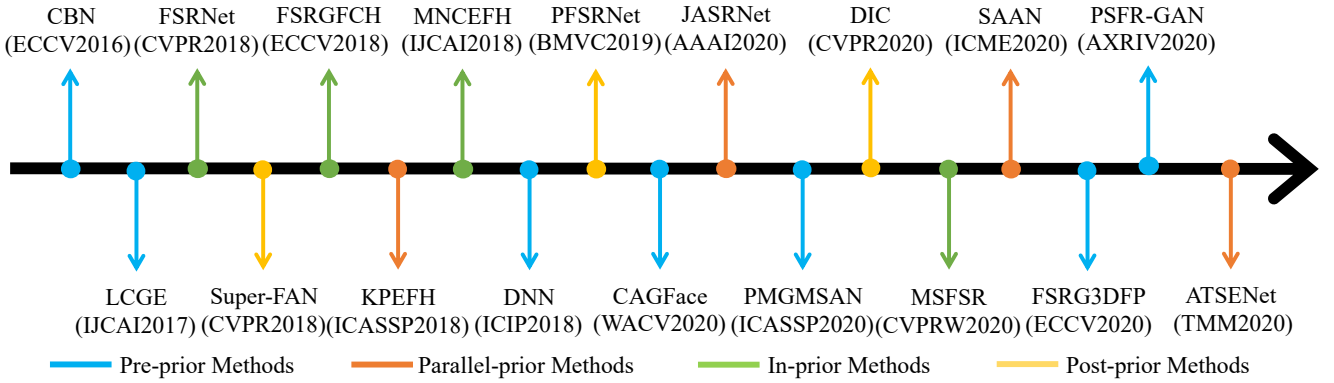


Fig. 7 Milestones of prior-guided face super-resolution methods. We simply list their names and venues.

PSFR-GAN. With constant vector C , LR and parsing map that is obtained from pre-trained extraction network as input, PSFR-GAN transfers semantic-aware style features from the concatenation of LR and parsing map to C by style transformation block (STB) in multi-scale way. Component semantic prior-guided generative adversarial network (CSPGAN) [104] also extracts facial prior from LR with the pre-trained model, and then utilizes prior to boost face reconstruction.

These methods take whole images and prior information as input, when LR is large (*i.e.*, 1024×1024), and the computational cost is expensive, which limits the application of the methods. Thus, component attention guided face super-resolution network (CAGFace) [72] designs a patch-based facial component-wise attention network. In common with PMGMSAN, CAGFace also develops a subnetwork to extract the parsing map. Instead of concatenation between LR and the parsing map,

CAGFace multiplies LR with this map, producing an attention mask which is concatenated with LR as input of SR subnetwork.

Instead of limiting to existing facial prior, Fan *et al.* [38] supposed that the depth map that provides geometric information also contributes to face reconstruction. Thus, they build a branch to learn the depth map from LR, and then import depth into the SR branch for improved face reconstruction. The prior used in above methods, including landmarks, heatmaps and parsing maps are all in 2D space, which may hinder network learning of 3D facial details. To solve this problem, face super-resolution guided by 3D facial priors (FSRG3DFP) [55] incorporates 3D priors into face super-resolution. In detail, FSRG3DFP builds a face reconstruction prior to generate 3D prior from LR, and a super-resolution model which takes advantage of 3D prior by spatial feature transforms blocks at different layers and

reconstructs SR results. Although the depth map and 3D prior are not contained in traditional facial prior, we also view them as a special prior.

Parallel-prior Methods: Although the above methods perform well, they ignore correlation between prior information estimation and face super-resolution. On the one hand, face prior estimation benefits from the enhancement of face resolution. On the other hand, much higher accurate prior estimated boosts face reconstruction. Earlier, key parts enhancement face hallucination (KPEFH) [90] built a common branch to extract features, thereafter feeds those features into five components branches which also consist of two branches for generating super-resolved faces and corresponding component mask (parsing map). Finally, KPEFH combines five different super-resolved results and component masks to obtain final super-resolved results. To take advantage of the correlation between prior estimation and face super-resolution, Yin *et al.* [177] developed joint alignment and super-resolution network (JASRNet). Specifically, JASRNet utilizes a shared decoder to extract features for super-resolution and prior estimation simultaneously. Through this design, a shared decoder can extract features with the amount of maximum information for both two tasks, which results in increased performance by 0.19 dB in Y channel. Later, semantic attention adaptation network (SAAN) [194] is proposed, consisting of two subnetworks: super-resolution network (SRN) and semantic parsing network (SPN). The SPN produces a parsing map of LR and parsing features from intermediate layers. To utilize the parsing features as facial prior, SAAN develops semantic attention adaption which takes the parsing features and LR features in SRN as input and generates adapted LR features that are further fed into the following layers of SRN. Finally, the reconstruction results are obtained. Facial semantic attribute transformation and self-attentive structure enhancement (ATSENet) [93] first takes advantage of attribute information to ensure right attribute (which is introduced in section 4.3), and then develops two parallel branches to generate facial prior and SR results.

In-prior Methods: Pre-prior methods directly extract prior information from LR, due to the low resolution of LR, detecting accurate prior information is difficult and challenging. To reduce the difficulty and improve the accuracy of prior estimation, researchers first coarsely enhanced LR and extracted prior information from enhanced results of LR (intermediate results or features).

Specifically, the first in-prior method is FSRNet [27], comprised of four subnetworks: coarse SR network to coarsely recover intermediate results, fine super-resolution encoder to extract informative features from intermediate results, prior estimation network that estimates prior information (including facial landmarks, heatmaps, and

parsing maps), and fine super-resolution decoder which constructs super-resolved faces with concatenation of features and prior information. The network is trained end-to-end with \mathcal{L}_2 loss between intermediate results or super-resolved results (estimated prior information) and HR (ground prior information). Analogously, Li *et al.* [92] designed a DNN that uses upsampling module to produce intermediate results, and then extract features and parsing maps from intermediate results, finally recovering the super-resolved results with a reconstruction module.

In contrast to extracting the prior from intermediate results, Face Super-resolution guided by facial component heatmaps (FSRGFCH) [179] embeds prior estimation branch (PEB) into super-resolution branch (SRB), and divides SR into two parts: former and latter. In the former SRB, LR is upsampled preliminarily. Without loss between intermediate features and HR, PEB directly estimates heatmaps from intermediate features for subsequent reconstruction. The latter concatenates intermediate features and heatmaps, further using heatmaps to repair high-quality face images. Compared with FSRNet, FSRGFCH extracts prior information from features instead of intermediate results, which may increase network flexibility. In contrast to prior information used in FSRNet and FSRGFCH, Multi-stage face super-resolution (MSFSR) [192] creates a novel facial prior (facial boundaries) and takes advantage of it to deal with LR progressively. For $8\times$ factor, three $2\times$ upsampling steps are adopted: coarse super-resolution followed by prior information extraction followed by fine super-resolution, which is applied on LR at every step.

Post-prior Methods: In contrast to pre- and in-prior methods, post-prior methods locate prior based on super-resolved results rather than LR or intermediate features.

One of the most popular post-prior methods is Super-FAN [12], which jointly addresses face super-resolution and facial landmarks detection tasks, including three parts: super-resolution network (SRN) which recovers HR face images, a discriminator which distinguishes between the real HR images and the super-resolved results generated by SRN, and the facial alignment network (FAN) which detects the facial landmarks from the super-resolved face images. Specifically, Super-FAN applies FAN to capture the facial structure information and constrain landmarks of super-resolved results and HR to be close. Similar to Super-FAN, progressive face super-resolution network (PFSRNet) [78] and fractal residual network [40] also employ prior information of super-resolved results and HR as loss constraint. However, PFSRNet is a progressive super-resolution model that generates multiple I_{SR} and

applies a unified FAN to acquire different-scale prior information.

Although Super-FAN achieves a performance breakthrough, it only applies prior information as a loss function, which fails to bring prior information into full play. To fully explore potential of prior information, deep-iterative-collaboration (DIC) [114] is proposed. Similar to Super-FAN, DIC also consists of a super-resolution model that generates SR results and a FAN model detecting SR prior information. In addition to using prior information as a loss function, DIC embeds prior information into a super-resolution model to boost reconstruction performance. Overall, DIC is a two-loop iteration architecture: i) inner loop, where a super-resolution model and FAN are both recurrent network, and ii) outer loop in which the recovery of LR and detection of landmarks is performed alternatively and iteratively. By iteratively collaborating between SR and prior estimation, the two tasks can promote each other.

4.2.2 Prior Fusion Methods

In addition to the order of prior extraction and face super-resolution, how to fuse prior information is also an important task. In this section, we introduce several commonly used prior information fusion methods.

Crop: Some methods only use prior information to crop facial components and deal with cropped facial components, *e.g.*, both LCGE [141] and MNCEFH [68] take advantage of landmarks to crop facial components.

Concatenation: Concatenation is the simplest and most direct way to fuse prior information into a super-resolution model, *e.g.*, CBN [201], FSRNet [27], PMGMSAN [147], FSRGFCH [179], and MSFSR [192]. In addition, a method multiplies LR with prior information to produce an attention map, and then feeds the concatenation of LR and the attention map into the following network, such as CAGFace [72].

Prior-based Loss: Another direct and commonly used way to utilize prior information is to design a prior-based loss function to constrain the super-resolution model. Specifically, post-prior methods always measure the distance between heatmaps or landmarks of SR and HR, which is called heatmap loss that is also contained in DIC [114], PFSRNet [78] (which applies heatmap loss to SR and HR at every scale), and others. The heatmap loss can be defined as:

$$\mathcal{L}_{\text{Heatmap}}(H_{\text{HR}}, H_{\text{SR}}) = \|H_{\text{HR}} - H_{\text{SR}}\|_F, \quad (24)$$

where $H_{\text{HR}}, H_{\text{SR}}$ are the heatmaps of HR and SR respectively, and F may be 1 or 2. Parsing map loss [90, 27], named as parsing map loss, is similar to heatmap loss, so we do not introduce it in detail.

Several researchers have proposed special loss function combined with prior information. PFSRNet [78] proposes facial attention loss based on heatmaps to enhance the adjacent area of the face as follows:

$$\mathcal{L}_{\text{Attention}} = M \cdot \|I_{\text{HR}} - I_{\text{SR}}\|_1, \quad (25)$$

where M is the attention map generated by selecting channel-wise max values of the ground truth heatmaps. To further enhance face details, PSFR-GAN [25] supposes semantic aware style loss that can improve FID [52] by 3.83:

$$\mathcal{L}_{\text{Semantic}} = \sum_{i=1}^5 \sum_{j=0}^{18} \|\mathcal{G}(\phi_i(I_{\text{SR}}), M_j) - \mathcal{G}(\phi_i(I_{\text{HR}}), M_j)\|_2, \quad (26)$$

where M_j is the j -th semantic label mask, ϕ means feature from the VGG [138] and \mathcal{G} is used to calculate Gram matrix.

Attention Fusion Module: Nevertheless, simple concatenation and prior-based loss cannot fully capture and explore prior information because different facial components contain different details and special information but they still treat these different facial components in the same way. To cope with this problem and maximize the use of prior information, attention fusion module (AFM) is developed in DIC [114]. First, based on facial five components, DIC groups heatmaps into five subsets. Then, DIC adds heatmaps in every group to produce the component map for every component, which is followed by softmax operation to activate and acquire attention map M_p . To recover facial components efficiently and pointedly, group convolutions are applied to generate individual feature F_p corresponding to different components. Finally, a fused feature is obtained by

$$\mathcal{F}_{\text{AFM}} = \sum_{p=1}^5 M_p \cdot F_p. \quad (27)$$

Compared with concatenation, AFM can integrate prior information better. Likewise, SAAN [194] also develops a semantic attention adoption (SAA) module to utilize prior information, and ATSENet builds a feature fusion unit (FFU) to combine prior information.

Spatial Feature Transform Inspired by SFT [165], FSRG3DFP employs a spatial feature transform block to efficiently use 3D prior formation for face reconstruction. For details, the 3D prior information are imported into three cascaded convolutional layers, extracting features F_{3D} which are fed into two parallel branches (two convolutional layers) to generate two transformation parameters μ, ν that transform features:

$$\mathcal{SFT}(F | \mu, \nu) = \mu F + \nu, \quad (28)$$

where \mathcal{SFT} denotes the function of SFT, and F is the input feature.

4.2.3 Comparisons of Prior-guided Face Super-resolution Methods

To have a global view of prior-guided face super-resolution methods, we compare the differences among all these methods from some aspects in Table 4. First, we list the prior information used in the methods, including commonly used landmarks, parsing maps and heatmaps, recently used depth prior [38], facial boundaries [192], and 3D prior [55]. Then, we show the pattern of prior information extraction, including joint extraction prior information and super-resolution, extraction prior information by pre-trained model. Finally, we represent the prior fusion pattern, for example, concatenation [201, 27, 192, 72, 179, 92, 147, 25], prior information based loss function [90, 177, 12] and special fusion modules [114, 55].

4.3 Attribute-constrained Face Super-resolution Methods

In addition to designing network architecture and utilizing prior information, facial attribute information is also exploited in face super-resolution, which is called attribute constrained face super-resolution methods. Facial attribute as a kind of semantic information provides semantic knowledge, *e.g.*, whether people wear glasses, which is useful to ensure the attributes of SR. The following are attribute-constrained face super-resolution methods. The overview is provided in Table 5.

4.3.1 Attribute Knowability

Different from prior information whose acquisition relies on the image itself, attribute information can be available but not acquire the existence of LR, such as in criminal cases where attribute information may not be clear in LR but accurately known by witnesses. Thus, some researchers construct network on condition that attribute information is given, while others relax this constraint. According to this concept, attribute constrained face super-resolution can be divided into two frameworks: given attribute methods and generated attribute methods, which are described as follows:

Given Attribute Methods: Based on the assumption that attribute information is available and given, many methods directly utilize attribute information to boost face repair. Given the attribute, how to integrate attribute into the super-resolution model is the key. For this problem, attribute-guided conditional CycleGAN (AGCycleGAN) [112] replicates attribute vector to match the size of LR, obtaining attribute maps that are then concatenated with LR as the input of the super-resolution model. The discriminator also takes SR and attribute as input to force super-resolution model to notice attribute

information. However, LR provides pixel-level information while attribute contains semantic-level information, thus direct concatenation of them is improper. Therefore, face super-resolution with supplementary attributes [180, 181] (also called FaceAttr) employs an encoder to encode LR into a high-level latent, which is then concatenated with attribute vector so that the super-resolution model incorporates with attribute information effectively. Moreover, FaceAttr also inserts attributes in the middle of a discriminator to prevent the super-resolution model from simply ignoring attributes. Overall, both AGCycleGAN and FaceAttr take advantage of concatenation to converge attribute information. However, Lee *et al.* [86] hold that information contained in LR and information is dissimilar, and direct concatenation is unsuitable for use and may hinder the performance. With regard to this view, Lee *et al.* construct attribute augmented convolutional neural network (AACNN), including a feature extraction network, a super-resolution model, and a discriminator. Different from FaceAttr and AGCycleGAN, the feature extraction module utilizes two branches to extract features from two different domains (LR and attribute), and then fuses them into a common branch by concatenation. To inject guidance into the super-resolution model, feature extraction module offers features in a common branch to the super-resolution model at multiple scales and assists it to learn information from different domains. Different from these methods, attribute transfer and enhancement network (ATENet) [91] and facial semantic attribute transformation and self-attentive structure enhancement (ATSENet) [93] (an improved version of ATENet) develops an attribute transfer network to first upsample LR and then fuses attribute information with the upsampled LR feature by concatenation, thereby generating upsampled LR with rational attribute, and then building an enhancement network for further detail enhancement.

Generated Attribute Methods: However, given attribute methods work on the condition that the attribute is available and given, making them limited in real-world scenes. Although the attribute vector can be set as unknown, such as 0 or random values, the performance may drop sharply. Many researchers build modules to generate attribute information for face super-resolution. Residual attribute attention network (RAAN) is built [161], which is based on cascaded residual attribute attention blocks (RAAB) that consists of a feature extraction and three branches including shape generation network, texture prediction network and attribute attention network generating shape, texture and attribute information respectively, and then fuses different level information for deep blocks. RAAN embeds generation and utilization of attribute information into the basic block, while facial attribute capsules network (FACN) [162] integrates

Table 4 Comparison of prior-guided face SR algorithms. Prior means the prior information used in the method. Extraction shows how the prior information is extracted, *i.e.*, extracting from pre-trained model (pre-trained), embedding subnetwork into the model and joint super-resolution and extraction (joint), or first pre-training the subnetwork then joint. Prior fusion lists the fusion pattern used in the method.

Frameworks	Methods	Prior	Extraction	Prior fusion
Pre-prior methods	CBN	Dense correspondence field	Joint	Concatenation
	LCGE [141]	Landmark	Pre-trai	Crop
	MNCEFH [68]	Landmark	Pre-trained	Crop
	PMGMSAN [147]	Parsing map	Pre-trained and then joint	Concatenation
	PSFR-GAN[25]	Parsing map	Pre-trained	Concatenation
	CAGFace [72]	Parsing map	Pre-trained	Concatenation
	AGFFSR [38]	Depth	Joint	Summation
	FSRG3DFP [55]	3D prior	Joint	SFT
Parallel-prior methods	KPEFH [90]	Parsing map	Joint	$\mathcal{L}_{\text{Parsing}}$
	JASRNet [177]	Heatmap	Joint	$\mathcal{L}_{\text{Heatmap}}$
	SAAN [194]	Parsing map	Joint	SAA
	ATSENet [93]	Facial boundary heatmaps	Joint	FFU
In-prior methods	FSRNet [27]	Landmark, parsing map, heatmaps	Joint	Concatenation
	DNN [92]	Parsing map	Joint	Concatenation
	FSRFCH [179]	Heatmaps	Joint	Concatenation
	MSFSR [192]	Facial boundaries	Joint	Concatenation
Post-prior methods	Super-FAN [12]	Heatmaps	Joint	$\mathcal{L}_{\text{Heatmap}}$
	PFSRNet [78]	Heatmaps	Pre-trained	$\mathcal{L}_{\text{Heatmap}}, \mathcal{L}_{\text{attention}}$
	DIC [114]	Heatmaps	Joint	$\mathcal{L}_{\text{Heatmap}}, \text{AFM}$

attributes in capsules. Specifically, FACN first encodes LR into encoded features which is fed into a capsule generation block that produces semantic capsules, probabilistic capsules, and facial attributes to represent the facial feature. With the combination of semantic capsules, probabilistic capsules and facial attributes as input, the decoder of FACN recovers the final SR results.

Concatenation: Similar to prior information, concatenation is the simplest way. From a macro perspective, all given attribute methods adopt concatenation to fuse attribute information. However, from a micro perspective, attribute information is concatenated in

different ways in these methods. On the one hand, every attribute vector is resized as an attribute map by replicating the value in AGCycleGAN [112]. Nevertheless, the attribute vector is directly used in FaceAttr [180,181], AACNN [86] and ATENet [91]. On the other hand, AGCycleGAN [112] concatenates the attributes with LR, while FaceAttr [180,181] concatenates attribute with encoded LR features. However, AACNN [86] concatenates features extracted from LR and attributes to obtain fused information, and then further concatenates features of fused information with LR features at different scales. We show the various concatenations of different models in Fig. 8.

Table 5 State-of-the-art attribute-constrained face super-resolution methods.

Methods	Algorithms	Network architectures
FaceAttr [180] (CVPR2018)	Concatenating attributes and feature of LR to embed attribute information for face super-resolution	Autoencoder
FaceAttr [181] (TPAMI2020)	Concatenating attributes and feature of LR to embed attribute information for face super-resolution	Autoencoder
AGCycleGAN [112] (ECCV2018)	Feeding the concatenation of LR and attributes into conditional CycleGAN to learn attribute-guided face super-resolution	CycleGAN
AACNN [86] (CVPRW2018)	Designing feature extraction network to extract features from LR and attributes, and Exploring attribute augmented convolutional neural network for face super-resolution	Feature extraction network
ATENet [91] (ICME2019)	Upsampling LR into high resolution space, and then embedding auto-encoder to combine attributes and features of LR for recovering high quality faces	Transfer network and enhancement network
RAAN [161] (AAAI2019)	Designing cascaded residual attribute attention block which can predict attributes and generate attribute attention for face super-resolution	Cascaded residual attribute attention blocks (RAAB)
FACN [162] (AAAI2020)	Embedding attribute into facial attribute capsules for face super resolution	Encoder, capsules network and decoder
ATSENet [93] (TMM2020)	Learning facial semantic attribute transformation, and then exploring facial prior information for face super-resolution	Attribute transformation network and structure enhancement network

Attribute Attention: However, simple concatenation cannot fully explore attribute information. Thus, the attention mechanism has attracted many researchers. Both RAAN [161] and FACN [162] build their unique attribute attention mechanism. RAAN [161] feeds generated attributes into two fully connected layers to obtain the attribute channel attention that is used to combine with other information (*e.g.*, shape information and texture information) through multiplication along the channel dimension. In contrast to RAAN [161] building attribute channel attention, FACN [162] takes attributes as a kind of mask in two ways. In detail, FACN [162] multiplies the output of semantic capsules with attributes to activate and select it. At the same time, the sum of the output of probabilistic capsules and attribute mask is used as an updated output version of probabilistic capsules.

Attribute-based Loss: Similar to prior-based loss, attribute-based loss is still designed by many methods. Some of them force the network to predict attribute information correctly, *e.g.*, RAAN [161], ATENet [91], and FACN [162], which can be described as

$$\mathcal{L}_{\text{Attribute}} = \|A_P - A_{GT}\|, \quad (29)$$

where A_P is predicted attribute and A_{GT} is the ground truth attribute. Some of them constrain that super-resolved results should have the same attribute as HR, *e.g.* FRN [40],

$$\mathcal{L}_{\text{Attribute}} = A_{GT} \log(A_P) + (1 - A_{GT}) \log(1 - A_P). \quad (30)$$

The others require a given attribute to match the given face images, *e.g.* AGCycleGAN [112] and FaceAttr [180, 181],

$$\begin{aligned} \mathcal{L}_{\text{Attribute}_D} = & -\log D(I_{GT}, A_{GT}) - \log(1 - D(I_{SR}, A)) \\ & - \log(1 - D(I_{GT}, \tilde{A})), \end{aligned} \quad (31)$$

where A is attribute matched with I_{GT} while \tilde{A} is a mismatched one. In this pattern, the loss can force the network to generate real face images that match a given attribute.

4.3.2 Comparisons of Attribute-constrained Face Super-resolution Methods

We compare the differences between all attribute-constrained face super-resolution methods and show them in Table. 6. In attribute-constrained face super-resolution task, an important task is to obtain

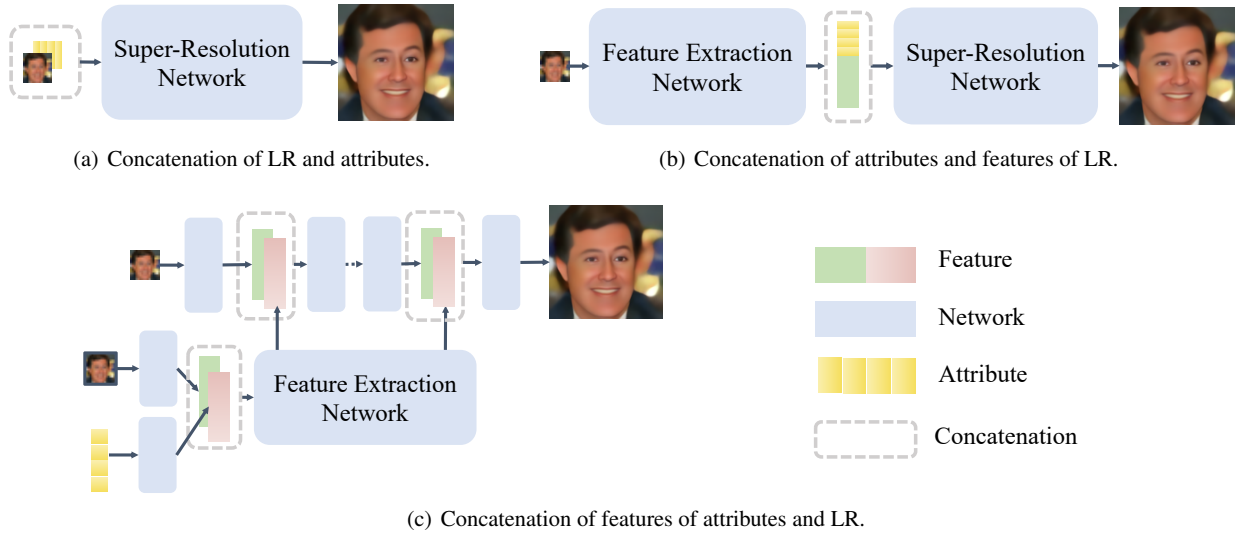


Fig. 8 Different concatenations in given attribute methods.

attribute information and tap the potential of attribute. First, we list the number of used attributes in every model. Then, we turn to the availability of attributes. On the one hand, the models suppose that attributes are given, such as FaceAttr [180,181], AGCycleGAN [112], and AACNN [86]. On the other hand, the attributes can be extracted and learned by designing a subnetwork, *e.g.*, RAAN [161] and FACN [162]. No matter how the attributes are obtained, the utilization is equally crucial. From the preceding review, we can know given attribute method AGCycleGAN [112] directly concatenates LR with attribute, FaceAttr [180,181] uses features of LR instead of LR, and AACNN [86] selects features of LR and attributes to perform concatenation. Without the known attribute, RAAN [161] and FACN [162] employ attention module with channel attention of attribute, summation and multiplication operation on attribute and other information.

4.4 Identity-preserving Face Super-resolution

Compared with prior and attribute information providing facial details, identity information containing identity-aware details is similarly important and identity preserving face super-resolution methods have received an increasing amount of attention in recent years. These methods aim to maintain the identity consistency between SR and LR. They also hope to improve the performance of face recognition. We show the overview of identity-preserving face super-resolution methods in Table 7. The following are methods that combine the identity information.

4.4.1 Identity Preserving Pattern

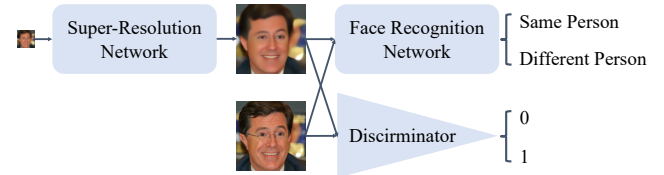


Fig. 9 Framework of method defining identity-based loss.

Simple Identity Loss: The introduction of identity information aims to maintain identity consistency between SR and HR. To this end, a commonly used design is to define identity-based loss, *e.g.*, deep joint face hallucination and recognition (DJFHR) [159], super-identity convolutional neural network (SICNN) [188], face hallucination generative adversarial network (FH-GAN) [8], optimizing super resolution for face recognition (FRSRResNet) [1], WaveSRGAN [60], low-resolution face recognition based on identity-preserved face hallucination (SRNet) [82], identity-preserving face hallucination via deep reinforcement learning (IPFH) [30] and identity-preserving face super-resolution (IPFSR) [16], ID preserving face super-resolution generative adversarial networks [89], cascaded super-resolution and identity priors (C-SRIP) [48] and facial semantic attribute transformation and self-attentive structure enhancement [93]. The frameworks of these methods consist of two main components: super-resolution model that generates SR results, and pre-trained face recognition

Table 6 Comparisons of attribute-constrained face SR algorithms. Number of attributes is the number of used attributes, and the attribute embedding method shows how to take advantage of attribute information. Given denotes whether attribute information is given, \checkmark means attribute is provided, and \times indicates that the attribute is not given and needs to be generated. A represents the attribute, and F_{LR} , F_A denote the features extracted from LR and attribute, respectively. F_{LR^\uparrow} is a feature of upsampled version of LR.

Frameworks	Methods	Number of attributes	Given	Attribute embedding methods
Given attribute methods	FaceAttr [180, 181]	18	\checkmark	Concatenation of A and F_{LR}
	AGCycleGAN [112]	18	\checkmark	Concatenation of A and LR
	AACNN [86]	38	\checkmark	Concatenation of F_{LR} and F_A
	ATENet [91]	-	\checkmark	Concatenation of F_{LR^\uparrow} and A
	ATSENet [93]	-	\checkmark	Concatenation of F_{LR^\uparrow} and A
Generated attribute methods	RAAN [161]	-	\times	Attribute channel attention
	FACN [162]	18	\times	Attribute attention mask

Table 7 State-of-the-art identity-preserving face super-resolution methods.

Methods	Algorithms	Network architectures
SICNN [188] (ECCV2018)	Designing super-resolution network followed by face recognition network and developing identity loss for face hallucination	Super-resolution network and face recognition network
SiGAN [53] (TIP2019)	Utilizing contrastive property that different LR corresponds to different identities to maintain identity consistency	Twin generator and discriminator
FH-GAN [8] (Arxiv2019)	Joint performing face super-resolution and face recognition to achieve identity-preserving face super-resolution	Face recognition network
WaveSRGAN [60] (IJCV2019)	Transforming image domain into wavelet domain for multi-scale face super-resolution and designing identity loss based on face recognition network to achieve identity-preserving	WaveletSRNet and face recognition network
IPFH [30] (TCSVT2019)	Recovering a coarse face with global network, and then using deep reinforcement learning and face recognition network to preserve facial identity	Face hallucination and face recognition networks
IADFH [77] (BTAS2019)	Feeding multi-scale super-resolved results of different LR into discriminator and using adversarial face verification (different LR corresponds to different identities) to achieve identity-aware deep face super-resolution	Three-way classification discriminator
C-SRIP [48] (TIP2020)	Recovering multi-scale faces progressively and designing multi-scale face recognition networks to explore identity priors	Multi-scale face recognition network blocks
SPGAN [189] (TMM2020)	Designing supervised pixel-wise GAN with identity embedding discriminator for face super-resolution	Identity embedding discriminator

network (FRN) for identity preserving, probably an additional discriminator. The pipeline is shown in Fig. 9. The super-resolution model super-resolves LR, generating I_{SR} which is fed into FRN to obtain its identity features. (Note that WaveSRGAN is based on WaveletSRNet, and the input of its discriminator is in wavelet domain instead of image.) At the same time, I_{HR} is also passed into FRN, acquiring its corresponding identity features. The identity-based loss is calculated by

$$\mathcal{L}_{\text{Identity}} = \|\text{FRN}(I_{HR}) - \text{FRN}(I_{SR})\|_F. \quad (32)$$

Towards F , when F is 1, the loss is used in WaveSRGAN [60] and IPFSR [16], F is 2 in FH-GAN [8] and FRSRResNet [1], or

$$\mathcal{L}_{\text{Identity}} = \left\| \frac{\text{FRN}(I_{HR})}{\|\text{FRN}(I_{HR})\|_2} - \frac{\text{FRN}(I_{SR})}{\|\text{FRN}(I_{SR})\|_2} \right\|_2, \quad (33)$$

in SICNN [188], and

$$\mathcal{L}_{\text{Identity}} = 1 - \frac{\text{FRN}(I_{SR}) \cdot \text{FRN}(I_{HR})}{\|\text{FRN}(I_{SR})\|_2 \|\text{FRN}(I_{HR})\|_2}, \quad (34)$$

in SRNet [82], and A-softmax loss in [106, 30].

Rather than directly extracting identity features from I_{SR} and I_{HR} , C-SRIP [48] feeds residual maps between I_{HR} or I_{SR} and I_{LR}^{\uparrow} (upsampled by bicubic interpolation) respectively into FRN (which is pre-trained with a residual map.), and applies cross-entropy loss on them. Moreover, C-SRIP generates multi-scale face images which are fed into different scale face recognition network.

Identity Embedding Discriminator: Building a face recognition network-based identity loss cannot fully explore identity information, as features in pre-trained FRN with informative facial identity prior are ignored. To cope with this limitation, based on SPGAN [189], an identity-based pixel-wise discriminator (IPD) that introduces facial identity prior from FRN is constructed. Specifically, a face recognition network (FRN) with the same architecture as IPD is designed to provide identity information for the discriminator at different locations. The way to provide identity information is to concatenate low-, middle-, and high-level features extracted from pre-trained FRN with corresponding features in the discriminator, as shown in Fig. 10.

In addition, SPGAN [189] designs attention-based identity GAN loss to further utilize identity prior as follows:

$$p^D = -\min(0, \mathcal{D}_{\text{IPD}}(I_{LR}, I_{HR}) - |I_{HR} - I_{SR}|), \quad (35)$$

$$p^G = \alpha * D_{\text{SP}}(I_{LR}, I_{SR}) + b, \quad (36)$$

$$\mathcal{L}_{\text{Identity(D)}} = \|\text{FRN}(I_{HR}) - \text{FRN}(I_{SR})\|_2 \otimes p^D, \quad (37)$$

$$\mathcal{L}_{\text{Identity(G)}} = \|\text{FRN}(I_{HR}) - \text{FRN}(I_{SR})\|_2 \otimes p^G, \quad (38)$$

where FRN denotes the function of FRN, $\mathcal{L}_{\text{Identity(G)}}$ and $\mathcal{L}_{\text{Identity(D)}}$ are loss functions of the generator and discriminator respectively, \mathcal{D}_{sp} is the final output of SPGAN (which has been introduced in section 4.1.2), \mathcal{D}_{IPD} denotes the function of IPD, b is weight, and $p_{i,j}$ denotes the element of p in the i -th row and j -th column.

Pairwise Data Contrastive Loss: Identity loss is based on pre-trained FRN, which needs well-labeled datasets. However, a large well-labeled dataset is very costly. To solve the problem of costly labeled datasets, one solution is to deal with weak-labeled instead of well-labeled datasets. In consideration of this, siamese generative adversarial network (SiGAN) [53] takes advantage of weak pairwise label (in which different LR corresponds to different identities) to achieve identity preservation. Specifically, SiGAN has twin GAN that share the same architecture but super-resolve different LR at the same time. As identities of different LR are different, the identities of super-resolved results corresponding to LR are also vary. Based on this observation, SiGAN designs an identity-preserving contrastive loss that intends to minimize the difference between same-identity pairs and maximize the difference between different-identity pairs. The contrastive loss is defined as

$$\mathcal{L}_{\text{Contrastive}} = (1 - y) \frac{1}{2} [\max(0, 0.5 - E_w)]^2 + y \frac{1}{2} (E_w)^2, \quad (39)$$

$$E_w = \|P(I_{SR}^1), P(I_{SR}^2)\|_1, \quad (40)$$

where P is a 128-neuron fully connected layer that connects with the end of the second resblock in G_1 and G_2 , y is 0 when recovering two SR face images that belong to the same identity, and y is set to 1 when the images belong to different identities.

Compared with SiGAN that designs twin generators for pair data, identity-aware deep face hallucination (IADFH) [77] feeds pair data into the discriminator. Its discriminator is a three-way classifier that generates fake, genuine and imposter. When the HR face of one identity and SR of the downsampled version of the same face (pair₁) or HR face of one identity and SR of the downsampled version of another identity face (pair₂) are imported into the discriminator, the output is fake, while two different HR with the same identity (pair₃) correspond to genuine, two HR with different identities (pair₄) correspond to imposter. In this pattern, the generator

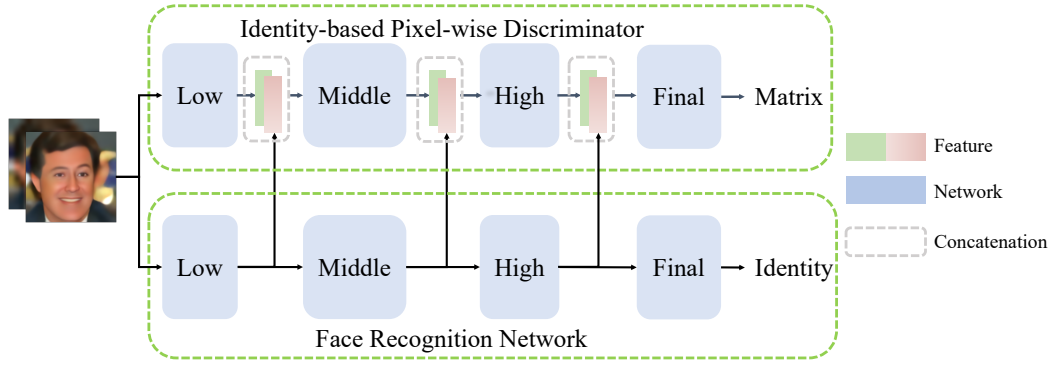


Fig. 10 Architecture of IPD.

considers identity information. The loss is called adversarial face verification loss (AFVL),

$$\mathcal{L}_{\text{AFVL(D)}}(\text{pair}) = \log d_f(\text{pair}_1) + \log d_f(\text{pair}_2) + \log d_{\text{gen}}(\text{pair}_3) + \log d_{\text{imp}}(\text{pair}_4), \quad (41)$$

$$\mathcal{L}_{\text{AFVL(G)}}(\text{pair}, D) = \log d_{\text{gen}}(\text{pair}_1) + \log d_{\text{imp}}(\text{pair}_2), \quad (42)$$

where $\mathcal{L}_{\text{AFVL(D)}}$ and $\mathcal{L}_{\text{AFVL(G)}}$ are loss functions for discriminator and generator respectively, and $d_f, d_{\text{gen}}, d_{\text{imp}}$ correspond to the output of the discriminator for fake, genuine and imposter pair.

4.4.2 Comparisons of Identity-preserving Face Super-resolution Methods

We compare the differences among all identity-preserving face super-resolution methods in Table 8. In the Table 8, we present whether the network includes FRN, and how these methods take advantage of identity information in detail. In essence, SICNN [188], FH-GAN [8], IPFSR [16] all take advantage of the distance between $\text{FRN}(I_{\text{HR}})$ and $\text{FRN}(I_{\text{SR}})$ as identity loss, although some differences between them, *e.g.*, MSE loss or \mathcal{L}_1 loss. Compared with them, C-SRIP's identity loss measures cross entropy distance between features of the residual map. SPGAN [189] proposes its unique identity-based pixel-wise discriminator, while SiGAN [53] and IADFH [77] utilize the fact that different LRs have different identities to construct pair-wise contrastive loss.

4.5 Reference Face Super-resolution

The face super-resolution networks discussed all exploit the face itself information. However, in some conditions, we may obtain the high-quality face image of the same identity as the LR, for example, the person of LR we want to recover may have other high-quality face images. In view

of this condition, reference face super-resolution methods are proposed, and we present an overview of these methods in Table 9. Different from single face super-resolution methods which recover corresponding HR face image from a LR face image, the problem is formulated as that of a high-quality reference (R) face image of the same identity as LR, which is used to guide the LR reconstruction. Out of the same identity of LR and R, R can provide many identity-aware face details to the network, which can contribute to the recovery of LR. The following methods are reference face super-resolution methods. In general, reference face super-resolution methods recover HR face from given LR and reference face image R . Obviously, R can be only one image, or multiple images, according to the number of R , a guided framework can be partitioned into single-face guided, multi-face guided, and dictionary-guided methods.

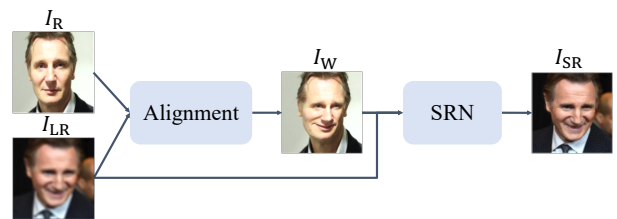


Fig. 11 Framework of single-face guided methods.

4.5.1 Single-face Guided Methods

At first, a high-quality face image which shares the same identity of LR serves as R , such as guided face restoration network (GFRNet) [97], guidance warper adversarial-loss, identity-loss network (GWAInet) [34] and guided cascaded super-resolution network (GCFSRNet) [17]. Since R and LR may have different poses and expressions, which may hinder the recovery of face images, single-face guided

Table 8 Comparison of identity-preserving face SR algorithms. FRN indicates whether the method needs a face recognition network. Specific method list the how identity information is used in different networks. Notably, I_{RH} is the residual map between I_{HR} and $I_{LR}^{\uparrow s}$, and I_{RS} is the residual map between I_{SR} and $I_{LR}^{\uparrow s}$.

Identity preserving pattern	Methods	FRN	Specific methods
Simple identity loss	SICNN [188]	✓	MSE loss on normalized $FRN(I_{SR})$ and $FRN(I_{HR})$
	FH-GAN [8]	✓	MSE loss on $FRN(I_{SR})$ and $FRN(I_{HR})$
	FRSRResNet [1]	✓	MSE loss on $FRN(I_{SR})$ and $FRN(I_{HR})$
	IPFSR [16]	✓	\mathcal{L}_1 loss on $FRN(I_{SR})$ and $FRN(I_{HR})$
	WaveSRGAN [60]	✓	\mathcal{L}_1 loss on $FRN(I_{SR})$ and $FRN(I_{HR})$
	C-SRIP [48]	✓	Cross entropy loss on $FRN(I_{RS})$ and $FRN(I_{RH})$
	IPFH [30]	✓	A-softmax loss on $FRN(I_{SR})$ and $FRN(I_{HR})$
Identity embedding discriminator	SPGAN [189]	✓	Identity-based discriminator
Pairwise data contrastive loss	SiGAN [53]	×	Pair contrastive loss on pair I_{SR}
	IADFH [77]	×	Adversarial face verification loss

Table 9 State-of-the-art reference face super-resolution methods.

Methods	Algorithms	Network architectures
GFRNet [97] (ECCV2018)	Warping reference to get aligned guidance, and then feeding the concatenation of aligned guidance and LR into super-resolution network to recover faces	WarpNet and RecNet
GWAInet [34] (CVPRW2018)	Warping reference without landmarks, then extracting features of warped reference to boost face super-resolution	GFENet and SRNet
JSRFC [95] (ICTAI2019)	Recovering a coarse face and warping the reference, and then learning similar facial components from multiple references	Face alignment network, parsing network
ASFFNet [96] (CVPR2020)	Selecting the most similar reference from multiple references and designing adaptive spatial feature fusion to fuse reference and LR for face super-resolution	Feature extraction, adaptive feature fusion and reconstruction
MEFSR [148] (ACCV2020)	Fusing multiple exemplars and LR, and then extracting similar features for face super-resolution	Weighted pixel average module
DFDNet [94] (ECCV2020)	Learning deep multi-scale component dictionaries and utilizing component dictionaries to guide face super-resolution	Transfer and enhancement networks

methods tend to perform the alignment between R and LR. After alignment, both LR and aligned R (we name it I_w) are fed into a reconstruction network to repair the SR. The pipeline of this category is shown in Fig. 11. The differences between GFRNet and GWAInet include two aspects: i) GFRNet realizes alignment by landmarks which is difficult to obtain from LR, while GWAInet employs flow field to finish the alignment; ii) in the reconstruction network, GFRNet directly concatenates LR and I_w as the input. Nevertheless, GWAInet believes that direct concatenation cannot bring I_w into play completely, thereby building a GFENet to extract features from I_w and transferring useful features of I_w to the reconstruction network to recover SR. Reference-based face super-resolution (RefSR-Face) [108] consists of an encoder to learn latent vectors of joint LR and R, a sampling generator that samples from the latent space, and a decoder that recovers face images from latent vectors.

4.5.2 Multi-face Guided Methods

Single-face guided methods set the problem as a LR only has one high-quality reference face image of the same person, but many high-quality face images are available in reality, and multiple high-quality guidance faces can further decrease the difficulty of face super-resolution, so why not exploit the potential of many high-quality reference images? Thus, multi-face guided face methods are proposed.

Adaptive spatial feature fusion network (ASFFNet) [96] is the first to explore multi-face guided face super-resolution methods. Instead of limiting to only one R, ASFFNet is constructed based on multiple high-quality guidance faces. ASFFNet simply selects the best reference image with the most similar pose and expression from all high-quality images of the same person for a LR finished by Guidance Selection Module. However, misalignment remains between R and LR, and thus, alignment is necessary. Meanwhile, R and LR always have different illumination backgrounds, which also have to be considered. To cope with these two problems, ASFFNet applies weighted least-square alignment [134] and AdaIN [62] for alignment and illumination translation on extracted F_R, F_{LR}, F_L by feeding R, LR and landmark of LR into three feature extraction subnetworks respectively, generating aligned and illumination consistent features $F_{R, w, a}$. After eliminating the illumination difference and misalignment, adaptive feature fusion block (AFFB) is proposed to combine $F_{R, w, a}, F_{LR}$ and F_L , and the process can be expressed as:

$$\begin{aligned} F_{\text{AFFB}}(F_{LR}, F_{R, w, a}) &= (1 - F_m) * C_{LR}(F_{LR}) + F_m * C_G(F_{R, w, a}) \\ &= C_{LR}(F_{LR}) + F_m * (C_R(F_{G, w, a}) - C_{LR}(F_{LR})), \end{aligned} \quad (43)$$

where C_{LR}, C_R is the convolutional layer with F_{LR} and F_R as input respectively, $*$ is the element-wise product, and the F_m is the attention mask generated by the convolutional layer taking the concatenation of F_{LR}, F_G and F_L , which aims to complement the information from G and LR. After several AFFBs, the features are fed into the reconstruction module to recover the final SR.

Similar to ASFFNet [96], multiple exemplar face super-resolution (MEFSR) [148] also assumes that there are a set of high-quality reference faces (R) of the same people to provide high-frequency information. Instead of aligning R with LR, MEFSR designs the weighted pixel average (PWAVE) module to extract a good combination of feature maps from the concatenation of all reference faces, and then adds the combination information from PWAVE to intermediate features in a super-resolution model at two different scales for fusing high-frequency information.

Different from ASFFNet and MEFSR that need reference images to have the same identity with LR, joint super-resolution and face composite (JSRFC) [95] only requires that reference images have similar components with LR (every reference face image is labeled with a vector to indicate which components are similar.). Specifically, JSRFC first builds a coarse super-resolution network for recovering global faces, then in order to enhance details, it warps reference images and extracts similar facial components from the warped reference images by a parsing map, which is fed into a refining network to refine coarse SR results.

4.5.3 Dictionary-guided Methods

Although single-face guided methods and multi-face guided methods perform excellently, guidance faces meeting identity constraints are also unknown when the identity of LR is unavailable, thereby resulting in unrealistic reference face methods. Li *et al.* [94] observed

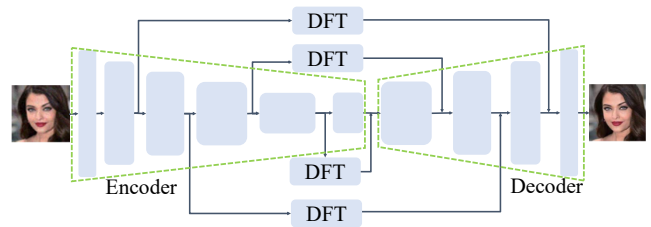


Fig. 12 The framework of DFDNet.

that different people can also have similar facial components, so reference images should not be limited to the same identity. According to this observation, they constructed a Deep Face Dictionary Network (DFDNet)

including two steps: i) component dictionary generation and ii) face reconstruction with component dictionary. They used pre-trained VGGFace [19] to extract features in different scales from high-quality faces, then crop and resample four components with landmarks, then cluster K classes for every component by K-means, generating $\text{RDic}_{s,c}^k$ where s is scale and c denotes component.

Given component dictionaries, DFDNet recovers SR results with dictionaries. DFDNet is an autoencoder with pre-trained VGGFace as encoder to keep feature consistency, and inserts several dictionary feature transfer (DFT) for transferring features from dictionaries to LR, as shown in Fig. 12. In DFT, facial components are cropped and then AdaIn is applied to eliminate illumination difference between facial components of every cluster and LR. Thereafter, an inner product is adopted to select the most similar cluster $\text{RDic}_{s,c}^*$ with $F_{s,c}^{\text{LR}}$. Then, transferred features are obtained by

$$\hat{F}_{s,c} = F_{s,c}^{\text{LR}} + \text{RDic}_{s,c}^* * F_{\text{Conf}}(\text{RDic}_{s,c}^* - F_{s,c}^{\text{LR}}), \quad (44)$$

where F_{Conf} is the block that generates confidence score producing a confidence score representing the difference degree between $F_{s,c}^{\text{LR}}$ and $\text{RDic}_{s,c}^*$. Then, inverse RoIAlign is used to put back facial components and feed transferred features into two parallel subnetworks, generating scale α and shift β similar to SFT [165] for the decoder to learn details for reconstruction. Through DFT at different scales, the details from dictionaries can be progressively transferred to degraded input.

4.5.4 Comparisons of Reference Face Super-resolution Methods

We compare the aforementioned guided face super-resolution methods in several aspects shown in Table 10. In reference face super-resolution task, the main pending problems include the following points: i) Does reference face image need to be of the same identity as LR? First, GFRNet [97], GWAInet [34], MEFSR [148], and ASFFNet [96] require reference images of the same people, while DFDNet [94] relaxes the limitation and does not need to keep identity consistency. ii) How many reference images are used? From GFRNet [97], GWAInet [34] to ASFFNet [96] and MEFSR [148], to DFDNet [94], the number increases progressively. iii) How can the misalignment between reference images and LR be addressed? To solve this problem, GFRNet [97] chooses flow field based on landmark to align, while GWAInet [34] also uses flow field but does not require landmark, ASFFNet [96] takes advantage of moving least-square. MEFSR [148] and DFDNet [94] take no account of this problem. iv) How can reference and LR features be combined? GFRNet [97] directly concatenates LR and

warped reference images for deep network. To fully exploit information from reference image, GWAInet [34] designs feature extraction network and fuses the features extracted. ASFFNet [96] and DFDNet [94] propose specialized AFF and DFT to compensate details from the reference images. MEFSR [148] applies weighted pixel average to merge the information.

4.6 Audio-guided Face Super-resolution

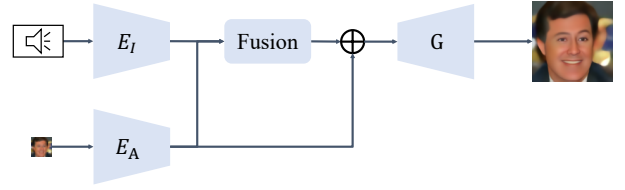


Fig. 13 The framework of audio-guided face super-resolution method.

In addition to aforementioned facial information, audio also carries facial attributes and identity information. Considering this assumption, Meishvili *et al.* [117] develop FSA that is the first audio-guided face super-resolution method, and the architecture of FSA is shown in Fig. 13. Due to the difference of multi-modal in two aspects, convergence and capacity, FSA builds two encoders, one E_I for images and the other E_A for audio. First, a generator mapping latent representation to HR is pre-trained, then a reference encoder (inverse generator) is trained for mapping HR to latent representation with auto-encoding constrain, ensuring that latent representation corresponds with HR. In this pattern, LR correspond with HR, HR corresponds with latent representation, thus LR corresponds with latent representation and E_I can be trained with LR-latent pair, and E_A is similar. Finally, a fusion network is trained and E_A is finetuned to merge audio and LR latent representations and feed acquired latent representation into the generator, thereby recovering the final SR results.

4.7 Loss Functions

In deep learning-based face super-resolution methods, loss function which measure the difference between SR and HR plays an important role to guide the network update for better performance. Initially, pixel \mathcal{L}_2 loss (also known as MSE loss) was popular, but later on, researchers found that \mathcal{L}_2 loss always guides model to generate smooth results, and many kinds of loss functions are employed (*e.g.*, pixel \mathcal{L}_1 loss, perceptual loss, adversarial loss and total variation

Table 10 Comparisons of reference face super-resolution algorithms. Number of R represents the number of reference images. Identity means whether the reference images keep the same identity as LR. Alignment shows what alignment implementation is based on. The utilization of R represents how the method takes advantage of information of R.

Frameworks	Methods	Number of R	Identity	Alignment	Utilization of R
Single-face guided methods	GFRNet [97]	1	✓	Landmark	Concatenation
	GWAInet [34]	1	✓	Flow field	GFENet
Multi-face guided methods	ASFFNet [96]	k	✓	Moving Least-Square	AFFB
	MEFSR [148]	k	✓	-	PWAve
	JSRFC [95]	k	×	landmark	Concatenation
Dictionary-guided methods	DFDNet [94]	-	×	-	DFT

loss). Among the particular faces, many methods always design multi-task networks (*e.g.*, face super-resolution network and face alignment network), which usually collaborate with heatmaps loss or identity loss. In this section, we study the loss functions in face super-resolution.

4.7.1 Pixel Loss

Pixel loss measures the distance between the two images at pixel level, which consists of two kinds: \mathcal{L}_1 loss that calculates the mean absolute error and \mathcal{L}_2 loss that calculates the mean square error, and can be expressed as follows:

$$\mathcal{L}_1(I_{HR}, I_{SR}) = \|I_{HR} - I_{SR}\|_1 = \frac{1}{hwc} \sum_{i,j,k} |I_{HR}^{i,j,k} - I_{SR}^{i,j,k}|, \quad (45)$$

$$\mathcal{L}_2(I_{HR}, I_{SR}) = \|I_{HR} - I_{SR}\|_2 = \frac{1}{hwc} \sum_{i,j,k} (I_{HR}^{i,j,k} - I_{SR}^{i,j,k})^2, \quad (46)$$

where h , w and c denote the height, width and channel of the image, and $I^{i,j,k}$ is the pixel value on location (i, j, k) . With the constrain of the pixel loss, the obtained I_{SR} can be close enough to the I_{HR} on the pixel value, and the decrease of the pixel loss leads to the increase of the PSNR as the definition of PSNR, making the pixel loss popular and widely used. Although \mathcal{L}_1 and \mathcal{L}_2 loss are both pixel loss, \mathcal{L}_1 loss has many advantages in improving the performance and convergence over \mathcal{L}_2 loss. From the view of the expression, \mathcal{L}_2 loss is sensitive to large errors and

indifferent to small errors, and \mathcal{L}_1 loss equally treats large and small errors.

In addition to \mathcal{L}_1 loss and \mathcal{L}_2 loss, Huber loss used in [72] is a kind of pixel loss in nature, which is a combination of \mathcal{L}_1 loss and \mathcal{L}_2 loss in terms of calculation and advantages. Scholars find that the \mathcal{L}_2 loss may cause instability of the training at the early training phase while \mathcal{L}_1 loss converges slowly at the later training phase. To address this problem, Huber loss combines them and uses \mathcal{L}_1 loss at the early stage and \mathcal{L}_2 loss at the late stage, which can be defined as:

$$L_{\text{Huber}}(d) = \begin{cases} \frac{1}{2}d^2 & \text{for } |d| \leq \delta \\ \delta|d| - \frac{\delta^2}{2} & \text{otherwise} \end{cases}, \quad (47)$$

where

$$d = I_{HR}^{i,j,k} - I_{SR}^{i,j,k}, \quad (48)$$

is pixel-wise difference between I_{SR} and I_{HR} . Likewise, Carbonnier penalty function is also designed in [83], which can be expressed as,

$$\mathcal{L}_{\text{Car}}(I_{HR}, I_{SR}) = \rho(I_{HR} - I_{SR}), \quad (49)$$

where

$$\rho(I_{HR} - I_{SR}) = \left((I_{HR} - I_{SR})^2 + \varepsilon^2 \right)^{1/2}, \quad (50)$$

where ε is a compensation parameter.

Although pixel loss can improve PSNR, the image quality is always ignored and over-smooth images lacking high-frequency detail are generated. Many novel loss functions are used to refine the perceptual quality.

4.7.2 SSIM Loss

Similar to pixel loss, SSIM loss is designed to improve SSIM of super-resolved images, which is based on the following equation:

$$\mathcal{L}_{\text{SSIM}}(I_{\text{HR}}, I_{\text{SR}}) = \frac{1}{2} (1 - \mathcal{SSIM}(I_{\text{HR}}, I_{\text{SR}})), \quad (51)$$

where \mathcal{SSIM} denotes the calculation of SSIM. Except SSIM loss, multi-scale SSIM loss is proposed and called MSSIM loss, which calculates SSIM loss at different scales. We do not introduce it in detail.

4.7.3 Perceptual Loss

To refine the perceptual quality of the image, one of the most popular loss functions is perceptual loss, which is expressed as

$$\mathcal{L}_{\text{Perceptual}}(I_{\text{HR}}, I_{\text{SR}}, \Phi, l) = \frac{1}{h_l w_l c_l} \sum_{i,j,k} (\Phi_{i,j,k}^l(I_{\text{HR}}) - \Phi_{i,j,k}^l(I_{\text{SR}}))^2, \quad (52)$$

where Φ is the pre-trained network to extract semantic, l is the l -th layer, and h_l, w_l, c_l represent the height, width, and number of channel of features in l -th layer respectively. In essence, the perceptual loss measures the distance between the image features extracted from pre-trained classification network Φ (e.g., VGG [138], face recognition network designed by the algorithm itself in [31]), it presents the difference at the semantic level. Different from pixel loss at the pixel level, perceptual loss usually encourages the network to generate I_{SR} more perceptually similar to I_{HR} . The I_{SR} with perceptual loss always has more high-frequency and lower PSNR than pixel loss.

4.7.4 Adversarial Loss

Except for perceptual loss and style loss, adversarial loss is also widely used in face super-resolution. Adversarial loss comes from generative adversarial network (GAN) [46]. For details, the GAN is comprised of two models: a generator (G) and a discriminator (D), where G is aimed at generating images (or text), while D takes both generated image (or text) and instance from the target domain as inputs to discriminate whether the images (or text) is from target domain. In face super-resolution, the GAN can be described as follows: G is the super-resolution model which generates the super-resolved face with LR as input, and D discriminates whether the image is generated or real. In GAN training, G and D are trained alternatively: first train D with fixed G, and then train G with fixed D. After G and D are converged, the image generated by G is consistent with the real image and can cheat D, while D cannot

distinguish generated images (SR in face super-resolution) from real images (HR in face super-resolution).

In face super-resolution, the super-resolution model acts as G and we only need to design an additional D with adversarial loss. Early methods URDGN [182] and TDAE [164] use cross entropy-based adversarial loss expressed as follows:

$$\mathcal{L}_G(I_{\text{SR}}) = -\log(\mathcal{D}(I_{\text{SR}})), \quad (53)$$

$$\mathcal{L}_D(I_{\text{SR}}, I_{\text{HR}}) = -\log(\mathcal{D}(I_{\text{HR}})) - \log(1 - \mathcal{D}(I_{\text{SR}})), \quad (54)$$

where \mathcal{L}_D denotes the loss function of D, \mathcal{L}_G denotes the loss function of G, \mathcal{D} denotes the function of D, and I_{HR} is sampled randomly from HR. However, the GAN trained with this loss is always unstable and may cause the model to collapse. Due to these problems, Wasserstein GAN (WGAN) [3] is proposed. On the one hand, WGAN removes the sigmoid or softmax layer in the output layer, and clips the weight to range $(-c, c)$, and uses RMSprop or SGD rather than Adam. On the other hand, WGAN redefines the loss function as:

$$\mathcal{L}_G(I_{\text{SR}}) = -\mathcal{D}(I_{\text{SR}}), \quad (55)$$

$$\mathcal{L}_D(I_{\text{SR}}, I_{\text{HR}}) = \mathcal{D}(I_{\text{SR}}) - \mathcal{D}(I_{\text{HR}}). \quad (56)$$

However, researchers find that the weight clipping may cause unreasonable phenomenon, then WAGN-GP [49] is proposed, which removes the weight clipping and uses gradient penalty to constrain the discriminator further:

$$\mathcal{L}_G(I_{\text{SR}}) = -\mathcal{D}(I_{\text{SR}}), \quad (57)$$

$$\mathcal{L}_D(I_{\text{SR}}, I_{\text{HR}}) = \mathcal{D}(I_{\text{SR}}) - \mathcal{D}(I_{\text{HR}}) + \lambda(\|\nabla_{I_{\text{SR}}} \mathcal{D}(I_{\text{SR}})\|_2 - 1)^2. \quad (58)$$

WGAN and WGAN-GP are both popular in GAN-based face super-resolution method, e.g., FH-GAN [8] and WGANFSR [29]. Furthermore, [13] takes hinge loss in [120] for faster training, and hinge loss can be defined as

$$\mathcal{L}_D(I_{\text{SR}}) = -\mathcal{D}(I_{\text{SR}}), \quad (59)$$

$$\mathcal{L}_G(I_{\text{SR}}, I_{\text{HR}}) = \min(0, \mathcal{D}(I_{\text{SR}}) - 1) + \min(0, -\mathcal{D}(I_{\text{HR}}) - 1). \quad (60)$$

4.7.5 Cycle Consistency Loss

Cycle consistency loss is proposed by CycleGAN [200] which is so popular that many CycleGAN-based face super-resolution methods (*e.g.*, LRGAN [13], CR [31], AGCycleGAN [112]) have been designed. In CycleGAN-based face SR, two cooperated models are used: one super-resolves the I_{LR} to recover the SR, and the other downsamples the SR back to $I_{LR'}$. The cycle consistent loss is aimed to keep the consistency between LR and $I_{LR'}$ because if the models are effective, then the NLR should be identical to LR. The cycle consistency loss can be expressed as

$$\mathcal{L}_{\text{Cycle}}(I_{LR}, I_{LR'}) = \frac{1}{hwc} \sum_{i,j,k} |I_{LR}^{i,j,k} - I_{LR'}^{i,j,k}|, \quad (61)$$

where the cycle consistency loss can be also used on I_{HR} and $I_{HR'}$.

4.7.6 Style Loss

In face super-resolution, style loss is commonly used for generating faces with fine details and better visual quality, *e.g.*, ASFFNet [96]. Style loss is first proposed in [43] and used for image style transfer. To a certain extent, this loss is similar to perceptual loss because they are both loss functions on feature level. Both SR and HR are fed into a pre-trained network (*e.g.*, VGGFace [126], and others.) to obtain their corresponding features F_{SR} , F_{HR} , and Gram matrices are calculated from both features. These matrices are used to calculate loss, which is defined as

$$\mathcal{L}_{\text{Style}}(F_{HR}, F_{SR}) = \|\mathcal{G}(F_{HR}) - \mathcal{G}(F_{SR})\|_2, \quad (62)$$

where \mathcal{G} denotes the operation to acquire the feature's Gram matrix.

4.7.7 Total Variation Loss

Total variation (TV) [130] loss is produced for noise reduction. In some face super-resolution methods [97], TV loss is used to constrain images (which are not the super-resolved results) or features to be smooth. We introduce it here. TV loss calculates the difference between neighboring pixels to impose images that are spatially smooth. TV loss is defined as:

$$\mathcal{L}_{\text{TV}}(I_{SR}) = \frac{1}{hwc} \sum_{i,j,k} \sqrt{\left(I_{SR}^{i,j+1,k} - I_{SR}^{i,j,k}\right)^2 + \left(I_{SR}^{i+1,j,k} - I_{SR}^{i,j,k}\right)^2}. \quad (63)$$

4.7.8 Feature Match Loss

Feature match loss [151] is an improved version of adversarial loss to facilitate stable training of GAN. In detail, feature match loss measures the distance between features extracted from multi-layer of the discriminator rather than its output, which can be expressed as

$$\mathcal{L}_{\text{Feature Match}} = \sum_{i=1}^T \left[\left\| \mathcal{D}_n^{(i)}(I_{HR}) - \mathcal{D}_n^{(i)}(I_{SR}) \right\|_1 \right], \quad (64)$$

where i is the i -th layer of the n -th discriminator \mathcal{D}_n , T is the total number of layers. In [151] there are n discriminators to distinguish SR and HR at different scales.

4.8 Comparisons of Different Methods

After introducing different face super-resolution methods and loss functions, we compare different models in some aspects and show the comparisons in Table 11, 12, 13, 14. To be specific, we compare the input and output of network corresponding to I and O respectively, LR denotes that the input is LR without any interpolation (*i.e.*, 16×16) while B indicates that the input is interpolated LR face image sharing the same resolution with ground truth. Notably, some models take noise vector as input, such as PULSE [118], and we use z to represent the noise vector. Towards O, we use D to denote direct output, which means the model directly recovers the final super-resolved results instead of progressively generating results at different scales (which is denoted as P). After listing the input and output of the model, we compare the upscale factor, framework (*e.g.*, PyTorch, TensorFlow, Caffe, Torch, and others), datasets used in methods, degradation process and loss function. Due to the space limitations, we use the abbreviation of loss function introduced above. Adversarial loss is denoted as \mathcal{L}_{Ad} , \mathcal{L}_{Pe} corresponds to perceptual loss, \mathcal{L}_{He} represents heatmap loss, \mathcal{L}_{Id} is identity loss, \mathcal{L}_{Pa} means parsing map loss, \mathcal{L}_{Hu} is Huber loss, \mathcal{L}_{At} is attribute loss, \mathcal{L}_{Cy} is cycle consistent loss, \mathcal{L}_{Ca} corresponds to Carbonnier function; \mathcal{L}_{MS} is MSSIM loss, \mathcal{L}_{La} is landmark loss, \mathcal{L}_{St} is style loss, and \mathcal{L}_{FM} is feature match loss.

4.9 Joint Face Super-resolution and Other Tasks

Joint face super-resolution and other tasks refer to methods that jointly consider multiple degradation and take advantage of their relationship to facilitate both tasks each other, such as joint face completion and super-resolution, deblurring and super-resolution, face frontalization and super-resolution. In the following, we introduce these multi-task methods.

Table 11 Comparisons of different face super-resolution models in which degradation process is: $I_{LR} = I_{HR} \downarrow_s$.

Methods	I	O	Upscale	Frameworks	Optimizer	Datasets	Loss functions
WaveletSRNet [59]	LR	D	$\times 4, \times 8, \times 16$	Caffe	SGD	CelebA [107], Helen [85]	$\mathcal{L}_{Wave}, \mathcal{L}_2$
RBPNet [153]	LR	D	$\times 8$	PyTorch	Adam	CelebA [107]	$\mathcal{L}_{Wave}, \mathcal{L}_2$
URDGN [182]	LR	D	$\times 8$	Torch	RMSProp	CelebA [107]	$\mathcal{L}_{Ad}, \mathcal{L}_2$
WGANFSR [29]	LR	D	$\times 4$	TensorFlow	RMSProp	CelebA [107]	\mathcal{L}_{Ad}
TDN [183]	LR	D	$\times 8$	Torch	RMSProp	CelebA [107]	$\mathcal{L}_{Ad}, \mathcal{L}_2$
TDAE [164]	LR	D	$\times 8$	Torch	RMSProp	CelebA [107]	$\mathcal{L}_{Ad}, \mathcal{L}_2$
FCGAN [10]	B	D	$\times 4$	-	Adam	CelebA [107]	$\mathcal{L}_{Ad}, \mathcal{L}_1$
Attention-FH [18]	LR	D	$\times 4, \times 8$	-	Adam	BioID [66], LFW [57]	-
PCA-SRGAN [36]	LR	D	$\times 4, \times 8$	-	Adam	CelebA [107], FFHQ [75]	$\mathcal{L}_{Ad}, \mathcal{L}_1$
SPARNet [24]	B	D	$\times 8$	PyTorch	Adam	CelebA [107], Helen [85], UMDFace [7]	$\mathcal{L}_{FM}, \mathcal{L}_1, \mathcal{L}_{Pe}$
FSRNet [27]	B	D	$\times 8$	Torch	RMSProp	CelebA [107], Helen [85]	$\mathcal{L}_{Ad}, \mathcal{L}_2, \mathcal{L}_{He}$
Super-FAN [12]	LR	D	$\times 8$	PyTorch	RMSProp	LS3D-W [11]	$\mathcal{L}_{Ad}, \mathcal{L}_2, \mathcal{L}_{Pe}, \mathcal{L}_{Pe}$
PFSRNet [78]	LR	P	$\times 8$	PyTorch	Adam	CelebA [107], AFLW [81]	$\mathcal{L}_2, \mathcal{L}_{He}, \mathcal{L}_{Ad}, \mathcal{L}_{Pe}$
PMGMSAN [147]	LR	D	$\times 8$	PyTorch	Adam	CelebA [107], CelebAMask-HQ [74]	$\mathcal{L}_{Pa}, \mathcal{L}_1$
DIC [114]	LR	D	$\times 8$	PyTorch	Adam	CelebA [107], Helen [85]	$\mathcal{L}_2, \mathcal{L}_{He}, \mathcal{L}_{Ad}, \mathcal{L}_{Pe}$
JASRNet [177]	LR	D	$\times 8$	PyTorch	Adam	CelebA [107], Helen [85]	$\mathcal{L}_1, \mathcal{L}_{He}$
CAGFace [72]	LR	D	$\times 4$	PyTorch	-	CelebAMask-HQ [74], FFHQ [75]	\mathcal{L}_{Hu}
FH-GAN [8]	LR	D	$\times 8$	PyTorch	-	CFP [135], LFW [57], VGGFace2 [19]	$\mathcal{L}_{Ad}, \mathcal{L}_1, \mathcal{L}_{Id}, \mathcal{L}_{Pe}$
FSRGRCH [179]	LR	D	$\times 8$	Torch	RMSProp	CelebA [107], Menpo [187]	$\mathcal{L}_{Ad}, \mathcal{L}_1, \mathcal{L}_{He}, \mathcal{L}_{Pe}$
KPEFH [90]	LR	D	$\times 4$	TensorFlow	SGD	CelebA [107], LFW [57]	$\mathcal{L}_2, \mathcal{L}_{Pa}, \mathcal{L}_{Pe}$
MSFSR [192]	LR	P	$\times 4, \times 8$	PyTorch	Adam	CelebAMask-HQ [74], Helen [85], WLF [160]	$\mathcal{L}_{Ad}, \mathcal{L}_2, \mathcal{L}_{He}$
FSRG3DFP [55]	LR	D	$\times 4, \times 8$		Adam	CelebA [107], Menpo [187]	-
FaceAttr [180, 181]	LR	D	$\times 8$	Torch	RMSProp	CelebA [107]	$\mathcal{L}_2, \mathcal{L}_{Pe}, \mathcal{L}_{At}$
AGCycleGAN [112]	LR	D	$\times 8$	-	SGD	CelebA [107]	$\mathcal{L}_{Cy}, \mathcal{L}_{At}$
AACNN [86]	LR	D	$\times 8$	Caffe	RMSProp	CelebA [107]	$\mathcal{L}_{Ad}, \mathcal{L}_2$
RAAN [161]	LR	D	$\times 8$	PyTorch	Adam	CelebA [107]	$\mathcal{L}_{Ad}, \mathcal{L}_1, \mathcal{L}_{At}$
FACN [162]	LR	D	$\times 8$	PyTorch	Adam	CelebA [107]	$\mathcal{L}_{Ad}, \mathcal{L}_1, \mathcal{L}_{At}, \mathcal{L}_{Pe}$
ATSENet [93]	LR	D	$\times 4, \times 8$	-	SGD	CelebA [107], LFW [158]	$\mathcal{L}_{Ad}, \mathcal{L}_1, \mathcal{L}_{He}, \mathcal{L}_{Id}, \mathcal{L}_{At}$
SICNN [188]	LR	D	$\times 8$	PyTorch	SGD	CASIA-WebFace [176], CACD2000 [23], CelebA [107], VGGFaces [126]	$\mathcal{L}_1, \mathcal{L}_{Id}$
SiGAN [53]	LR	D	$\times 4, \times 8$	TensorFlow	Adam, SGD	CASIA-WebFace [176], CelebA [107], LFW [57]	$\mathcal{L}_{Ad}, \mathcal{L}_2, \mathcal{L}_{Id}$
PULSE [118]	z	D	$\times 32$	TensorFlow	Adam	FFHQ [75]	$\mathcal{L}_1, \mathcal{L}_2$
SPGAN [189]	LR	D	$\times 8$	PyTorch	Adam	VGGFace2, CelebA [107], Helen [85], LFW [57], CFP-FP [136], AgeDB-30 [121]	$\mathcal{L}_{Id}, \mathcal{L}_2, \mathcal{L}_{SPGAN}$

Table 12 Comparisons of different face super-resolution models in which degradation process is learned by the models.

Methods	I	O	Upscale	Frameworks	Optimizer	Datasets	Loss functions
LRGAN [13]	LR	D	$\times 8$	PyTorch	Adam	CelebA [107], AFLW [81], LS3D-W [11], VGGFace2 [19], WiderFace [173]	$\mathcal{L}_{Cy}, \mathcal{L}_2$
CR [31]	LR	D	$\times 8$	TensorFlow	Adam	LRPIPA [82], LR-DukeMTMC [82], WiderFace [173], CelebA [107]	$\mathcal{L}_{Cy}, \mathcal{L}_2$

Table 13 Comparisons of different face super-resolution models in which degradation process is: $I_{LR} = ((I_{HR} \otimes k) \downarrow_s + n)_{JPEG}$.

Methods	I	O	Upscale	Frameworks	Optimizer	Datasets	Loss functions
PSFR-GAN [25]	B	D	$\times 32$	PyTorch	Adam	CelebAMask-HQ [74], FFHQ [75], PSFR-RealTest [25]	$\mathcal{L}_{Ad}, \mathcal{L}_2, \mathcal{L}_{FM}$
GFRNet [97]	B	D	$\times 4, \times 8$	Torch	Adam	CASIA-WebFace [176], VGGFace2 [19]	$\mathcal{L}_{Ad}, \mathcal{L}_2, \mathcal{L}_{TV}, \mathcal{L}_{La}$
GWAInet [34]	B	D	$\times 8$	Torch	Adam	CelebA [107], CASIA-WebFace [176], VGGFace2 [19]	$\mathcal{L}_{Ad}, \mathcal{L}_1, \mathcal{L}_{Id}$
ASFFNet [96]	B	D	$\times 4, \times 8$	PyTorch	Adam	CelebA [107], CASIA-WebFace [176], VGGFace2 [19]	$\mathcal{L}_{Ad}, \mathcal{L}_2, \mathcal{L}_{Pe}, \mathcal{L}_{St}$
DFDNet [94]	B	D	$\times 4, \times 8$	PyTorch	Adam	CelebA [107], FFHQ [75], VGGFace2 [19]	$\mathcal{L}_{Ad}, \mathcal{L}_2, \mathcal{L}_{Pe}$

Table 14 Comparisons of different face super-resolution models in which degradation process is: $I_{LR} = (I_{HR} \otimes k) \downarrow_s$.

Methods	I	O	Upscale	Frameworks	Optimizer	Datasets	Loss functions
BCCNN [198]	LR	D	$\times 2$ to $\times 5$	-	SGD	Collect from web [198]	\mathcal{L}_2
SRCNN-IBP [56]	IN	D	$\times 4$	Caffe	SGD	LFW [57]	\mathcal{L}_2
GLN [144]	LR	D	$\times 4, \times 8$	Caffe	SGD	LFW-a [158], FRGC [127]	$\mathcal{L}_{Ad}, \mathcal{L}_2$
DPDFN [69]	LR	D	$\times 4, \times 8$	TensorFlow	-	CelebA [107], Helen [85], LFW [57]	\mathcal{L}_{Ca}
C-SRIP [48]	LR	P	$\times 2, \times 4, \times 8$	-	Adam	CASIA-WebFace [176], LFW [57], Helen [85], CelebA [107]	$\mathcal{L}_{Id}, \mathcal{L}_{MS}$

4.9.1 Joint Face Completion and Super-resolution

Face super-resolution recovers low-quality faces, while face completion aims to fill in the missing area. Both of them are only effective at single task, but in the real world both LR and occluded images always coexist. Thus, restoration of faces that degraded by these two modes is important. The

most straightforward way is to first complete the occluded part with face completion model and then super-resolve the completed LR with super-resolution model, *e.g.*, obscured face hallucination network (OFHNet) [170] that builds an inpainting subnetwork for repair and then a super-resolution network for super-resolution, or switches the order. However, the results of applying two models

Table 15 Comparisons of different face super-resolution models in which degradation process is: $I_{LR} = (I_{HR}) \downarrow_s + n$, $I_{LR} = (I_{HR} \downarrow_s)_{JPEG}$ and $I_{LR} = (I_{HR} \otimes k) \downarrow_s + n$ respectively.

Methods	I	O	Upscale	Frameworks	Optimizer	Datasets	Loss functions
SGEN [28]	LR	D	$\times 4$	PyTorch	Adam	CelebA [107]	$\mathcal{L}_{Cy}, \mathcal{L}_2$
ATFMN [70]	-	D	$\times 8$	PyTorch	Adam	CelebA [107], Helen [85]	\mathcal{L}_{Ca}
CDFH [101]	LR	D	$\times 4$	TensorFlow	Adam	CelebA [107]	\mathcal{L}_2

successively always contain large artifacts. Cai *et al.* [15, 14] assume that applying face super-resolution and face completion successively does not benefit from multitask learning, leading to artifacts. Thus, FCSR-GAN, an end-to-end generative model for face completion and face super-resolution, is proposed. FCSR-GAN pre-trains a face completion (FC) model, then combines FC with super-resolution (SR) model, and trains SR with fixed FC, and finally finetunes the whole network. Then, Liu *et al.* [109] proposed a graph convolution pyramid blocks and design MFG-GAN which only needs one step to be trained rather than multiple steps of FCSR-GAN.

4.9.2 Joint Face Deblurring and Super-resolution

Blurry low-resolution faces always arise in real surveillance and sports videos, which cannot be repaired effectively by single task model, *e.g.*, super-resolution or deblurring model. Thus, deblurring and super-resolution of blurry low-quality faces simultaneously emerged. In the literature, Yu *et al.* [166] developed SCGAN to jointly deblur and super-resolve the input. Then, Song *et al.* [140] found that restored faces lack high-frequency details and ignore the utilization of facial prior information. Thus, FSGN first utilizes a parsing map and LR to recover a basic result which is then fed into detail enhancement module to compensate high-frequency detail from high-quality exemplar. Later on, FSDGAN [186] built two parallel branches: one generates a high-quality blurry face that is fed into the other branch to generates a high-quality clean face. DGFAN [168] develops two feature extraction modules for different tasks to extract features that are fed into well-designed gated fusion modules, and then generates the deblurring high-quality results.

4.9.3 Joint Illumination Compensation and Face Super-resolution

Abnormal illumination face super-resolution has also attracted the attention of many scholars. Ding *et al.* [33] built a pipeline of face detection in surroundings, and then

recovered detected faces with landmarks. Zhang *et al.* [191] formulated the problem as a normal illumination external HR guidance guide abnormal illumination LR to compensate illumination and enhance. Specifically, they developed a copy-and-paste GAN (CPGAN), including two key components: an internal copy-and-paste network to utilize face intern information for reconstruction, and an external copy-and-paste network to compensate illumination.

4.9.4 Joint Face Frontalization and Super-resolution

Faces in the real world have various poses, some of which may not be frontal. When existing face super-resolution methods are applied to non-frontal faces, the reconstruction performance drops sharply and has poor visual quality. Artifacts exist even when apply face super-resolution and face frontalization are applied in sequence or inverse order. To alleviate this problem, the method in [190] first takes advantage of STN and CNN to coarsely frontalize and hallucinate the faces, and designs fine upsampling network for refining face details. Yu *et al.* [185] proposed transformative adversarial neural network for joint face frontalization and hallucination; this method consists of a generator that contains a transformer network, an upsampling network, and a discriminator. The transformer network encodes non-frontal LR and frontal LR into latent space and constrain the non-frontal one to be close to the frontal one, and then according to encoded latent extract features that are imported into upsampling network to repair the final results.

4.10 Related Applications

4.10.1 Face Video Super-resolution

Faces usually appear in low-resolution video sequences, such as surveillance, multi-frame faces. The correlation and dependency between frames provide much more facial details. Thus, the face video super-resolution methods always pay attention to how to take advantage of the

dependency of interframes. One solution is to fuse multi-frame information and exploit inter-frame dependency. Evgeniya *et al.* [145] designed a deep multi-frame face super-resolution network that extracts features and warps the adjacent features of frames to central frames, and finally recovers the reconstruction of the central frame with a reconstruction block from concatenation of warped features. Identity-guided generative adversarial networks [88] considers identity information, which not only constrain super-resolved results and HR in pixel space but also in features extracted by face recognition network. Multi-input-single-output [4] employs a generator to generate the super-resolved results for every frame, and a fusion module to converge the super-resolved frames and estimate the reconstruction of the central frame. Deshmukh *et al.* [32] built a pipeline from video frame extraction, to face detection, to recover frames with detected faces. Considering that the aforementioned methods lack the ability to model the complex temporal dependency, Xin *et al.* [163] proposed a motion-adaptive feedback cell which captures inter-frame motion information and updates the current frames adaptively, and a novel network named Motion-Adaptive Feedback Network. Fang *et al.* [39] assumed that multiple super-resolved frames are crucial for reconstruction of subsequent frame, and thus they designed the ConvLSTM based recurrence strategies to make better use of inter-frame information.

4.10.2 Old Photo Restoration

Restoration of old pictures is important in the real world, and the degradation process is complex. Naturally, one solution is to learn the map from real LR (viewing real old images as real LR) to artificial LR (such as LRGAN and CR), which is the same as BOPBL [146]. However, BOPBL transforms images at latent space rather than image space. Specifically, a LR VAE encodes real LR and artificial LR into the same latent space by constraining their latent z_r, z_a belong to the same distribution and close to Gaussian prior, and another HR VAE encodes HR into its corresponding latent space z_y . Then, a mapping network, designed for mapping z_a, z_r to Z_Y (latent of Y), is trained solely with fixed two VAEs and feature match loss and L_1 loss between output and z_y . Finishing these, $z_r \rightarrow y$ is achieved by the encoder of LR VAE, mapping network, and the decoder of HR VAE.

4.10.3 Face 3D Reconstruction

Face 3D reconstruction is also a face-related challenging task. Richardson *et al.* [129] trained a model with synthetic data. Then unsupervised autoencoder-based model is built

in [143]. However, [64] believed that these models need a complex and inefficient pipeline, and then it constructs an efficient volumetric regression network for reconstruction. Considering facial structure, [142] combines facial prior to boost face 3D reconstruction. [42] finishes 3D reconstruction and face alignment jointly instead of relying on other prior models. [44] built a GAN-based network and proposed a novel loss based on face recognition network. Since depth map provides 3D information, there is still a method that uses depth for face 3D reconstruction [197]. Nevertheless, when these methods are applied on wild images, the results are limited. To solve this problem, AvatarMe [84] captured a large dataset and assumed the first methodology that reconstructing 3D faces from arbitrary images. [152] built a lightweight network that combines sparse photometric stereo and facial prior information for face 3D reconstruction.

4.10.4 Cross Spectral Face Super-resolution

Synthesizing photo-realistic visible faces (VIS) from near-infrared (NIR) images is important in face recognition, which is challenging especially when paired training data are unavailable. Zhang *et al.* [87] took advantage of low-rank embedding to maintain identity consistency of cross spectral hallucination. Song *et al.* [139] built a generative adversarial face super-resolution network. Given VIS image and misaligned NIR, Yu *et al.* [178] combined CycleGAN and face recognition to keep generated VIS and original VIS identity consistency, and constructed an attention warping module to achieve pose and expression-preserving VIS from NIR. Duan *et al.* propose PACH [37], which first aligns NIR to VIS with the UV map, and then transfers the texture from text prior to aligned NIR and generates a corresponding VIS.

5 Conclusion and Future Directions

In this review, we presented a taxonomy of deep learning-based face super-resolution methods according to face-specific information. This field can be divided into six categories: general face super-resolution methods, prior-guided face super-resolution methods, attribute-constrained face super-resolution methods, identity-preserving face super-resolution methods, reference face super-resolution methods, and audio-guided super-resolution method. Depending on the design of network architecture or the specific utilization of face-specific information, every category is further divided. In particular, general face super-resolution methods are further divided into basic CNN-based methods, GAN-based methods, reinforcement learning-based methods and ensemble-learning based methods. For other methods that

combine face-specific information, the categorization is executed according to the specific utilization pattern of face-specific information. Furthermore, we introduce commonly used loss function in this field. Of course, face super-resolution technique is not limited to the methods we presented. However, a fully panoramic view of this fast-expanding field is challenging, thereby resulting in possible omissions. Therefore, this review serves as a pedagogical tool to provide researchers with insights into typical methods of face super-resolution. In practice, researchers could use these general guidelines to develop the most suitable technique for their specific studies.

Despite great breakthroughs, face super-resolution still presents many challenges and it is expected to continue its fast growth. In this section, we simply provided an outlook on problems to be solved and trends to expect in the future.

5.1 High-dimensional Facial Prior Information

High-dimensional facial prior information will gain increasing research attention. The special information used in face super-resolution methods is increasingly complex and has higher than higher dimensions, from 2D images (facial landmarks, facial heatmaps, parsing maps) to 3D prior in [55], which means higher-dimensional prior provides richer information. Thus, higher-dimensional prior can significantly enhance face reconstruction.

5.2 Balance between Subjective and Objective Quality

Researchers are likely to pay more attention to the balance between subjective and objective quality. Similar to general super-resolution, face super-resolution tends to recover SR with higher PSNR but poorer visual quality and vice versa. However, the balance between subjective and objective quality is important. In general super-resolution, approaches are designed to achieve the balance, but existing face super-resolution methods ignore how to find a balance between them. This trend is expected to diffuse in face super-resolution and make a difference.

5.3 Lightweight Face Super-resolution Models

Lightweight face super-resolution models will attract more and more attention. Although deep learning-based face super-resolution methods have achieved great breakthroughs, they have difficulty in deploying real-world applications, which is caused by a mass of parameters and high computation cost. Hence, developing models with more lightweight and lower computation cost is still a major challenge. In the future, it is expected that this trend

will diffuse and play an important role in face super-resolution field.

5.4 More Challenging Scales

More challenging scales, such as $\times 32$, $\times 64$, will be explored. Existing face super-resolution methods mainly focus on the case of the magnification factors $\times 8$. However, with fast developing high resolution display devices, face super-resolution with larger scales come into existence and will become more and more popular. Naturally, larger scales must increase the difficulty of the problems. Thus, more efficient and powerful face super-resolution methods should be explored.

5.5 Unsupervised or Self-supervised Methods

Unsupervised or self-supervised methods will become mainstream face super-resolution methods. Due to the observation that the degradation process in the real world is too complex to be simulated, thereby resulting in a large gap existing between synthesized LR-HR pair and real-world data. When applying models trained by synthesized pair to real-world LR, the performance of methods drops dramatically. Therefore, unsupervised or self-supervised methods will become a general design concept in face super-resolution.

5.6 Utilization of Cross Modal Information

The utilization of cross modal information (including audio, depth, NIR (near infrared)) will be increasingly promoted. Evidently, different modals provide different information. In this field, researchers always explore face-related information, such as attribute, identity and others. Nevertheless, the appearance of audio guided face super-resolution [38] inspires us to take advantage of information belonging to different modals. This trend will undoubtedly continue and diffuse into every category in this field. The introduction of cross-modal information will also spur the development of face super-resolution.

Acknowledgements The research was supported by the National Natural Science Foundation of China (61971165, 61922027, 61773295); in part by Natural Science Foundation of Heilongjiang Province (YQ2020F004); and in part by the Fundamental Research Funds for the Central Universities.

References

1. Abello, A.A., Hirata, R.: Optimizing super resolution for face recognition. In: Proceedings of the SIBGRAPI Conference on Graphics, Patterns and Images, pp. 194–201 (2019)

2. Anwar, S., Khan, S., Barnes, N.: A deep journey into super-resolution: A survey. *ACM Computing Surveys* **53**(3), 1–34 (2020)
3. Arjovsky, M., Chintala, S., Bottou, L.: Wasserstein generative adversarial networks. pp. 214–223 (2017)
4. Ataer-Cansizoglu, E., Jones, M.: Super-resolution of very low-resolution faces from videos. In: *Proceedings of the British Machine Vision Conference* (2018)
5. Autee, M.P., Mehta, M.S., Desai, M.S., Sawant, V., Nagare, A.: A review of various approaches to face hallucination. *Procedia Computer Science* **45**, 361–369 (2015)
6. Baker, S., Kanade, T.: Hallucinating faces. In: *Proceedings of the IEEE International Conference on Automatic Face and Gesture Recognition*, pp. 83–88 (2000)
7. Bansal, A., Nanduri, A., Castillo, C.D., Ranjan, R., Chellappa, R.: Umdfaces: An annotated face dataset for training deep networks. *CoRR* **abs/1611.01484** (2016)
8. Bayramli, B., Ali, U., Qi, T., Lu, H.: Fh-gan: Face hallucination and recognition using generative adversarial network. In: *Neural Information Processing*, pp. 3–15 (2019)
9. Berthelot, D., Schumm, T., Metz, L.: BEGAN: boundary equilibrium generative adversarial networks. *CoRR* **abs/1703.10717** (2017)
10. Bin, H., Chen, W., Wu, X., Chun-Liang, L.: High-quality face image SR using conditional generative adversarial networks. *CoRR* **abs/1707.00737** (2017)
11. Bulat, A., Tzimiropoulos, G.: How far are we from solving the 2d 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In: *Proceeding of the IEEE International Conference on Computer Vision*, pp. 1021–1030 (2017)
12. Bulat, A., Tzimiropoulos, G.: Super-fan: Integrated facial landmark localization and super-resolution of real-world low resolution faces in arbitrary poses with gans. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 109–117 (2018)
13. Bulat, A., Yang, J., Tzimiropoulos, G.: To learn image super-resolution, use a gan to learn how to do image degradation first. In: *Proceeding of the European conference on computer vision*, pp. 187–202 (2018)
14. Cai, J., Han, H., Shan, S., Chen, X.: Fcsr-gan: Joint face completion and super-resolution via multi-task learning. *IEEE Transactions on Biometrics, Behavior, and Identity Science* **2**, 109–121 (2020)
15. Cai, J., Hu, H., Shan, S., Chen, X.: Fcsr-gan: End-to-end learning for joint face completion and super-resolution. In: *Proceeding of the IEEE International Conference on Automatic Face Gesture Recognition*, pp. 1–8 (2019)
16. Cansizoglu, E., Jones, M., Zhang, Z., Sullivan, A.: Verification of very low-resolution faces using an identity-preserving deep face super-resolution network. *ArXiv* **abs/1903.10974** (2019)
17. Cao, L., Liu, J., Du, K., Guo, Y., Wang, T.: Guided cascaded super-resolution network for face image. *IEEE Access* **8**, 173387–173400 (2020)
18. Cao, Q., Lin, L., Shi, Y., Liang, X., Li, G.: Attention-aware face hallucination via deep reinforcement learning. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 690–698 (2017)
19. Cao, Q., Shen, L., Xie, W., Parkhi, O.M., Zisserman, A.: Vg-gface2: A dataset for recognising faces across pose and age. In: *Proceeding of the IEEE International Conference on Automatic Face and Gesture Recognition*, pp. 67–74 (2018)
20. Chakrabarti, A., Rajagopalan, A., Chellappa, R.: Super-resolution of face images using kernel pca-based prior. *IEEE Transactions on Multimedia* **9**(4), 888–892 (2007)
21. Chan, K.C., Wang, X., Xu, X., Gu, J., Loy, C.C.: Glean: Generative latent bank for large-factor image super-resolution. *arXiv preprint arXiv:2012.00739* (2020)
22. Chang, H., Yeung, D.Y., Xiong, Y.: Super-resolution through neighbor embedding. In: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 1, pp. 1–275 (2004)
23. Chen, B., Chen, C., Hsu, W.H.: Face recognition and retrieval using cross-age reference coding with cross-age celebrity dataset. *IEEE Transactions on Multimedia* **17**(6), 804–815 (2015)
24. Chen, C., Gong, D., Wang, H., Li, Z., Wong, K.Y.K.: Learning spatial attention for face super-resolution. *IEEE Transactions on Image Processing* (2020)
25. Chen, C., Li, X., Yang, L., Lin, X., Zhang, L., Wong, K.: Progressive semantic-aware style transformation for blind face restoration. *ArXiv* **abs/2009.08709** (2020)
26. Chen, Y., Phoneyilay, V., Tao, J., Chen, X., Xia, R., Zhang, Q., Yang, K., Xiong, J., Xie, J.: The face image super-resolution algorithm based on combined representation learning. *Multimedia Tools and Applications* pp. 1–23 (2020)
27. Chen, Y., Tai, Y., Liu, X., Shen, C., Yang, J.: Fsrnet: End-to-end learning face super-resolution with facial priors. In: *Proceeding of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2492–2501 (2018)
28. Chen, Z., Lin, J., Zhou, T., Wu, F.: Sequential gating ensemble network for noise robust multiscale face restoration. *IEEE Transactions on Cybernetics* pp. 1–11 (2019)
29. Chen, Z., Tong, Y.: Face super-resolution through wasserstein gans. *ArXiv* **abs/1705.02438** (2017)
30. Cheng, X., Lu, J., Yuan, B., Zhou, J.: Identity-preserving face hallucination via deep reinforcement learning. *IEEE Transactions on Circuits and Systems for Video Technology* pp. 1–1 (2019)
31. Cheng, Z., Zhu, X., Gong, S.: Characteristic regularisation for super-resolving face images. In: *Proceeding of the IEEE Winter Conference on Applications of Computer Vision*, pp. 2424–2433 (2020)
32. Deshmukh, A.B., Rani, N.U.: Face video super resolution using deep convolutional neural network. In: *Proceeding of the International Conference On Computing, Communication, Control And Automation*, pp. 1–6 (2019)
33. Ding, X., Hu, R.: Learning to see faces in the dark. In: *Proceeding of the IEEE International Conference on Multimedia and Expo*, pp. 1–6 (2020)
34. Dogan, B., Gu, S., Timofte, R.: Exemplar guided face image super-resolution without facial landmarks. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pp. 0–0 (2019)
35. Dong, C., Loy, C.C., He, K., Tang, X.: Image super-resolution using deep convolutional networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **38**(2), 295–307 (2016)
36. Dou, H., Chen, C., Hu, X., Xuan, Z., Hu, Z., Peng, S.: Pca-srgan: Incremental orthogonal projection discrimination for face super-resolution. In: *Proceedings of the ACM International Conference on Multimedia*, pp. 1891–1899 (2020)
37. Duan, B., Fu, C., Li, Y., Song, X., He, R.: Cross-spectral face hallucination via disentangling independent factors. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7930–7938 (2020)
38. Fan, Z., Hu, X., Chen, C., Wang, X., Peng, S.: Facial image super-resolution guided by adaptive geometric features. *EURASIP Journal on Wireless Communications and Networking* **2020**(1), 1–15 (2020)
39. Fang, C., Li, G., Han, X., Yu, Y.: Self-enhanced convolutional network for facial video hallucination. *IEEE Transactions on Image Processing* **29**, 3078–3090 (2020)
40. Fang, Y., Ran, Q., Li, Y.: Fractal residual network for face image super-resolution. In: I. Farkaš, P. Masulli, S. Wermter (eds.) *Proceeding of the Artificial Neural Networks and Machine Learning*, pp. 15–26 (2020)

41. Farrugia, R.A., Guillemot, C.: Face hallucination using linear models of coupled sparse support. *IEEE Transactions on Image Processing* **26**(9), 4562–4577 (2017)
42. Feng, Y., Wu, F., Shao, X., Wang, Y., Zhou, X.: Joint 3d face reconstruction and dense alignment with position map regression network. In: *Proceedings of the European Conference on Computer Vision*, pp. 534–551 (2018)
43. Gatys, L.A., Ecker, A.S., Bethge, M.: Image style transfer using convolutional neural networks. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2414–2423 (2016)
44. Gecer, B., Ploumpis, S., Kotsia, I., Zafeiriou, S.: Ganfit: Generative adversarial network fitting for high fidelity 3d face reconstruction. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1155–1164 (2019)
45. Geng, C., Chen, L., Zhang, X., Zhou, P., Gao, Z.: A wavelet-based learning for face hallucination with loop architecture. In: *Proceeding of the IEEE Visual Communications and Image Processing*, pp. 1–4 (2018)
46. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: *Proceeding of the Advances in Neural Information Processing Systems*, vol. 2, pp. 2672–2680 (2014)
47. Grm, K., Pernus, M., Cluzel, L., Scheirer, W.J., Dobriska, S., Struc, V.: Face hallucination revisited: An exploratory study on dataset bias. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 0–0 (2019)
48. Grm, K., Scheirer, W.J., Štruc, V.: Face hallucination using cascaded super-resolution and identity priors. *IEEE Transactions on Image Processing* **29**, 2150–2165 (2020)
49. Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., Courville, A.C.: Improved training of wasserstein gans. In: *Proceeding of the Advances in Neural Information Processing Systems*, pp. 5767–5777 (2017)
50. Gunturk, B.K., Batur, A.U., Altunbasak, Y., Hayes, M.H., Mersereau, R.M.: Eigenface-domain super-resolution for face recognition. *IEEE Transactions on Image Processing* **12**(5), 597–606 (2003)
51. Guo, J., Chen, J., Han, Z., Liu, H., Wang, Z., Hu, R.: Adaptive aggregation network for face hallucination. In: R. Hong, W.H. Cheng, T. Yamasaki, M. Wang, C.W. Ngo (eds.) *Advances in Multimedia Information Processing*, pp. 190–199 (2018)
52. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. In: *Proceeding of the Advances in Neural Information Processing Systems*, pp. 6626–6637 (2017)
53. Hsu, C., Lin, C., Su, W., Cheung, G.: Sigan: Siamese generative adversarial network for identity-preserving face hallucination. *IEEE Transactions on Image Processing* **28**(12), 6225–6236 (2019)
54. Hu, X., Ma, P., Mai, Z., Peng, S., Yang, Z., Wang, L.: Face hallucination from low quality images using definition-scalable inference. *Pattern Recognition* **94**, 110–121 (2019)
55. Hu, X., Ren, W., Lamaster, J., Cao, X., Li, X., Li, Z., Menze, B., Liu, W.: Face super-resolution guided by 3d facial priors. In: *Proceeding of the European Conference on Computer Vision*, pp. 763–780 (2020)
56. Huang, D., Liu, H.: Face hallucination using convolutional neural network with iterative back projection. In: *Proceedings of the Chinese Conference on Biometric Recognition*, pp. 167–175 (2016)
57. Huang, G.B., Ramesh, M., Berg, T., Learned-Miller, E.: Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Tech. Rep. 07-49, University of Massachusetts, Amherst (2007)
58. Huang, H., He, H., Fan, X., Zhang, J.: Super-resolution of human face image using canonical correlation analysis. *Pattern Recognition* **43**(7), 2532–2543 (2010)
59. Huang, H., He, R., Sun, Z., Tan, T.: Wavelet-srnet: A wavelet-based cnn for multi-scale face super resolution. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1689–1697 (2017)
60. Huang, H., He, R., Sun, Z., Tan, T.: Wavelet domain generative adversarial network for multi-scale face hallucination. *International Journal of Computer Vision* **127**(6-7), 763–784 (2019)
61. Huang, W., Chen, Y., Mei, L., You, H.: Super-resolution reconstruction of face image based on convolution network. In: *Proceeding of the Advances in Intelligent Systems and Computing*, pp. 288–294 (2018)
62. Huang, X., Belongie, S.: Arbitrary style transfer in real-time with adaptive instance normalization. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1501–1510 (2017)
63. Indradi, S.D., Arifianto, A., Ramadhani, K.N.: Face image super-resolution using inception residual network and gan framework. In: *Proceeding of the International Conference on Information and Communication Technology*, pp. 1–6 (2019)
64. Jackson, A.S., Bulat, A., Argyriou, V., Tzimiropoulos, G.: Large pose 3d face reconstruction from a single image via direct volumetric cnn regression. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1031–1039 (2017)
65. Jaderberg, M., Simonyan, K., Zisserman, A., Kavukcuoglu, K.: Spatial transformer networks. In: *Proceeding of the Advances in Neural Information Processing Systems*, pp. 2017–2025 (2015)
66. Jesorsky, O.: Robust face detection using the hausdorff distance. *Audio- and Video-Based Biometric Person Authentication* pp. 90–95 (2001)
67. Jiang, J., Hu, R., Wang, Z., Han, Z.: Noise robust face hallucination via locality-constrained representation. *IEEE Transactions on Multimedia* **16**(5), 1268–1281 (2014)
68. Jiang, J., Yu, Y., Hu, J., Tang, S., Ma, J.: Deep cnn denoiser and multi-layer neighbor component embedding for face hallucination. In: *Proceedings of the International Joint Conference on Artificial Intelligence*, p. 771–778 (2018)
69. Jiang, K., Wang, Z., Yi, P., Lu, T., Jiang, J., Xiong, Z.: Dual-path deep fusion network for face image hallucination. *IEEE Transactions on Neural Networks and Learning Systems* pp. 1–14 (2020)
70. Jiang, K., Wang, Z., Yi, P., Wang, G., Gu, K., Jiang, J.: Atmfn: Adaptive-threshold-based multi-model fusion network for compressed face hallucination. *IEEE Transactions on Multimedia* pp. 1–1 (2019)
71. Jung, C., Jiao, L., Liu, B., Gong, M.: Position-patch based face hallucination using convex optimization. *IEEE Signal Processing Letters* **18**(6), 367–370 (2011)
72. Kalarot, R., Li, T., Porikli, F.: Component attention guided face super-resolution network: Cagface. In: *Proceeding of the IEEE Winter Conference on Applications of Computer Vision*, pp. 359–369 (2020)
73. Kanakaraj, S., V K, G., Kalady, S.: Face super resolution: A survey. *International Journal of Image, Graphics and Signal Processing* **9**, 54–67 (2017)
74. Karras, T., Aila, T., Laine, S., Lehtinen, J.: Progressive growing of gans for improved quality, stability, and variation. *CoRR abs/1710.10196* (2017)
75. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: *Proceeding of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4396–4405 (2019)
76. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: *Proceeding of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4396–4405 (2019)

77. Kazemi, H., Taherkhani, F., Nasrabadi, N.M.: Identity-aware deep face hallucination via adversarial face verification. In: *Proceeding of the IEEE 10th International Conference on Biometrics Theory, Applications and Systems*, pp. 1–10 (2019)
78. Kim, D., Kim, M., Kwon, G., Kim, D.S.: Progressive face super-resolution via attention to facial landmark. In: *Proceedings of the British Machine Vision Conference*, pp. I–I (2019)
79. Kim, J., Lee, J.K., Lee, K.M.: Accurate image super-resolution using very deep convolutional networks. In: *Proceeding of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1646–1654 (2016)
80. Ko, W., Chien, S.: Patch-based face hallucination with multitask deep neural network. In: *Proceeding of the IEEE International Conference on Multimedia and Expo*, pp. 1–6 (2016)
81. Köstinger, M., Wohlhart, P., Roth, P.M., Bischof, H.: Annotated facial landmarks in the wild: A large-scale, real-world database for facial landmark localization. In: *Proceeding of the IEEE International Conference on Computer Vision Workshops*, pp. 2144–2151 (2011)
82. Lai, S., He, C., Lam, K.: Low-resolution face recognition based on identity-preserved face hallucination. In: *Proceeding of the IEEE International Conference on Image Processing*, pp. 1173–1177 (2019)
83. Lai, W.S., Huang, J.B., Ahuja, N., Yang, M.H.: Deep laplacian pyramid networks for fast and accurate super-resolution. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 624–632 (2017)
84. Lattas, A., Moschoglou, S., Gecer, B., Ploumpis, S., Triantafyllou, V., Ghosh, A., Zafeiriou, S.: Avatarme: Realistically renderable 3d facial reconstruction. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 760–769 (2020)
85. Le, V., Brandt, J., Lin, Z., Bourdev, L., Huang, T.S.: Interactive facial feature localization. In: *Proceeding of the European conference on computer vision*, pp. 679–692 (2012)
86. Lee, C.H., Zhang, K., Lee, H.C., Cheng, C.W., Hsu, W.: Attribute augmented convolutional neural network for face hallucination. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 721–729 (2018)
87. Lezama, J., Qiu, Q., Sapiro, G.: Not afraid of the dark: Nir-vis face recognition via cross-spectral hallucination and low-rank embedding. In: *Proceeding of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6807–6816 (2017)
88. Li, D., Wang, Z.: Face video super-resolution with identity guided generative adversarial networks. In: *CCF Chinese Conference on Computer Vision*, pp. 357–369 (2017)
89. Li, J., Zhou, Y., Ding, J., Chen, C., Yang, X.: Id preserving face super-resolution generative adversarial networks. *IEEE Access* **8**, 138373–138381 (2020)
90. Li, K., Bare, B., Yan, B., Feng, B., Yao, C.: Face hallucination based on key parts enhancement. In: *Proceeding of the IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 1378–1382 (2018)
91. Li, M., Sun, Y., Zhang, Z., Xie, H., Yu, J.: Deep learning face hallucination via attributes transfer and enhancement. In: *Proceeding of the IEEE International Conference on Multimedia and Expo*, pp. 604–609 (2019)
92. Li, M., Sun, Y., Zhang, Z., Yu, J.: A coarse-to-fine face hallucination method by exploiting facial prior knowledge. In: *Proceeding of the IEEE International Conference on Image Processing*, pp. 61–65 (2018)
93. Li, M., Zhang, Z., Yu, J., Chen, C.W.: Learning face image super-resolution through facial semantic attribute transformation and self-attentive structure enhancement. *IEEE Transactions on Multimedia* pp. 1–1 (2020)
94. Li, X., Chen, C., Zhou, S., Lin, X., Zuo, W., Zhang, L.: Blind face restoration via deep multi-scale component dictionaries. In: *Proceedings of the European Conference on Computer Vision*, pp. 399–415 (2020)
95. Li, X., Duan, G., Wang, Z., Ren, J., Zhang, Y., Zhang, J., Song, K.: Recovering extremely degraded faces by joint super-resolution and facial composite. In: *Proceeding of the IEEE International Conference on Tools with Artificial Intelligence*, pp. 524–530 (2019)
96. Li, X., Li, W., Ren, D., Zhang, H., Wang, M., Zuo, W.: Enhanced blind face restoration with multi-exemplar images and adaptive spatial feature fusion. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2706–2715 (2020)
97. Li, X., Liu, M., Ye, Y., Zuo, W., Lin, L., Yang, R.: Learning warped guidance for blind face restoration. In: *Proceeding of the European Conference on Computer Vision*, pp. 272–289 (2018)
98. Liang, Y., Lai, J.H., Zheng, W.S., Cai, Z.: A survey of face hallucination. In: *Proceedings of the Chinese Conference on Biometric Recognition*, pp. 83–93 (2012)
99. Liang, Y., Xie, X., Lai, J.H.: Face hallucination based on morphological component analysis. *Signal Processing* **93**(2), 445–458 (2013)
100. Liu, C., Shum, H.Y., Zhang, C.S.: A two-step approach to hallucinating faces: global parametric model and local nonparametric model. In: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 1, pp. I–I (2001)
101. Liu, H., Han, Z., Guo, J., Ding, X.: A noise robust face hallucination framework via cascaded model of deep convolutional networks and manifold learning. In: *Proceeding of the IEEE International Conference on Multimedia and Expo*, pp. 1–6 (2018)
102. Liu, H., Ruan, Z., Zhao, P., Shang, F., Yang, L., Liu, Y.: Video super resolution based on deep learning: A comprehensive survey. *arXiv preprint arXiv:2007.12928* (2020)
103. Liu, H., Zheng, X., Han, J., Chu, Y., Tao, T.: Survey on gan-based face hallucination with its model development. *IET Image Processing* **13**(14), 2662–2672 (2019)
104. Liu, L., Wang, S., Wan, L.: Component semantic prior guided generative adversarial network for face super-resolution. *IEEE Access* **7**, 77027–77036 (2019)
105. Liu, Q., Jia, R., Zhao, C., Liu, X., Sun, H., Zhang, X.: Face super-resolution reconstruction based on self-attention residual network. *IEEE Access* **8**, 4110–4121 (2020)
106. Liu, W., Wen, Y., Yu, Z., Li, M., Raj, B., Song, L.: Sphereface: Deep hypersphere embedding for face recognition. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 212–220 (2017)
107. Liu, Z., Ping, L., Wang, X., Tang, X.: Deep learning face attributes in the wild. In: *IEEE International Conference on Computer Vision*, pp. 3730–3738 (2016)
108. Liu, Z., Siu, W., Chan, Y.: Reference based face super-resolution. *IEEE Access* **7**, 129112–129126 (2019)
109. Liu, Z., Wu, Y., Li, L., Zhang, C., Wu, B.: Joint face completion and super-resolution using multi-scale feature relation learning. *arXiv preprint arXiv:2003.00255* (2020)
110. Lu, T., Hao, X., Zhang, Y., Liu, K., Xiong, Z.: Parallel region-based deep residual networks for face hallucination. *IEEE Access* **7**, 81266–81278 (2019)
111. Lu, T., Wang, H., Xiong, Z., Jiang, J., Zhang, Y., Zhou, H., Wang, Z.: Face hallucination using region-based deep convolutional networks. In: *Proceeding of the IEEE International Conference on Image Processing*, pp. 1657–1661 (2017)
112. Lu, Y., Tai, Y.W., Tang, C.K.: Attribute-guided face generation using conditional cyclegan. In: *Proceedings of the European Conference on Computer Vision*, pp. 282–297 (2018)
113. Luo, Y., Huang, K.: Super-resolving tiny faces with face feature vectors. In: *Proceeding of the International Conference on Information Science and Technology*, pp. 145–152 (2020)

114. Ma, C., Jiang, Z., Rao, Y., Lu, J., Zhou, J.: Deep face super-resolution with iterative collaboration between attentive recovery and landmark estimation. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5569–5578 (2020)
115. Ma, X., Zhang, J., Qi, C.: Hallucinating face by position-patch. *Pattern Recognition* **43**(6), 2224–2236 (2010)
116. Majdabadi, M.M., Ko, S.: Msg-capsgan: Multi-scale gradient capsule gan for face super resolution. In: *Proceeding of the International Conference on Electronics, Information, and Communication*, pp. 1–3 (2020)
117. Meishvili, G., Jenni, S., Favaro, P.: Learning to have an ear for face super-resolution. In: *Proceeding of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1364–1374 (2020)
118. Menon, S., Damian, A., Hu, M., Ravi, N., Rudin, C.: Pulse: Self-supervised photo upsampling via latent space exploration of generative models. In: *Proceeding of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2223–2232 (2020)
119. Mittal, A., Soundararajan, R., Bovik, A.C.: Making a “completely blind” image quality analyzer. *IEEE Signal Processing Letters* **20**(3), 209–212 (2013)
120. Miyato, T., Kataoka, T., Koyama, M., Yoshida, Y.: Spectral normalization for generative adversarial networks. *CoRR* **abs/1802.05957** (2018)
121. Moschoglou, S., Papaioannou, A., Sagonas, C., Deng, J., Kotsia, I., Zafeiriou, S.: Agedb: The first manually collected, in-the-wild age database. In: *Proceeding of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 1997–2005 (2017)
122. Nguyen, K., Fookes, C., Sridharan, S., Tistarelli, M., Nixon, M.: Super-resolution for biometrics: A comprehensive survey. *Pattern Recognition* **78**, 23–42 (2018)
123. Ni, P., Zhang, D., Hu, K., Jing, C., Yang, L.: Single face image super-resolution using local training networks. In: *Proceeding of the International Conference on Systems and Informatics*, pp. 1277–1281 (2017)
124. Nie, H., Lu, Y., Ikram, J.: Face hallucination via convolution neural network. In: *Proceeding of the IEEE 28th International Conference on Tools with Artificial Intelligence*, pp. 485–489 (2016)
125. Park, J.S., Lee, S.W.: An example-based face hallucination method for single-frame, low-resolution facial images. *IEEE Transactions on Image Processing* **17**(10), 1806–1816 (2008)
126. Parkhi, O.M., Vedaldi, A., Zisserman, A.: Deep face recognition. In: *Proceedings of the British Machine Vision Conference*, pp. 41.1–41.12 (2015)
127. Phillips, P.J., Flynn, P.J., Scruggs, T., Bowyer, K.W., Jin Chang, Hoffman, K., Marques, J., Jaesik Min, Worek, W.: Overview of the face recognition grand challenge. In: *Proceeding of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 1, pp. 947–954 (2005)
128. Rajput, S.S., Arya, K.V., Singh, V., Bohat, V.K.: Face hallucination techniques: A survey. In: *Proceeding of the Conference on Information and Communication Technology*, pp. 1–6 (2018)
129. Richardson, E., Sela, M., Kimmel, R.: 3d face reconstruction by learning from synthetic data. In: *Proceeding of the International Conference on 3D Vision*, pp. 460–469 (2016)
130. Rudin, L.I., Osher, S., Fatemi, E.: Nonlinear total variation based noise removal algorithms. *Physica D: Nonlinear Phenomena* **60**(1), 259–268 (1992)
131. Sagonas, C., Tzimiropoulos, G., Zafeiriou, S., Pantic, M.: A semi-automatic methodology for facial landmark annotation. In: *Proceeding of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 896–903 (2013)
132. Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., Chen, X.: Improved techniques for training gans. In: *Proceeding of the Advances in Neural Information Processing Systems*, pp. 2234–2242 (2016)
133. Saurabh, G., Aakanksha, Rajagopalan, N.: Robust super-resolution of real faces using smooth features. In: *Proceeding of the European Conference on Computer Vision Workshop*, pp. I–I (2020)
134. Schaefer, S., McPhail, T., Warren, J.: Image deformation using moving least squares. In: *Proceeding of the ACM Special Interest Group on Computer Graphics*, p. 533–540 (2006)
135. Sengupta, S., Chen, J., Castillo, C., Patel, V.M., Chellappa, R., Jacobs, D.W.: Frontal to profile face verification in the wild. In: *Proceeding of the IEEE Winter Conference on Applications of Computer Vision*, pp. 1–9 (2016)
136. Sengupta, S., Chen, J., Castillo, C., Patel, V.M., Chellappa, R., Jacobs, D.W.: Frontal to profile face verification in the wild. In: *Proceeding of the IEEE Winter Conference on Applications of Computer Vision*, pp. 1–9 (2016)
137. Shi, Y., Guanbin, L., Cao, Q., Wang, K., Lin, L.: Face hallucination by attentive sequence optimization with reinforcement learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence* pp. 2809–2824 (2019)
138. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. *CoRR* **abs/1409.1556** (2015)
139. Song, L., Zhang, M., Wu, X., He, R.: Adversarial discriminative heterogeneous face recognition. *arXiv preprint arXiv:1709.03675* (2017)
140. Song, Y., Zhang, J., Gong, L., He, S., Bao, L., Pan, J., Yang, Q., Yang, M.H.: Joint face hallucination and deblurring via structure generation and detail enhancement. *International Journal of Computer Vision* **127**(6–7), 785–800 (2019)
141. Song, Y., Zhang, J., He, S., Bao, L., Yang, Q.: Learning to hallucinate face images via component generation and enhancement. In: *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence*, pp. 4537–4543 (2017)
142. Tewari, A., Zollhöfer, M., Garrido, P., Bernard, F., Kim, H., Pérez, P., Theobalt, C.: Self-supervised multi-level face model learning for monocular reconstruction at over 250 Hz. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2549–2559 (2018)
143. Tewari, A., Zollhofer, M., Kim, H., Garrido, P., Bernard, F., Perez, P., Theobalt, C.: Mofa: Model-based deep convolutional face autoencoder for unsupervised monocular reconstruction. In: *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pp. 1274–1283 (2017)
144. Tuzel, O., Taguchi, Y., Hershey, J.: Global-local face upsampling network. *ArXiv* **abs/1603.07235** (2016)
145. Ustinova, E., Lempitsky, V.: Deep multi-frame face super-resolution. *arXiv preprint arXiv:1709.03196* (2017)
146. Wan, Z., Zhang, B., Chen, D., Zhang, P., Chen, D., Liao, J., Wen, F.: Bringing old photos back to life. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2747–2757 (2020)
147. Wang, C., Zhong, Z., Jiang, J., Zhai, D., Liu, X.: Parsing map guided multi-scale attention network for face hallucination. In: *Proceeding of the IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 2518–2522 (2020)
148. Wang, K., Oramas, J., Tuytelaars, T.: Multiple exemplars-based hallucination for face super-resolution and editing. In: *Proceedings of the Asian Conference on Computer Vision* (2020)
149. Wang, M., Chen, Z., Wu, Q.M.J., Jian, M.: Improved face super-resolution generative adversarial networks. *Machine Vision and Applications* **31**(4), 22 (2020)
150. Wang, N., Tao, D., Gao, X., Li, X., Li, J.: A comprehensive survey to face hallucination. *International Journal of Computer Vision* **106**(1), 9–30 (2014)
151. Wang, T.C., Liu, M.Y., Zhu, J.Y., Tao, A., Kautz, J., Catanzaro, B.: High-resolution image synthesis and semantic manipulation

- with conditional gans. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 8798–8807 (2018)
152. Wang, X., Guo, Y., Deng, B., Zhang, J.: Lightweight photometric stereo for facial details recovery. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 740–749 (2020)
 153. Wang, X., Lu, Y., Chen, X., Li, W., Wang, Z.: Rbpnet: An asymptotic residual back-projection network for super resolution of very low resolution face image. In: Neural Information Processing, pp. 175–186 (2019)
 154. Wang, X., Tang, X.: Hallucinating face by eigentransformation. *IEEE Transactions on Systems, Man, and Cybernetics, Part C* **35**(3), 425–434 (2005)
 155. Wang, Y., Lu, T., Wang, Y., Zhang, Y.: Face hallucination using split-attention in split-attention network. *arXiv preprint arXiv:2010.11575* (2020)
 156. Wang, Z., Chen, J., Hoi, S.C.: Deep learning for image super-resolution: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2020)
 157. Wang, Z., Simoncelli, E.P., Bovik, A.C.: Multiscale structural similarity for image quality assessment. In: The Thirty-Seventh Asilomar Conference on Signals, Systems Computers, 2003, vol. 2, pp. 1398–1402 (2003)
 158. Wolf, L., Hassner, T., Taigman, Y.: Effective unconstrained face recognition by combining multiple descriptors and learned background statistics. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **33**(10), 1978–1990 (2011)
 159. Wu, J., Ding, S., Xu, W., Chao, H.: Deep joint face hallucination and recognition. *arXiv preprint arXiv:1611.08091* (2016)
 160. Wu, W., Qian, C., Yang, S., Wang, Q., Cai, Y., Zhou, Q.: Look at boundary: A boundary-aware face alignment algorithm. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 2129–2138 (2018)
 161. Xin, J., Wang, N., Gao, X., Li, J.: Residual attribute attention network for face image super-resolution. In: Proceeding of the Association for the Advancement of Artificial Intelligence, pp. 9054–9061 (2019)
 162. Xin, J., Wang, N., Jiang, X., Li, J., Gao, X., Li, Z.: Facial attribute capsules for noise face super resolution. *Proceeding of the Association for the Advancement of Artificial Intelligence* **34**(7), 12476–12483 (2020)
 163. Xin, J., Wang, N., Li, J., Gao, X., Li, Z.: Video face super-resolution with motion-adaptive feedback cell. *Proceeding of the Association for the Advancement of Artificial Intelligence* **34**(7), 12468–12475 (2020)
 164. Xin, Y., Porikli, F.: Hallucinating very low-resolution unaligned and noisy face images by transformative discriminative autoencoders. In: Proceeding of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3760–3768 (2017)
 165. Xintao Wang, Ke Yu, C.D., Loy, C.C.: Recovering realistic texture in image super-resolution by deep spatial feature transform. In: Proceeding of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 606–615 (2018)
 166. Xu, X., Sun, D., Pan, J., Zhang, Y., Pfister, H., Yang, M.H.: Learning to super-resolve blurry face and text images. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 251–260 (2017)
 167. Xu, Y., Li, X.: Research on face super resolution reconstruction based on dbgan. In: Proceeding of the International Conference on Mechanical, Control and Computer Engineering, pp. 8070–8075 (2019)
 168. Yang, C.H., Chang, L.W.: Deblurring and super-resolution using deep gated fusion attention networks for face images. In: Proceeding of the IEEE International Conference on Acoustics, Speech and Signal Processing, pp. 1623–1627 (2020)
 169. Yang, C.Y., Liu, S., Yang, M.H.: Structured face hallucination. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1099–1106 (2013)
 170. Yang, L., Shao, B., Sun, T., Ding, S., Zhang, X.: Hallucinating very low-resolution and obscured face images. *CoRR abs/1811.04645* (2018)
 171. Yang, L., Wang, P., Gao, Z., Wang, S., Ren, P., Ma, S., Gao, W.: Implicit subspace prior learning for dual-blind face restoration. *arXiv preprint arXiv:2010.05508* (2020)
 172. Yang, L., Wang, S., Ma, S., Gao, W., Liu, C., Wang, P., Ren, P.: Hifacegan: Face renovation via collaborative suppression and replenishment. In: Proceedings of the ACM International Conference on Multimedia, pp. 1551–1560 (2020)
 173. Yang, S., Luo, P., Loy, C.C., Tang, X.: Wider face: A face detection benchmark. In: Proceeding of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 5525–5533 (2016)
 174. Yang, W., Zhang, X., Tian, Y., Wang, W., Xue, J.H., Liao, Q.: Deep learning for single image super-resolution: A brief review. *IEEE Transactions on Multimedia* **21**(12), 3106–3121 (2019)
 175. Yang, X., Lu, T., Wang, J., Zhang, Y., Wu, Y., Wang, Z., Xiong, Z.: Enhanced discriminative generative adversarial network for face super-resolution. In: Advances in Multimedia Information Processing, pp. 441–452 (2018)
 176. Yi, D., Lei, Z., Liao, S., Li, S.Z.: Learning face representation from scratch. *CoRR abs/1411.7923* (2014)
 177. Yin, Y., Robinson, J.P., Zhang, Y., Fu, Y.: Joint super-resolution and alignment of tiny faces. *Proceeding of the Association for the Advancement of Artificial Intelligence* **34**(07), 2693–12700 (2019)
 178. Yu, J., Cao, J., Li, Y., Jia, X., He, R.: Pose-preserving cross spectral face hallucination. In: Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, pp. 1018–1024 (2019)
 179. Yu, X., Fernando, B., Ghanem, B., Porikli, F., Hartley, R.: Face super-resolution guided by facial component heatmaps. In: Proceedings of the European Conference on Computer Vision, pp. 217–233 (2018)
 180. Yu, X., Fernando, B., Hartley, R., Porikli, F.: Super-resolving very low-resolution face images with supplementary attributes. In: Proceeding of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 908–917 (2018)
 181. Yu, X., Fernando, B., Hartley, R., Porikli, F.: Semantic face hallucination: Super-resolving very low-resolution face images with supplementary attributes. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **42**(11), 2926–2943 (2020)
 182. Yu, X., Porikli, F.: Ultra-resolving face images by discriminative generative networks. In: Proceeding of the European conference on computer vision, pp. 318–333 (2016)
 183. Yu, X., Porikli, F.: Face hallucination with tiny unaligned images by transformative discriminative neural networks. In: Proceeding of the Association for the Advancement of Artificial Intelligence, pp. 4327–4333 (2017)
 184. Yu, X., Porikli, F., Fernando, B., Hartley, R.: Hallucinating unaligned face images by multiscale transformative discriminative networks. *International Journal of Computer Vision* **128**(2), 500–526 (2020)
 185. Yu, X., Shiri, F., Ghanem, B., Porikli, F.: Can we see more? joint frontalization and hallucination of unaligned tiny faces. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **42**(9), 2148–2164 (2020)
 186. Yun, J.U., Jo, B., Park, I.K.: Joint face super-resolution and deblurring using generative adversarial network. *IEEE Access* **8**, 159661–159671 (2020)
 187. Zafeiriou, S., Trigeorgis, G., Chrysos, G., Deng, J., Shen, J.: The menpo facial landmark localisation challenge: A step towards the solution. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, pp. 170–179 (2017)

188. Zhang, K., Zhang, Z., Cheng, C.W., Hsu, W.H., Qiao, Y., Liu, W., Zhang, T.: Super-identity convolutional neural network for face hallucination. In: Proceedings of the European Conference on Computer Vision, pp. 183–198 (2018)
189. Zhang, M., Ling, Q.: Supervised pixel-wise gan for face super-resolution. *IEEE Transactions on Multimedia* pp. 1–1 (2020)
190. Zhang, Y., Tsang, I.W., Li, J., Liu, P., Lu, X., Yu, X.: Face hallucination with finishing touches. *arXiv preprint arXiv:2002.03308* (2020)
191. Zhang, Y., Tsang, I.W., Luo, Y., Hu, C.H., Lu, X., Yu, X.: Copy and paste gan: Face hallucination from shaded thumbnails. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 7355–7364 (2020)
192. Zhang, Y., Wu, Y., Chen, L.: Msfsr: A multi-stage face super-resolution with accurate facial representation via enhanced facial boundaries. In: Proceeding of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, pp. 2120–2129 (2020)
193. Zhanxiang Feng, Lai, J., Xie, X., Dakun Yang, Ling Mei: Face hallucination by deep traversal network. In: Proceeding of the International Conference on Pattern Recognition, pp. 3276–3281 (2016)
194. Zhao, T., Zhang, C.: Saan: Semantic attention adaptation network for face super-resolution. In: Proceeding of the IEEE International Conference on Multimedia and Expo, pp. 1–6 (2020)
195. Zheng, W., Yan, L., Zhang, W., Gou, C., Wang, F.: Guided cyclegan via semi-dual optimal transport for photo-realistic face super-resolution. In: Proceeding of the IEEE International Conference on Image Processing, pp. 2851–2855 (2019)
196. Zheng, X., Liu, H., Han, J., Hou, S.: Deep feature-preserving based face hallucination: Feature discrimination versus pixels approximation. In: Proceedings of Chinese Conference on Pattern Recognition and Computer Vision, pp. 114–125 (2019)
197. Zhong, Y., Pei, Y., Li, P., Guo, Y., Ma, G., Liu, M., Bai, W., Wu, W.H., Zha, H.: Depth-based 3d face reconstruction and pose estimation using shape-preserving domain adaptation. *IEEE Transactions on Biometrics, Behavior, and Identity Science* pp. 1–1 (2020)
198. Zhou, E., Fan, H., Cao, Z., Jiang, Y., Yin, Q.: Learning face hallucination in the wild. In: Proceeding of the Association for the Advancement of Artificial Intelligence, pp. 3871–3877 (2015)
199. Zhou Wang, Bovik, A.C.: A universal image quality index. *IEEE Signal Processing Letters* **9**(3), 81–84 (2002)
200. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 2223–2232 (2017)
201. Zhu, S., Liu, S., Chen, C.L., Tang, X.: Deep cascaded bi-network for face hallucination. In: Proceeding of the European conference on computer vision, vol. 9909, pp. 614–630 (2016)
202. Zhuang, Y., Zhang, J., Wu, F.: Hallucinating faces: Lph super-resolution and neighbor reconstruction for residue compensation. *Pattern Recognition* **40**(11), 3178–3194 (2007)