Recognizing Tiny Faces

Siva Chaitanya Mynepalli

July 22, 2019

The Robotics Institute Carnegie Mellon University Pittsburgh, Pennsylvania 15213

> **Thesis Committee:** Deva Ramanan, Chair Simon Lucey Peiyun Hu

Thesis proposal submitted in partial fulfillment of the requirements for the degree of Master of Science in Robotics

©Siva Chaitanya Mynepalli, 2019

Abstract

Objects are naturally captured over a continuous range of distances, causing dramatic changes in appearance, especially at low resolutions. Recognizing such small objects at range is an open challenge in object recognition. In this paper, we explore solutions to this problem by tackling the fine-grained task of face recognition. State-of-the-art embeddings aim to be *scale-invariant* by extracting representations in a canonical coordinate frame (by resizing a face window to a resolution of say, 224x224 pixels). However, it is well known in the psychophysics literature that human vision is decidedly scale *variant*: humans are much less accurate at lower resolutions. Motivated by this, we explore *scale-variant* multiresolution embeddings that explicitly disentangle factors of variation across resolution and scale. Importantly, multiresolution embeddings can adapt in size and complexity to the resolution of input image *on-the-fly* (e.g., high resolution input images produce more detailed representations that result in better recognition performance). Compared to state-of-the-art "one-size-fits-all" approaches, our embeddings dramatically reduce error for small faces by at least **70%** on standard benchmarks (i.e. IJBC, LFW and MegaFace).

Acknowledgement

I would like to sincerely thank my adviser Prof. Deva Ramanan. His guidance, encouragement, and support were pivotal to achieve my targets from this Masters program.

To my labmates, thank you for the amazing environment, and the weekly discussions. I learned a great deal from everyone of you. A special shout out to Peiyun Hu and Aayush Bansal for their belief, and words of encouragement at all times.

My main objective from this program was to learn how to conduct research. Thanks to all of you, I am definitely leaving this place wiser.

To my Masters cohort, thank you for making my time at Carnegie Mellon University much more enjoyable. I thank you for all the special, much cherished moments and the amazing discussions.

To my friends elsewhere, thank you for existing. You are all a constant source of motivation and inspiration.

To my parents, words fall short to express my gratitude. Thank you for your endless support. I love you.

Contents

1	Introduction					
2	Related Work2.1CNN based face recognition2.2Human vision2.3Multi scale representations in neural networks2.4Low resolution face recognition	4 4 4 5				
3	Method3.1Resolution-specific models3.2Multiple resolution-specific embeddings3.3Multi-resolution (MR) embeddings3.4Adaptive inference3.5Training details	6 8 8 10 12				
4	Experiments4.1Single image verification4.2Identification4.3Image set-based verification4.4Megaface4.5Qualitative Results4.6Off-the-shelf super-resolution	14 14 16 16 17 19 20				
5	Additional Experiments5.1Multi-resolution pooling5.2IJB-C overall performance5.3Improvement is orthogonal to training loss	24 24 25 25				
6	Conclusion	29				

List of Figures

1.1 Traditional approaches for matching compare embedding vectors of a query and reference image. We introduce multi-resolution embeddings with several desirable properties (1) they adapt in complexity to the resolution of the input, such that larger embeddings are produced when additional high-res information is available (**bottom**). (2) they produce disentangled representations where frequency-specific components can be "switched off" when not present in the input (**top**). (3) they can adapted *on-thefly* to any desired resolution by "zero'ing out" certain frequencies (the **bottom-right** embedding).

1

2

2

7

7

- 1.2 We illustrate the drop in recognition performance with resolution. The numbers at the bottom of each probe image is the similarity score obtained by comparing a probe of specified resolution with the reference image using a state-of-the-art face recognition model [7]. However, humans can make accurate inferences on these pairs of images by comparing resolution-specific features. For example, we rely on hairstyle, face shape etc. to accurately compare the very low resolution probe image with the reference image, and on finer details like eyebrows when verifying high res images.
- 1.3 To explore how resolution affects recognition performance, we evaluate a state-of-theart face embedding (VGGFace2 [7]) on resolution-constrained subsets of a standard face recognition dataset IJBC [25]. Note the significant drop in performance as resolution decreases (i.e. 20 pixels). At a false-positive rate of 10^{-3} , the true positive rate for small (20 pixel) faces drops by 60%.
- 3.1 Impact of resolution- specific models We demonstrate the massive improvement in the performance of our resolution-specific models compared to the baseline VGG2 embedding (trained for 224x224) on the task of low-res face verification. On the left, we test our resolution-specific model tuned for images of height 16 (LFW-16), SR+FT(16). On the right, we test a resolution-specific model tuned for images of height 20 (LFW-20), SR+FT(20). We show that super-resolving the low res image back to 224x224 (SR+FT) performs better than basic bicubic upsampling (Bicubic), and VGG2. We also show that SR+FT(20) performs better than SR+FT(16) on LFW-20. It shows that we need to train resolution-specific models at multiple resolutions for best performance. Full plots shown in supp. material.
- 3.2 We describe different strategies for learning embedding networks tuned for a particular resolution of r pixels, that make use of pre-training. **Bicubic** interpolates training images of size r to a canonical resolution (224x224), and then fine-tunes a pre-trained embedding network. **Super-resolution(SR)** replaces bicubic interpolation with an off-the-shelf super-resolution network (not shown in figure). **SR+Finetuning(SR+FT)** fine-tunes both the front-end super-resolution and the embedding network.

3.3	We illustrate the difference in upsampling strategies, on images of height 16 (left), and images of height 20 (right). Bicubic interpolated images are shown in the top row, while SR+FT upsampled images are shown in the central row. We can observe that the SR+FT upsampled images are sharper near the edges from the difference images in the bottom row. Zoom in to contrast between the sets of images.	8
3.4	Jointly trained multi-resolution embedding MR-J . Each low-res image is super-resolved by SR . The figure shows that certain parts of the network are designed to only operate on images of specific resolution. These parts output embeddings tuned to images of those resolutions. As discussed earlier, (1)they adapt in complexity to the resolution of the input, such that larger embeddings are produced when additional high-res in- formation is available (bottom). (2)they produce disentangled representations where frequency-specific components can be "switched off" when not presenting the input (top/centre). (3) they can be adapted on-the-fly to any desired resolution	11
3.5	Multi-resolution pooling. The embeddings corresponding to each scale are pooled separately to generate the final multi-resolution representation of a template, as described by Eqn.3.4	12
4.1	Performance on Single image verification . The plots show ROC curves generated by verifying pairs from IJBC Covariate- <i>X</i> . We observe that MR(16) almost doubles VGG2's performance and easily outperforms FT-all for IJBC Covariate-20, and IJBC Covariate-25. Similarly, SR+FT surpasses VGG2's performance by 45% on IJBC Covariate-35, and 6% on IJBC Covariate-50. Remarkably, MR models outperform VGG2 by 70% on IJBC Covariate-35 and 11% on IJBC Covariate-50. Notice that MR-I models outperform MR-J models at both these resolutions. It is interesting to observe that the difference between our best models and FT-all increases with decrease in probe resolution. Numbers in the legend refer to true positive rate (TPR on y-axis) at 10^{-3} false positive rate (FPR on x-axis). The plots also show the number of comparisons evaluated to generate them (top).	15
4.2	Performance on Identification. The plots show CMC curves generated by classifying single image probes from IJBC- <i>X</i> , to one of a defined gallery of 3531 classes, each represented by a <i>set</i> of images. The plots for IJBC-20, and IJBC-25 show that MR at least doubles VGG2's performance. The plots for IJBC-35, and IJBC-50 show that SR+FT models perform much better VGG2. They also demonstrate that MR models surpass VGG2's performance by 66 % and 22 % respectively. Numbers in the legend show the percentage of probes with ground truth in the top-10 predictions. Number of probes in each subset are shown at the top of each plot.	17
4.3	Image set based verification. These plots show TPR (y-axis) at specified FPRs (x-axis), for probe sets with varying ratios of low-resolution images. SR+FT outperforms VGG2 and FT-all at higher ratios (0.8, 0.9). MR models (particularly MR-J) outperform all other approaches with increasingly larger margins for higher ratios. Numbers in the legend refer to TPR (y-axis) at 10^{-3} FPR (x-axis). The plots also show the number of comparisons evaluated to generate them (top).	18
1.4	We visualize the salient features captured by a 16px embedding by plotting both the low-res image pairs and their high-res counterparts. The top-left quadrant show face pairs of the same identity with high cosine similarity (shown at the bottom of each pair). The top right shows face of same identity with low similarity (due to expression or makeup changes). The bottom left mistakes suggest that the low res model relies heavily on racial and face shape cues, as different people from the same race are predicted to have high similarity. The bottom right suggests that gender appears to be an easily distinguishable feature at low resolution, previously observed in [33]	19
	pair). The top right shows face of same identity with low similarity (due to express or makeup changes). The bottom left mistakes suggest that the low res model re- heavily on racial and face shape cues, as different people from the same race are p dicted to have high similarity. The bottom right suggests that gender appears to be easily distinguishable feature at low resolution, previously observed in [33]	ion lies >re- e an

4.5	Joint training forces multi-resolution models to learn complementary features. The figure shows an input image (Input) and its nearest neighbors of a different identity determined using, (a)LR:the low resolution part of a jointly trained multi-resolution embedding (MR-J), (b)HR: high resolution part of MR-J, (c)MR: full multi-resolution embedding, and (d)Normal: a normal baseline embedding. We find that the low resolution embeddings capture features like shape of face, race, gender, and distinguishable features like high cheek bones. On the other hand, the high resolution part downplays these features and emphasizes on complementary ones like shape of nose. The multi resolution embeddings combine both these features to return a set of nearest neighbors similar to those of a normal embedding.	20
4.6	Off-the-shelf super resolution networks fails to preserve identity. The figure shows that recognizing low resolution images by super-resolving them with an off-the-shelf super resolution network, and processing them with the baseline face recognition model (FSR+VGG2) performs worse than just using the baseline (VGG2). This indicates that super-resolution networks lose identity information when operating on real images. Also, we observe that our SR-16 massively outperforms the baseline and FSR+VGG2, demonstrating that our SR models generalize better to real low resolution images. We use OTS to denote 'off-the-shelf'.	21
4.7	We present high res images from the LFW dataset along with their low res (16x16) counterparts. The figure also shows super-resolved outputs of these low res images when passed through a deep off-the-shelf super-resolution network FSR-OTS [37] and our SR , SR+FT models. We observe that FSR-OTS hallucinates high resolution images but fails to preserve the identity of the person even with these artificially resized images.	21
4.8	The figure shows more examples of real low resolution images of different identities super-resolved with various approaches. It emphasizes poor performance of the off-the-shelf super resolution networks (FSR-OTS) on real images.	22
4.9	The figure shows <i>real</i> low res images of the same identity from the IJBC dataset. The outputs of FSR-OTS show that (a)deep super-resolution networks generalize poorly to real low res images, and (b)the outputs do not retain identity information. In contrast, SR performs similarly on both real and artificially downsampled images (Fig. 4.7). From the outputs of SR+FT, and SR we observe that SR+FT sacrifices its ability to output sharper images to aid recognition. SR+FT primarily focuses on eliminating JPEG artifacts, and noise in the original image.	22
5.1	Impact of multiresolution pooling The figure shows true positive rate (TPR) at 10^{-3} false positive rate for probe sets with at least 60% low resolution images. The results illustrate that naive pooling, MR-J(NaPo), and MR-I(NaPo), MR features perform worse than multi-resolution pooling, MR-J and MR-I. Interestingly, naive pooling performs worse than SR+FT for probes with large fraction of low resolution images (0.7, 0.8, 0.9), showing that it is sub-optimal to corrupt 'high resolution features' from high res images with 'high resolution features' extracted from low resolution images	24
5.2	The figure shows a relative comparison between our multi-resolution model and the baseline VGG2 on all comparisons in the IJBC single image verification protocol (Covariate verification). We find that our model induces a small overall improvement. Note that the improvement is small because of the smaller number of comparisons involving low resolution images.	25

5.3 The figures show 2D embeddings output by a CNN trained to distinguish between 8 face identities (classes). Each dot represents an embedding of a training image, and its color represents the class of the image. The lines from the center of the circle to the cluster centres depict the direction of column weights W_{y_i} . On the left, we can observe that the embeddings output by a CNN trained on softmax loss have large intra-class variance and low inter-class separation. On the right, we can observe that the ArcFace loss explicitly optimizes for these two desirable properties.

26

List of Tables

3.1	Given a fixed embedding dimension (say 128), does MR embedding perform better	
	than its fixed counterparts? The table shows that MR embeddings, both joint and in-	
	dependent, composed of two 64 dimensional embeddings perform much better than	
	a single resolution embeddings of same size SR+FT, and also the 2048 dimensional	
	baseline model VGG2 on real low resolution images (height ; 40px.). We use real im-	
	ages to better visualize the difference between the models. Full plots are shown in the	
	supp. material.	10
3.2	Given a probe and gallery image pair of different resolutions, what should be the	
	resolution of the embeddings used to compare them? The table shows that in case of a	
	large mismatch in resolution of the probe and the gallery image: the best performance	
	is achieved by resizing the higher resolution image (25 px) to the lower resolution (16	
	px), and employing lower-resolution (16 px) embedding (left). If the mismatch is not	
	large, we can use either representation (right). Full plots are shown in the supp. material.	11
4.1	The table shows that our multiresolution models continue to outperform the baseline	
	models (VGG2, FT-all), and also the SR+FT models. However, note that the difference	

4.1 The table shows that our infutnesolution models continue to outperform the baseline models (VGG2, FT-all), and also the SR+FT models. However, note that the difference between SR+FT and MR-X is not high because the test images are artificially down-sampled and the models may overfit to this downsampling method. More results are shown in the supplementary material.

Chapter 1 Introduction

Objects are visually captured at a continuous range of distances in the real world. One of the remaining open challenges in object recognition is recognition of small objects at range [21]. We focus on the illustrative task of recognizing faces across a wide range of scales, a crucial task in surveillance [6]. This is a well-known challenge because distinctive features (such as eyebrows [29]) may not be resolvable in low resolution. Contemporary face recognition systems, which now outperform the average forensic examiner on high quality images [17], perform dramatically worse for lower resolutions (Fig. 1.2 and 1.3).



Figure 1.1: Traditional approaches for matching compare embedding vectors of a query and reference image. We introduce multi-resolution embeddings with several desirable properties (1) they adapt in complexity to the resolution of the input, such that larger embeddings are produced when additional high-res information is available (**bottom**). (2) they produce disentangled representations where frequency-specific components can be "switched off" when not present in the input (**top**). (3) they can adapted *on-the-fly* to any desired resolution by "zero'ing out" certain frequencies (the **bottom-right** embedding).

Scale: Recognition is often cast as an image retrieval task, where the central challenge is learning an embedding for matching image queries (probes) to a stored library (gallery). Virtually all contemporary retrieval systems learn a scale-*in*variant embedding, by first canonicalizing a given image crop to a standard resolution (of say, 224x224 pixels) before feature extraction [19]. However, recognition accuracy for human vision is decidedly scale *variant*. Humans are much more accurate at higher resolutions, and moreover, tend to rely on resolution-specific features to make inferences at particular resolutions [32]. Fig. 1.2 shows a reference image and candidate probe matches at varying resolutions. At low resolutions, coarse features such as the hairline and jaw shape seem to reveal the identity. At



Figure 1.2: We illustrate the drop in recognition performance with resolution. The numbers at the bottom of each probe image is the similarity score obtained by comparing a probe of specified resolution with the reference image using a state-of-the-art face recognition model [7]. However, humans can make accurate inferences on these pairs of images by comparing resolution-specific features. For example, we rely on hairstyle, face shape etc. to accurately compare the very low resolution probe image with the reference image, and on finer details like eyebrows when verifying high res images.



Figure 1.3: To explore how resolution affects recognition performance, we evaluate a state-of-theart face embedding (VGGFace2 [7]) on resolution-constrained subsets of a standard face recognition dataset IJBC [25]. Note the significant drop in performance as resolution decreases (i.e. 20 pixels). At a false-positive rate of 10^{-3} , the true positive rate for small (20 pixel) faces drops by 60%.

high resolutions, subtle features such as the eyebrow and nose shape appear to play an important role. Such resolution-specific identity cues cannot be captured by a scale-invariant embedding.

Mulitresolution embeddings: We begin by showing that a conceptually simple solution is to train *multiple* fixed-resolution embeddings, and use the appropriate one depending on the resolution of the query (probe) and reference (gallery) face to be compared. Moreover, one can significantly improve accuracy by combining these resolution-specific embeddings into a *single* multiresolution representation that explicitly disentangles factors of identity into frequency-specific components. For example, certain dimensions of the embedding vector are trained to encode low-frequency cues such as hairlines, while other dimensions are trained to encode high-frequency cues such as nose shape. In the limit, one can interpret our embeddings as a "fourier" decomposition of identity into frequency-specific components. Importantly, because the resolution of an input image is known, missing frequencies for low-res inputs can be "switched off". Moreever, even when present in high-res input, they can be "zero'd out" on-the-fly to facilitate comparisons to low-res images (Fig. 1.1).

Disentangled representations: We illustrate two applications that specifically exploit disentangled embeddings. The first is *adapation*: given a probe at a particular resolution, we adapt the gallery embedding *on-the-fly* by selecting the appropriate frequency-specific components in the embedding (Fig. 1.1). The second is *aggregation*: practical face recognition methods often match *sets* of faces (say, extracted from a video sequence). Such methods typically produce an aggregate template representation by pooling embeddings from faces in the set [7,28]. We show that multiresolution pooling, that uses only high-resolution faces to produce the high-frequency components in the final embedding, is considerably more accurate.

Evaluation: Evaluating our method is hard because most benchmarks provide faces only at highresolution. This reveals the inherent bias of the community for scale invariance! It is tempting to create artificial scale variation by resizing such images [18]. In fact, we do so for diagnostic experiments, resizing the well-known LFW datset [15] into different resolutions. However, recent work has shown that downsampling is not a good model for natural scale degradation [5]. As such, we present final results on the IJBC [25] benchmark, which is unique in that it includes the raw images on which faces were extracted, and so contains natural scale variation. We also compare our algorithm on *resized* versions of the popular Megaface dataset to showcase our algorithm on a larger scale. Additionally, we compare the performance of our approach with more recent face recognition networks in the supplement.

Chapter 2

Related Work

2.1 CNN based face recognition

Recent methods for face recognition aim to learn an nonlinear embedding through a variety of loss functions, including triplet loss [30], softmax loss [26], contrastive loss [9], center loss [36]. We use the well-known VGG face network [7] as our backbone for fine-tuning in most of the experiments.

More recently, [10, 24, 34] propose angular softmax as minor modifications to the traditional softmax loss to learn an embedding that is more suitable for open-set recognition (i.e. identities in the test set do not overlap with those in the training set). We also test our method on models trained with ArcFace [10], to show that the improvement caused by our method is orthogonal to the improvement caused by modifications to the training loss. This shows that our method of training resolution-specific embeddings is necessary for optimal performance over a large range of scales. Instead of learning an "one-size-fits-all" embedding, we learn a multiresolution representation that can be adapted to different resolutions. Our approach to scale-invariance is inspired by previous work on pose invariance [16], which learns separate models for frontal and profile faces.

2.2 Human vision

Extensive studies on human vision show that human are surprisingly good at recognizing low-res faces [32]. [11] shows that human accurately recognize familiar faces even as small as 16x16. [6] points out the familiarity is the key – the more human are familiar with the face subject the more they can tolerate the poor quality of imagery. Perhaps the closest analogy to familiarity is learning-based recognition methods. Contemporary face recognition approaches train face embeddings on millions of images for many iterations. In some sense, given any new face image, it must have seen faces that feel familiar.

2.3 Multi scale representations in neural networks

Using representations drawn from multiple scales has been integral to computer vision tasks ever since the seminal work on gaussian pyramids [1]. More recently, researchers have been using deep representations drawn from multiple scales to include greater context for Semantic Segmentation [39], Object Detection [20] and other vision tasks. Our work is inspired by such approaches, but differs in its execution because the *dimensionality* of our underlying embedding depends on the image resolution.

2.4 Low resolution face recognition

Recent works on low-resolution face recognition can be classified into two categories [35]. The first category can be referred to as super-resolution based [2, 3, 13, 14, 18, 22, 23, 38, 40] approaches. Given a low-res probe, these methods first hallucinate the high-res version, and then verify/classify the high-res version. Alternatively, one might learn a feature representation that is designed to work at low resolutions [4, 8]. Such representations are often based on handcrafted features (such as color). In our approach, we learn resolution-specific features instead of hand-crafting them. Additionally, we employ super-resolution networks as a pre-processing stage that is trained end-to-end with the resolution-specific embedding.

Perhaps the most relevant work to ours is [38], which learns a *fixed-resolution* deep network to regress a high-res embedding from low-res images using a L2 loss. In comparison, we learn *multi-resolution* embeddings that are directly trained to minimize (categorical) identity mis-classifications.

Chapter 3

Method

As argued above traditional face recognition models suffer a massive drop in performance on low-resolution images (Fig. 1.3). In this section, we explore various simple strategies to remedy this. We make use of an artificially-resized LFW dataset where all images are sized to a target resolution of X pixels (denoted as LFW-X) to support design decisions.

3.1 Resolution-specific models

The most intuitive way to alleviate the impact of resolution is to train separate models for specific resolutions. But, how does one train an embedding for say, a 16x16 image?

3.1.1 Training images

Ideally, we should train these models with *real* low-resolution images of size 16x16, but in general, there may not be enough in a given training set. An attractive alternative is to augment the training set with resized images, a common practice in multi-scale training. We find that upsampling images may introduce blurry artifacts, but downsampling is a relatively benign form of augmentation (even given the caveats of [5]). In practice, we downsample images from VGGFace2 to the resolution of interest to train resolution-specific models for all resolutions ; 60.

3.1.2 Pre-training

Armed with a training set of 16x16 images, which network architecture do we use to learn an embedding? One option is training a custom architecture from scratch for that resolution. But this makes it hard to take advantage of pretrained backbone networks trained on faces resized to a fixed input size (finetuning networks pretrained on high-res images was shown to perform better than training them from scratch on low-res images [27, 31]). So, we *upsample* the downsampled images *back* to 224x224 with **Bicubic** interpolation, and fine-tune a ResNet-50 (pretrained on VGGFace2 at full-resolution) on such training images. To evaluate this approach, we train and test a face verification model on LFW-X. Fig. 3.1 demonstrates that simple resolution-specific models results in a dramatic relative improvement over an off-the-shelf embedding (**VGG2**): 60% for LFW-16 and 15% for LFW-20.

3.1.3 Super-resolution (SR)

We posit that the specific method for upsampling the input image might have a large effect on recognition performance. Fig. 3.2 replaces the bicubic upsampler with a (lightweight) super-resolution (**SR**) network. Interestingly, Fig. 3.1 demonstrates that super-resolution networks may lose identity



Figure 3.1: **Impact of resolution- specific models** We demonstrate the massive improvement in the performance of our resolution-specific models compared to the baseline **VGG2** embedding (trained for 224x224) on the task of low-res face verification. On the left, we test our resolution-specific model tuned for images of height 16 (LFW-16), **SR+FT(16)**. On the right, we test a resolution-specific model tuned for images of height 20 (LFW-20), **SR+FT(20)**. We show that super-resolving the low res image back to 224x224 (SR+FT) performs better than basic bicubic upsampling (**Bicubic**), and **VGG2**. We also show that **SR+FT(20**) performs better than **SR+FT(16)** on LFW-20. It shows that we need to train resolution-specific models at multiple resolutions for best performance. Full plots shown in supp. material.



Figure 3.2: We describe different strategies for learning embedding networks tuned for a particular resolution of r pixels, that make use of pre-training. **Bicubic** interpolates training images of size r to a canonical resolution (224x224), and then fine-tunes a pre-trained embedding network. **Super-resolution(SR)** replaces bicubic interpolation with an off-the-shelf super-resolution network (not shown in figure). **SR+Finetuning(SR+FT)** fine-tunes both the front-end super-resolution at the embedding network.

relevant information (also observed in [18]). In supplementary material, we show that this effect is even more pronounced with deeper state-of-the-art super-res networks operating on *real* images.



Figure 3.3: We illustrate the difference in upsampling strategies, on images of height 16 (left), and images of height 20 (right). **Bicubic** interpolated images are shown in the top row, while **SR+FT** upsampled images are shown in the central row. We can observe that the **SR+FT** upsampled images are sharper near the edges from the difference images in the bottom row. Zoom in to contrast between the sets of images.

3.1.4 Super-resolution with Fine-tuning (SR+FT)

Finally, we finetune the lightweight super-resolution network along with the backbone face embedding model, to guide the SR model to retain identity information. Fig. 3.1 shows that **SR+FT** outperforms bicubic interpolation. Fig. 3.3 visualizes images generated by the fine-tuned super-resolution network, which are sharper than the bicubic result.

3.2 Multiple resolution-specific embeddings

Fig. 3.1 suggests models tuned for particular resolutions (16px) might outperform models tuned for similar but distinct sizes (20px). To avoid training an exorbitant number of models, we choose a fixed number of 'anchor resolutions' r spaced along a linear scale of 16px, 35px, and 50px. We found this to provide a good tradeoff of memory and performance. Please see the Experiments section for additional details.

3.3 Multi-resolution (MR) embeddings

The above results suggest that one should train a set of resolution-specific models to improve recognition performance. It is natural to ask if these different resolution-specific embeddings could be *ensembled* together to improve performance. In order to apply a different network to a given input image, we would need to upsample or downsample it. As previously argued, downsampling an image is less prone to introducing artifacts, unlike upsampling. This suggests that given an image at a fixed resolution, one can ensemble together embeddings tuned for lower resolutions by downsampling.

3.3.1 Independent MR (MR-I)

A reasonable solution is to concatenate these embeddings together to produce a 'multi-resolution' embedding.

$$\Phi(x) = \begin{bmatrix} \frac{\phi_1(x_1)}{\|\phi_1(x_1)\|} & \frac{\phi_2(x_2)}{\|\phi_2(x_2)\|} & \dots & \frac{\phi_n(x_n)}{\|\phi_n(x_n)\|} \end{bmatrix}$$
(3.1)

where, x_i is a lower resolution version of a given image x resized to anchor height r_i , and ϕ_i denotes a resolution-specific model tuned for that specific resolution. We find that normalizing each resolution-specific embedding is necessary to match the relative scales of the embeddings.

MR-I inference: Given an input image at a particular resolution, we create its downsampled versions corresponding to anchor resolutions of equal or smaller size. This collection of blurred images are processed with resolution-specific streams that produce embeddings that are concatenated together to produce the final multi-resolution vector given by Equation 3.1 for **MR-I** models. With such a representation, the similarity score between an image pair (x, y) downsampled to the same anchor resolution is evaluated as follows,

$$s(\Phi(x), \Phi(y)) = \frac{\Phi(x)^T \Phi(y)}{\|\Phi(x)\| \|\Phi(y)\|}$$
(3.2)

The similarity score is equal to the mean cosine similarity of the resolution-specific embeddings. Qualitatively, this is equivalent to comparing probe and reference images at multiple scales.

3.3.2 Jointly-trained MR (MR-J)

Because the above approach naively concatenates together independently-trained embeddings, they might contain redundant information. To truly *disentangle* features across scale, we would like to jointly train all constituent resolution-specific embedding "streams" of a network. Following the grand tradition of *residual networks* [12], jointly training the resolution-specific streams would force the streams tuned for higher resolutions to learn residual complementary information.

3.3.3 Embedding dimension

We test this hypothesis by examining the embeddings generated by our multiresolution networks. Particularly, we ask the salient question: given a limited budget to store information by constraining the target dimension for an embedding (say, 128d), are multi-resolution embeddings able to store more information? The answer is yes! As shown in Table. 4.1, multi-resolution embeddings composed of two 64 dimensional embeddings (**MR-I,128-dim** and **MR-J,128-dim**), are better at face recognition than single-res embeddings of equal size, **SR+FT,128-dim** which are trained on the same data (other than benign blurring) with the same loss function.

Also, we would like embeddings with small memory footprints. Our multi-res embeddings might generate large memory footprints if implemented naively. This experiment demonstrates that given limited storage space, it is better to store multi-resolution embeddings of the same size.

3.3.4 MR-J inference

Fig.3.4 demonstrates the operation of a joint multi-resolution model. It shows that certain parts of an MR-J network are designed to only operate on inputs of certain resolutions, while other parameters are shared. For example, given a low resolution image (r_1xr_1) , the network outputs only a part of the overall embedding (blue), while it outputs the full embedding for a higher resolution image (r_3xr_3) . Given an input image, the outputs of these resolution specific streams are concatenated together to output a true multi-resolution embedding as discussed earlier. We show in the supplementary material that joint training forces higher resolution streams to learn to ignore low resolution features like gender [33] etc., demonstrating that they encode disentangled features.

Method	Embed dim.	TPR at 1e-3 FPR		
MR-J	128	61.2		
MR-I	128	60.7		
SR+FT	128	54.3		
VGG2	2048	38.7		

Table 3.1: Given a fixed embedding dimension (say 128), does MR embedding perform better than its fixed counterparts? The table shows that MR embeddings, both joint and independent, composed of two 64 dimensional embeddings perform much better than a single resolution embeddings of same size SR+FT, and also the 2048 dimensional baseline model VGG2 on real low resolution images (height ; 40px.). We use real images to better visualize the difference between the models. Full plots are shown in the supp. material.

3.3.5 Parameter sharing

How does one decide the optimal policy to share parameters between the resolution-specific streams? We experiment with two extreme strategies to help us identify the ideal approach. (a) we test a model in which no parameters are shared across different resolutions, i.e. each stream operates independently till the final output stage. We refer to this model as **MR-J(W)** or MR-J(Wide). (b) at the other extreme, we test another model in which a small 3-layered resolution-specific streams operate on an embedding output by a fully shared network. We refer to this model as **MR-J**. As a consequence of aggressively sharing parameters across different resolutions, **MR-J** much more efficient than **MR-J(W)**. Its memory footprint and computational complexity are comparable to a single ResNet-50 model (25M vs 23M params). We direct the reader to the supplementary material for a detailed description of the training scheme.

In the Experiments section, we show that multi-resolution embeddings significantly outperform VGG2, and also our resolution-specific models SR+FT.

3.4 Adaptive inference

3.4.1 Choosing the ideal representation

Thus far, our results indicate that when comparing two images at a particular resolution, we should use **MR** embeddings tuned for that resolution. Now, what about comparing two faces at *different* resolutions? Two natural options are (a) downsample the larger image to the smaller size, and use a model tuned for the smaller resolution or (b) upsample the smaller image and use a model tuned for the larger image. We analyze these strategies along with the baseline approach on dissimilarly resized LFW datasets for a clean evaluation. Table. 3.2 shows that when the two resolutions are similar (20px vs 25px), it doesn't quite matter. But for a large mismatch (16px vs 25px), (a) using a representation tuned for the lower resolution image is more effective.

3.4.2 Adaptive multi-resolution inference

Assume we are given a gallery of high-resolution face images. Our model produces a multi-resolution embedding that is stored for all gallery images. Given a probe of a particular size r, our prior experiments suggest that we should tune the gallery to the closest-anchor resolution, r_i . This is trivial to do with a *disentangled* multi-resolution embedding. Simply tune the gallery embeddings "on-the-fly" with array indexing:

$$\Phi(x)[1:i] \tag{3.3}$$



Figure 3.4: Jointly trained multi-resolution embedding **MR-J**. Each low-res image is super-resolved by **SR**. The figure shows that certain parts of the network are designed to only operate on images of specific resolution. These parts output embeddings tuned to images of those resolutions. As discussed earlier, (1)they adapt in complexity to the resolution of the input, such that larger embeddings are produced when additional high-res information is available (bottom). (2)they produce disentangled representations where frequency-specific components can be "switched off" when not presenting the input (top/centre). (3) they can be adapted on-the-fly to any desired resolution

TPR at FPR 1e-3						
LFW-16 vs LFW-25			LFW-20 vs LFW-25			
SR+FT(16)	SR+FT(25)	VGG2	SR+FT(20)	SR+FT(25)	VGG2	
94.1	89.0	74.5	96.7	96.7	88.5	

Table 3.2: Given a probe and gallery image pair of different resolutions, what should be the resolution of the embeddings used to compare them? The table shows that in case of a large mismatch in resolution of the probe and the gallery image: the best performance is achieved by resizing the higher resolution image (25 px) to the lower resolution (16 px), and employing lower-resolution (16 px) embedding (left). If the mismatch is not large, we can use either representation (right). Full plots are shown in the supp. material.

3.4.3 Multi-resolution pooling

Practical face recognition methods often operate on *sets* of faces (say, extracted from a video sequence). Such methods generate an aggregate template representation by pooling embeddings of face images in the set. The templates are then used to efficiently compare these sets with a single image or with an other set of faces. In our supplementary material, we show that naive pooling of our multi-resolution embeddings is not optimal. Intuitively, naive pooling mixes information across scales. Rather, we should use only high-resolution faces to construct the pooled high-frequency feature. We operate on the *i*th anchor resolution as follows:

$$\bar{\phi}_i = \frac{1}{|S_i|} \sum_{x_i \in S_i} \phi_i(x_i) \tag{3.4}$$



Figure 3.5: **Multi-resolution pooling.** The embeddings corresponding to each scale are pooled separately to generate the final multi-resolution representation of a template, as described by Eqn.3.4

where S_i is the set of images in the set that are of at least the resolution of r_i , and $\bar{\phi}_i$ is the pooled feature for anchor resolution r_i (Fig. 3.5). These features are concatenated to output a multi-resolution template embedding, as done earlier.

3.5 Training details

In this section, we go over the exact training procedure for our models in detail.

3.5.1 Architecture

We use the ResNet-50 network that has been pretrained on VGGFace2 dataset as the backbone architecture for each resolution-specific component of our models. By default, the network is trained to classify input images of size 224x224 as one of over 8000 unique face identities.

3.5.2 Preprocessing

We conform to the procedure recommended in [7] to obtain face image patches of size 224x224. We then blur these image patches with a Gaussian filter before resizing to the target "anchor resolution" to prevent aliasing. We then subtract the mean and normalize it to be within the range of [-1,1]. All networks described in the paper are trained on similarly downsampled versions of the VGG2 dataset.

3.5.3 Training SR networks

We pre-train a separate instance of a small 5-layer CNN for each anchor resolution r. We first train the network for each anchor resolution to super-resolve downsized images of resolution r back to 224x224 with L1 loss. We then append these pre-trained CNNs to the recognition model to further fine-tune them together, as shown in Fig.5 as **SR+FT** in the original paper(5K parameters).

3.5.4 Training wide joint multi-resolution models (MR-J(W))

In practice, we found it easier and effective to train a multi-resolution network in a coarse-to-finemanner. In this scheme, we train the coarse-low resolution image stream first, and then train resolutionspecific streams for progressively higher resolutions while feeding these streams concatenated embeddings of lower resolution streams as an auxiliary input. Same as earlier, our final multiresolution embedding in this scenario is smaller in size to the default embedding of our baseline, as we use three resolution-specific streams each of size 512 (compared to the baseline embedding size of 2048).

3.5.5 Training joint multi-resolution models (MR-J)

As argued earlier, the optimal strategy to train a multi-resolution network is to have disjoint resolutionspecific streams. However, storing such a huge model is not optimal. We propound an alternate endto-end approach to reduce the memory footprint of the model. In this approach, the network consists of a single CNN which acts as a feature extractor. The embedding output by this CNN is processed by resolution-specific streams which output embeddings which exclusively encode information of specific resolutions. We force the various resolution-specific components of the multi-resolution embedding to encode disentangled information by using appropriate subsets of the full multi-resolution embedding to classify an image of a specific resolution. For example, we use only the lowest resolution embedding for a 16x16 face. For a higher resolution face, say 35x35, we pass both the lowest resolution embedding and the mid-resolution embedding as input to a softmax layer. This operation can be intuitively observed in Fig. 7 of the original submission.

Chapter 4

Experiments

As argued earlier, we focus our final results on the IJB-C dataset because it uniquely includes *real* low resolution images. We create 4 resolution constrained subsets of low resolution faces (height ; 60) from the IJBC dataset to test the effectiveness of our algorithm at various scales. Each of these subsets, named **IJBC-***X*, contains faces of height close to $X \in \{20, 25, 35, 50\}$. For example, a face of height 28 px is placed in the IJBC-25 subset. We will make these subsets publicly available.

In the following subsections, we discuss the results of our algorithms on probe images drawn from these splits when tested under various protocols of the IJB-C dataset and compare them with the baseline VGG2. Additionally, we compare our results with a VGG2 model finetuned with artificially downsampled images of *all* resolutions, **FT-all**.

4.1 Single image verification

4.1.1 Setup

The simplest IJB-C protocol is 1:1 covariate verification, where a single probe image is compared with a single reference image. The protocol specifies over 48M verification pairs from which we sample those pairs with at least one low resolution image (height ; 60). We bin verification pairs into one of 4 groups, **IJBC Covariate-***X*, when the lower resolution image in the pair belongs to IJBC-*X*.

4.1.2 Results

Fig. 4.1 shows the true positive rate (TPR) at 1e-3 false positive rate (FPR). The plot shows that a simple resolution-specific model tuned for images of height 16, (both MR-J, SR+FT) almost **doubles** the performance of VGG2 on both IJBC Covariate-20, IJBC Covariate-25. Note that for the lowest anchor resolution (16x16), MR-J is same as SR+FT. Similarly, resolution-specific models SR+FT, exceed the baseline's performance by 45% on IJBC Covariate-35, and 6% on IJBC Covariate-50 respectively. More importantly, we draw attention to the remarkable performance of multi-resolution embeddings, MR-J(W), MR-J and MR-I. We find that the MR models outperform VGG2 by 70% on IJBC Covariate-35, and 11% on IJBC Covariate-50. They also easily surpass the resolution-specific models and FT-all. All relative improvements are reported at 10^{-3} False Positive Rate.

4.1.3 Discussion

(a)Why do MR models massively outperform other models? Disentangling resolution-specific features forces models to learn to encode scale-specific features which were ignored when trained on



Figure 4.1: **Performance on Single image verification**. The plots show ROC curves generated by verifying pairs from IJBC Covariate-*X*. We observe that MR(16) almost **doubles** VGG2's performance and easily outperforms FT-all for IJBC Covariate-20, and IJBC Covariate-25. Similarly, SR+FT surpasses VGG2's performance by **45%** on IJBC Covariate-35, and **6%** on IJBC Covariate-50. Remarkably, MR models outperform VGG2 by **70%** on IJBC Covariate-35 and **11%** on IJBC Covariate-50. Notice that MR-I models outperform MR-J models at both these resolutions. It is interesting to observe that the difference between our best models and FT-all increases with decrease in probe resolution. Numbers in the legend refer to true positive rate (TPR on y-axis) at 10^{-3} false positive rate (FPR on x-axis). The plots also show the number of comparisons evaluated to generate them (top).

higher resolution images. Also, verifying faces by comparing them at multiple scales seems to help recognition.

(b)In particular, we demonstrate that although FT-all and MR-J are trained on same images, with the same loss, and similar size (25M vs 23M params.), the small resolution-specific streams operating at the top of MR-J greatly improve its recognition performance at all low resolutions. FT-all also allows us to show that an unmodified single ResNet model cannot optimally encode both low and high resolution features.

(c)MR-J(W) models slightly outperform MR-I models. This shows that joint training of multiresolution models enjoys an advantage over training independently, as they do not encode redundant information. MR-J(W) also slightly outperform MR-J. We propose that, apart from model complexity (3 times larger), the inability of a single network to optimally model scale variation is also a contributing factor.

(d) In our experiments, we observed that a model tuned for images of height 16 alone performs better than tuning multiple resolution-specific models for images of height ; 30. This is surprising, as we would expect appropriately tuned resolution-specific models to perform better! One probable reason is that the *effective resolution* of real images is influenced by other factors such as JPEG compression, motion blur etc., and the additional blur created by using a model tuned for a lower resolution assists in dealing with them. Check supp. material for more experiments.

(e) The difference between our best models, and FT-all increases with a drop in resolution. Also the performance of our MR-J model which shares parameters across all resolution drops in comparison to MR-J(W). This observation validates our method, as it shows that lower resolution images need separate models for optimal performance.

4.2 Identification

4.2.1 Setup

Given a face image from IJBC-*X*, this protocol asks which one of N (3531) identities does it belong to? Each of the N subjects in the gallery is represented by a set of high quality images. It is an important protocol resembling the operational work of law enforcement [25]. Moreover, it allows us to test test *multi-resolution pooling*, and *adaptive inference* for multi-resolution embeddings.

4.2.2 Results

Fig. 4.2 presents the percentage of probe images which had the ground truth (GT) in one of their top-10 predictions for each of our models and the baseline over various IJBC-*X*. From the figure, we observe that the resolution-specific embeddings MR-J(W) **quadruples** the performance of VGG2 for probes from IJBC-20, and **double** the baseline's performance for probes from IJBC-25. Similar to earlier experiment, SR+FT surpasses VGG2's performance by 44% and 13.5% for IJBC-35, IJBC-50 respectively.

We can observe that MR-I, and MR-J again outperform the baseline by 66% on probes from IJBC-35, and 22% on IJBC-50. Also, MR models' significantly better performance validates adaptive multiresolution inference.

4.3 Image set-based verification

4.3.1 Setup

This is the more common 1:1 verification protocol defined in IJB-C dataset [25]. In this setting, probe sets are compared with gallery sets. We sample relevant probe sets with more than 60% images of very low resolution (height;30) to perform this experiment.

4.3.2 Results

In the plots of Fig. 4.3, we show our results with probes containing increasing fractions of very low resolution images. The figure shows that the SR+FT outperforms VGG2, and FT-all, by 11%, 30% respectively, on probe sets with larger fraction of very low res images (0.8, 0.9). Their performances are comparable for probe sets with lower fractions (0.6, 0.7) of low res images, as SR+FT is unable to capitalize on the additional high-res information in the probe set. We show that both MR models outperform all other approaches with increasingly larger margins on probe sets with increasing fractions of low resolution images. Particularly, the MR-J, MR-J(W) models beat the baseline by 11.1%, 11.9%, 28.8%, 47.1% for probe sets with fraction of low resolution images greater than 0.6, 0.7, 0.8, and 0.9



Figure 4.2: **Performance on Identification.** The plots show CMC curves generated by classifying single image probes from IJBC-*X*, to one of a defined gallery of 3531 classes, each represented by a *set* of images. The plots for IJBC-20, and IJBC-25 show that MR at least **doubles** VGG2's performance. The plots for IJBC-35, and IJBC-50 show that SR+FT models perform much better VGG2. They also demonstrate that MR models surpass VGG2's performance by **66**% and **22**% respectively. Numbers in the legend show the percentage of probes with ground truth in the top-10 predictions. Number of probes in each subset are shown at the top of each plot.

respectively, proving that the MR models optimally combine both high-resolution and low-resolution features of images in the probe and reference sets.

4.4 Megaface

Megaface is a popular large-scale testing benchmark for face recognition. However, the dataset does not contain images of low resolutions. To test our method at this large scale, we resize all images in the Megaface dataset to specific sizes before evaluating our methods on these resized images. The table shows the Rank-1 accuracy obtained by our models and the baseline at various such sizes. All results are obtained by using a distractor set of 100K images. Table 4.1 shows that our multiresolution models continue to outperform the baseline models (VGG2, FT-all), and also the SR+FT models. However, note that the difference between SR+FT and MR-X is not high because the test images are artificially downsampled and the models may overfit to this downsampling method. More results are shown in



Figure 4.3: **Image set based verification.** These plots show TPR (y-axis) at specified FPRs (x-axis), for probe sets with varying ratios of low-resolution images. SR+FT outperforms VGG2 and FT-all at higher ratios (0.8, 0.9). MR models (particularly MR-J) outperform all other approaches with increasingly larger margins for higher ratios. Numbers in the legend refer to TPR (y-axis) at 10^{-3} FPR (x-axis). The plots also show the number of comparisons evaluated to generate them (top).

the supplementary material.

	Rank 1. Acc.					
Face height	MR-J(W)	MR-J	SR+FT	FT-all	VGG2	
20	40.1	38.9	40.1	32.0	15.9	
35	71.5	70.7	70.2	58.0	51.0	
50	79.2	77.5	77.4	64.3	65.2	

Table 4.1: The table shows that our multiresolution models continue to outperform the baseline models (VGG2, FT-all), and also the SR+FT models. However, note that the difference between SR+FT and MR-X is not high because the test images are artificially downsampled and the models may overfit to this downsampling method. More results are shown in the supplementary material.



Figure 4.4: We visualize the salient features captured by a 16px embedding by plotting both the lowres image pairs and their high-res counterparts. The top-left quadrant show face pairs of the same identity with high cosine similarity (shown at the bottom of each pair). The top right shows face of same identity with low similarity (due to expression or makeup changes). The bottom left mistakes suggest that the low res model relies heavily on racial and face shape cues, as different people from the same race are predicted to have high similarity. The bottom right suggests that gender appears to be an easily distinguishable feature at low resolution, previously observed in [33]

4.5 Qualitative Results

4.5.1 What features are captured by the embedding trained at lowest resolution (16px.)?

First, we begin with qualitative results for resolution-specific verification (Fig. 4.4). We visualize the salient features captured by a 16px embedding by plotting both the low-res image pairs and their high-res counterparts. The top-left quadrant show face pairs of the same identity with high cosine similarity (shown at the bottom of each pair). The top right shows face of same identity with low similarity (due to expression or makeup changes). The bottom left mistakes suggest that the low res model relies heavily on racial and face shape cues, as different people from the same race are predicted to have high similarity. The bottom right suggests that gender appears to be an easily distinguishable feature at low resolution, previously observed in [33].

4.5.2 What do multi-resolution embeddings encode in their resolution-specific parts?

In Section 3.2 of the submission, we state that jointly training multi-resolution models *disentangles* the features across scale, and prevents encoding redundant information. We support that claim with an experiment to identify which visual features are encoded in different resolution-specific parts of a jointly trained multi-resolution embedding (MR-J). For this experiment, we determine the nearest neighbors (cosine similarity) for faces from the LFW dataset using different resolution-specific components of MR-J and show these results in Fig. 4.5. In the figure, we only show nearest neighbors



Figure 4.5: Joint training forces multi-resolution models to learn complementary features. The figure shows an input image (Input) and its nearest neighbors of a different identity determined using, (a)LR:the low resolution part of a jointly trained multi-resolution embedding (MR-J), (b)HR: high resolution part of MR-J, (c)MR: full multi-resolution embedding, and (d)Normal: a normal baseline embedding. We find that the low resolution embeddings capture features like shape of face, race, gender, and distinguishable features like high cheek bones. On the other hand, the high resolution part downplays these features and emphasizes on complementary ones like shape of nose. The multi resolution embeddings combine both these features to return a set of nearest neighbors similar to those of a normal embedding.

from a different identity to best illustrate our results.

The figure shows that the different resolution-specific components of MR-J learn to encode complementary features. For example, the low resolution part of MR-J encodes features like face shape, race, gender, and easily distinguishable features like high cheek bones (see faces labeled LR). On the other hand, the high resolution part downplays these features and emphasizes complementary features observed in a high res image like shape of nose (see faces labeled HR). The multi-resolution embedding combines both these features to return a set of nearest neighbors (MR) similar to those of a normal embedding (Normal).

4.6 Off-the-shelf super-resolution

In Section 3.1 of the paper, we use a simpler super-resolution model for better generalization to real low resolution images. In this section, we substantiate this choice with both quantitative and qualitative results that compare our models to an off-the-shelf super-resolution network [37].

4.6.1 Quantitative

We support our choice of a light weight CNN for super-resolution by comparing the performance of our **SR** approach with an off-the-shelf face super-resolution network, **FSR** [37], on real low resolution images drawn from the IJBC dataset. The poor performance of the **FSR+VGG2** curve in Fig.4.6 shows that traditional deep super-resolution methods lose identity information when operating on low-resolution images. Further, it shows that these approaches generalize poorly to real images, as has been previously established by [5], and supported by our qualitative experiments.



Figure 4.6: **Off-the-shelf super resolution networks fails to preserve identity.** The figure shows that recognizing low resolution images by super-resolving them with an off-the-shelf super resolution network, and processing them with the baseline face recognition model (FSR+VGG2) performs worse than just using the baseline (VGG2). This indicates that super-resolution networks lose identity information when operating on real images. Also, we observe that our SR-16 massively outperforms the baseline and FSR+VGG2, demonstrating that our SR models generalize better to real low resolution images. We use OTS to denote 'off-the-shelf'.



Figure 4.7: We present high res images from the LFW dataset along with their low res (16x16) counterparts. The figure also shows super-resolved outputs of these low res images when passed through a deep off-the-shelf super-resolution network **FSR-OTS** [37] and our **SR**, **SR+FT** models. We observe that **FSR-OTS** hallucinates high resolution images but fails to preserve the identity of the person even with these artificially resized images.



Figure 4.8: The figure shows more examples of real low resolution images of different identities superresolved with various approaches. It emphasizes poor performance of the off-the-shelf super resolution networks (FSR-OTS) on real images.



Figure 4.9: The figure shows *real* low res images of the same identity from the IJBC dataset. The outputs of FSR-OTS show that (a)deep super-resolution networks generalize poorly to real low res images, and (b)the outputs do not retain identity information. In contrast, SR performs similarly on both real and artificially downsampled images (Fig. 4.7). From the outputs of SR+FT, and SR we observe that SR+FT sacrifices its ability to output sharper images to aid recognition. SR+FT primarily focuses on eliminating JPEG artifacts, and noise in the original image.

4.6.2 Qualitative

Fig. 4.7 illustrates that FSR-OTS fails to preserve the identity of the person even in artificially downsized images (evident from the George Bush, Jennifer Aniston images). In Figs. 4.9 and 4.8, we observe that off-the-shelf super-resolution networks fail to generalize to real images, and their inability to preserve identity becomes more evident. This leads to a drop in recognition performance as shown earlier. In comparison, our **SR** model performs similarly on both artificially down-sampled and real low-resolution images.

Also, comparing SR and SR+FT we see that SR+FT generates images which are ideally suited for our face recognition models tuned to low resolutions. Specifically, SR+FT focuses on eliminating noise and JPEG artifacts from real low resolution images.

Chapter 5

Additional Experiments

5.1 Multi-resolution pooling

In Section 3.3 of the submission, we argue that only high-resolution faces should be used to construct the pooled high-frequency feature, and proposed multi-resolution pooling as a solution (Eqn.4 and Fig.11 in the submission). We validate that claim in this section by comparing the performance of templates constructed with multi-resolution pooling, MR-I and MR-J, with templates constructed by naive pooling, MR-I(NaPo) and MR-J(NaPo) on the image set based verification described in Section 4.3 of the submission and show the results in Fig. 5.1.



Figure 5.1: **Impact of multiresolution pooling** The figure shows true positive rate (TPR) at 10^{-3} false positive rate for probe sets with at least 60% low resolution images. The results illustrate that naive pooling, MR-J(NaPo), and MR-I(NaPo), MR features perform worse than multi-resolution pooling, MR-J and MR-I. Interestingly, naive pooling performs worse than SR+FT for probes with large fraction of low resolution images (0.7, 0.8, 0.9), showing that it is sub-optimal to corrupt 'high resolution features' from high res images with 'high resolution features' extracted from low resolution images.

The figure shows the true positive rate (TPR) at 10^{-3} false positive rate for probe sets with at least 60% low resolution images. The results illustrate that naive pooling of our multi-resolution embeddings, MR-J(NaPo) and MR-I(NaPo), performs worse than multi-resolution pooling, MR-J and



Figure 5.2: The figure shows a relative comparison between our multi-resolution model and the baseline VGG2 on all comparisons in the IJBC single image verification protocol (Covariate verification). We find that our model induces a small overall improvement. Note that the improvement is small because of the smaller number of comparisons involving low resolution images.

MR-I. Also, we find that naive pooling performs worse than SR+FT for probes with larger fraction of low resolution images (0.7, 0.8, 0.9), showing that it is sub-optimal to corrupt 'high resolution features' from high res images with 'high resolution features' extracted from low resolution images.

5.2 IJB-C overall performance

In this section, we show a relative comparison between our model and the baseline over all comparisons in IJB-C single image verification protocol. Figure 5.2 shows that our model induces a small improvement in the overall performance. Please note that the improvement is small because of the small number of comparisons involving low resolution images.

5.3 Improvement is orthogonal to training loss

In this section, we show the improvement caused by our method is orthogonal to modifications in training loss by studying the effect of scale variation on state-of-the-art loss functions. Primarily, we use ArcFace loss [10](AF) for our experiments.

ArcFace, or Additive Angular Margin Loss was introduced by Deng et al. [10]. This paper is one of a group of works [10, 24, 34], which try to design appropriate loss functions to enhance the discriminative power of CNNs employed for large scale face recognition.

All results till this point were generated with CNNs trained with the widely used categorical cross entropy loss (softmax loss) function, given by

$$CE = -\frac{1}{N} \sum_{i=1}^{N} \log \frac{e^{W_{y_i}^T x_i + b_i}}{\sum_{j=1}^{n} e^{W_j^T x_i + b_j}}$$
(5.1)

where $x_i \in \mathbf{R}^d$ is the embedding output by the CNN for the *i*-th sample, belonging to the y_i -th class. $W_j \in \mathbf{R}^d$ is the *j*-th column of the weight $W \in \mathbf{R}^d$ and $b_j \in \mathbf{R}^n$. N is the bacth size, and n is the number of classes

In the 2D embeddings in Fig. 5.3a, training a CNN with this loss results in large intra-class variance and almost no inter-class separation. This is not a problem for classification where we directly predict the class of an input image. But, for face recognition the CNN operates on previously unseen faces and the cosine distance between the output embeddings is used to predict the class of the input faces.

ArcFace loss was proposed as a modification of softmax loss Eqn.5.1 to explicitly minimize intraclass variance and maximize the inter-class separation. Specifically the bias b_j is fixed to 0. Now, we write the logit $W_j^T x_i = ||W_j|| ||x_i|| \cos\theta_j$, where θ_j is the angle between the weight W_j and x_i . Here we set both $||W_j|| = 1$ and $||x_i|| = 1$ by simple L_2 normalization. Therefore,

$$W_j^T x_i = \cos\theta_j \tag{5.2}$$

By enforcing these constraints, we force the embeddings of training images x_{y_i} belonging to class y_i to be distributed around its corresponding column weight W_{y_i} on a hyper sphere of radius 1.

In this setting, the ArcFace loss adds an additive angular margin penalty m between x_{y_i} and W_{y_i} . This margin penalty reduces the intra-class variance while also boosting the inter-class separation. This loss function is given by,

$$AF = -\frac{1}{N} \sum_{i=1}^{N} \log \frac{e^{s(\cos(\theta_{y_i} + m))}}{e^{s(\cos(\theta_{y_i} + m))} + e^{\sum_{j=1, j \neq y_i}^{n} s\cos(\theta_j)}}$$
(5.3)

In Fig. **??**, we can observe that the ArcFace loss forces the 2D embeddings to be much more compact (reducing intra-class variance), while enhancing the inter-class separation.



Figure 5.3: The figures show 2D embeddings output by a CNN trained to distinguish between 8 face identities (classes). Each dot represents an embedding of a training image, and its color represents the class of the image. The lines from the center of the circle to the cluster centres depict the direction of column weights W_{y_i} . On the left, we can observe that the embeddings output by a CNN trained on softmax loss have large intra-class variance and low inter-class separation. On the right, we can observe that the ArcFace loss explicitly optimizes for these two desirable properties.

Now, Figs. 5.4 and 5.5 demonstrate show a relative comparison between our models and baselines trained with the two loss functions. The translucent bars (hatchet with positive slope) show the performance of our models trained with ArcFace loss [10], while the solid bars show the performance of our models trained with traditional cross-entropy loss.



Figure 5.4: The figure shows a relative comparison between our models and baselines trained with different loss functions. The translucent bars (hatchet with positive slope) show the performance of our models trained with ArcFace loss (AF), while the solid bars show the performance of our models trained with traditional cross-entropy loss(CE). We use the same training dataset to train all models. The results show that even state-of-the-art models suffer a marked drop in performance at low resolutions (Base, trained with ArcFace). The figure shows that our method massively improves the performance of even these state-of-the-art models, indicating that this improvement is caused by novel modifications to the network architecture and training scheme. Specifically, the figure compares the performance of various models under the Single image verification protocol described earlier. It also shows that training our multi-resolution models with ArcFace improves their performance compared to cross-entropy loss.



Figure 5.5: Similar to Fig. 5.4, this figure shows a relative comparison between our multi-resolution models and baselines trained with different loss functions. This figure shows the performance of all these models under the Identification protocol described earlier. We had observed earlier that multi-resolution models massively impact the performance of networks trained with traditional cross-entropy loss. This figure shows that they have a similar impact on even state-of-the-art models (compare Base ArcFace, MR-J ArcFace). As argued earlier, this demonstrates that the improvement caused by our models is orthogonal to the influence of modifications in loss function

The results in Fig. 5.4 show that even state-of-the-art models suffer a marked drop in performance at low resolutions (Base, trained with ArcFace). The figure shows that our method massively improves the performance of even these state-of-the-art models, indicating that this improvement is caused by novel modifications to the network architecture and training scheme. Specifically, the figure compares the performance of various models under the Single image verification protocol described earlier. It also shows that training our multi-resolution models with ArcFace improves their performance compared to cross-entropy loss.

We show similar results in Fig. 5.5, in which we compare the performance of the various models in the Identification protocol described earlier.

From these results, we can clearly observe that even state-of-the-art models suffer a marked drop in performance at low resolutions. Both figures show that our method to train models to output multi-resolution embeddings massively improves the performance of such models. This indicates that the improvement caused by our model at low resolutions is caused by novel modifications to the traditional way of training CNNs.

Chapter 6

Conclusion

We propose a simple yet effective approach for recognizing faces at low resolution. We first point out that state-of-the-art face recognizers, which use fixed-resolution embeddings, perform dramatically worse as face resolution drops below 30 pixels. We then show that by simply tuning resolution-specific embedding we can significantly improve the recognition accuracy. We further explore multi-resolution embedding that efficiently adapts in size and complexity to the resolution of test image *on-the-fly*. Finally, comparing to state-of-the-art fixed-resolution embeddings, our proposed embedding dramatically reduces recognition error on small faces on standard benchmarks.

In the future, we would like to test this approach on more recently proposed QMUL-SurvFace dataset. This dataset was deveised to test the performance of face recognition models on low resolution face images extracted from surveillance footage, and to facilitate more research in this direction. Positive results on this dataset would go a long way in cementing the effectiveness of our approach.

Further, we would also like to extend our approach to more generic computer vision tasks like object detection, and image retrieval.

Object detection: Small object detection is one of the major challenges facing the vision community. Despite the significant progress caused by CNN-based object detection models, there is a large disparity in their performance between large and small objects. This could be due to , 1. they do not occur frequently, and 2. obviously, their small size.

We face the same two issues in recognizing tiny faces. As such, we believe that our data augmentation scheme, and multi-resolution embeddings will help solve the challenging problem of small object detection.

2. **Image retrieval:** Face recognition, as performed today, is very similar to the task of image retrieval and few-shot recognition tasks. We can effectively use our scheme to retrieve other instances of an object captured at low resolution using our approach.

Also, in our experiments, we observed that different parts of a multi-resolution embedding capture different scale-specific features. This property can be useful to retrieve images with specific attributes.

Bibliography

- E. H. Adelson, C. H. Anderson, J. R. Bergen, P. J. Burt, and J. M. Ogden. Pyramid methods in image processing. *RCA engineer*, 1984.
- [2] S. Baker and T. Kanade. Hallucinating faces. In Automatic Face and Gesture Recognition, 2000. Proceedings. Fourth IEEE International Conference on, pages 83–88. IEEE, 2000.
- [3] S. Baker and T. Kanade. Limits on super-resolution and how to break them. IEEE Transactions on Pattern Analysis and Machine Intelligence, 24(9):1167–1183, 2002.
- [4] S. Biswas, K. W. Bowyer, and P. J. Flynn. Multidimensional scaling for matching low-resolution face images. IEEE transactions on pattern analysis and machine intelligence, 34(10):2019–2030, 2012.
- [5] A. Bulat, J. Yang, and G. Tzimiropoulos. To learn image super-resolution, use a gan to learn how to do image degradation first. *arXiv preprint arXiv:1807.11458*, 2018.
- [6] A. M. Burton, S. Wilson, M. Cowan, and V. Bruce. Face recognition in poor-quality video: Evidence from security surveillance. *Psychological Science*, 10(3):243–248, 1999.
- [7] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman. Vggface2: A dataset for recognising faces across pose and age. arXiv preprint arXiv:1710.08092, 2017.
- [8] J. Y. Choi, Y. M. Ro, and K. N. Plataniotis. Color face recognition for degraded face images. IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics), 39(5):1217–1230, 2009.
- [9] S. Chopra, R. Hadsell, Y. LeCun, et al. Learning a similarity metric discriminatively, with application to face verification. In CVPR (1), pages 539–546, 2005.
- [10] J. Deng, J. Guo, N. Xue, and S. Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. arXiv preprint arXiv:1801.07698, 2018.
- [11] L. D. Harmon. The recognition of faces. *Scientific American*, 229(5):70–83, 1973.
- [12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 770–778, 2016.
- [13] P. H. Hennings-Yeomans, S. Baker, and B. V. Kumar. Simultaneous super-resolution and feature extraction for recognition of low-resolution faces. In *Computer Vision and Pattern Recognition*, 2008. CVPR 2008. IEEE Conference on, pages 1–8. IEEE, 2008.
- [14] P. H. Hennings-Yeomans, B. V. Kumar, and S. Baker. Robust low-resolution face identification and verification using high-resolution features. In *Image Processing (ICIP)*, 2009 16th IEEE International Conference on, pages 33–36. IEEE, 2009.
- [15] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical report, Technical Report 07-49, University of Massachusetts, Amherst, 2007.
- [16] M. Iacopo, S. Rawls, G. Medioni, and P. Natarajan. Pose-aware face recognition in the wild. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016.
- [17] P. Jonathon, A. Yates, Y. Hu, C. Hahn, E. Noyes, K. Jackson, and J. Cavazos. Face recognition accuracy of forensic examiners, superrecognizers, and face recognition algorithms. *Proceedings of the National Academy* of Sciences, 2018.
- [18] Z. Kaipeng, Z. Zhang, C.-W. Cheng, W. H. Hsu, Y. Qiao, W. Liu, and T. Zhang. Super-identity convolutional neural network for face hallucination. *ECCV*, 2018.

- [19] G. Koch, R. Zemel, and R. Salakhutdinov. Siamese neural networks for one-shot image recognition. In *ICML Deep Learning Workshop*, volume 2, 2015.
- [20] T.-Y. Lin, P. DollÂąr, R. B. Girshick, K. He, B. Hariharan, and S. J. Belongie. Feature pyramid networks for object detection. CVPR, 2017.
- [21] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [22] C. Liu, H.-Y. Shum, and W. T. Freeman. Face hallucination: Theory and practice. International Journal of Computer Vision, 75(1):115–134, 2007.
- [23] C. Liu, H.-Y. Shum, and C.-S. Zhang. A two-step approach to hallucinating faces: global parametric model and local nonparametric model. In *Computer Vision and Pattern Recognition*, 2001. CVPR 2001. Proceedings of the 2001 IEEE Computer Society Conference on, volume 1, pages I–I. IEEE, 2001.
- [24] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song. Sphereface: Deep hypersphere embedding for face recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, 2017.
- [25] B. Maze, J. Adams, J. A. Duncan, N. Kalka, T. Miller, C. Otto, A. K. Jain, W. T. Niggel, J. Anderson, J. Cheney, et al. Iarpa janus benchmark-c: Face dataset and protocol.
- [26] O. M. Parkhi, A. Vedaldi, A. Zisserman, et al. Deep face recognition. In *BMVC*, volume 1, page 6, 2015.
- [27] X. Peng, J. Hoffman, X. Y. Stella, and K. Saenko. Fine-to-coarse knowledge transfer for low-res image classification. In *Image Processing (ICIP), 2016 IEEE International Conference on*, pages 3683–3687. IEEE, 2016.
- [28] R. Ranjan, C. Castillo, and R. Chellappa. L2-constrained softmax loss for discriminative face verification. arXiv preprint arXiv:1703.09507, 2017.
- [29] J. Sadr, I. Jarudi, and P. Sinha. The role of eyebrows in face recognition. *Perception*, 32(3):285–293, 2003.
- [30] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015.
- [31] B. Singh and L. S. Davis. An analysis of scale invariance in object detection snip. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 3578–3587, 2018.
- [32] P. Sinha, B. Balas, Y. Ostrovsky, and R. Russell. Face recognition by humans: Nineteen results all computer vision researchers should know about. *Proceedings of the IEEE*, 94(11):1948–1962, 2006.
- [33] S. Tamura, H. Kawai, and H. Mitsumoto. Male/female identification from 8x6 very low resolution face images by neural network. *Pattern recognition*, pages 331–335, 1996.
- [34] H. Wang, Y. Wang, Z. Zhou, X. Ji, D. Gong, J. Zhou, Z. Li, and W. Liu. Cosface: Large margin cosine loss for deep face recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5265–5274, 2018.
- [35] Z. Wang, Z. Miao, Q. J. Wu, Y. Wan, and Z. Tang. Low-resolution face recognition: a review. *The Visual Computer*, 30(4):359–386, 2014.
- [36] Y. Wen, K. Zhang, Z. Li, and Y. Qiao. A discriminative feature learning approach for deep face recognition. In *European conference on computer vision*, pages 499–515. Springer, 2016.
- [37] C. Yu, Y. Tai, X. Liu, C. Shen, and J. Yang. Fsrnet: End-to-end learning face super-resolution with facial priors. pages 2492–2501, 2018.
- [38] E. Zangeneh, M. Rahmati, and Y. Mohsenzadeh. Low resolution face recognition using a two-branch deep convolutional neural network architecture. *arXiv preprint arXiv:1706.06247*, 2017.
- [39] H. Zhao, S. Jianping, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. CVPR, 2017.
- [40] W. W. Zou and P. C. Yuen. Very low resolution face recognition problem. IEEE Transactions on Image Processing, 21(1):327–340, 2012.