# Deep Feature Interpolation for Image Content Changes

Paul Upchurch[1,*]    Jacob Gardner[1,*]    Geoff Pleiss[1]    Robert Pless[2]    Noah Snavely[1]    Kavita Bala[1]
Kilian Weinberger[1]

[1]Cornell University
[2]George Washington University
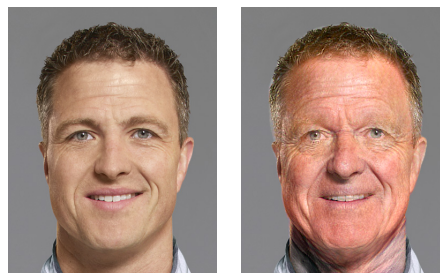[*]Authors contributed equally

## Abstract

*We propose* Deep Feature Interpolation (DFI)*, a new data-driven baseline for automatic high-resolution image transformation. As the name suggests, DFI relies only on simple linear interpolation of deep convolutional features from pre-trained convnets. We show that despite its simplicity, DFI can perform high-level semantic transformations like "make older/younger", "make bespectacled", "add smile", among others, surprisingly well—sometimes even matching or outperforming the state-of-the-art. This is particularly unexpected as DFI requires no specialized network architecture or even any deep network to be trained for these tasks. DFI therefore can be used as a new baseline to evaluate more complex algorithms and provides a practical answer to the question of which image transformation tasks are still challenging after the advent of deep learning.*

## 1. Introduction

Generating believable changes in images is an active and challenging research area in computer vision and graphics. Until recently, algorithms were typically hand-designed for individual transformation tasks and exploited task-specific expert knowledge. Examples include transformations of human faces [41, 17], materials [2, 1], color [50], or seasons in outdoor images [23]. However, recent innovations in deep convolutional auto-encoders [33] have produced a succession of more versatile approaches. Instead of designing each algorithm for a specific task, a conditional (or adversarial) generator [21, 13] can be trained to produce a set of possible image transformations through supervised learning [48, 43, 52]. Although these approaches can perform a variety of seemingly impressive tasks, in this paper we show that a surprisingly large set of them can be solved via linear interpolation in deep feature space and may not require specialized deep architectures.

How can linear interpolation be effective? In pixel space, natural images lie on an (approximate) non-linear manifold [44]. Non-linear manifolds are locally Euclidean, but



Input          Older

Figure 1. Aging a face with DFI.

globally curved and non-Euclidean. It is well known that in pixel space linear interpolation between images introduces ghosting artifacts, a sign of departure from the underlying manifold, and linear classifiers between image categories perform poorly.

On the other hand, deep convolutional neural networks (convnets) are known to excel at classification tasks such as visual object categorization [38, 14, 15]—while relying on a simple linear layer at the end of the network for classification. These linear classifiers perform well because networks map images into new representations in which image classes are *linearly* separable. In fact, previous work has shown that neural networks that are trained on sufficiently diverse object recognition classes, such as VGG [38] trained on ImageNet [22], learn surprisingly versatile feature spaces and can be used to train linear classifiers for additional image classes. Bengio *et al.* [3] hypothesize that convnets linearize the manifold of natural images into a (globally) Euclidean subspace of deep features.

Inspired by this hypothesis, we argue that, in such deep feature spaces, some image editing tasks may no longer be as challenging as previously believed. We propose a simple framework that leverages the notion that in the right feature space, image editing can be performed simply by linearly interpolating between images with a certain attribute and images without it. For instance, consider the task of adding facial hair to the image of a male face, given two sets of images: one set with facial hair, and one set without. If con-
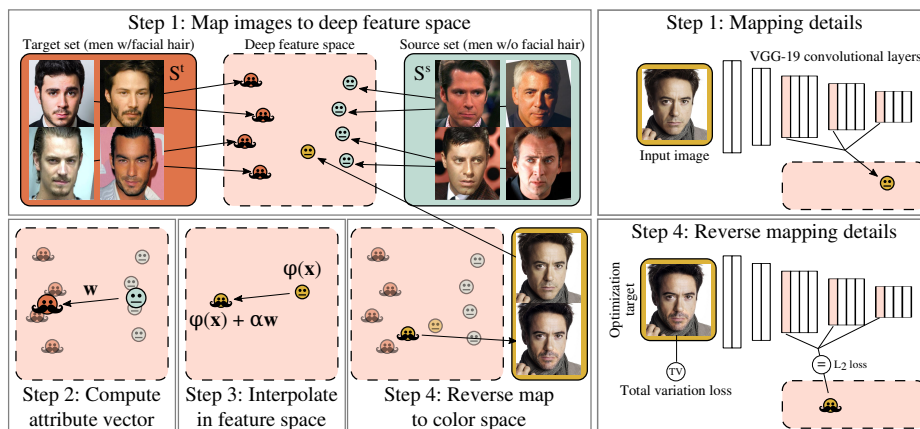
1

Deep Feature Interpolation



Figure 2. A schematic outline of the four high-level DFI steps.

vnets can be trained to distinguish between male faces with facial hair and those without, we know that these classes must be linearly separable. Therefore, motion along a single linear vector should suffice to move an image from deep features corresponding to "no facial hair" to those corresponding to "facial hair". Indeed, we will show that even a simple choice of this vector suffices: we average convolutional layer features of each set of images and take the difference.

We call this method Deep Feature Interpolation (DFI). Figure 1 shows an example of a facial transformation with DFI on a $390 \times 504$ image.

Of course, DFI has limitations: our method works best when all images are aligned, and thus is suited when there are feature points to line up (e.g. eyes and mouths in face images). It also requires that the sample images with and without the desired attribute are otherwise similar to the target image (e.g. in the case of Figure 2, the other images should contain Caucasian males).

However, these assumptions on the data are comparable to what is typically used to train generative models, and in the presence of such data DFI works surprisingly well. We demonstrate its efficacy on several transformation tasks commonly used to evaluate generative approaches. Compared to prior work, it is much simpler, and often faster and more versatile: It does not require re-training a convnet, is not specialized on any particular task, and it is able to process much higher resolution images. Despite its simplicity we show that on many of these image editing tasks it outperforms state-of-the-art methods that are substantially more involved and specialized.

## 2. Related Work

Probably the generative methods most similar to ours are [24] and [32], which similarly generate data-driven at-

tribute transformations using deep feature spaces. We use these methods as our primary points of comparison; however, they rely on specially trained generative auto-encoders and are fundamentally different from our approach to learning image transformations. Works by Reed *et al*. [33, 34] propose content change models for challenging tasks (identity and viewpoint changes) but do not demonstrate photo-realistic results. A contemporaneous work [4] edits image content by manipulating latent space variables. However, this approach is limited by the output resolution of the underlying generative model. An advantage of our approach is that it works with pre-trained networks and has the ability to run on much higher resolution images. In general, many other uses of generative networks are distinct from our problem setting [13, 5, 51, 37, 30, 6, 8], as they deal primarily with generating novel images rather than changing existing ones.

Gardner *et al*. [9] edits images by minimizing the witness function of the Maximum Mean Discrepancy statistic. The memory needed to calculate the transformed image's features by their method grows linearly whereas DFI removes this bottleneck.

Mahendran and Vedaldi [28] recovered visual imagery by inverting deep convolutional feature representations. Gatys *et al*. [11] demonstrated how to transfer the artistic style of famous artists to natural images by optimizing for feature targets during reconstruction. Rather than reconstructing imagery or transferring style, we edit the content of an existing image while seeking to preserve photo-realism and all content unrelated to the editing operation.

Many works have used vector operations on a learned generative latent space to demonstrate transformative effects [7, 32, 12, 46]. In contrast, we suggest that vector operations on a discriminatively-trained feature space can achieve similar effects.

In concept, our work is similar to [41, 10, 42, 19, 17] that use video or photo collections to transfer the personality and

character of one person's face to a different person (a form of puppetry [39, 45, 20]). This difficult problem requires a complex pipeline to achieve high quality results. For example, Suwajanakorn *et al.* [41] combine several vision methods: fiducial point detection [47], 3D face reconstruction [40] and optical flow [18]. Our method is less complicated and applicable to other domains (e.g., product images of shoes).

While we do not claim to cover all the cases of the techniques above, our approach is surprisingly powerful and effective. We believe investigating and further understanding the reasons for its effectiveness would be useful for better design of image editing with deep learning.

## 3. Deep Feature Interpolation

In our setting, we are provided with a test image $\mathbf{x}$ which we would like to change in a believable fashion with respect to a given attribute. For example, the image could be a man without a beard, and we would like to modify the image by adding facial hair while preserving the man's identity. We further assume access to a set of *target* images *with* the desired attribute $\mathcal{S}^t = \{\mathbf{x}_1^t, ..., \mathbf{x}_n^t\}$ (e.g., *men with facial hair*) and a set of *source* images *without* the attribute $\mathcal{S}^s = \{\mathbf{x}_1^s, ..., \mathbf{x}_m^s\}$ (e.g., *men without facial hair*). Further, we are provided with a pre-trained convnet trained on a sufficiently rich object categorization task—for example, the openly available VGG network [38] trained on ImageNet [35]. We can use this convnet to obtain a new representation of an image, which we denote as $\mathbf{x} \rightarrow \phi(\mathbf{x})$. The vector $\phi(\mathbf{x})$ consists of concatenated activations of the convnet when applied to image $\mathbf{x}$. We refer to it as the *deep feature representation* of $\mathbf{x}$.

**Deep Feature Interpolation** can be summarized in four high-level steps (illustrated in Figure 2):

1. We map the images in the target and source sets $\mathcal{S}^t$ and $\mathcal{S}^s$ into the deep feature representation through the pre-trained convnet $\phi$ (e.g., VGG-19 trained on ILSVRC2012).

2. We compute the mean feature values for each set of images, $\bar{\phi}^t$ and $\bar{\phi}^s$, and define their difference as the *attribute vector*

$$\mathbf{w} = \bar{\phi}^t - \bar{\phi}^s. \quad (1)$$

3. We map the test image $\mathbf{x}$ to a point $\phi(\mathbf{x})$ in deep feature space and move it along the attribute vector $\mathbf{w}$, resulting in $\phi(\mathbf{x}) + \alpha\mathbf{w}$.

4. We can reconstruct the transformed output image $\mathbf{z}$ by solving the reverse mapping into pixel space w.r.t. $\mathbf{z}$

$$\phi(\mathbf{z}) = \phi(\mathbf{x}) + \alpha\mathbf{w}. \quad (2)$$

Although this procedure may appear deceptively simple, we show in Section 4.2 that it can be surprisingly effective. In

the following we will describe some important details to make the procedure work in practice.

**Selecting $\mathcal{S}^t$ and $\mathcal{S}^s$.** DFI assumes that the attribute vector $\mathbf{w}$ isolates the targeted transformation, i.e., it points towards the deep feature representation of image $\mathbf{x}$ with the desired attribute change. If such an image $\mathbf{z}$ was available (e.g., the same image of Mr. Robert Downey Jr. with beard), we could compute $\mathbf{w} = \phi(\mathbf{z}) - \phi(\mathbf{x})$ to isolate exactly the difference induced by the change in attribute. In the absence of the exact target image, we estimate $\mathbf{w}$ through the target and source sets. It is therefore important that both sets are as similar as possible to our test image $\mathbf{x}$ and there is no systematic attribute bias across the two data sets. If, for example, all target images in $\mathcal{S}^t$ were images of more senior people and source images in $\mathcal{S}^s$ of younger individuals, the vector $\mathbf{w}$ would unintentionally capture the change involved in aging. Also, if the two sets are too different from the test image (e.g., a different race) the transformation would not look believable. To ensure sufficient similarity we restrict $\mathcal{S}^t$ and $\mathcal{S}^s$ to the $K$ nearest neighbors of $\mathbf{x}$. Let $\mathcal{N}_K^t$ denote the $K$ nearest neighbors of $\mathcal{S}^t$ to $\phi(\mathbf{x})$; we define

$$\bar{\phi}^t = \frac{1}{K} \sum_{\mathbf{x}^t \in \mathcal{N}_K^t} \phi(\mathbf{x}^t) \text{ and } \bar{\phi}^s = \frac{1}{K} \sum_{\mathbf{x}^s \in \mathcal{N}_K^s} \phi(\mathbf{x}^s). \quad (3)$$

These neighbors can be selected in two ways, depending on the amount of information available. When attribute labels are available, we find the nearest images by counting the number of matching attributes (e.g., matching gender, race, age, hair color). When attribute labels are unavailable, or as a second selection criterion, we take the nearest neighbors by cosine distance in deep feature space.

**Deep feature mapping.** There are many choices for a mapping into deep feature space $\mathbf{x} \rightarrow \phi(\mathbf{x})$. We use the convolutional layers of the normalized VGG-19 network pre-trained on ILSVRC2012, which has proven to be effective at artistic style transfer [11]. The deep feature space must be suitable for two very different tasks: (1) linear interpolation and (2) reverse mapping back into pixel space. For the interpolation, it is advantageous to pick deep layers that are further along the linearization process of deep convnets [3]. In contrast, for the reverse mapping, earlier layers capture more details of the image [28]. The VGG network is divided into five pooling regions (with increasing depth). As an effective compromise we pick the first layers from the last three regions, `conv3_1`, `conv4_1` and `conv5_1` layers (after ReLU activations), flattened and concatenated. As the pooling layers of VGG reduce the dimensionality of the input image, we *increase* the image resolution of small images to be at least $200 \times 200$ before applying $\phi$.

**Image transformation.** Due to the ReLU activations used in most convnets (including VGG), all dimensions in $\phi(\mathbf{x})$

are non-negative and the vector is sparse. As we average over $K$ images (instead of a single image as in [3]), we expect $\bar{\phi}^t, \bar{\phi}^s$ to have very small components in most features. As the two data sets $\mathcal{S}^t$ and $\mathcal{S}^s$ only differ in the target attribute, features corresponding to visual aspects unrelated to this attribute will be averaged to very small values and approximately subtracted away in the vector $\mathbf{w}$.

**Reverse mapping.** The final step of DFI is to reverse map the vector $\phi(\mathbf{x}) + \alpha\mathbf{w}$ back into pixel space to obtain an output image $\mathbf{z}$. Intuitively, $\mathbf{z}$ is an image that corresponds to $\phi(\mathbf{z}) \approx \phi(\mathbf{x}) + \alpha\mathbf{w}$ when mapped into deep feature space. Although no closed-form inverse function exists for the VGG mapping, we can obtain a color image by adopting the approach of [28] and find $\mathbf{z}$ with gradient descent:

$$\mathbf{z} = \arg\min_{\mathbf{z}} \frac{1}{2}\|(\phi(\mathbf{x})+\alpha\mathbf{w})-\phi(\mathbf{z})\|_2^2 + \lambda_{V^\beta} R_{V^\beta}(\mathbf{z}), \quad (4)$$

where $R_{V^\beta}$ is the Total Variation regularizer [28] which encourages smooth transitions between neighboring pixels,

$$R_{V^\beta}(\mathbf{z}) = \sum_{i,j} \left( (z_{i,j+1} - z_{i,j})^2 + (z_{i+1,j} - z_{i,j})^2 \right)^{\frac{\beta}{2}} \quad (5)$$

Here, $z_{i,j}$ denotes the pixel in location $(i,j)$ in image $\mathbf{z}$. Throughout our experiments, we set $\lambda_{V^\beta} = 0.001$ and $\beta = 2$. We solve (4) with the standard hill-climbing algorithm L-BFGS [26].

# 4. Experimental Results

We evaluate DFI on a variety of tasks and data sets. For perfect reproducibility our code is available at https://github.com/paulu/deepfeatinterp.

## 4.1. Changing Face Attributes

We compare DFI to AEGAN [24], a generative adversarial autoencoder, on several face attribute modification tasks. Similar to DFI, AEGAN also makes changes to faces by vector operations in a feature space. We use the Labeled Faces in the Wild (LFW) data set, which contains 13,143 images of faces with predicted annotations for 73 different attributes (e.g., SUNGLASSES, GENDER, ROUND FACE, CURLY HAIR, MUSTACHE, etc.). We use these annotations as attributes for our experiments. We chose six attributes for testing: SENIOR, MOUTH SLIGHTLY OPEN, EYES OPEN, SMILING, MOUSTACHE and EYEGLASSES. (The negative attributes are YOUTH, MOUTH CLOSED, NARROW EYES, FROWNING, NO BEARD, NO EYEWEAR.) These attributes were chosen because it would be plausible for a single person to be changed into having each of those attributes. Our test set consists of 38 images that did not have any of the six target attributes, were not WEARING HAT, had MOUTH CLOSED, NO BEARD and NO EYEWEAR. As LFW is highly gender imbalanced,

| older | mouth open | eyes open | smiling | moustache | glasses |
|-------|------------|-----------|---------|-----------|---------|
| 4.57  | 7.09       | 17.6      | 20.6    | 24.5      | 38.3    |

Table 1. Perceptual study results. Each column shows the ratio at which workers preferred DFI to AEGAN on a specific attribute change (see Figure 3 for images).

we only used images of the more common gender, men, as target, source, and test images.

Matching the approach of [24], we align the face images and crop the outer pixels leaving a $100 \times 100$ face image, which we resize to $200 \times 200$. Target (source) collections are LFW images which have the positive (negative) attributes. From each collection we take the $K = 100$ nearest neighbors (by number of matching attributes) to form $\mathcal{S}^t$ and $\mathcal{S}^s$.

We empirically find that scaling $\mathbf{w}$ by its mean squared feature activation makes the free parameter somewhat more consistent across multiple attribute transformations. If $d$ is the dimensionality of $\phi(\mathbf{x})$ and $pow$ is applied element-wise then we define

$$\alpha = \frac{\beta}{\frac{1}{d}pow(\mathbf{w}, 2)}. \quad (6)$$

We set $\beta = 0.4$.

Comparisons are shown in Figure 3. Looking down each column, we expect each image to express the target attribute. Looking across each row, we expect to see that the identity of the person is preserved. Although AEGAN often produces the right attributes, it does not preserve identity as well as the much simpler DFI.

**Perceptual Study.** Judgments of visual image changes are inherently subjective. To obtain an objective comparison between DFI and AEGAN we conducted a blind perceptual study with Amazon Mechanical Turk workers. We asked workers to pick the image which best expresses the target attribute while preserving the identity of the original face. This is a nuanced task so we required workers to complete a tutorial before participating in the study. The task was a forced choice between AEGAN and DFI (shown in random order) for six attribute changes on 38 test images. We collected an average of 29.6 judgments per image from 136 unique workers and found that DFI was preferred to AEGAN by a ratio of 12:1. The least preferred transformation was Senior at 4.6:1 and the most preferred was Eyeglasses at 38:1 (see Table 1).

## 4.2. High Resolution Aging and Facial Hair

One of the major benefits of DFI over many generative models is the ability to run on high resolution images. However, there are several challenges in presenting results on high resolution faces.

First, we need a high-resolution dataset from which to select $\mathcal{S}^s$ and $\mathcal{S}^t$. We collect a database of 100,000

Figure 3. **(Zoom in for details.)** Adding different attributes to the same person (random test images). **Left.** Original image. **Middle.** DFI. **Right.** AEGAN. The goal is to add the specified attribute while preserving the identity of the original person. For example, when adding a moustache to Ralf Schumacher (3rd row) the hairstyle, forehead wrinkle, eyes looking to the right, collar and background are all preserved by DFI. No foreground mask or human annotation was used to produce these test results.

high resolution face images from existing computer vision datasets (CelebA, MegaFace, and Helen) and Google image search [27, 29, 25]. We augment existing datasets, selecting only clear, unobstructed, front-facing high-resolution faces. This is different from many existing datasets which may have noisy and low-resolution images.

Next, we need to learn the attributes of the images present in the face dataset to properly select source and target images. Because a majority of images we collect do not have labels, we use face attribute classifiers developed using labeled data from LFW and CelebA.

Finally, the alignment of dataset images to the input image needs to be as close as possible, as artifacts that result from poor alignment are more obvious at higher resolutions. Instead of aligning our dataset as a preprocessing step, we

use an off-the-shelf face alignment tool in DLIB [16] to align images in $\mathcal{S}^s$ and $\mathcal{S}^t$ to the input image at test time.

We demonstrate results on editing megapixel faces for the tasks of aging and adding facial hair on three different faces. Due to the size of these images, selected results are shown in Figure 5. For full tables of results on these tasks, please see the supplementary materials.

### 4.3. Inpainting Without Attributes

Inpainting fills missing regions of an image with plausible pixel values. There can be multiple correct answers. Inpainting is hard when the missing regions are large (see Figure 4 for our test masks). Since attributes cannot be predicted (e.g., eye color when both eyes are missing) we use distance in feature space to select the nearest neighbors.

Figure 4. (**Zoom in for details.**) Filling missing regions. **Top.** LFW faces. **Bottom.** UT Zappos50k shoes. Inpainting is an interpolation from masked to unmasked images. Given any dataset we can create a source and target pair by simply masking out the missing region. DFI uses $K = 100$ such pairs derived from the nearest neighbors (excluding test images) in feature space. The face results match wrinkles, skin tone, gender and orientation (compare noses in 3rd and 4th images) but fail to fill in eyeglasses (3rd and 11th images). The shoe results match style and color but exhibit silhouette ghosting due to misalignment of shapes. Supervised attributes were not used to produce these results. For the curious, we include the source image but we note that the goal is to produce a plausible region filling—not to reproduce the source.

Inpainting may seem like a very different task from changing face attributes, but it is actually a straightforward application of DFI. All we need are source and target pairs which differ only in the missing regions. Such pairs can be generated for any dataset by taking an image and masking out the same regions that are missing in the test image. The images with mask become the source set and those without the target set. We then find the $K = 100$ nearest neighbors in the masked dataset (excluding test images) by cosine distance in VGG-19 pool5 feature space. We experiment on two datasets: all of LFW (13,143 images, including male and female images) and the Shoes subset of UT Zappos50k (29,771 images) [49, 31]. For each dataset we find a single $\beta$ that works well (1.6 for LFW and 2.8 for UT Zappos50k).

We show our results in Figure 4 on 12 test images (more in supplemental) which match those used by disCVAE [48] (see Figure 6 of their paper). Qualitatively we observe that the DFI results are plausible. The filled face regions match skin tone, wrinkles, gender, and pose. The filled shoe regions match color and shoe style. However, DFI failed to produce eyeglasses when stems are visible in the input and some shoes exhibit ghosting since the dataset is not perfectly aligned. DFI performs well when the face is missing (i.e., the central portion of each image) but we found it performs worse than disCVAE when half of the image is missing (Figure 8). Overall, DFI works surprisingly well on these

inpainting tasks. The results are particularly impressive considering that, in contrast to disCVAE, it does not require attributes to describe the missing regions.

## 4.4. Varying the free parameters

Figure 6 illustrates the effect of changing $\beta$ (strength of transformation) and $K$ (size of source/target sets). As $\beta$ increases, task-related visual elements change more noticeably (Figure 7). If $\beta$ is low then ghosting can appear. If $\beta$ is too large then the transformed image may jump to a point in feature space which leads to an unnatural reconstruction. $K$ controls the variety of images in the source and target sets. A lack of variety can create artifacts where changed pixels do not match nearby unchanged pixels (e.g., see the lower lip, last row of Figure 6). However, too much variety can cause $\mathcal{S}^t$ and $\mathcal{S}^s$ to contain distinct subclasses and the set mean may describe something unnatural (e.g., in the first row of Figure 6 the nose has two tips, reflecting the presence of left-facing and right-facing subclasses). In practice, we pick an $\beta$ and $K$ which work well for a variety of images and tasks rather than choosing per-case.

## 5. Discussion

In the previous section we have shown that Deep Feature Interpolation is surprisingly effective on several image

Figure 5. (**Zoom in for details.**) Editing megapixel faces. **First column.** Original image. **Right columns.** The top 3 rows show aging ($\beta = \{0.15, 0.25\}$) and the bottom 3 rows show the addition of facial hair ($\beta = \{0.21, 0.31\}$). High resolution images are challenging since artifacts are more perceivable. We find DFI to be effective on the aging and addition of facial hair tasks.
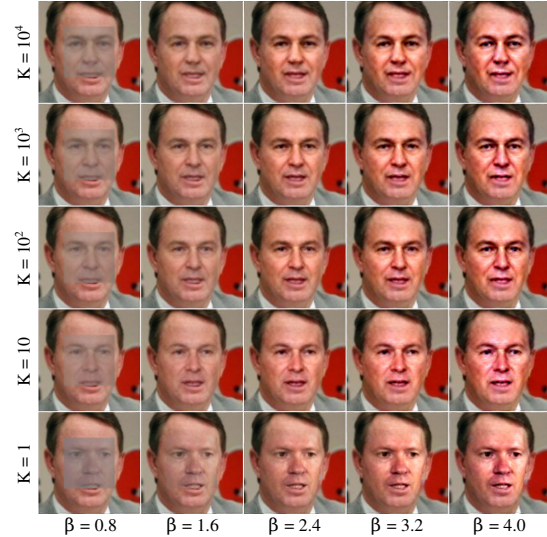


Figure 6. Inpainting and varying the free parameters. **Rows:** $K$, the number of nearest neighbors. **Columns:** $\beta$, higher values correspond to a larger perturbation in feature space. When $K$ is too small the generated pixels do not fit the existing pixels as well (the nose, eyes and cheeks do not match the age and skin tone of the unmasked regions). When $K$ is too large a difference of means fails to capture the discrepancy between the distributions (two noses are synthesized). When $\beta$ is too small or too large the generated pixels look unnatural. We use $K = 100$ and $\beta = 1.6$.

transformation tasks. This is very promising and may have implications for future work in the area of automated image transformations. However, DFI also has clear limitations and requirements on the data. We first clarify some of the aspects of DFI and then focus on some general observations.

**Image alignment** is a necessary requirement for DFI to work. We use the difference of means to cancel out the contributions of convolutional features that are unrelated to the attribute we wish to change, particularly when this attribute is centered in a specific location (adding a mustache, opening eyes, adding a smile, etc). For example, when adding a mustache, all target images contain a mustache and therefore the convolutional features with the mustache in their receptive field will not average out to zero. While max-pooling affords us some degree of translation invariance, this reasoning breaks down if mustaches appear in highly varied locations around the image, because no specific subset of convolutional features will then correspond to "mustache features". Image alignment is a limitation but not for faces, an important class of images. As shown in Section 4.2, existing face alignment tools are sufficient for DFI.

**Time and space complexity.** A significant strength of DFI is that it is very lean. The biggest resource footprint is
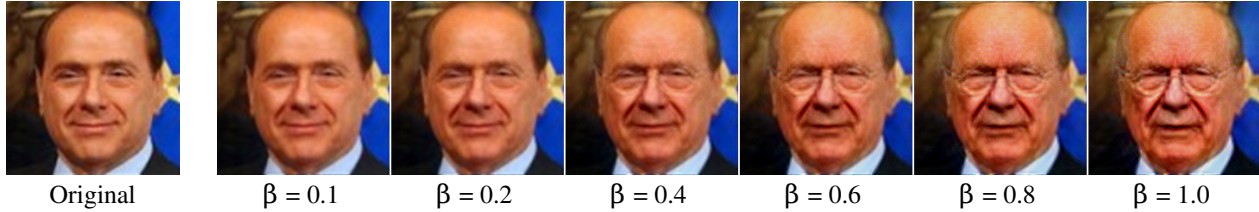
|  |  |  |  |  |  |  |
| Original | β = 0.1 | β = 0.2 | β = 0.4 | β = 0.6 | β = 0.8 | β = 1.0 |

Figure 7. Morphing a face to make it appear older. The transformation becomes more pronounced as the value of $\beta$ increases.

GPU memory for the convolutional layers of VGG-19 (the large fully-connected layers are not needed). A $1280 \times 960$ image requires 4 GB and takes 5 minutes to reconstruct. A $200 \times 200$ image takes 20s to process. The time and space complexity are linear. In comparison, many generative models only demonstrate $64 \times 64$ images. Although DFI does not require the training of a specialized architecture, it is also fair to say that during test-time it is significantly slower than a trained model (which, typically, needs sub-seconds) As future work it may be possible to incorporate techniques from real-time style-transfer [36] to speed-up DFI in practice.

**DFI's simplicity.** Although there exists work on high-resolution style transfer [11, 28, 36], to our knowledge, DFI is the first algorithm to enable automated high resolution content transformations. The simple mechanisms of DFI may inspire more sophisticated follow-up work on scaling up current generative architectures to higher resolutions, which may unlock a wide range of new applications and use cases.

**Generative vs. Discriminative networks.** To our knowledge, this work is the first cross-architectural comparison of an AE against a method that uses features from a discriminatively trained network. To our great surprise, it appears that a discriminative model has a latent space as good as an AE model at editing content. A possible explanation is that the AE architecture could organize a better latent space if it were trained on a more complex dataset. AE are typically trained on small datasets with very little variety compared to the size and richness of recognition datasets. The richness of ImageNet seems to be an important factor: in early experiments we found that the convolutional feature spaces of VGG-19 outperformed those of VGG-Face on face attribute change tasks.

**Linear interpolation as a baseline.** Linear interpolation in a pre-trained feature space can serve as a first test for determining if a task is interesting: problems that can easily be solved by DFI are unlikely to require the complex machinery of generative networks. Generative models can be much more powerful than linear interpolation, but the current problems (in particular, face attribute editing) which are used to showcase generative approaches are too simple. Indeed, we



Figure 8. Example of a hard task for DFI: inpainting an image with the right half missing.

do find many problems where generative models outperform DFI. In the case of inpainting we find DFI to be lacking when the masked region is half the image (Figure 8). DFI is also incapable of shape [53] or rotation [34] transformations since those tasks require aligned data. Finding more of these difficult tasks where generative models outshine DFI would help us better evaluate generative models. We propose DFI to be the linear interpolation baseline because it is very easy to compute, it will scale to future high-resolution models, it does not require supervised attributes, and it can be applied to nearly any aligned class-changing problems.

# 6. Conclusion

We have introduced DFI which interpolates in a pretrained feature space to achieve a wide range of image transformations like aging, adding facial hair and inpainting. Overall, DFI performs surprisingly well given the method's simplicity. It is able to produce high quality images over a variety of tasks, in many cases of higher quality than existing state-of-the-art methods. This suggests that, given the ease with which DFI can be implemented, it should serve as a highly competitive baseline for certain types of image transformations on aligned data. Given the performance of DFI, we hope that this spurs future research into image transformation methods that outperform this approach.

# References

[1] M. Aittala, T. Aila, and J. Lehtinen. Reflectance modeling by neural texture synthesis. *ACM Transactions on Graphics (TOG)*, 35(4):65, 2016. 1

[2] R. Bellini, Y. Kleiman, and D. Cohen-Or. Time-varying weathering in texture space. *ACM Transactions on Graphics (TOG)*, 35(4):141, 2016. 1

[3] Y. Bengio, G. Mesnil, Y. Dauphin, and S. Rifai. Better mixing via deep representations. In *ICML (1)*, pages 552–560, 2013. 1, 3, 4

[4] A. Brock, T. Lim, J. Ritchie, and N. Weston. Neural photo editing with introspective adversarial networks. *arXiv preprint arXiv:1609.07093*, 2016. 2

[5] E. L. Denton, S. Chintala, R. Fergus, et al. Deep generative image models using a laplacian pyramid of adversarial networks. In *Advances in Neural Information Processing Systems*, pages 1486–1494, 2015. 2

[6] J. Donahue, P. Krähenbühl, and T. Darrell. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*, 2016. 2

[7] A. Dosovitskiy, J. Tobias Springenberg, and T. Brox. Learning to generate chairs with convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1538–1546, 2015. 2

[8] V. Dumoulin, I. Belghazi, B. Poole, A. Lamb, M. Arjovsky, O. Mastropietro, and A. Courville. Adversarially learned inference. *arXiv preprint arXiv:1606.00704*, 2016. 2

[9] J. R. Gardner, P. Upchurch, M. J. Kusner, Y. Li, K. Q. Weinberger, K. Bala, and J. E. Hopcroft. Deep Manifold Traversal: Changing labels with convolutional features. *arXiv preprint arXiv:1511.06421*, 2015. 2

[10] P. Garrido, L. Valgaerts, O. Rehmsen, T. Thormahlen, P. Perez, and C. Theobalt. Automatic face reenactment. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4217–4224, 2014. 2

[11] L. A. Gatys, A. S. Ecker, and M. Bethge. A neural algorithm of artistic style. *arXiv preprint arXiv:1508.06576*, 2015. 2, 3, 8

[12] R. Girdhar, D. F. Fouhey, M. Rodriguez, and A. Gupta. Learning a predictable and generative vector representation for objects. *arXiv preprint arXiv:1603.08637*, 2016. 2

[13] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2672–2680, 2014. 1, 2

[14] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Computer Vision and Pattern Recognition (CVPR), 2016 IEEE Conference on*. IEEE, 2016. 1

[15] G. Huang, Z. Liu, and K. Q. Weinberger. Densely connected convolutional networks. *arXiv preprint arXiv:1608.06993*, 2016. 1

[16] V. Kazemi and J. Sullivan. One millisecond face alignment with an ensemble of regression trees. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1867–1874, 2014. 5

[17] I. Kemelmacher-Shlizerman. Transfiguring portraits. *ACM Transactions on Graphics (TOG)*, 35(4):94, 2016. 1, 2

[18] I. Kemelmacher-Shlizerman and S. M. Seitz. Collection flow. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 1792–1799. IEEE, 2012. 3

[19] I. Kemelmacher-Shlizerman, E. Shechtman, R. Garg, and S. M. Seitz. Exploring photobios. In *ACM Transactions on Graphics (TOG)*, volume 30, page 61. ACM, 2011. 2

[20] N. Kholgade, I. Matthews, and Y. Sheikh. Content retargeting using parameter-parallel facial layers. In *Proceedings of the 2011 ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, pages 195–204. ACM, 2011. 3

[21] D. P. Kingma, S. Mohamed, D. J. Rezende, and M. Welling. Semi-supervised learning with deep generative models. In *Advances in Neural Information Processing Systems*, pages 3581–3589, 2014. 1

[22] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012. 1

[23] P.-Y. Laffont, Z. Ren, X. Tao, C. Qian, and J. Hays. Transient attributes for high-level understanding and editing of outdoor scenes. *ACM Transactions on Graphics (TOG)*, 33(4):149, 2014. 1

[24] A. B. L. Larsen, S. K. Sønderby, H. Larochelle, and O. Winther. Autoencoding beyond pixels using a learned similarity metric. *arXiv preprint arXiv:1512.09300*, 2015. 2, 4

[25] V. Le, J. Brandt, Z. Lin, L. Bourdev, and T. Huang. Interactive facial feature localization. *Computer Vision–ECCV 2012*, pages 679–692, 2012. 5

[26] D. C. Liu and J. Nocedal. On the limited memory BFGS method for large scale optimization. *Mathematical programming*, 45(1-3):503–528, 1989. 4

[27] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015. 5

[28] A. Mahendran and A. Vedaldi. Understanding deep image representations by inverting them. In *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2015. 2, 3, 4, 8

[29] A. Nech and I. Kemelmacher-Shlizerman. Megaface 2: 672,057 identities for face recognition. 2016. 5

[30] A. Odena, C. Olah, and J. Shlens. Conditional image synthesis with auxiliary classifier GANs. *arXiv preprint arXiv:1610.09585*, 2016. 2

[31] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros. Context encoders: Feature learning by inpainting. *arXiv preprint arXiv:1604.07379*, 2016. 6

[32] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *International Conference on Learning Representations*, 2016. 2

[33] S. Reed, K. Sohn, Y. Zhang, and H. Lee. Learning to disentangle factors of variation with manifold interaction. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 1431–1439, 2014. 1, 2

[34] S. E. Reed, Y. Zhang, Y. Zhang, and H. Lee. Deep visual analogy-making. In *Advances in Neural Information Processing Systems*, pages 1252–1260, 2015. 2, 8

[35] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. 3

[36] F. Sadeghi, C. L. Zitnick, and A. Farhadi. Visalogy: Answering visual analogy questions. In *Advances in Neural Information Processing Systems*, pages 1873–1881, 2015. 8

[37] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training GANs. *arXiv preprint arXiv:1606.03498*, 2016. 2

[38] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015. 1, 3

[39] R. W. Sumner and J. Popović. Deformation transfer for triangle meshes. In *ACM Transactions on Graphics (TOG)*, volume 23, pages 399–405. ACM, 2004. 3

[40] S. Suwajanakorn, I. Kemelmacher-Shlizerman, and S. M. Seitz. Total moving face reconstruction. In *Computer Vision–ECCV 2014*, pages 796–812. Springer, 2014. 3

[41] S. Suwajanakorn, S. M. Seitz, and I. Kemelmacher-Shlizerman. What makes tom hanks look like tom hanks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3952–3960, 2015. 1, 2, 3

[42] J. Thies, M. Zollhöfer, M. Nießner, L. Valgaerts, M. Stamminger, and C. Theobalt. Real-time expression transfer for facial reenactment. *ACM Transactions on Graphics (TOG)*, 34(6):183, 2015. 2

[43] X. Wang and A. Gupta. Generative image modeling using style and structure adversarial networks. *arXiv preprint arXiv:1603.05631*, 2016. 1

[44] K. Q. Weinberger and L. K. Saul. Unsupervised learning of image manifolds by semidefinite programming. *International Journal of Computer Vision*, 70(1):77–90, 2006. 1

[45] T. Weise, H. Li, L. Van Gool, and M. Pauly. Face/off: Live facial puppetry. In *Proceedings of the 2009 ACM SIGGRAPH/Eurographics Symposium on Computer animation*, pages 7–16. ACM, 2009. 3

[46] J. Wu, C. Zhang, T. Xue, W. T. Freeman, and J. B. Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. *arXiv preprint arXiv:1610.07584*, 2016. 2

[47] X. Xiong and F. Torre. Supervised descent method and its applications to face alignment. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 532–539, 2013. 3

[48] X. Yan, J. Yang, K. Sohn, and H. Lee. Attribute2Image: Conditional image generation from visual attributes. In *Computer Vision–ECCV 2016*. 2016. 1, 6

[49] A. Yu and K. Grauman. Fine-grained visual comparisons with local learning. In *Computer Vision and Pattern Recognition (CVPR)*, June 2014. 6

[50] R. Zhang, P. Isola, and A. A. Efros. Colorful image colorization. *arXiv preprint arXiv:1603.08511*, 2016. 1

[51] J. Zhao, M. Mathieu, and Y. LeCun. Energy-based generative adversarial network. *arXiv preprint arXiv:1609.03126*, 2016. 2

[52] T. Zhou, S. Tulsiani, W. Sun, J. Malik, and A. A. Efros. View synthesis by appearance flow. *arXiv preprint arXiv:1605.03557*, 2016. 1

[53] J.-Y. Zhu, P. Krähenbühl, E. Shechtman, and A. A. Efros. Generative visual manipulation on the natural image manifold. In *Proceedings of European Conference on Computer Vision (ECCV)*, 2016. 8