# Predicting Image Differences Based on Image-Difference Features

*Ingmar Lissner, Jens Preiss, and Philipp Urban;*
*Institute of Printing Science and Technology, Technische Universität Darmstadt;*
*Magdalenenstr. 2, 64289 Darmstadt, Germany*

## Abstract

*An accurate image-difference measure would greatly simplify the optimization of imaging systems and image processing algorithms. The prediction performance of existing methods is limited because the visual mechanisms responsible for assessing image differences are not well understood. This applies especially to the cortical processing of complex visual stimuli.*

*We propose a flexible image-difference framework that models these mechanisms using an empirical data-mining strategy. A pair of input images is first normalized to specific viewing conditions by an image appearance model. Various image-difference features (IDFs) are then extracted from the images. These features represent assumptions about visual mechanisms that are responsible for judging image differences. Several IDFs are combined in a blending step to optimize the correlation between image-difference predictions and corresponding human assessments.*

*We tested our method on the Tampere Image Database 2008, where it showed good correlation with subjective judgments. Comparisons with other image-difference measures were also performed.*

## Introduction

An image difference-measure (IDM) that accurately predicts human judgments is the Holy Grail of perception-based image processing. An IDM takes two images and parameters that specify the viewing conditions (e.g., viewing distance, illuminant, and luminance level). It returns a prediction of the perceived difference between the images under the specified viewing conditions. An accurate IDM could supersede tedious psychophysical experiments that are required to optimize imaging systems and image processing algorithms.

In the past decades many attempts were made to create increasingly sophisticated IDMs. Unfortunately, evaluations show that they cannot replace human judgments for a wide range of distortions and arbitrary images so far [1, 2]. How an observer perceives a distortion depends on his interpretation of the image content — for example, changing a person's skin color is likely to cause a larger perceived difference than changing the color of a wall by the same amount. It is therefore improbable that IDMs will perfectly predict human perception before the cortical visual processing is comprehensively understood. However, IDMs could provide a reasonable median prediction of human judgments for only a few selected distortions, e.g., lossy compression or gamut mapping.

### The Role of Image Appearance Models

Many IDMs use image appearance models such as S-CIELAB [3], Pattanaik's multiscale model [4], or iCAM [5, 6] to transform the input images into an opponent color space defined for specific viewing conditions (e.g., 10° observer, illuminant D65, and average viewing distance). This can be seen as a normalization of the images to the given viewing conditions. Advanced models also consider various appearance phenomena to adjust pixel values to human perception. Typically, they account for spatial properties of the visual system by convolving the images with the chromatic and achromatic contrast sensitivity functions. This allows a meaningful pixelwise comparison of, e.g., halftone and continuous-tone images. For instance, S-CIELAB has been used as an IDM [7] in combination with the CIEDE2000 [8] color-difference formula.

Note that image appearance models are still an active research area and have room for improvement. Ideally, they normalize an input image to specific viewing conditions and remove imperceptible content. The result is an image in an opponent color space from which color attributes (lightness, chroma, and hue) can be obtained for each pixel. This space is referred to as the *working color space* in the following.

### The Role of the Color Space

It is advantageous for image-difference analysis if the working color space is highly perceptually uniform, meaning that Euclidean distances correlate well with perceived color differences. Note that a color space cannot be perfectly perceptually uniform because of geometrical issues and the effect of *diminishing returns in color-difference perception* [9]. In addition, color-difference data is obtained using color patches instead of complex visual stimuli. Nevertheless, image gradients and edges require perceptually meaningful normalization, i.e., their perceptual magnitudes should be reflected by the corresponding values as closely as possible. Analyzing such image features in a highly non-uniform color space may cause an over- or underestimation of their perceptual significance.

### Image-Difference Features

Many IDMs create *image-difference maps* showing perceived pixel deviations between two input images. For image-difference evaluation, these maps are transformed into a single characteristic value, such as the mean or the 95th percentile. However, psychophysical experiments show that the degree of difference visibility is not well correlated with perceived overall image difference [10]. For example, global intensity changes are generally less objectionable than compression artifacts [10]. It is therefore likely that the prediction performance of IDMs that only operate on image-difference maps can be improved.

Our approach uses hypotheses of perceptually significant image differences. We call these hypotheses *image-difference features* (IDFs). Various examples can be found in the literature [10, 11, 12]. Fig. 1 outlines the normalization and feature-extraction steps of our proposed image-difference framework.

We assess the relevance of our IDFs using data that relate image distortions (e.g., noise, lossy compression) to perceived image differences. A vector of IDFs is computed for each image pair (reference image and distorted image). This allows us to determine the

correlations of individual IDFs with the perceived differences of the image pairs, which are expressed by *mean opinion scores* (MOS).
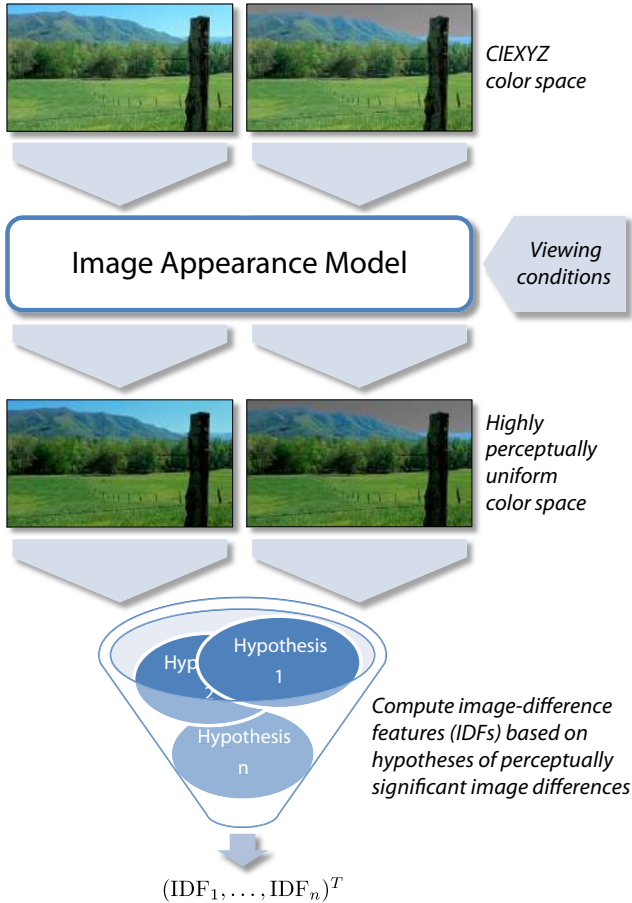


**Figure 1.** *Normalization and feature-extraction steps of the proposed image-difference framework. The example image is part of Mark Fairchild's HDR Photographic Survey [13].*

## Extracting Image-Difference Features

Even though the human visual system is not comprehensively understood, we assume that several superposed mechanisms contribute to the assessment of image differences. Our method reflects this assumption: it considers not just one but multiple hypotheses of perceptually important image differences. These hypotheses are mathematically described as image-difference features (IDFs). Although the features have low complexity, the combination of several IDFs allows us to model complex mechanisms of the visual system with a structurally simple algorithm. The IDFs are combined into an image-difference measure in a blending step, which is described in the next section.

Image-difference features are computed for image pairs, which in most cases consist of a reference image and a distorted image. As already mentioned, all images are processed by an image appearance model and transformed into a working color space before the IDF computation. In this paper, we use S-CIELAB in combination with the LAB2000 [14] color space. This space is approximately perceptually uniform with respect to the CIEDE2000 color-difference formula. Note that LAB2000 is optimized for the D65/10° white point,

which causes some error when transforming sRGB images (D65/2°) into the space. MATLAB implementations of S-CIELAB [15] and LAB2000 [16] are available online.

The IDF computation consists of three steps (see Fig. 2 for details):

1. **Per-pixel feature computation**
   (e.g., low-pass filtering, image-gradient computation)
2. **Image-difference merging**
   (e.g., Euclidean/chroma/hue differences)
3. **Characteristic-value computation**
   (e.g., mean, median, 95th percentile)

An example of an IDF computation is shown in Fig. 3. Our approach is strictly modular — the operations at each processing step can be changed without difficulty. In addition, operations can be added to incorporate new hypotheses into the framework.

Note that the low-pass filters at the first processing step represent *suprathreshold* filtering operations, whereas *threshold* filtering is performed at the initial normalization step (using S-CIELAB).

## Combining Image-Difference Features

Each image-difference feature (IDF) represents a simple hypothesis of perceptually significant image differences. To model the complex mechanisms of the visual system without increasing the complexity of our IDFs, we can combine several IDFs into an image-difference measure (IDM) in a blending step. Blending was the key to winning the Netflix Grand Prize [17, 18], a competition based on data from a recommender system.

The blending model should be simple and only depend on a few parameters to avoid overfitting the data. Thus, the number $p$ of IDFs in a blending model should be small. In this paper, we consider three models:

*1. Linear model:*

$$\text{IDM}_{\text{Lin}}(I,J) = \sum_{i=1}^{p} \lambda_i \cdot \text{IDF}_i(I,J), \tag{1}$$

where $I, J$ denote the input images and $\lambda_i$ are the model parameters. These parameters can be determined by linear regression on a set of training image pairs $I_k, J_l$ with experimentally determined mean opinion scores $\text{MOS}(I_k, J_l)$, $k = 1, \ldots, m$ and $l = 1, \ldots, n$.

*2. Polynomial model:*

$$\text{IDM}_{\text{Poly}}(I,J) = \sum_{i=1}^{p} \lambda_i \cdot \text{IDF}_i(I,J) + \sum_{i<j} \lambda_{ij} \cdot \text{IDF}_i(I,J) \cdot \text{IDF}_j(I,J). \tag{2}$$

Equation (2) shows a second-order polynomial without the quadratic terms. The polynomial order should be small ($\leq 3$) to avoid oscillation and overfitting. The model parameters $\lambda_i$ can be determined by linear regression on a training image database.

*3. Factorial model:*

$$\text{IDM}_{\text{Fac}}(I,J) = \prod_{i=1}^{p} \text{IDF}_i(I,J)^{\lambda_i}. \tag{3}$$

The parameters $\lambda_i$ of the factorial model can be optimized on a training image database using linear regression of log-transformed data.
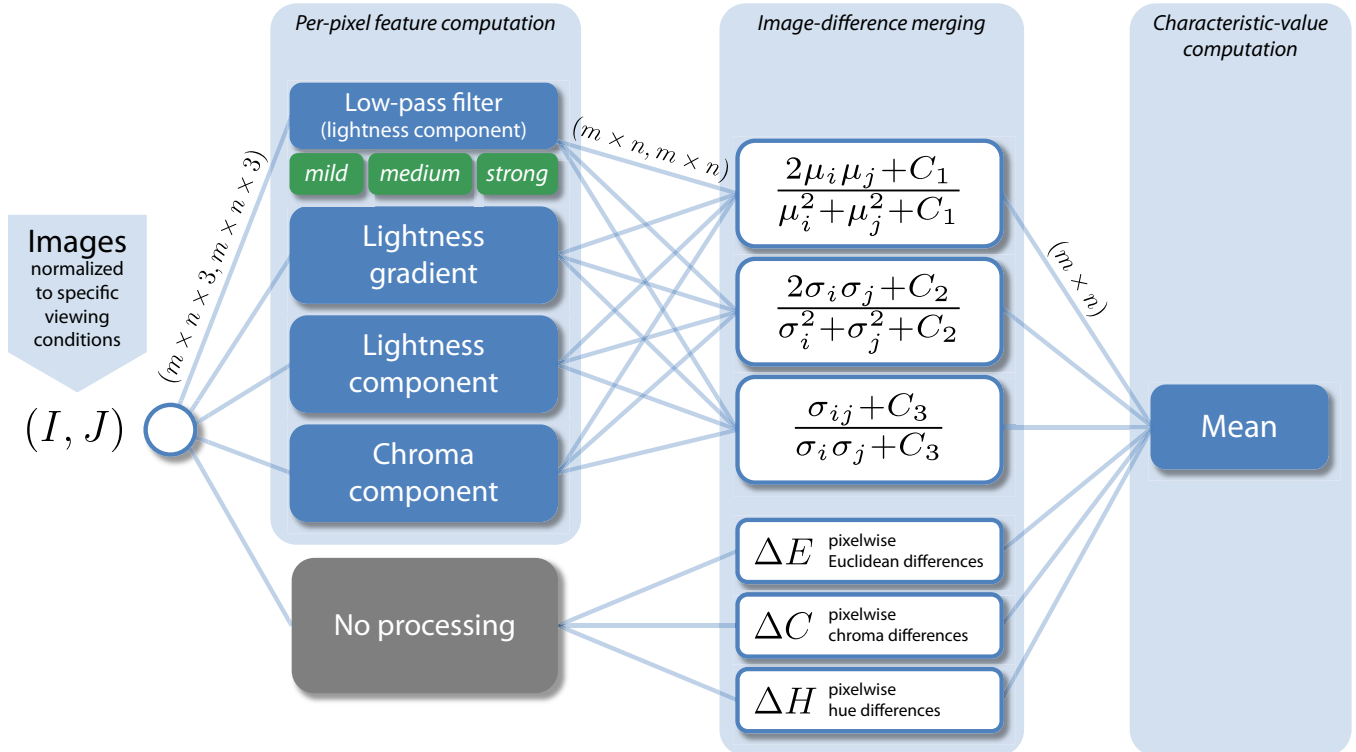
**Figure 2.** *Processing steps of the image-difference feature (IDF) computation. Each path from left to right represents a different IDF. The dimensions of the input and output data at each step are indicated at the connecting lines. The upper three operations of the "image-difference merging" step represent image comparison functions proposed by Wang et al. [10]. For reasons of simplicity we only used the mean in the characteristic-value-computation step of our implementation.*

The performance of the resulting image-difference measure strongly depends on the choice of combined IDFs. Combining two IDFs that reflect the same hypothesis does not improve the prediction accuracy, even if they both correlate well with the subjective assessments expressed by the mean opinion scores (MOS). It is therefore advisable to sort the IDFs according to their impact on the prediction performance. This sorting is performed using a set of training images with corresponding MOS. To avoid overfitting the training data, only the few most important IDFs should be considered.

A sorting algorithm based on the Spearman rank-order correlation is outlined below. Note that the result depends on the selected blending model.

---

**Algorithm 1** IDF SORTING

---

```
INPUT: MOS for M image pairs, N IDFs

IDF₁ = IDF with highest Spearman correlation to MOS
FOR i = 2 : N ITERATIONS
    FOR EACH IDF ∉ {IDF₁,..., IDF_{i-1}}
        Optimize blending model parameters based
            on {IDF₁,..., IDF_{i-1}, IDF} with respect to MOS
        Compute Spearman correlation between
            blending model predictions and MOS
    END FOR
    IDF_i = IDF resulting in highest Spearman
            correlation between predictions and MOS
END FOR

OUTPUT: Sorted IDFs {IDF₁,..., IDF_N}
```

---

In summary, we obtain image-difference measures by:

1. **Feature extraction:** Computing a large number of IDFs for a set of training images.
2. **Sorting:** Selecting the most important IDFs considering redundancies and prediction performance of individual IDFs.
3. **Blending:** Optimizing the parameters of the selected blending model on the training images.

Steps 2 and 3 are performed simultaneously (see Algorithm 1).

## Image Database

We trained and tested our method on the Tampere Image Database 2008 [1, 19]. It contains 1700 distorted images derived from 25 reference images and more than 256 000 quality judgments from more than 800 observers.

Seventeen image distortions in four intensities were applied to each reference image. Some examples are shown in Fig. 4. The distortions can be divided into the following categories: noise; lossy compression; miscellaneous (e.g., blur, denoising, intensity shifts).

Subjective scores were obtained through pair comparisons of two distorted images with the corresponding reference image. Applying the "Swiss competition principle" [1], a mean opinion score (MOS) between 0 (worst) and 9 (best) was determined for each distorted image. Since we designed our IDFs to return values within $[0,1]$, we scaled the MOS to the same range to fit the parameters of our blending models.

As our method uses S-CIELAB to normalize the input images, we had to specify the image resolution in samples per degree (spd) of
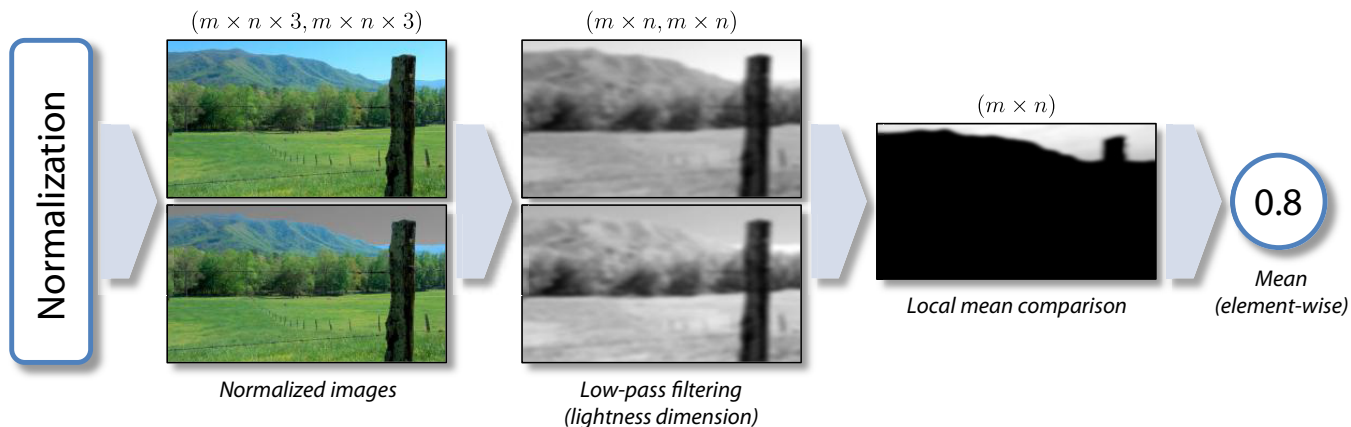
*Figure 3.* *Example of an image-difference feature (IDF) computation. The two input images are first normalized using an image appearance model and transformed into an opponent color space. In this example of an IDF, the following processing steps are then performed: 1. Low-pass filtering of the lightness component (the chromatic components are discarded); 2. Computation of a difference image using local mean comparison (as proposed by Wang et al. [10]); 3. Computation of the element-wise mean of the resulting difference image. The example image is part of Mark Fairchild's HDR Photographic Survey [13].*

the visual field. Assuming a viewing distance of two screen heights at a resolution of $1152 \times 864$ pixels and a $19''$ display [1] yields an image resolution of approximately 30 spd; this was our S-CIELAB parameter.

## Results and Discussion

To evaluate our image-difference framework, we first compiled a set of image-difference features (IDFs) as shown in Fig. 2. Following the principle of cross-validation, we divided the image database into two disjoint sets — a training set and a test set — and computed the IDFs for all training image pairs.

Using the IDF sorting algorithm (Algorithm 1) we determined the five most significant IDFs for each of the three blending models [Eqs. (1)–(3)]. The model parameters were optimized on the training set. We then computed the predictions of the resulting image-difference measures (IDMs) for the test set images and compared them with those of various quality assessment methods: MSE, SNR, PSNR, SSIM [10], MSSIM [21], VSNR [22], VIF [23], VIFP [23], UQI [24], IFC [25], NQM [26], and WSNR [26]. The predictions of these methods were obtained using the MeTriX MuX Visual Quality Assessment Package [27].

We used the Spearman and the Kendall rank-order correlations to compare the image-difference predictions with the corresponding mean opinion scores (MOS). The results are shown in Table 1. Higher Spearman and Kendall correlations indicate better prediction accuracy. We used these common rank correlations instead of, e.g., the linear Pearson correlation, because rank correlations are not affected by nonlinear relations between the input variables. However, rank correlations do not show the absolute deviations between perceived and predicted image differences.

Following common practice, we computed overall correlations between all predictions and MOS. This does not properly reflect the psychophysical data. The subjective scores of our test database are based on comparisons between images derived from the same original (all compared images show the same scene). This yields *local* instead of *global* scores, i.e., MOS of images derived from different originals are not comparable. For a meaningful overall correlation, individual correlations between MOS and predictions for each scene

| Method | Correlation | |
| --- | --- | --- |
| | Spearman | Kendall |
| WSNR | 0.490 | 0.395 |
| SNR | 0.534 | 0.382 |
| MSE | 0.554 | 0.401 |
| PSNR | 0.554 | 0.401 |
| IFC | 0.571 | 0.430 |
| UQI | 0.609 | 0.450 |
| NQM | 0.620 | 0.460 |
| SSIM | 0.624 | 0.454 |
| VIFP | 0.635 | 0.476 |
| VSNR | 0.704 | 0.530 |
| VIF | 0.735 | 0.570 |
| MSSIM | 0.862 | 0.665 |
| IDF-based (linear) | 0.877 | 0.687 |
| IDF-based (factorial) | 0.880 | 0.698 |
| IDF-based (polynomial) | 0.890 | 0.714 |

**Table 1. Correlations between perceived images differences and corresponding predictions for a set of test image pairs. Different blending models were used to create the three IDF-based IDMs that consist of five IDFs each. The results are sorted by increasing Spearman correlation.**

should be averaged.

It is evident from Table 1 that the IDF-based image-difference predictions correlate well with perceived image differences based on the Spearman and Kendall correlations. The polynomial blending method leads to slightly better predictions than the linear and factorial methods. Note that some of our IDFs include parts of the SSIM quality measure (see Fig. 2), which is also included in our comparison. Since training and test images are part of the same database, the corresponding subjective assessments were collected under similar viewing conditions. Consequently, the correlation may decrease for judgments obtained under different viewing conditions.

A comparison of IDF-based image-difference predictions with corresponding MOS is shown in Fig. 5. The predictions seem to be linearly related to the perceived differences up to a MOS of approximately 0.6, but not for higher MOS. It is generally preferable if the
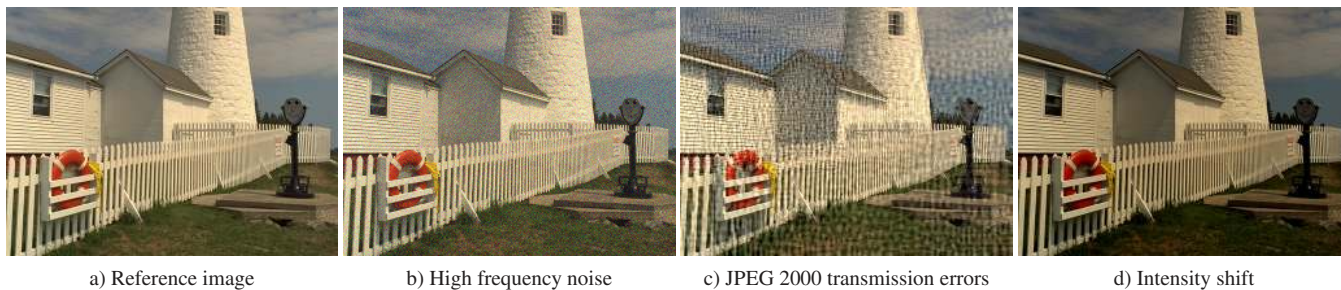
| a) Reference image | b) High frequency noise | c) JPEG 2000 transmission errors | d) Intensity shift |

**Figure 4.** *Example images from the Tampere Image Database 2008 [1, 19]. The original image is available online [20].*
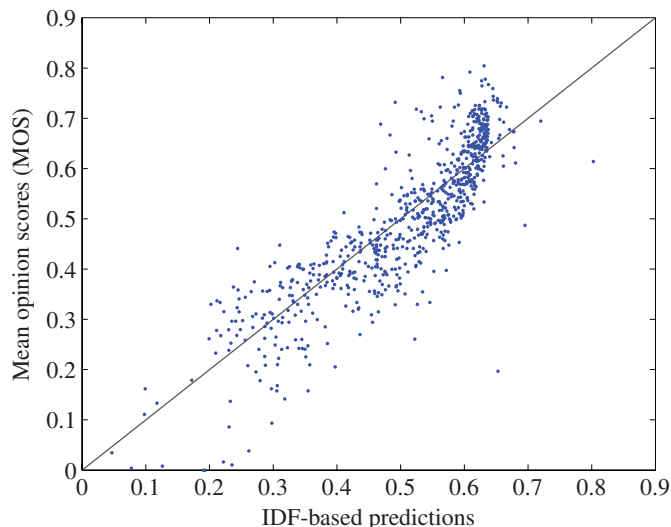


**Figure 5.** *Comparison of mean opinion scores with the predictions of an IDF-based IDM, which was created from five IDFs and a polynomial blending model. The computations are based on test images that were not used to train the IDM.*



**Figure 6.** *Spearman rank-order correlations between perceived and predicted image differences depending on the number of combined IDFs. The correlations are based on test images that were not used to train the IDMs.*

predictions of the IDM are proportional to the perceived image differences. For example, if the perceived difference of an image pair is twice as large as that of another pair, this should be reflected by the corresponding predictions. In addition, a linear relationship enables the use of linear correlation measures to calculate and optimize the correlation between MOS and predictions.

Fig. 6 illustrates how the prediction performance of our IDMs depends on the number of combined IDFs considering different blending methods. Note that the number of parameters for the linear and factorial blending models increases linearly with the number of IDFs, whereas it grows rapidly for the polynomial model. With the linear and factorial models, the prediction performance decreases if more than four IDFs are combined. With the polynomial model, the performance drops considerably if more than eight IDFs are used. In all three cases the decreasing prediction performance indicates that the training data are overfitted.

## Conclusions

We proposed a method to create image-difference measures (IDMs) based on image-difference features (IDFs), which represent simple hypotheses of how the visual system assesses image differences. Several IDFs are combined in a blending step. Their weights are optimized using training images. The resulting set of weighted IDFs
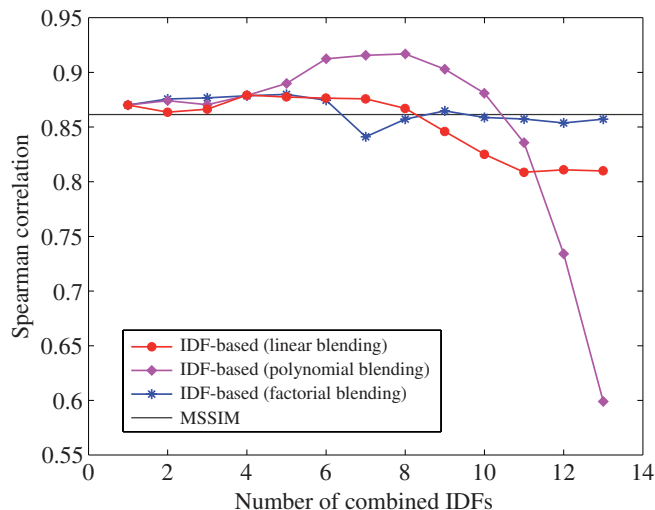
is an image-difference measure.

We created three IDMs using linear, polynomial, and factorial blending. Their predictions showed good correlation with human judgments for a set of test images. To avoid overfitting, only the few most important IDFs should be combined.

The prediction performance could be improved by adding non-redundant IDFs (reflecting new hypotheses) to the pool of IDFs. Including so-called saliency maps [28, 29] in the feature computation may also increase the prediction accuracy.

Some common image distortions were not part of our test image database, e.g., gamut mapping and HDR tone mapping. If we trained our IDMs on such distortions, entirely different IDFs would be selected for an optimal prediction performance.

## Acknowledgments

## References

[1] N. Ponomarenko, V. Lukin, A. Zelensky, K. Egiazarian, M. Carli, and F. Battisti. TID2008 — A database for evaluation of full-reference visual quality assessment metrics. *Advances of Modern Radioelectronics*, 10:30–45, 2009.

[2] N. Bonnier, F. Schmitt, H. Brettel, and S. Berche. Evaluation of spatial gamut mapping algorithms. In *14th Color Imaging Conference*, pages 56–61, Scottsdale, AZ, 2006.

[3] X. Zhang and B. A. Wandell. A spatial extension of CIELAB for digital color image reproduction. *Society for Information Display Symposium Technical Digest*, 27:731–734, 1996.

[4] S. N. Pattanaik, J. A. Ferwerda, M. D. Fairchild, and D. P. Greenberg. A multiscale model of adaptation and spatial vision for realistic image display. In *Proceedings of the 25th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH 98, pages 287–298, Orlando, FL, 1998.

[5] M. D. Fairchild and G. M. Johnson. iCAM framework for image appearance, differences, and quality. *Journal of Electronic Imaging*, 13(1):126–138, 2004.

[6] J. Kuang, G. M. Johnson, and M. D. Fairchild. iCAM06: A refined image appearance model for HDR image rendering. *Journal of Visual Communication and Image Representation*, 18(5):406–414, 2007.

[7] G. M. Johnson and M. D. Fairchild. A top down description of S-CIELAB and CIEDE2000. *Color Research and Application*, 28(6):425–435, 2003.

[8] CIE Publication No. 142. Improvement to industrial colour-difference evaluation. Technical report, Central Bureau of the CIE, Vienna, Austria, 2001.

[9] D. B. Judd. Ideal color space. *Color Engineering*, 8(2):37–52, 1970.

[10] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: From error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.

[11] G. Hong and M. R. Luo. Perceptually based colour difference for complex images. In *Proceedings of the 9th Congress of the International Colour Association*, pages 618–621, Rochester, NY, 2001.

[12] P.-L. Sun and J. Morovic. What differences do observers see in colour image reproduction experiments? In *First European Conference on Color in Graphics, Imaging, and Vision*, pages 181–186, Poitiers, France, 2002.

[13] M. D. Fairchild. The HDR Photographic Survey: `http://www.cis.rit.edu/fairchild/HDR.html`.

[14] P. Urban, D. Schleicher, M. R. Rosen, and R. S. Berns. Embedding non-Euclidean color spaces into Euclidean color spaces with minimal isometric disagreement. *Journal of the Optical Society of America A*, 24(6):1516–1528, 2007.

[15] X. Zhang. MATLAB implementation of S-CIELAB: `http://white.stanford.edu/~brian/scielab/`.

[16] MATLAB implementation of the LAB2000 space: `http://www.idd.tu-darmstadt.de/color/papers/`.

[17] Y. Koren. The BellKor solution to the Netflix Grand Prize, 2009.

[18] M. Jahrer, A. Töscher, and R. Legenstein. Combining predictions for accurate recommender systems. In *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 693–701, Washington, DC, 2010.

[19] N. Ponomarenko, F. Battisti, K. Egiazarian, J. Astola, and V. Lukin. Metrics performance comparison for color image database. In *Fourth International Workshop on Video Processing and Quality Metrics for Consumer Electronics*, Scottsdale, AZ, 2009.

[20] Kodak Lossless True Color Image Suite: `http://r0k.us/graphics/kodak/`.

[21] Z. Wang, E. P. Simoncelli, and A. C. Bovik. Multi-scale structural similarity for image quality assessment. In *Proceedings of the 37th IEEE Asilomar Conference on Signals, Systems and Computers*, pages 1398–1402, Pacific Grove, CA, 2003.

[22] D. M. Chandler and S. S. Hemami. VSNR: A wavelet-based visual signal-to-noise ratio for natural images. *IEEE Transactions on Image Processing*, 16(9):2284–2298, 2007.

[23] H. R. Sheikh and A. C. Bovik. Image information and visual quality. *IEEE Transactions on Image Processing*, 15(2):430–444, 2006.

[24] Z. Wang and A. C. Bovik. A universal image quality index. *IEEE Signal Processing Letters*, 9(3):81–84, 2002.

[25] H. R. Sheikh, A. C. Bovik, and G. de Veciana. An information fidelity criterion for image quality assessment using natural scene statistics. *IEEE Transactions on Image Processing*, 14(12):2117–2128, 2005.

[26] N. Damera-Venkata, T. D. Kite, W. S. Geisler, B. L. Evans, and A. C. Bovik. Image quality assessment based on a degradation model. *IEEE Transactions on Image Processing*, 9(4):636–650, 2000.

[27] MeTriX MuX Visual Quality Assessment Package: `http://foulard.ece.cornell.edu/gaubatz/metrix_mux/`.

[28] L. Itti and C. Koch. Computational modelling of visual attention. *Nature Reviews Neuroscience*, 2(3):194–203, 2001.

[29] Y. Tong, H. Konik, F. A. Cheikh, and A. Tremeau. Full reference image quality assessment based on saliency map analysis. *Journal of Imaging Science and Technology*, 54(3):030503–(14), 2010.

## Author Biographies

*Ingmar Lissner received his degree in Computer Science and Engineering from the Hamburg University of Technology (Germany) in 2009. He is currently working toward the PhD degree with the Institute of Printing Science and Technology, Technische Universität Darmstadt (Germany). His research interests include color perception, uniform color spaces, and image-difference measures for color images.*

*Jens Preiss received his diploma in Physics (equivalent to a M.S.) from the University of Freiburg (Germany) in 2010. He is currently a research assistant at the Institute of Printing Science and Technology, Technische Universität Darmstadt (Germany), where he works as a doctoral candidate in the area of color and imaging science.*

*Philipp Urban has been head of an Emmy-Noether research group at the Technische Universität Darmstadt (Germany) since 2009. His research focuses on color science and spectral imaging. From 2006–2008 he was a visiting scientist at the RIT Munsell Color Science Laboratory. He holds a MS in mathematics from the University of Hamburg and a PhD from the Hamburg University of Technology (Germany).*