

# The Evaluation of Stylized Facial Expressions

Christian Wallraven<sup>1</sup>, Jan Fischer<sup>2</sup>, Douglas W. Cunningham<sup>1,2</sup>, Dirk Bartz<sup>2</sup>, Heinrich H. Bühlhoff<sup>1</sup>

<sup>1</sup> Max Planck Institute for Biological Cybernetics, Tübingen, Germany

<sup>2</sup> WSI-GRIS, University of Tübingen, Germany



Figure 1: Stylized happy expressions: a) original, b) Brush Stroke, c) Cartoon, and d) Illustrative stylization.

## Abstract

Stylized rendering aims to abstract information in an image making it useful not only for artistic but also for visualization purposes. Recent advances in computer graphics techniques have made it possible to render many varieties of stylized imagery efficiently. So far, however, few attempts have been made to characterize the perceptual impact and effectiveness of stylization. In this paper, we report several experiments that evaluate three different stylization techniques in the context of dynamic facial expressions. Going beyond the usual questionnaire approach, the experiments compare the techniques according to several criteria ranging from introspective measures (subjective preference) to task-dependent measures (recognizability, intensity). Our results shed light on how stylization of image contents affects the perception and subjective evaluation of facial expressions.

**CR Categories:** I.3.7 [Computer Graphics]: Three-Dimensional Graphics and Realism—Animation J.4 [Computer Application]: Social and Behavioural Sciences—Psychology;

**Keywords:** facial animations, stylized rendering, recognition, evaluation, psychophysics, perceptual graphics

## 1 Introduction

Stylized or non-photo-realistic (NPR) rendering has attracted much interest in the computer graphics community over the last decade and has established itself firmly alongside the quest for increasing realism. Techniques that allow automatic creation of images that convey complex meaning and support high degrees of abstraction “with a few brush strokes” have applications ranging from illustration to information visualization to artistic expression.

One of the major challenges in designing stylization algorithms lies in identifying *principled ways* for creating such images. These principled ways depend, of course, crucially on the task at hand: creating an image so that it conveys specific information efficiently or so that it conforms to particular aesthetic principles are two vastly different goals and require two very different solutions. Even when focusing on a particular task, it is often unknown which visual information is needed to support this task: If faced with the task of rendering a facial expression such that it is easily recognizable, no one can clearly describe exactly what information is necessary or sufficient in order to perceive a thoughtful smile<sup>1</sup>. Finally, in some

<sup>1</sup>This study focuses on stylized renderings of complex, natural images rather than on the highly abstracted semiotics of iconography.

cases it is also difficult to evaluate or measure the success of a particular technique. Although questionnaires and other introspective measures are commonly used and offer quick and easy answers, they provide only rather indirect insights. For example, one might simply ask observers “Is this an effective technique for rendering facial expressions?”, and get very valid data about what the observers *think* about the effectiveness of a technique. Reflections about the potential effectiveness of a technique, however, is not necessarily the same thing as actually measuring its effectiveness. For that, one needs a more direct measure.

In this paper, we conduct a detailed evaluation of stylized rendering. More specifically, we examine the degree with which three stylization techniques - which greatly differ in terms of their visual impression - are able to render animated facial expressions effectively. In contrast to earlier studies, we employ several measures ranging from introspective ratings to task-specific performance characteristics. This allows us not only to contrast the evaluation criteria but also to paint a more complete picture of the impact of the different stylization techniques in the context of facial expressions.

## 2 Related Work and Motivation

In this section, we briefly discuss related work in stylized rendering techniques, evaluation of those techniques, as well as psychophysical research on perception of facial expressions. In addition, we also state how this paper aims to advance on issues raised in each context.

### 2.1 Stylization Techniques

Artistic and illustrative stylization have been areas of very active research for several years. Strothotte and Schlechtweg have published a good survey of methods used in the field [Strothotte and Schlechtweg 2002]. While artistic and illustrative stylization techniques usually remove some detail from the original image, in some applications they can actually convey the relevant information better than unprocessed data. This principle was dubbed “functional realism” in [Ferwerda 2003]. The following brief discussion lists

Copyright © 2006 by the Association for Computing Machinery, Inc. Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from Permissions Dept, ACM Inc., fax +1 (212) 869-0481 or e-mail [permissions@acm.org](mailto:permissions@acm.org).

APGV 2006, Boston, Massachusetts, July 28–29, 2006.

© 2006 ACM 1-59593-429-4/06/0007 \$5.00

three selected classes of stylized rendering algorithms that have inspired the work in this study.

In [DeCarlo and Santella 2002], a technique for *cartoon-like stylization* of photographs is presented that uses a combination of color segmentation and edge detection. The areas in which the stylization is applied are selected by frequency of fixations determined through eye tracking of users looking at an image. The resulting images consist of uniformly colored image regions on top of which edges are painted to emphasize contours. A second class of algorithms creates *Brush Stroke* stylization of images and videos [Haerberli 1990; Litwinowicz 1997]. These are often used in painterly rendering as they give the impression that the images were painted using a paintbrush. A third class of algorithms is based on *halftoning*, where the goal is to transform a grayscale or color continuum into black-and-white hatching [Freudenberg et al. 2002]. Images created using halftoning resemble sketches done with a pen by capturing the underlying shading in the image using (cross-)hatching.

*Aim:* The aim of this paper is not to develop additional stylization techniques, but to evaluate the effectiveness of three existing algorithms – one from each class. Such a perceptual evaluation will provide valuable information for the areas of computer animation and NPR on how to create easily recognizable, stylized animations.

## 2.2 Evaluation Approaches

Evaluations of stylized or NPR techniques have been conducted using three main approaches. In the first approach, which is based on introspection, users (or experts) are asked for their impression of some aspect (e.g., effectiveness [Agrawala and Stolte 2001]) of the abstracted imagery. This approach is easy to apply (often using questionnaires) and analyze, which might explain why it is perhaps the most common evaluation method not only for stylized techniques but also for computer graphics in general (see, e.g., [Stokes et al. 2004]). Its main drawback is that it has limited validity and generalizability as the desired information may not be readily and reliably accessible by introspection.

The second approach - most often used in human factors studies - evaluates the performance of users in a task-dependent context (for examples, see [Fischer et al. 2006a; Gooch and Willemsen 2002; Gooch et al. 2004; Wallraven et al. 2005]). [Gooch et al. 2004], for example, evaluated the impact of a line-drawing stylization method on identification and learning of faces. They found that the stylization faces were just as easily identified as photographs. Critically, they also found that users learned novel faces *faster* when they were stylized than when they were real photographs. This shows that abstracting the right information not only results in a more efficient data representation but also in more effective processing. [Fischer et al. 2006a], on the other hand, investigated the use of stylization for creating a consistent augmented reality environment. Currently, the placement of virtual objects in real video generally results in very noticeable visual artifacts. Psychophysical experiments showed, however, that if both the virtual objects and the real scene were stylized, participants failed to distinguish between real and virtual objects, thus demonstrating the usefulness of abstraction. The major disadvantage of this second evaluation approach is that the large number of potential tasks makes it near impossible to measure performance on every level. For techniques that are designed with a specific task in mind, however, such a direct, task-specific evaluation approach is, of course, to be preferred.

In a recent paper, [Santella and DeCarlo 2004] presented an interesting third approach to evaluation. This approach is based on eye movements that are known to reflect not only overt but also covert cognitive processes. In their study, statistical analyses of

fixation clusters were conducted to show how different NPR techniques guide and capture the users' gaze. Although the results seem promising, it is unclear exactly what the method is measuring and how it compares to the other approaches. Furthermore, data acquisition (which requires eye tracking equipment), analysis, and interpretation are difficult and less than straightforward.

*Aim:* In this study, we will take an *integrative* approach to evaluating stylized imagery by collecting both introspective *and* task-specific data in order to paint a more complete picture. More specifically, we will investigate effectiveness of stylized facial expressions through evaluating a battery of measures: these include introspective questionnaires, direct comparisons, recognition performance, perceived intensity, and perceived sincerity.

## 2.3 Perception of Facial Expressions

Facial expressions have been extensively studied in the cognitive sciences over the last few decades (for a recent review, see [Adolphs 2002]). Studies by [Ekman 1972], for example, suggest that there are seven universally recognized facial expressions (anger, contempt, disgust, fear, happiness, sadness, and surprise). In addition, facial expressions have been shown to provide a rich non-verbal communication channel that is able to alter the meaning of what is being said, to provide emphasis to spoken words, or to control the flow of a conversation (see [Bull 2001]). Recently, a series of papers ([Cunningham et al. 2003; Cunningham et al. 2004; Cunningham et al. 2005; Wallraven et al. 2004]) has started to characterize the visual information that drives the recognizability, intensity and believability of conversational facial expressions. In [Wallraven et al. 2005] this research was used in an initial set of experiments to evaluate the perceptual realism of several 3D animation methods. Since these animation methods allow full control over important information channels used in facial expressions (such as internal motion of the face, rigid head motion, shape, and texture), they provide an ideal tool for highly controlled psychophysical experiments. After determining how perceptually realistic the animations are, one of the animation methods was used to investigate the relative contribution of shape, texture, and motion to the recognizability and perceived sincerity of conversational expressions. In all experiments, a strong influence of dynamic information (both rigid head motion and non-rigid facial motion) was found. Whereas shape and texture manipulations showed only little influence on recognizability, intensity and sincerity were more affected by these dimensions. Overall, the results determined the differential contribution of a variety of perceptual characteristics and animation methods to the perception of facial expressions as well as the benefits of a close coupling of psychophysical and computer animation methods.

*Aim:* Analyses of the visual information that is emphasized by different stylization techniques will not only help to evaluate the effectiveness of the stylization techniques, but will shed further light on the processing of facial expressions.

## 3 Stylized Facial Expressions

In the following, we first briefly review the facial animation system that will be used in this study. We then discuss the three different stylization techniques that were applied to these animated expressions in order to create stylized sequences. These specific techniques were selected as each enhances or decreases the importance of quite different image characteristics (such as color, edges, or motion continuity) to a different degree. In addition, for each technique, we also determined a suitable parameter that allowed us to manipulate the level of detail contained in the image. This was done in order to investigate the impact of increasing or coarser abstraction on the effectiveness of each technique.

### 3.1 The Avatar

The avatar that was used in this paper is based on the design by [Breidt et al. 2003] and was introduced in [Wallraven et al. 2005] (see Figure 1a). It is based on a combination of high *spatial* resolution 3D-scans of peak expressions together with high *temporal* resolution motion capture of the facial deformation during these expressions. Both scan and motion capture data for these expressions are taken from a trained actor using a method-acting protocol that ensures very “natural” expressions. Scans of peak expressions are first put into correspondence using a manually designed control mesh in order to create a basis set of *morphable meshes*. From the motion capture data (captured with 72 markers), non-rigid motion is extracted and used to specify linear detectors for the expression-specific deformations in the face. The detectors provide the weights that drive the morph channels. Finally, eyes and teeth geometry are added to the scans and anchored to the rigid head motion. The movements of the eyes are created by fixating them on the virtual camera throughout the expression sequence. This corresponds closely to the real eye movements made by the actor during the recordings. The whole pipeline results in a realistic animation based on the amplitude and timing of marker motion in the motion capture data. The expression sequences used in this study were the same as in [Wallraven et al. 2005] (except for added eyes and teeth).

### 3.2 Brush Stroke Stylization

The Brush Stroke stylization used in this study (Figure 1b) is a painterly style where the output images are composed of a number of small brush strokes. The algorithm used for achieving this effect was presented in [Fischer et al. 2005a]. Briefly, a two-dimensional sampling grid is generated in a one-time preprocessing step. The grid remains fixed throughout the processing of consecutive input images and consists of an array of sampling point records. Each sampling point record contains the 2D position of the point and additional information about the brush stroke that is to be painted there. The point position is based on a regular grid with a horizontal and vertical grid spacing. Each brush stroke location is randomly displaced from this initial regular grid position. Additionally, the radius of the brush stroke is randomly generated, with a user-definable random number range. Finally, a random color offset is computed for each brush stroke. The image stylization process samples the input image by reading pixel colors at the sampling point positions in a random order (this random order is determined in a preprocessing step and remains constant for all images). The color offset is then added to each pixel color, and the resulting RGB components are clamped to the valid color number range. Each brush stroke is drawn as a textured square with a side length of the stored stroke radius, centered at the brush stroke location - this radius introduces a natural resolution scale. During brush stroke rendering, alpha blending is enabled to achieve partial transparency for overlapping brush strokes.

**Characteristics:** The Brush Stroke stylization preserves local colors in the image, with only a limited random color offset added to the input pixel. Depending on the selected brush stroke radius, however, small or medium-sized regions are masked in the output image. The discrete sampling of input pixels and the typically rather large sampling point distance result in limited frame coherence or motion continuity for animated image sequences.

### 3.3 Cartoon-like Stylization

In the cartoon-like stylization technique (see Figure 1c) used here, each input image is processed so that the resulting image consists of mostly uniformly colored areas enclosed by black silhouette

lines. The algorithm, which was described in [Fischer and Bartz 2005], is designed as a post-processing filter in a real-time rendering pipeline. The implementation of the algorithm uses vertex and fragment shaders, which run on the programmable graphics processing units (GPUs) of recent graphics cards.

The stylization filter consists of two steps. In the first step, a simplified color image is computed from the input image. The basis of this computation is a non-linear filter, which is inspired by bilateral filtering, as described in [Tomasi and Manduchi 1998]. The non-linear filter performs a repeated, photometric weighting of the pixels, taking into account only their chrominance components. The repetition of the filter operation is necessary in order to achieve a sufficiently good color simplification. The second stage of the image stylization filter is an edge detection step based on the simplified color image. Thus, the silhouette lines are located between similarly colored regions in the image, which is an approximation of a cartoon-like rendering style. Finally, the simplified color image is combined with the edge detection results. The responses of the edge detection filter are drawn over the output image as black lines. A specific weight function is used for computing a transparency for the detected edge pixels, which produces a smooth blending over the color image. As an increasing number of filtering iterations in the first step results in highly simplified, blurred image as well as less distinct edges, this parameter was chosen to create different resolution levels for this technique.

**Characteristics:** The cartoon-like stylization stresses high contrast edges in the image and preserves the dominant color in larger image regions. It does, however, remove small details as well as low-contrast details as an effect of the non-linear filter.

### 3.4 Illustrative Stylization

The Illustrative stylization used here (see Figure 1d) generates output images which reproduce the brightness of the input image with black-and-white hatching. Moreover, high contrast edges are rendered as black lines. This algorithm is based on aspects of the illustrative rendering method described in [Fischer et al. 2005b].

In order to creating the hatching, a procedural halftoning technique similar to the one described by Strothotte and Schlechtweg is applied [Strothotte and Schlechtweg 2002]. Parameters of this algorithm are the orientation of the main hatching direction, the minimum intensity required for the addition of perpendicular cross-hatching strokes, as well as the size of the hatching pattern in pixels. As with the Brush Stroke technique, the size of the pattern was used to determine the resolution levels. In addition to the black-and-white representation of the input image, silhouette lines are added to the stylized output. A Sobel edge detection filter delivers the location of high contrast edges in the image. These locations are then overlaid as black lines over the output image. As can be seen in Figure 1d, these lines contribute to the final image mainly in high-contrast regions such as the eyes.

**Characteristics:** The illustrative stylization emphasizes intensities in the image. These are computed as the Y component of the YUV color space representation of each input pixel. Moreover, high-contrast edges are stressed by the edge detection step used in this stylization method. The illustrative stylization removes all color information from the image, rendering a purely black-and-white representation of the input image. Small details in the image are not preserved, depending on the selected size of the halftoning pattern.

## 4 Experiments

The stylization techniques were evaluated by:

**Direct preference:** By showing two stylized sequences side-by-side and asking “which sequence captures the essence of the expression better?”, we allow participants to compare and contrast two stylization techniques at the same time. One of the advantages of this method is that although the question asks for a very subjective evaluation, participants are forced to choose one sequence which allows for a clean analysis of the data.

**Introspective questionnaire:** The questionnaire asks participants to rank the three stylization techniques. The techniques are ranked three times: Once according to aesthetic principles, once according to effectiveness in rendering facial expressions, and once according to subjective preference. All three criteria are evaluated by introspection.

**Recognizability:** It is of course crucial that the different techniques support the recognition of the facial expressions. In particular, stylization should neither decrease recognition accuracy nor increase recognition time compared to the non-stylized version. In contrast to the two previous criteria, recognition accuracy constitutes an objective, quantifiable criterion.

**Intensity and sincerity:** These two criteria constitute higher-level characteristics of facial expressions. The ratings can, for example, be of interest if the goal is not only to create recognizable but also convincing facial expressions. This will be important in areas such as virtual sales and kiosk applications.

### 4.1 Experiment 1 - Direct Comparison

In the first experiment, participants directly compared two sequences in terms of their effectiveness. Additionally, participants had to fill out a standard introspective questionnaire.

#### 4.1.1 Stimulus Sequences

In this experiment, we used seven expression sequences from the avatar: confusion, disgust, fear, happy, sad, surprise, and thinking. Each sequence was contrast normalized in order to provide a consistent input to each of the three stylization techniques. We chose default parameters for rendering for each technique that were derived from their standard use [Fischer et al. 2006b; Fischer et al. 2005a; Fischer et al. 2005b]. Three different resolution levels were obtained by a) setting the diameter of the element size for the Brush Stroke algorithm to 2, 8, and 16 pixels, b) treating the texture map in the Cartoon algorithm by blurring the image 2, 8, and 16 times, and c) setting the basic element size for the Illustrative stylization to 2, 8, and 16 pixels. Finally, we created three different resolution levels for the avatar by blurring it with an equivalent blurring filter as for the Cartoon stylization (see Figure 2a-d for example images).

#### 4.1.2 Design

The two video sequences were presented side-by-side at a resolution of 1024x768 pixels on a CRT monitor (each sequence was shown at 512x512 pixels). Participants viewed these sequences using a chin-rest at a distance of 0.5 meters (each face on the monitor subtended a visual angle of 11.4°). A single trial of the experiment consisted of the video sequences being shown repeatedly in the center of the screen. A 200ms blank screen was inserted between repetitions of the video clip. In each trial, participants had to indicate by a timed button press whether they thought the left or the right sequence captured the essence of the expression better (no name or description of expression itself was given). The experiment compared all methods and resolution levels within each expression, leaving out same-same comparisons - the total number of trials was thus  $(7 \text{ expressions}) \cdot [(12 \text{ combinations of method-resolution$

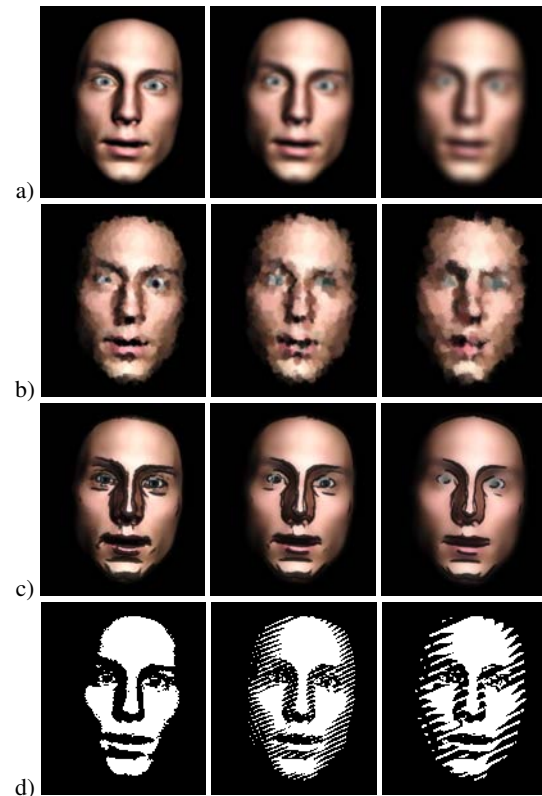


Figure 2: Stylization techniques used in this paper for a fearful expression and all resolution levels. a) Standard Avatar, b) Brush Stroke, c) Cartoon, d) Illustrative stylization.

level)-(12 - 12 same-same comparisons)/2] = 462 trials - where the order of the trials was fully randomized. Participants could take a break after half of the trials in order to avoid fatigue effects.

After the experiment (which lasted around an hour), we showed participants high-quality printouts of the different techniques and resolution levels that were taken from a frame of the happy expression. We then asked them to fill out a questionnaire in which they had to rank these 12 images according to three different criteria. The first criterion asked how artistic participants thought the different stylization techniques to be. The second criterion asked which of the techniques was the most effective in rendering facial expressions - the same question as in the direct comparison task. Finally, we asked participants to rank the techniques according to which one they liked best.

#### 4.1.3 Results & Discussion

The direct preference data from ten participants were analyzed as frequency histograms using  $\chi^2$ -tests with the factors “stylization technique”, “resolution level”, and “expression” for between technique comparisons. The analysis of the trials in which both image sequences used the same stylization technique was done separately in order to look for effects of “resolution level”, and “expression” within each technique.

As can be seen in Figure 3a, when faced with two different stylization techniques, participants most often choose the original Avatar animation, followed by the Illustrative and Cartoon techniques. The Brush Stroke method got chosen a mere 4 percent of the time. Figure 3b shows that there is a clear trend in preferences as a function of resolution level - the most detailed level is preferred over the

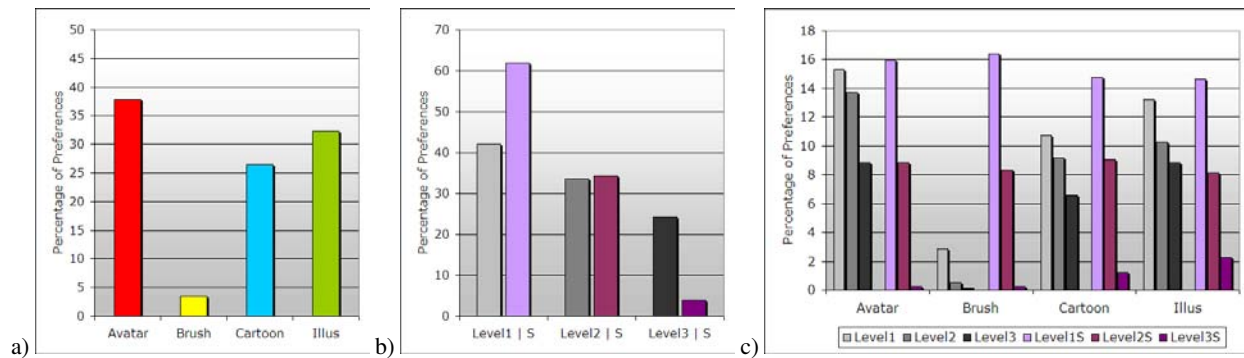


Figure 3: Experiment 1. Preference results for a) stylization techniques, b) resolution level, c) both techniques and levels. In b) and c), the grey-shaded bars represent the between-technique trials, while the colored bars represent the within-technique trials.

medium level, which in turn is preferred over the least detailed level with a 10 percent drop in preference for each step. Figure 3c shows that this effect depends on the stylization technique - for the Avatar and Cartoon conditions, there is virtually no difference between the first and the second level, but a large drop for the third. For both Illustrative and Brush Stroke techniques, the second and third levels are chosen equally often; both levels are chosen significantly less than the first, most detailed level.

Participants clearly thought that the Avatar captured the essence of the expressions best. Of the three stylization techniques, the Illustrative style seems to be preferred and the Brush Stroke seems to be not seen as not very effective. For all sequences, the highest level of detail is preferred - a result that might be expected as the least detailed levels contain only severely reduced visual information that masks much of the facial motion.

**Preference for within technique comparisons:** The analysis revealed a highly significant main effect for resolution level ( $p < 0.001$ ,  $df=2$ ,  $\chi^2=422.12$ ) as well as an interaction between level and method ( $p < 0.001$ ,  $df=6$ ,  $\chi^2=25.67$ ). These effects are plotted in Figure 3b-c using colored bars.

Figure 3b shows that for within technique comparisons, the first level is preferred 60 percent of the time, whereas the second level drops sharply to 35 percent and the third, least detailed, level is rarely chosen. Again, this pattern depends on the stylization technique - for both Illustrative and Cartoon stylization the least detailed level was occasionally chosen, whereas for both Avatar and Brush Stroke stylization it was almost never chosen. Interestingly, the preference of the highest level of detail is much more pronounced for within technique than for between technique comparisons.

**Response times:** Participants tended to respond more slowly in trials where at least one of the two image sequences was rendered with Brush Stroke stylization than in other trials (2.5s versus 2.1s,  $p < 0.05$ ). There were no other statistically significant effects for reaction time. Overall, however, the general lack of a reaction time effect as well as the absolute numbers jointly suggest that the decision was not dependent on the stylization technique and that it was not particularly hard for participants to reach a decision.

**Summary:** When directly comparing two image sequences, participants clearly felt that the original Avatar animation captured the essence of the expressions best. Among the stylization techniques, the Illustrative method was preferred. Additionally, the most detailed resolution was judged as more effective than lower resolution levels, with the exact pattern of decrease depending on the stylization technique used. This indicates that some techniques are less prone to degradation than others - a result that is confirmed by the pattern of preferences found for within-technique comparisons. In

summary, the results from this experiment indicate that stylization would actually *harm* the effective depiction of expressions.

**Preference for between technique comparisons:** We found highly significant main effects for technique ( $p < 0.001$ ,  $df=3$ ,  $\chi^2=1036.43$ ), resolution level ( $p < 0.001$ ,  $df=2$ ,  $\chi^2=178.76$ ) as well as an interaction between level and method ( $p < 0.001$ ,  $df=6$ ,  $\chi^2=101.84$ ). These effects are plotted in Figures 3.

## 4.2 Experiment 2 - Recognizability

The second experiment provides a different and complementary view to the same issue examined in the first experiment. Here, in the context of a specific task, several perceptual measures (including recognition, reaction time, and the perception of intensity and sincerity) are investigated.

### 4.2.1 Design

The setup and design of this experiment followed closely that of [Wallraven et al. 2005]: The first task was to *identify* the expression by selecting the name of the expression from a list displayed on the side of the screen. The list of choices included all seven expressions as well as “none of the above” (an eight-alternative-non-forced-choice task, see [Cunningham et al. 2003] for a detailed discussion of this paradigm). The second task was to rate the *intensity* of the expressions on a scale from 1 (not intense) to 7 (very intense). In the third task, participants were to rate the *sincerity* of the expressions, with a rating of 1 indicating that the actor was clearly pretending and a value of 7 indicating that the actor really meant the underlying emotion. Participants were explicitly instructed to anchor the scales at a value of 4 (normal intensity and sincerity) and to try and use the whole range of the scale during the experiment. The experiment used 3 repetitions of each sequence, yielding a total of (7 expressions)·(4 stylization techniques)·(3 resolution levels)·(3 repetitions) = 252 trials - again, trials were fully randomized. After the experiment, we asked participants to fill out the same questionnaire as in Experiment 1.

### 4.2.2 Results & Discussion

Data were collected from ten participants who had not taken part in the previous experiment. The results were analyzed using standard “analysis of variance” (ANOVA) methods which analyze statistical significances for each factor (expression, stylization technique, resolution level) for the different measures (recognition, reaction time, intensity, and sincerity).

**Recognition:** The ANOVA found main effects of expression ( $F(6,54)=12.201$ ,  $p < 0.001$ ) stylization method ( $F(3,27)=3.27$ ,

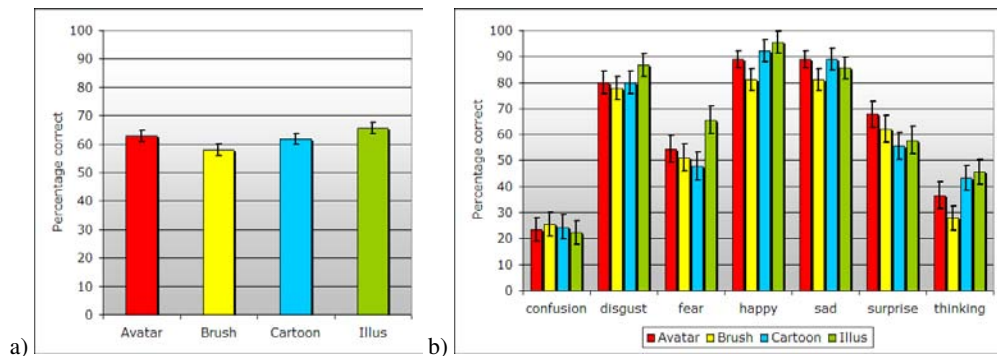


Figure 4: Experiment 2. Recognition results broken down by a) stylization technique, b) stylization technique and expression.

$p < 0.001$ ) as well as interactions between expression and method ( $F(18,162)=1.555$ ,  $p=0.05$ ) and expression and resolution level ( $F(12,108)=1.998$ ,  $p < 0.05$ ).

As was expected - and in accordance with the pattern of results obtained in [Cunningham et al. 2005; Wallraven et al. 2005], we found that some expressions were more easily recognized than others. In particular, “thinking” and “confusion” are hard to recognize and are often confused with each other - these expressions also cause the overall “low” level of performance of 60 percent. More interesting is the main effect of stylization technique (Figure 4a) - here we found that the Brush Stroke stylization was significantly worse than the remaining three rendering methods (t-test: all  $p < 0.05$ ). In addition, the Illustrative stylization was slightly better (t-test:  $p=0.05$ ,  $p=0.07$ , marginally significant) than both Avatar and Cartoon conditions. Interestingly, a closer analysis of the incorrect and “none of the above” responses for all techniques showed that the Illustrative stylization technique had significantly less incorrect responses than the other techniques (t-test: all  $p < 0.05$ ), whereas the Brush Stroke stylization technique had a significantly higher amount of “none of the above” responses than the other three techniques (t-test: all  $p < 0.05$ ). These results can be summarized by ranking the four different techniques according to their *discriminatory performance*: in this case, the Illustrative stylization technique supports the most discriminative recognition performance, followed by the original Avatar and the Cartoon stylization at the same rank, followed by the Brush Stroke stylization.

As suggested by the significant interaction between expression and method (see also Figure 4b), which stylization produces the best recognition performance depends on what the expression is. The Illustrative stylization technique produced superior performance for some expressions (“fear”, to a lesser degree also “disgust”). For “surprise”, the original Avatar animation is most easily recognizable. For “thinking”, both Cartoon and Illustrative stylization provide increased performance, whereas the “sad”, “confused”, and “happy” expressions show no clear trend favoring one single technique. This pattern of results suggests that different techniques emphasize different types of information that is relevant for the recognition of certain expressions.

Finally, an examination of the interaction between expression and resolution level shows that across all techniques “thinking” is recognized significantly better in the most detailed level, whereas “fear” is recognized much better in the lowest resolution. For the remaining expressions, no clear dependence on resolution level could be found. The result for “thinking” is due to the fact that this expression requires a close analysis of the eye-motion in order to be reliably recognized [Cunningham et al. 2005] - this eye-motion, however, is masked by the coarser, blurred visual information in the lower resolution levels. In contrast, “fear” is recognized much bet-

ter in the lowest resolution level - this expression is driven mainly by the large amount of rigid head motion [Wallraven et al. 2005], which is more visible at the lowest level.

**Response times:** Overall, response times in this experiment showed no significant effects. Restricting the analysis to just the correctly answered trials, we found a small but significant increase in response times for the Brush stylization method (2.5s as opposed to 2.2s for the other three methods). This small increase most probably mirrors the impaired recognition performance observed for this method. In general, however, stylization incurred no additional cost in processing time. In other words, we did not find a speed-accuracy tradeoff as in [Fischer et al. 2006a]. In addition, no effect of resolution level on reaction times was found which provides additional support for the data gathered in [Wallraven et al. 2005], who also found no effect of resolution level on response times.

**Intensity:** For intensity ratings<sup>2</sup>, the ANOVA found main effects of expression ( $F(6,54)=7.882$ ,  $p < 0.001$ ), stylization technique ( $F(3,27)=7.02$ ,  $p=0.001$ ), resolution level ( $F(2,18)=9.192$ ,  $p < 0.01$ ) as well as interactions of expression and method ( $F(18,162)=2.285$ ,  $p < 0.01$ ) and stylization technique and resolution level ( $F(6,54)=2.596$ ,  $p < 0.05$ ).

Similar to the results in [Wallraven et al. 2005], we found a large main effect of expressions - emotional expressions such as “disgust”, and “fear” were rated as more intense. We again found a main effect of stylization technique (Figure 5a) - in this case, Brush stylization was rated as much less intense than the other three methods (t-tests, all  $p < 0.05$ ). One reason for this is that the brush stroke pattern masks both rigid and non-rigid head motion, which are highly correlated with ratings of perceived intensity [Wallraven et al. 2005]. Analysis of the interaction between expression and method revealed in particular that “happy”, “sad”, “surprise”, and “thinking” were rated as much less intense for the Brush stroke technique than the remaining three expressions. Finally, the main effect and interaction for resolution level showed a large decrease in intensity for the Avatar and the Brush stroke technique at the lowest resolution level, whereas this decrease was less pronounced for the Cartoon technique, and virtually absent for the Illustrative technique. It thus seems that the Illustrative technique provides a very stable impression of intensity even at low resolutions.

**Sincerity:** We found main effects of expression ( $F(6,54)=3.602$ ,  $p < 0.01$ ), stylization technique ( $F(3,27)=2.92$ ,  $p=0.05$ ), resolution level ( $F(2,18)=6.736$ ,  $p < 0.01$ ) as well as an interaction of expression and method ( $F(18,162)=2.090$ ,  $p < 0.01$ ).

Both main effects of expressions and resolution level provide additional support for the data found in [Wallraven et al. 2005].

<sup>2</sup>Both intensity and sincerity ratings were analyzed only for correct answers.

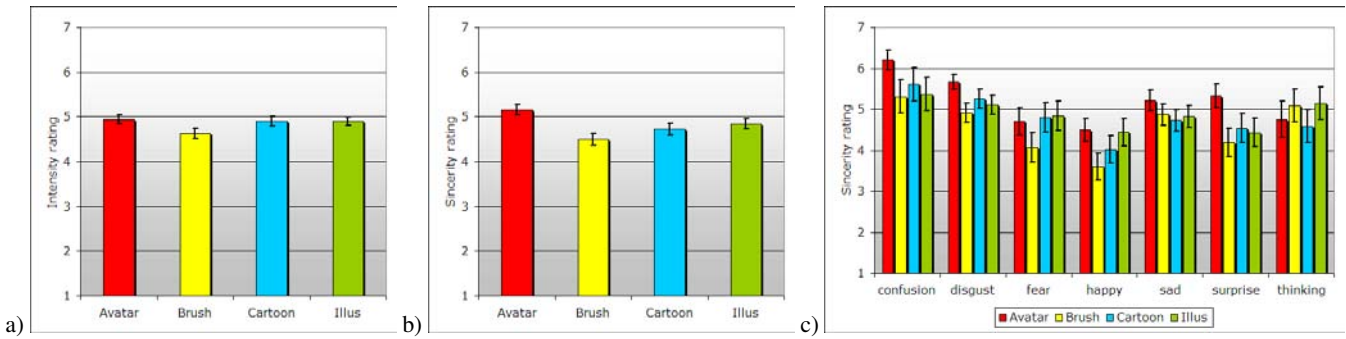


Figure 5: Experiment 2. a) Intensity ratings and b) Sincerity ratings broken down by stylization technique, c) Sincerity ratings broken down by technique and expression

Most importantly, lower resolutions provide less sincere expressions across all techniques (Figure 5b). The interaction of expression and method is shown in Figure 5c - in particular, “confusion”, “disgust”, and “surprise” are rated as most sincere in the original Avatar animation, whereas for “fear”, and “happy”, the Brush stroke technique is rated as the least sincere of all techniques.

**Summary:** While we found an effect of stylization across all measures, the most important result of this experiment is that the pattern of expression recognition did not match the subjective judgments measured in Experiment 1. In Experiment 1, participants thought that the Avatar captured the essence of the expressions best. In Experiment 2, however, Illustrative Stylization resulted in the highest level of discriminative recognition, demonstrating a small but significant advantage of abstraction for recognition performance. We found no effect of response times, which shows that abstraction of information induces no time penalty. Finally, analysis of the sincerity ratings revealed that stylization might also have an adverse effect as the original avatar animation had the highest degree of sincerity.

### 4.3 Questionnaires

Both sets of questionnaires provided similar trends, suggesting that the difference in task did *not* influence the introspective rankings of the different techniques. We therefore pooled the answers to the questionnaires for the final analysis. Analysis of this data was done by determining for each of the possible 12 ranks the winning technique (a combination of rendering technique and a specific resolution level). These results are summarized in Table 1 (note that double entries can occur using this analysis).

**Aesthetic preference:** The clear winner in terms of aesthetic preference is the original avatar animation - of the stylization techniques, the Illustrative stylization was judged as most aesthetic, followed by Cartoon and Brush stylization. One of the reasons why we did not find a clearer preference for one of the stylized techniques is probably that participants regarded all of techniques as equal, rather than judging them as one non-stylized and three stylized versions.

**Effectiveness preference:** For effectiveness, the avatar animation was ranked highest, followed closely by the Illustrative and Cartoon techniques, whereas the Brush stroke method was judged as least effective. This pattern mirrors more the one found in Experiment 1 rather closely, although the degree to which the techniques are separated in terms of their preference was much more pronounced in Experiment 1.

**Subjective preference:** For subjective preference, the ranking ordering changes - here, the Illustrative style clearly wins, followed by the Avatar and Cartoon renderings. As with the previous measures, the Brush stroke technique comes in last. It is interesting that

Rank	aesthetic	effectiveness	subjective
1	Ava1	Ava1	III1
2	Ava2	III1	Ava1
3	III1	Ava2,Car1	III1
4	Ava1	Ava2	III2
5	III2	Car2	III3
6	Car2	III2	Ava2,Car2
7	Ava2,Car2,III3	III2	Ava3
8	Car1	Car3	Ava2,Car3,III2
9	Car3	Ava3,Bru1	Ava3
10	Bru1	Bru1	Bru1
11	Bru2	Bru2	Bru2
12	Bru3	Bru3	Bru2

Table 1: Questionnaires. Results for aesthetic, effectiveness, and subjective preference judgments. Abbreviations indicate stylization technique and resolution level, respectively.

this measure clearly differs from the aesthetic preference - at least for subjective ratings, it seems that stylization is preferred much more than the original animation.

Finally, for all measures, responses show a clear preference ordering of the resolution levels from high to middle to low - a pattern that was seen throughout this study. Overall, the pattern seems to correspond quite well to the one found in Experiment 1. The main disadvantage of the questionnaire analysis, however, as can be seen from Table 1, is that it allows only for a rather coarse interpretation of the results. For an in-depth analysis, other approaches - such as done in the previous two experiments - are required.

## 5 Conclusion & Outlook

In this paper, we presented a series of evaluations of three different stylization techniques in the context of facial expressions and found effects of stylization on almost all measures. The first experiment investigated the question of effectiveness in a direct comparison task. The results indicated that stylization would potentially reduce effectiveness compared to the original avatar animation. A similar pattern of results was found for the introspective evaluation of effectiveness using questionnaires. In the second experiment, we collected several task-specific measures that were centered on recognizability as well as perceived intensity and sincerity. The results did not correlate with the effectiveness measures - the Illustrative stylization provided the most discriminatory performance.

The most obvious explanation for the difference between the recognition results and the introspective evaluation and direct comparison of “effectiveness” is that they do not measure the same thing. This explanation is, in part, contradicted by the fact that almost all participants mentioned during de-briefing that recognizability was one

the central criterion they used in determining “effectiveness”.

In Experiment 2, we found no effect of stylization on response times. This is similar to a study on face identification by [Gooch and Willemsen 2002] in which stylization also did not affect response times. Taken together with the fact that recognizability also was not affected by resolution level, our results demonstrate that the processing of facial expressions is based on mechanisms that operate very robustly even under severe changes in stimulus detail.

The systematic investigation of the visual information that drives the pattern of results observed in these experiments (especially the interaction effects) will need to be done in future studies - nevertheless, we might speculate that for the Illustrative technique, one of the reasons for its comparatively good performance lies in the emphasis of shape through hatching in connection with silhouette lines that highlight small details of the face. Previous studies have shown that the loss of color does not impact recognition of identity [Yip and Sinha 2002] - our study has shown the same for expression recognition in the data for the Illustrative technique.

In terms of practical applications, the results of our experiments can be summarized as preliminary guidelines for effective rendering: On the one hand, if the goal is to convey a facial expression most effectively, choosing a stylized rendering method (such as illustrative rendering) might help - apart from offering other dimensions such as aesthetics, sparse representation, etc. On the other hand, if the goal is to provide subjective certainty about the conveyed expression, one needs to resort to a “realistic” rendering method.

In summary, our study has evaluated the effectiveness of three different stylization techniques across multiple perceptual and introspective dimensions. Our results have provided further insight into the robustness of expression recognition as well as demonstrated critical differences of evaluation methodology.

## References

- ADOLPHS, R. 2002. Recognizing Emotions from Facial Expressions: Psychological and Neurological Mechanisms. *Behavioral and Cognitive Neuroscience Reviews* 1, 1, 21–61.
- AGRAWALA, M., AND STOLTE, C. 2001. Rendering Effective Route Maps: Improving Usability Through Generalization. In *Proc. of ACM SIGGRAPH*, ACM Press, New York, NY, USA, 241–249.
- BREIDT, M., WALLRAVEN, C., CUNNINGHAM, D. W., AND BÜLTHOFF, H. H. 2003. Facial Animation Based on 3D Scans and Motion Capture. *SIGGRAPH '03 Sketches & Applications*.
- BULL, P. 2001. State of the art: Nonverbal communication. *The Psychologist* 14, 644–647.
- CUNNINGHAM, D. W., BREIDT, M., KLEINER, M., WALLRAVEN, C., AND BÜLTHOFF, H. H. 2003. How Believable are Real Faces: Towards a Perceptual Basis for Conversational Animation. In *Computer Animation and Social Agents 2003*, 23–39.
- CUNNINGHAM, D., NUSSECK, M., WALLRAVEN, C., AND BÜLTHOFF, H. 2004. The Role of Image Size in the Recognition of Conversational Facial Expressions. *Computer Animation & Virtual Worlds* 15, 3-4 (07), 305–310.
- CUNNINGHAM, D., KLEINER, M., WALLRAVEN, C., AND BÜLTHOFF, H. 2005. Manipulating Video Sequences to Determine the Components of Conversational Facial Expressions. *ACM Transactions on Applied Perception* 2, 3 (07), 251–269.
- DECARLO, D., AND SANTELLA, A. 2002. Stylization and Abstraction of Photographs. In *Proc. of ACM SIGGRAPH*, 769–776.
- EKMAN, P. 1972. *Universal and Cultural Differences in Facial Expressions of Emotion*. University of Nebraska Press, 207–283.
- FERWERDA, J. 2003. Three Varieties of Realism in Computer Graphics. In *Proc. of SPIE Human Vision and Electronic Imaging*, 290–297.
- FISCHER, J., AND BARTZ, D. 2005. Real-time Cartoon-like Stylization of AR Video Streams on the GPU. Technical Report WSI-2005-18, Wilhelm Schickard Institute for Computer Science, Graphical-Interactive Systems (WSI/GRIS), University of Tübingen, September.
- FISCHER, J., BARTZ, D., AND STRASSER, W. 2005. Artistic Reality: Fast Brush Stroke Stylization for Augmented Reality. In *Proc. of ACM Symposium on Virtual Reality Software and Technology (VRST)*, 155–158.
- FISCHER, J., BARTZ, D., AND STRASSER, W. 2005. Illustrative Display of Hidden Iso-Surface Structures. In *Proc. of IEEE Visualization*, 663–670.
- FISCHER, J., CUNNINGHAM, D., BARTZ, D., WALLRAVEN, C., BÜLTHOFF, H., AND STRASSER, W. 2006. Measuring the Discernability of Virtual Objects in Conventional and Stylized Augmented Reality. In *Eurographics Symposium on Virtual Environments (EGVE)*.
- FISCHER, J., EICHLER, M., BARTZ, D., AND STRASSER, W. 2006. Model-based Hybrid Tracking for Medical Augmented Reality. In *Eurographics Symposium on Virtual Environments (EGVE)*.
- FREUDENBERG, B., MASUCH, M., AND STROTHOTTE, T. 2002. Real-Time Halftoning: A Primitive for Non-Photorealistic Shading. In *Proc. of Eurographics Workshop on Rendering*, 227–231.
- GOOCH, A. A., AND WILLEMSSEN, P. 2002. Evaluating Space Perception in NPR Immersive Environments. In *NPAR '02: Proceedings of the 2nd international symposium on Non-photorealistic animation and rendering*, ACM Press, New York, NY, USA, 105–110.
- GOOCH, B., REINHARD, E., AND GOOCH, A. 2004. Human Facial Illustrations: Creation and Psychophysical Evaluation. *ACM Trans. Graph.* 23, 1, 27–44.
- HAEBERLI, P. 1990. Paint By Numbers: Abstract Image Representations. In *Proc. of ACM SIGGRAPH*, 207–214.
- LITWINOWICZ, P. 1997. Processing Images and Video for an Impressionist Effect. In *Proc. of ACM SIGGRAPH*, 407–414.
- SANTELLA, A., AND DECARLO, D. 2004. Visual Interest and NPR: An Evaluation and Manifesto. In *NPAR '04: Proceedings of the 3rd international symposium on Non-photorealistic animation and rendering*, ACM Press, New York, NY, USA, 71–150.
- STOKES, W. A., FERWERDA, J. A., WALTER, B., AND GREENBERG, D. P. 2004. Perceptual illumination components: a new approach to efficient, high quality global illumination rendering. *ACM Trans. Graph.* 23, 3, 742–749.
- STROTHOTTE, T., AND SCHLECHTWEIG, S. 2002. *Non-Photorealistic Computer Graphics - Modelling, Rendering, and Animation*. Morgan Kaufmann Publishers.
- TOMASI, C., AND MANDUCHI, R. 1998. Bilateral Filtering for Gray and Color Images. In *Proc. of IEEE International Conference on Computer Vision (ICCV)*, 839–846.
- WALLRAVEN, C., CUNNINGHAM, D., BREIDT, M., AND BÜLTHOFF, H. 2004. View Dependence of Complex Versus Simple Facial Motions. *ACM SIGGRAPH*, H. H. Blthoff and H. Rushmeier, Eds., 181.
- WALLRAVEN, C., BREIDT, M., CUNNINGHAM, D. W., AND BÜLTHOFF, H. H. 2005. Psychophysical Evaluation of Animated Facial Expressions. In *Proc. of APGV '05*, ACM Press, New York, NY, USA, 17–24.
- YIP, A. W., AND SINHA, P. 2002. Contribution of Color to Face Recognition. *Perception* 31, 8, 995–1003.



# The Evaluation of Stylized Facial Expressions

Christian Wallraven<sup>1</sup>, Jan Fischer<sup>2</sup>, Douglas W. Cunningham<sup>1,2</sup>, Dirk Bartz<sup>2</sup>, Heinrich H. Bühlhoff<sup>1</sup>

<sup>1</sup> Max Planck Institute for Biological Cybernetics, Tübingen, Germany

<sup>2</sup> WSI-GRIS, University of Tübingen, Germany

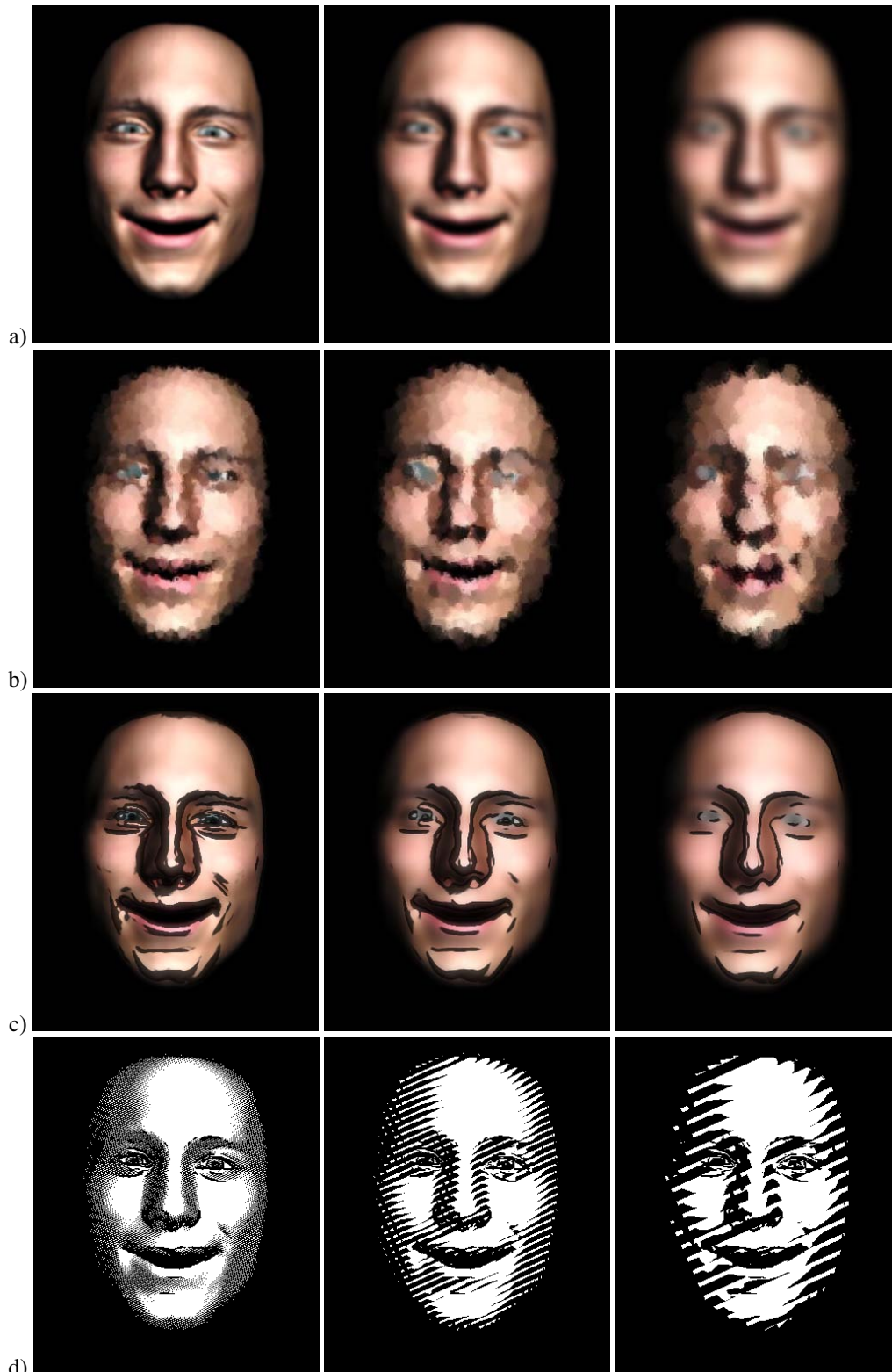


Figure 1: Stylization techniques used in the paper for a happy expression and all three resolution levels. a) Standard Avatar, b) Brush Stroke, c) Cartoon, d) Illustrative stylization - see paper for more details.