

Online Video Stream Stylization



Figure 1: Stylization example. 1: Original video frame; 2 - 4: Three frames of the stylized video with color scheme replacement.

Abstract

This paper gives an automatic method for online video stream stylization, producing a temporally coherent output video stream. Our system transforms video into an abstract style with large regions of constant color and highlighted bold edges. Our system includes two novel components. Firstly, to provide coherent and simplified output, we segment frames, and use optical flow to propagate segmentation information from frame to frame; an error control strategy is used to help ensure that the propagated information is reliable. Secondly, to achieve coherent and attractive coloring of the output, we use a color scheme replacement algorithm specifically designed for an online video stream. We demonstrate real-time performance, allowing our approach to be used for live communication, video games, and related applications.

1 Introduction

Video stream processing has many applications to areas such as live broadcast and communications, video games, and entertainment. Stylization and abstraction of video streams are popular special effects in such live applications [Winnemöller et al. 2006], and can also reduce the perceptual and cognitive effort required to understand the content.

We give here an online, automatic stylization method which performs abstraction to generate cartoon-like animation from an input video stream. Little user interaction is required, other than selection of a reference video which determines preferred output colors (or alternatively a hue histogram), and setting of preferences which control the amount of detail retained in colored regions and line drawing before processing.

Various potential applications exist for such a method which can work in real time. In live communications, stylization may be used to save bandwidth. It may hide someone’s identity or location, for example during an investigatory interview. It may be used to apply artistic effects to live scenes for tourist information. It may simply be used for fun in games or online chatting.

Processing live video streams requires efficient algorithms in order to cope with real-time data. Offline algorithms can make use of future as well as past information, but this is not available in live video streams. Furthermore, because of computational limits, such video has to be processed frame by frame, or at most make use of just a few past frames or an averaged history.

The particular goal of our video stylization approach is to produce an output style which is cartoon-like, having simplified regions

with user-guided colors, and high contrast. In comparison, Winnemöller [2006] produces a different artistic style based on simplified but smoothly shaded contents. Our cartoon-like style means that temporal coherence requirements are particularly strict.

Current stylization methods fall into three categories, each having limitations. Some methods focus on image processing and do not readily generalize to video. Others use simple image filters to achieve real-time performance, producing simplified and smoothly shaded contents, but the output may lack temporal coherence if input video streams are of low-quality. Yet others require significant user interaction to produce high-quality artistic results, and have a high computational cost.

We present a real-time system for a particular style of video stylization while providing good coherence. Our approach benefits from two novel aspects:

- a segmentation strategy which uses optical flow to propagate segmentation information from frame to frame, with an error control strategy to help ensure that the propagated information is reliable,
- a video stream color scheme replacement method that does not require complete knowledge of the source color distribution (i.e. does not need to know future frames), and which applies the color scheme from a reference video (or image, or a user designed histogram of hue) to the input video, while keeping color consistency between frames.

2 Related work

A number of significant papers have addressed stylization, but most of them concentrate on images rather than video [DeCarlo and Santella 2002; Yang and Yang 2008; Pang et al. 2008; Wen et al. 2006; Hertzmann 1998]. Much research has also considered the generation of cartoon-like video from real-world captured video with the aim of reducing the amount of work required for cartoon production, and the need for skilled artists. Examples include papers on ‘Snaketoonz’ [Agarwala 2002], keyframe based rotoscoping [Agarwala et al. 2004], and videotooning [Wang et al. 2004]. Such methods produce high quality results, but still require intensive user interaction, and generally restart video processing every few tens of frames to avoid error accumulation. Such video-processing tools inherently work in an offline manner: to process any given frame they consider future as well as past information. Such methods typically require considerable artistic skill for good results.

DeCarlo [2002; 2004] proposed an *image* stylization approach relying on eye-tracking data to guide image simplification or stylization, using a hierarchical segmentation structure. We pursue an

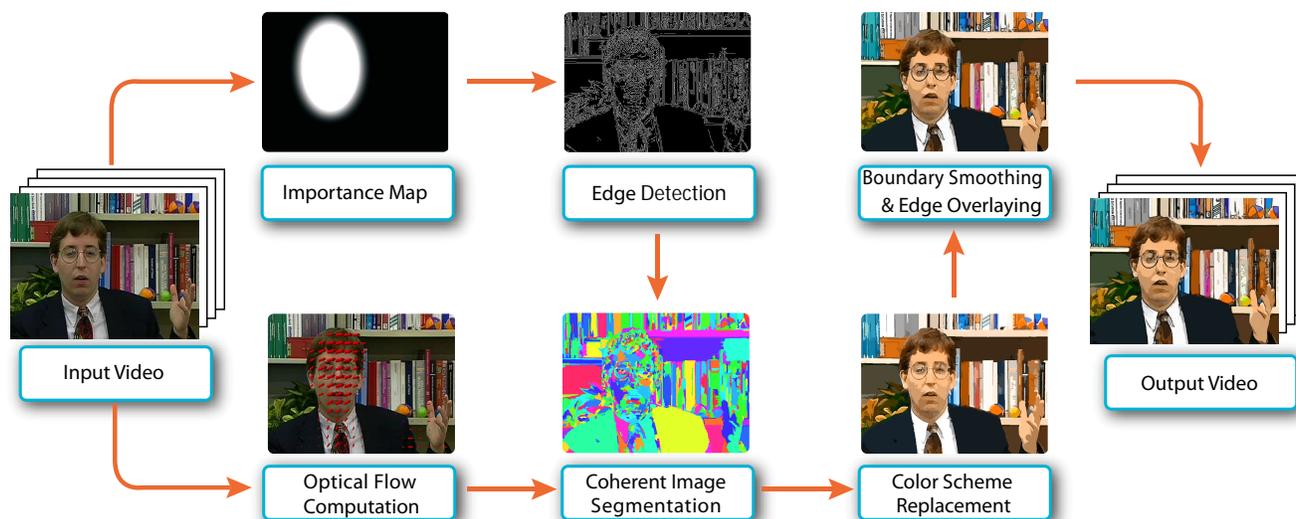


Figure 2: Pipeline. Frame segmentation uses optical flow and edges as inputs, tuned by an importance map, to produce coherent regions. The color scheme replacement step transforms region colors according to a reference video.

87 artistic style similar to DeCarlo [2002], and also use explicit image
88 structure, but produce coherent cartoon-like *animations*.

89 Several video stylization systems have been designed with the main
90 aim of helping artists with labor-intensive procedures [Wang et al.
91 2004; Collomosse et al. 2005]. Such systems use three-dimensional
92 video volume segmentation, and are so computationally expensive
93 that they must be used offline. Our method uses frames pairwise
94 during segmentation, rather than costly multi-frame analysis, while
95 being able to provide coherent online video stream stylization.

96 Other online video stylization methods exist. Fischer et al. [2009]
97 use an automatic stylization method for augmented reality, applying
98 stylization to both virtual and real inputs to fuse virtual objects into
99 a live video stream. Winnemöller [2006] and Kyprianidis [2008]
100 used simple image filters in real-time to process each frame into an
101 abstract style. Neither method takes any particular steps to achieve
102 temporal coherence, simply relying on the implicit image structure
103 itself, and stability of filter outputs. Our method uses an explicit
104 image structure representation which allows for stronger, better de-
105 fined stylization in terms of region shapes, and explicitly propa-
106 gates information from frame to frame to help achieve coherence.
107 Kang [2009] proposed a flow based image abstraction algorithm
108 which also considered shape smoothing, but it cannot be readily
109 applied to video.

110 Paris [2008] give a coherent approach to analyzing and filtering
111 video streams. Combining concepts of isotropic diffusion and
112 Gaussian convolution, this method achieves temporal coherence via
113 bilateral filtering and mean-shift segmentation in real time. How-
114 ever, as Paris notes, using optical flow can provide better results at
115 moving boundaries.

116 Litwinowicz [1997], Hertzmann [2000; 2001], Kovács [2002],
117 Hays [2004] and Vanderhaeghe [2007] use optical flow to direct
118 brush stroke movement or point placement during frames, but they
119 do not take into account errors in optical flow, so limited coherence
120 is achieved in image sequences. Bousseau [2007] also uses optical
121 flow in a watercolor generation pipeline, but relies on temporally bi-
122 directional optical flow interpolation to reduce optical flow errors,
123 an option which is not available in a live video system as future
124 information is unavailable.

125 Several works consider the problem of color scheme extrac-
126 tion [Greenfield and House 2003] and color transfer [Chang et al.
127 2003; Hertzmann et al. 2001; Reinhard et al. 2001]. These mainly
128 focus on image processing, and analyze source and reference fea-
129 tures to determine the color transformation. However, we cannot
130 extend such a method directly to live video stream color transfer,
131 again because of the lack of future information. On the other hand,
132 performing image color transfer for each frame independently can-
133 not guarantee color consistency in the output as source colors vary
134 in each frame. Wang [2006] gives a color transfer method for
135 still images forming a sequence; parameters are adjusted to present
136 gradually changing effects. Image analogies have also been used
137 for color transfer [Hertzmann et al. 2001]. These take into account
138 local similarities in the image, and process it pixel by pixel. The
139 latter is most suited to producing complex textures, but our out-
140 put style uses large regions of slowly varying color. Thus, instead,
141 we have devised a new efficient color replacement method which
142 captures the color style from a given reference video, and applies it
143 with temporal coherence to produce a stylized video stream without
144 the need for source video feature analysis.

145 Our method processes the video stream frame by frame, using a
146 region based representation to achieve a cartoon-like style, and
147 an optical flow based strategy to ensure coherence. Our method
148 is designed to produce relatively highly stylized video, while still
149 achieving real-time performance on a multi-core CPU. Compared
150 to earlier methods such as videotooning, our method is much faster
151 and needs much less user interaction.

152 3 Overview

153 Our goal is to replace incoming live video by stylized animation, in
154 real time. We must thus avoid time consuming non-local analysis
155 of the video (e.g. background reconstruction). During stylization,
156 our three goals are to

- 157 • simplify the content, without oversimplifying important areas,
- 158 • modify the coloring in a controllable manner,
- 159 • retain temporal coherence.

160 Taking these into account, our pipeline is as shown in Figure 2,
161 and includes six key elements: importance map computation, op-

Algorithm 1 Pseudocode for online video stylization

```

while (! end of video stream) do
  //—————Perform coherent image segmentation—————
  Compute importance map;
  Compute Canny edges;
  if (not first frame) then
    //—————Update segmentation using optical flow—————
    Compute optical flow;
    Propagate labels from previous frame;
    Apply morphological filtering to propagated labels;
    Grow regions;
  end if
  for ( $r = \text{initial radius}$ ;  $r \geq 1$ ;  $r \neq 2$ ) do
    Do trapped ball filling with ball radius  $r$ ;
    Grow regions;
  end for
  //—————End of image segmentation—————
  Apply color scheme replacement
  Smooth regions
  Compute DoG edges and overlay them
end while

```



Figure 3: Edge detection controlled by importance map. From Left to right: input, importance map, Canny edge detection, detected edges controlled by importance map.

162 tical flow computation, edge detection, coherent image segmen-
 163 tation, color scheme replacement, and boundary smoothing with edge
 164 overlay.

165 Of these six steps, importance map computation, optical flow com-
 166 putation and edge detection provide inputs for coherent image seg-
 167 mentation. The importance map is used to ensure greater detail in
 168 key areas. As optical flow information is inherently unreliable, we
 169 must be careful to avoid propagating and accumulating optical flow
 170 errors during segmentation. We use a careful error control strategy
 171 when propagating segmentation labels from frame to frame, to en-
 172 sure temporal coherence of region representations between frames.

173 For simple output, we could simply replace each region by its mean
 174 color. Better results are obtained by using a color scheme replace-
 175 ment process based on a desired ‘color style’ learnt offline from a
 176 reference video (or image, or provided by the user as a histogram).
 177 We use a mean-shift based color transformation to move input col-
 178 ors closer to a reference color distribution. We are careful to pre-
 179 serve color consistency between frames, while still not needing to
 180 know the potential colors in future frames.

181 Finally, we smooth boundary curves using a low pass filter to pro-
 182 duce more artistic effects, and then use a difference-of-Gaussians
 183 (DoG) operator on the input frames to detect and overlay edges, as
 184 in [Winnemöller et al. 2006].

185 Algorithm 1 gives pseudocode for the whole algorithm, and details
 186 of image segmentation are described in the following section.

187 4 Coherent image segmentation

188 As in [DeCarlo and Santella 2002], we perform segmentation in
 189 order to simplify image content. Several approaches exist for tem-
 190 porally coherent video segmentation [Zitnick et al. 2005; Kumar
 191 et al. 2005; Xiao and Shah 2005]. Most solve the problem via op-
 192 timization or a mean-shift approach, and thus perform a non-local
 193 analysis of the video. To process a *live* video stream, we need a
 194 method which only uses past frames, and to keep processing re-
 195 quirements low, only a few previous frames can be considered. We
 196 perform segmentation by propagating information from one frame
 197 to the next.

198 Our approach uses optical flow but in a way which avoids accu-

199 mulation of errors. Many optical flow computation methods ex-
 200 ist [Baker et al. 2007]; we use Zach’s [Zach et al. 2007] as it pro-
 201 vides good results in real time. The first frame is segmented using
 202 the trapped ball algorithm [Zhang et al. 2009]; this algorithm is ex-
 203 plained later. In each successive frame, optical flow is used to prop-
 204 agate segmentation labels for each pixel. In the result, some pixels
 205 may have *no* segmentation label (e.g. because they have been newly
 206 exposed from an occluded area), or may have *multiple* labels (e.g.
 207 because of optical flow errors). Labeling such pixels is the key to
 208 coherent image segmentation. Those which are sufficiently close in
 209 distance and color to an existing region are given the label of that
 210 region, while the remainder are segmented into *new* regions again
 211 using the trapped ball algorithm. Algorithm 1 gives pseudocode for
 212 image segmentation; we now consider these steps in detail.

213 4.1 Importance map and edge detection

214 Edges provide strong hints as to where region boundaries should ex-
 215 ist. In image and video stylization, *important areas*, such as faces,
 216 should be depicted in more detail than other areas, which requires
 217 segmenting them into more, smaller regions. Eye tracking [DeCarlo
 218 and Santella 2002] can find important areas, but is of limited appli-
 219 cability. We detect and track faces [Lienhart and Maydt 2002]: an
 220 ellipse is found inside which pixels are assigned high importance,
 221 and outside, low importance; a smoothed real-valued map is used,
 222 as the face position is not entirely robust. Alternative strategies
 223 could also be used to compute the importance map, e.g. taking into
 224 account saliency and motion information [Zhai and Shah 2006], but
 225 at a higher computational cost.

226 Importance map I values lie between 0 and 1, with 1 being the most
 227 important. The important map is used to control the number of
 228 edges found during edge detection: the hysteresis thresholds of a
 229 Canny edge detector [Canny 1986] are set inversely proportional to
 230 importance map values. In areas of greater importance, more edges
 231 are detected, as shown in Figure 3.

232 4.2 Region label propagation by optical flow

233 The first frame is segmented by the trapped ball method based on
 234 the detected edges. For subsequent frames, the previous frame has
 235 been segmented, and optical flow has been computed; these are
 236 used to provide an initial segmentation for the current frame. The
 237 segmentation is represented by labeling pixels with positive inte-
 238 gers corresponding to regions.

239 We start by labeling all pixels of the current frame as 0, i.e. unas-
 240 signed to any region. The optical flow is rounded to determine a
 241 spatial mapping of each source pixel from the previous frame to a
 242 target pixel in the current frame. This target pixel is given the same
 243 label as the source in the previous frame, except as below, when it
 244 retains its 0 label:

- 245 • the optical flow error is larger than a threshold,
- 246 • *more than one* source pixel ends up at this target pixel, and
 247 the labels of these source pixels are different,



Figure 4: Left: top, bottom: input frames 1, 2; Middle: top, bottom: segmentation results for frames 1, 2; Right: top: residual map before morphological filter; Right: bottom: region label propagation result after morphological filter.



Figure 5: Color scheme replacement. Top, left: a frame of the reference video, right: histogram of its H channel; bottom, left: segmented frame before replacement, right: frame after replacement.

- this target pixel has no corresponding source pixel.

To determine if an optical flow error has occurred, we compute the color difference in CIELAB color space between the source and a linear interpolation of the values of the 4-neighboring pixels of the target position. If this exceeds a threshold T_f , the target is labeled 0. T_f is chosen according to the video quality; values in the range 5–20 are typically appropriate. This prevents obvious optical flow failures. Fig. 4(top, right) shows the residual unlabeled pixels (red) after label propagation. Note that we do not explicitly need to detect scene changes in the incoming video: in such cases, the optical flow will have large errors, causing many pixels to be labeled 0, and will hence be resegmented.

After pixel label propagation, region boundaries are typically rather interlaced. These are smoothed, and other mislabelings due to errors in the optical flow corrected, using a morphological filter: if 5 or more of the 8-connected neighbors of the target pixel are the same, but differ from the target pixel, and if its label is 0, then it is changed to the most frequent label in the neighborhood, otherwise it is changed to 0. Fig. 4(bottom, right) shows the propagated labels after morphological filtering, where red pixels are unlabeled.

We must now complete the segmentation by labeling the unlabeled pixels. This is a problem of minimizing label errors for those unlabeled pixels which should be allocated to existing regions, and giving new labels to those pixels which are insufficiently similar to existing regions. At the same time, the labels assigned must respect the edge information, so that the regions are in agreement with any superimposed edges drawn later. These requirements, taken together with efficiency considerations, make it difficult to solve the problem by traditional optimization or graph cut methods. Instead, we improve the region growing and trapped ball filling [Zhang et al. 2009] to finish the segmentation.

4.3 Region growing

For each labeled region, a single, simple, color model is assumed to be an adequate fit. For speed, we use a constant color model, and assume pixels belonging to a region are near to its mean color.

We assign labels while respecting edges taking into account that (i) edge pixels may not be added to a region and (ii) pixels added to a region should agree with its color to within perceptual limits. Let the *reconstruction error* of a labeled pixel be the difference between its actual color and that predicted by the corresponding color model. All unlabeled pixels adjacent to the boundary of any labeled region are put into a priority queue, sorted by reconstruction error with respect to the adjacent region. We then pop the pixel with

minimum reconstruction error from the queue. If the reconstruction error is sufficiently small (in practice, below $20(1.2 - I(p))$ units in CIELab space), the pixel’s label is updated, and it is removed from the priority queue; at the same time, its unlabeled 4-connected neighbors are added to the queue. We repeatedly pop pixels until the queue is empty, or the least reconstruction error is too large.

4.4 Trapped-ball filling

Any pixels still remaining unlabeled belong to some *new* region. In general, multiple new regions may exist, so these unlabeled pixels still need to be segmented.

The basic idea of trapped-ball filling is to move a ball around, limited by a mask, formed here by the already labeled pixels and the detected edges. Each separate area in which the ball may move determines a new region of the image. The trapped ball has the advantage over floodfilling that it cannot leak out between short gaps in the detected edges. Initially, a large ball is used to find each region’s core, and then balls of successively smaller radii are used down to a radius of 1 to add more detailed parts of each region. In practice, the trapped ball approach is implemented using morphological erosion and dilation operations with a circular structuring element having the current ball radius.

The initial ball radius is chosen according to the number of unlabeled pixels N_0 . We set it to $\max(2, \sqrt{N_0}/30)$, which is typically in the range 2–32. In the first frame, *all* pixels are unlabeled, and an original segmentation for the whole frame is obtained using the trapped ball method with a *large* initial ball size. In subsequent frames, unlabeled regions are typically small, and a *smaller* initial ball radius is used. If a scene change occurs, many pixels are now unlabeled, so a larger initial ball size is again used. The extra processing entailed at a scene change can cause a slight lag, but as long as frames are processed faster than in real time, this lag can soon be overcome.

5 Color scheme replacement

Color selection is an important aspect of stylization, and different palettes can produce different emotional responses. Much work has been done on color transfer for images—applying the color scheme

327 from a reference image to a source image. Clearly, in this case,
 328 the color distribution for the *entire* source and reference are known.
 329 However, when processing *live* video, we do not have complete
 330 information: future video content is unknown. Thus our recoloring
 331 problem has two specific requirements: the transformation should
 332 cover the entire color space to allow for *any* possible future colors,
 333 and as we have video rather than images, we must also ensure inter-
 334 frame color stability for each region in the video.

335 Our color scheme replacement method uses the concept of *mean*
 336 *shift* [Comaniciu and Meer 2002]. Given a reference image or cartoon
 337 video clip, we count the number of pixels of each color to
 338 determine the color distribution. We must also choose the desired
 339 number of iterations and kernel radius for the mean shift procedure.
 340 Then, to determine the mapping of each input color to an output
 341 color, we determine a mean shift vector towards a local maximum
 342 in the color distribution, which shifts each input pixel’s color to-
 343 wards the nearest peak in the reference distribution. Overall, given
 344 an input color distribution, the mapping transforms it towards the
 345 reference color distribution.

346 It would be costly to compute the reference distribution and map-
 347 ping directly in CIELab color space. Instead, we simply consider
 348 the hue value in HSV space, and compute the reference distribu-
 349 tion and mean shift vector using H alone. In an analogous manner,
 350 [Cohen-Or et al. 2006] and [Sawant and Mitra 2008] also use hue
 351 as the most important feature of color when performing color har-
 352 monization, i.e. mapping a source color distribution to a template
 353 distribution.

354 As these authors note, there is discontinuity in new color when two
 355 nearby colors are transformed to different local maxima, which may
 356 causing flickering between frames. Thus, we also use an interframe
 357 blending method which takes into account correspondence between
 358 regions in successive frames; this information is provided by the
 359 image segmentation results. Thus, after first computing an ideal
 360 *target color* by the mean shift procedure, we modify it to give the
 361 *actual output color* by interframe blending.

362 To find the color mapping, first, we compute the color distribution
 363 of the reference video in HSV color space (which is more useful
 364 than RGB from a perceptual point of view), considering the three
 365 channels separately. In each frame, for the H channel, we compute a
 366 color histogram for all pixels having $S > 0.125$ and $V > 0.125$. Also
 367 we compute the mean μ_S and standard deviation σ_S for S channel
 368 pixels having $V > 0.125$, and the mean pixel value μ_V of the V
 369 channel. (Alternatively, instead of providing a reference video, the
 370 user could provide a reference image, or even manually design a
 371 color histogram represented in this way).

For each frame of the incoming video, all pixels in each segmented
 region are given the same target color. The *source color* (h, s, v) for
 a region is computed as the average color of all pixels in that region.
 This is mapped to a target color (h', s', v') . First, we compute μ_s ,
 σ_s and μ_v for this *frame* in an analogous manner to how they were
 computed for the reference video. We then compute (h', s', v') as
 follows. h' is set to the mean shift of h in the H channel histogram,
 by iterating the formula below k times; typically $k = 3$:

$$h_{i+1} = \frac{\sum_{c \in N(h_i)} c D(c)}{\sum_{c \in N(h_i)} D(c)}.$$

372 Here h_0 is h , and h_4 is the desired h' . $N(h_i)$ represents a 30° neigh-
 373 borhood in the histogram and $D(c)$ represents the histogram value
 374 for c . (This function is precomputed for all possible h values and
 375 stored in a look-up table for speed). Simple formulae are used for
 376 s' and v' with the aim of adjusting the brightness and contrast to be



Figure 6: Boundary smoothing and edge overlaying. Top, left: segmentation result, right: after boundary smoothing; bottom, left: detected lines, right: smoothed result after edge overlay.

377 closer to those of the reference video:

$$\begin{aligned} s' &= \mu_S + (s - \mu_S) \sigma_S / \sigma_s, \\ v' &= v(2 + \mu_V / \mu_v) / 3. \end{aligned}$$

378 This approach ensures that a given hue is always mapped to the
 379 same new hue in different frames. Such a mapping is not appropri-
 380 ate for s and v if the input video varies in brightness.

After obtaining the *target color*, we now compute the *actual output color*. Suppose some region R_i has a target color c_i , and the corresponding region R_{i-1} in the previous frame had actual output color o_{i-1} . We compute its actual output color o_i by color blending to improve color coherence:

$$o_i = (1 - \alpha(R_i, R_{i-1})) o_{i-1} + \alpha(R_i, R_{i-1}) c_i$$

381 Here $\alpha(R_i, R_{i-1})$ measures the correlation of the two correspond-
 382 ing regions: $\alpha(R_i, R_{i-1})$ is the ratio of: the number of pixels in R_{i-1}
 383 whose destination lies in R_i after optical flow, to the geometric av-
 384 erage of the areas of R_{i-1} and R_i . If region R_i is a *new* region in this
 385 frame, we simply set $o_i = c_i$.

386 An H channel color histogram extracted from a reference video and
 387 an example of color scheme replacement are shown in Fig. 5 (final
 388 results after boundary smoothing and edge overlaying are shown in
 389 Fig. 6).

6 Results and discussion

391 We have implemented our framework using a combination of CPU
 392 and GPU, on an Intel Core2 Quad Q9300 CPU at 2.5GHz, with
 393 4GB memory, and a GeForce 8600GT, using CUDA, OpenCV and
 394 Visual C++ 2008’s parallelizing compiler. Performance depends on
 395 image size and framework parameters. For a CIF (352×288) video
 396 stream, commonly used for live communication, face detection and
 397 edge detection is done on one CPU core, taking 40ms per frame,
 398 while simultaneously, optical flow computation takes 20ms on the
 GPU. Image segmentation takes from 20–30ms, depending on the
 number of residual unlabeled pixels. Color scheme replacement
 and edge overlay take under 10ms. Boundary smoothing takes 20–
 50ms, depending on the number of boundary pixels. Typically we
 can process a CIF video stream with default parameter settings at



Figure 7: Video stylization results. Columns 1, 3: original video frames; columns 2, 4: corresponding stylization results.



Figure 8: Effect of importance map. Left: segmentation with importance map, right: segmentation without importance map.

about 9–12 frames per second using the CPU for everything but optical flow (we directly use Zach’s GPU optical flow code [Zach et al. 2007]). We note that edge detection, morphological operations and the DoG filter could also benefit from GPU implementation, permitting faster frame rates or higher video resolution. The current Nvidia GTX295 graphics card offers perhaps 10× the performance of the older GeForce 8600GT card.

Figs. 1 and 7 show various stylization results. The goals of stylization are subjective. Ours are to achieve cartoon-like effects with simplified content comprising well-defined regions with bold boundaries, while using a chosen color style, and retaining temporal coherence. This definition is consistent with that in [DeCarlo and Santella 2002]. In comparison, Winnemöller aims to produce soft, simplified content, achieved by smoothing, as shown in Fig. 9. Our method has the ability to perform object shape simplification. The side-by-side comparison is provided in the supplemental material. Different styles can be produced by the user, but objective evaluation of the level of success is difficult. Our results have been independently evaluated by a company who randomly chose 20 employees to subjectively evaluate 5 video clips each. The average score awarded was 85/100, with 60/100 being considered acceptable; 97% of evaluations were at a level of 60/100 or higher.

Fig. 8 shows the effect of importance map in the segmentation. With importance map, the segmentation produces more regions but well maintains the detail of the faces, while other areas are simplified. The incoherence of importance map will also decrease the coherence of final stylized video. So if using importance map, better and faster temporal coherent saliency map is important to our method.

Optical flow is the sole information used to guide interframe correspondence, so its accuracy is very important. Zach’s method works

well in most cases, as demonstrated by the number of residual pixels (Figs. 4 and 10). When the scene changes slowly, few residual pixels remain (Fig. 10, above), and most are labeled, according to color, during the region growing step. When the video content changes rapidly, the number of residual pixels increases. Indeed, essentially the whole image is composed of residual pixels at scene changes (Fig. 10, below), as useful optical flow information no longer exists. Thus, we do not need to explicitly detect scene changes: in cases with large numbers of residual pixels, trapped-ball filling provides resegmentation of the scene. While newly generated regions may cause certain inconsistencies, viewers find lack of coherence much more noticeable in slowly changing scenes than in fast moving scenes. Currently we use the rounded optical flow for label propagation, which may cause flickering in long and thin regions. And image noise and motion blur may also decrease the quality of optical flow, any ultimately decrease the coherence of final results.

Fig. 11 illustrates the results of color scheme replacement using different reference images (these have been chosen to exaggerate the effects, and are atypical). The resulting hue distributions have been transformed towards the reference distribution so that peaks of hue agree with the reference image’s hue peaks: output pixels exhibit colors which tend to follow the reference distribution. To test color scheme replacement, we randomly choose 20 pictures as input, and applied our method using the 3 reference images in Figs. 5 and 11, giving 60 output pictures. These were evaluated by 10 participants, who were each shown 10 pictures randomly chosen from the above 60 results, and asked which of the three reference images had the most similar color style. A correct correspondence was suggested in 93/100 cases. The others contained many green pixels whose hue was almost unchanged in color scheme replacement, as they were far from peaks in the hue histogram.

Our method has other limitations. Like Winnemöller’s bilateral filtering method, our method needs to estimate visual salience. Certain high-contrast features such as specular highlights that may be emphasized by our framework are actually deemphasized in human vision, and repeated texture causes unwanted edge responses. Another limitation lies in our handling of occluded regions. When such regions are gradually uncovered by an overlying object, they typically initially merge with the overlying object, then suddenly separate from it when they have grown to a sufficient size. This may cause jumps in color, although our approach in Section 5 alleviates such effects to some extent.



Figure 9: Comparison of stylization effects. From left to right: source image, results using Winnemöller’s method, results of color quantization, results using our method.



Figure 10: Label propagation by optical flow: two examples. From left to right: previous frame, current frame, residual map before morphological filtering, propagation map after morphological filtering.

7 Conclusions

Our real-time video stream stylization method uses a segmentation strategy guided by optical flow, with care taken to appropriately allow for optical flow errors, so that we can consistently segment each frame while preserving temporal coherence. To achieve attractive stylization, we use a color scheme replacement method which applies colors learnt from a video clip or an image.

Various areas exist for future work. Firstly, we intend to investigate more general ways of computing the importance map while providing temporal coherence, yet with sufficient performance for a real-time system. One obvious consideration is that moving objects are typically more important than the static background. Secondly, we also intend to consider vectorization of the stylized video. Currently our methods directly produce vectorized output on a *per-frame basis* in the form of a boundary representation of each region together with its color. However, determining *inter-frame* correspondences in terms of affine transformation of coherent regions, and region overlap, would allow more semantically useful vectorization, as well as enabling greater compression. Thirdly, we would like to extend this work to other styles of rendering, such as oil paintings. Our experiments so far have shown that the more structured the appearance produced by the rendering method, the more objectionable any distortions due to optical flow errors appear. Finally, we only use constant color models during the segmentation, and higher order color models (e.g. quadratic models) would provide greater scope for stylization.

References

AGARWALA, A., HERTZMANN, A., SALESIN, D. H., AND SEITZ, S. M. 2004. Keyframe-based tracking for rotoscoping and animation. In *ACM SIGGRAPH '04*, 584–591.

AGARWALA, A. 2002. Snaketoonz: a semi-automatic approach to creating cel animation from video. In *NPAP '02*, 139–ff.

BAKER, S., ROTH, S., SCHARSTEIN, D., BLACK, M., LEWIS, J., AND SZELISKI, R. 2007. A database and evaluation methodol-

ogy for optical flow. In *International Conference on Computer Vision*, 1–8.

BOUSSEAU, A., NEYRET, F., THOLLOT, J., AND SALESIN, D. 2007. Video watercolorization using bidirectional texture advection. *ACM Transactions on Graphics* 26, 3, 104.

CANNY, J. 1986. A computational approach to edge detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 8, 6, 679–698.

CHANG, Y., SAITO, S., AND NAKAJIMA, M. 2003. A framework for transfer colors based on the basic color categories. *Computer Graphics International Conference* 0, 176.

CHUNG-MING WANG, Y.-H. H., AND HUANG, Y.-H. 2006. An effective algorithm for image sequence color transfer. *Mathematical and Computer Modelling* 44, 7-8, 608–627.

COHEN-OR, D., SORKINE, O., GAL, R., LEYVAND, T., AND XU, Y.-Q. 2006. Color harmonization. *ACM Transactions on Graphics (Proceedings of ACM SIGGRAPH)* 25, 3, 624–630.

COLLOMOSSE, J. P., ROWNTREE, D., AND HALL, P. M. 2005. Stroke surfaces: Temporally coherent artistic animations from video. *IEEE Transactions on Visualization and Computer Graphics* 11, 5, 540–549.

COMANICIU, D., AND MEER, P. 2002. Mean shift: a robust approach toward feature space analysis. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 24, 5, 603–619.

DECARLO, D., AND SANTELLA, A. 2002. Stylization and abstraction of photographs. In *ACM SIGGRAPH '02*, 769–776.

FISCHER, J., HALLER, M., AND THOMAS, B. 2009. Stylized Depiction in Mixed Reality. *International Journal of Virtual Reality*.

GREENFIELD, G. R., AND HOUSE, D. H. 2003. Image recoloring induced by palette color associations. *International Conferences in Central Europe on Computer Graphics, Visualization and Computer Vision* 11, 189–196.

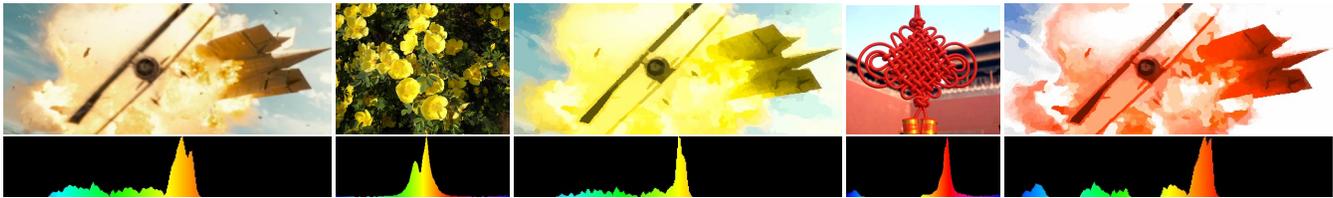


Figure 11: Color scheme replacement examples. Left to right, above: source image, reference image 1, color transferred from reference image 1, reference image 2, color transferred from reference image 2. Below: corresponding hue histograms in HSV space.

- 545 HAYS, J., AND ESSA, I. 2004. Image and video based painterly
546 animation. In *NPAR '04*, 113–120.
- 547 HERTZMANN, A., AND PERLIN, K. 2000. Painterly rendering for
548 video and interaction. In *NPAR '00*, 7–12.
- 549 HERTZMANN, A., JACOBS, C. E., OLIVER, N., CURLESS, B.,
550 AND SALESIN, D. H. 2001. Image analogies. In *ACM SIG-*
551 *GRAPH '01*, 327–340.
- 552 HERTZMANN, A. 1998. Painterly rendering with curved brush
553 strokes of multiple sizes. In *ACM SIGGRAPH '98*, 453–460.
- 554 HERTZMANN, A. 2001. Paint by relaxation. In *CGI '01: Computer*
555 *Graphics International 2001*, IEEE Computer Society, Washing-
556 ton, DC, USA, 47–54.
- 557 KANG, H., LEE, S., AND CHUI, C. K. 2009. Flow-based image
558 abstraction. *IEEE Transactions on Visualization and Computer*
559 *Graphics 15*, 1, 62–76.
- 560 KOVÁCS, L., AND SZIRÁNYI, T. 2002. Creating animations combin-
561 ing stochastic paintbrush transformation and motion detec-
562 tion. *International Conference on Pattern Recognition 2*, 1090–
563 1093.
- 564 KUMAR, M. P., TORR, P. H. S., AND ZISSERMAN, A. 2005.
565 Learning layered motion segmentation of video. In *International*
566 *Conference on Computer Vision*, I: 33–40.
- 567 KYPRIANIDIS, J. E., AND D"OLLNER, J. 2008. Image abstrac-
568 tion by structure adaptive filtering. In *Proc. EG UK Theory and*
569 *Practice of Computer Graphics*, 51–58.
- 570 LIENHART, R., AND MAYDT, J. 2002. An extended set of haar-like
571 features for rapid object detection. vol. 1, 1–900–I–903 vol.1.
- 572 LITWINOWICZ, P. C. 1997. Processing images and video for an
573 impressionist effect. In *ACM SIGGRAPH '97*, 407–414.
- 574 PANG, W.-M., QU, Y., WONG, T.-T., COHEN-OR, D., AND
575 HENG, P.-A. 2008. Structure-aware halftoning. *ACM Trans-*
576 *actions on Graphics (SIGGRAPH 2008 issue) 27*, 3, 89:1–89:8.
- 577 PARIS, S. 2008. Edge-preserving smoothing and mean-shift seg-
578 mentation of video streams. In *ECCV08*, II: 460–473.
- 579 REINHARD, E., ASHIKHMIN, M., GOOCH, B., AND SHIRLEY, P.
580 2001. Color transfer between images. *IEEE Comput. Graph.*
581 *Appl. 21*, 5, 34–41.
- 582 SANTELLA, A., AND DECARLO, D. 2004. Visual interest and
583 npr: an evaluation and manifesto. In *NPAR '04*, 71–150.
- 584 SAWANT, N., AND MITRA, N. 2008. Color harmonization for
585 videos. In *Sixth Indian Conference on Computer Vision, Graph-*
586 *ics & Image Processing, 2008. ICVGIP '08.*, 576–582.
- 587 VANDERHAEGHE, D., BARLA, P., THOLLOT, J., AND SILLION,
588 F. 2007. Dynamic point distribution for stroke-based rendering.
In *Rendering Techniques 2007 (Proceedings of the Eurographics*
589 *Symposium on Rendering)*, 139–146.
- 590
- 591 WANG, J., XU, Y., SHUM, H.-Y., AND COHEN, M. F. 2004.
592 Video tooning. *ACM Transactions on Graphics 23*, 3, 574–583.
- 593 WEN, F., LUAN, Q., LIANG, L., XU, Y.-Q., AND SHUM, H.-Y.
594 2006. Color sketch generation. In *NPAR '06*, 47–54.
- 595 WINNEMÖLLER, H., OLSEN, S. C., AND GOOCH, B. 2006. Real-
596 time video abstraction. *ACM Transactions on Graphics 25*, 3,
597 1221–1226.
- 598 XIAO, J., AND SHAH, M. 2005. Motion layer extraction in
599 the presence of occlusion using graph cuts. *IEEE Transactions*
600 *on Pattern Analysis and Machine Intelligence 27*, 10 (October),
601 1644–1659.
- 602 YANG, C.-K., AND YANG, H.-L. 2008. Realization of seurat's
603 pointillism via non-photorealistic rendering. *The Visual Com-*
604 *puter 24*, 5, 303–322.
- 605 ZACH, C., POCK, T., AND BISCHOF, H. 2007. A duality based
606 approach for realtime tv-l1 optical flow. In *Pattern Recognition*
607 *(Proc. DAGM)*, 214–223.
- 608 ZHAI, Y., AND SHAH, M. 2006. Visual attention detection in video
609 sequences using spatiotemporal cues. In *ACM MULTIMEDIA*
610 *'06*, 815–824.
- 611 ZHANG, S.-H., CHEN, T., ZHANG, Y.-F., HU, S.-M., AND
612 MARTIN, R. 2009. Vectorizing cartoon animations. *IEEE Trans-*
613 *actions on Visualization and Computer Graphics 15*, 4, 618–629.
- 614 ZITNICK, C. L., JOJIC, N., AND KANG, S. B. 2005. Consis-
615 tent segmentation for optical flow estimation. In *International*
616 *Conference on Computer Vision*, II: 1308–1315.