

Graphic Narrative with Interactive Stylization Design

I. Garcia-Dorado and P. Getreuer and M. Le and R. Debreuil and A. Kauffmann and P. Milanfar

Google Research



Figure 1: Our system converts a set of input images into a series of storyboards with varying layouts, images, and stylization.

Abstract

We present a system to convert any set of images (e.g., a video clip or a photo album) into a storyboard. We aim to create multiple pleasing graphic representations of the content at interactive rates, so the user can explore and find the storyboard (images, layout, and stylization) that best suits their needs and taste. The main challenges of this work are: selecting the content images, placing them into panels, and applying a stylization. For the latter, we propose an interactive design tool to create new stylizations using a wide range of filter blocks. This approach unleashes the creativity by allowing the user to tune, modify, and intuitively design new sequences of filters. In parallel to this manual design, we propose a novel procedural approach that automatically assembles sequences of filters for innovative results. We aim to keep the algorithm complexity as low as possible such that it can run interactively on a mobile device. Our results include examples of styles designed using both our interactive and procedural tools, as well as their final composition into interesting and appealing storyboards.

CCS Concepts

•Computing methodologies → Image manipulation; Image processing; •General and reference → Design; •Human-centered computing → Interactive systems and tools;

1. Introduction

A *storyboard* is a series of image panels representing a sequence of actions that will be performed and/or captured in motion, such as in a film, an animation, or onstage. Storyboards were first developed at Walt Disney Studios during the early 1930s [Whi04]. The idea was, during early stages of brainstorming and development, to draw key scenes on separate sheets of paper and pin them to create a narrative bulletin board story. This allowed them to design and iterate faster by interchanging the order and by replacing panels with

a more interesting ones at any point. Films have relied on this approach since the days of silent cinema. Nowadays, storyboards are used to illustrate, summarize, and pre-visualize many kinds of dynamic content including TV shows, movies, plays, novels, and even software. Storyboards represent dynamic content in a concise and explanatory manner, similar in many ways to comics. The main distinction is while storyboards focus mainly on the action and aim to convey information, comics focus on the conversation (with speech bubbles) with the goal of entertainment.

Our approach uses storyboarding as an intermediate representation between videos and comics, extracting representative frames

and laying them out on a stylized single page. We do not seek to optimize the layout, the content, and the stylization to create a unique output. Our goal is to allow a user to quickly explore dozens or even hundreds of alternatives and to select the representation(s) that best fit his or her requirements (See Fig. 1). Quick exploration and filtering require algorithms that run at near real-time speeds on mobile devices, necessarily limiting their complexity. Moreover, our approach applies beyond video to any set of input images such as photo albums and images from multiple different sources. In summary, our goal is to translate videos and images into storyboards that present the scenes within them in a visually pleasing, effective, and interesting manner.

The stylization filters and several other aspects developed here are components of Google Research’s “Storyboard” app [sto]. A variant of the described pipeline is used in this application.

Central to the visual impact of our storyboards is stylization. Most recent work in stylization has focused on convolutional neural network style transfer algorithms (e.g., [GEB16a]). However, as our results will show, the resulting effects are hard to control, despite recent efforts to make them more tunable [GEB*16b]. Alternatives such as abstraction and filtering focus on creating one appealing stylization that would satisfy the requirements (e.g., abstraction simplifies an image by removing details). Our goal is to create a filter block framework that comprises both simple and advanced filters that enable creative control and fine tuning. We will also explore a procedural approach, inspired by [MWH*06], for generating novel filters.

The main contributions of this work can be summarized as follows:

- We propose an automatic system to convert an album or video into storyboards, i.e., interesting and visually pleasing visual stories. Our method focuses on speed so it can interactively run on a mobile device.
- We present an interactive design tool for stylization. This tool allows tuning, modifying, and designing styles on the fly, unleashing creativity.
- We propose a novel procedural style generation technique built on our design tool. Using simple rules we create original stylization effects with a small number of filter blocks.

The rest of the paper is organized as follows: Section 2 reviews previous work. Section 3 is an overview of our system. Section 4 describes the basic process for discarding blurry and near-duplicate images. Section 6 explains the process for selecting and framing input images. Section 7 introduces our design application and the filter blocks. Section 8 presents our results. Finally, Section 9 contains our conclusions and future work.

2. Previous Work

We review work in comic generation, summarization, re-targeting, and stylization.

2.1. Movie to Comic Generation

The papers most related to our work are those focused on generating a comic from a movie sequence.

Ryu et al. [RPL*08] present a semi-automatic method for converting a film into a comic. The user manually selects the frames, relevant stickers, and word balloons. Then, this is integrate into a comic using a black and white cartoonization effect. Chu et al. [CYW15] propose a method to transform a temporal image sequence into a comic-based presentation using a genetic algorithm. This work focuses on the page allocation, layout, and speech balloon placement. Wang et al. [WHY*12] present a method for automatically generating a comic representation of a movie. The aim of the work is to keep all the information (including the script). Thus, this work focuses on finding the shots where a character talks so they can draw the appropriate speech bubble. As a result, a twenty-minute video becomes a seventy-page comic. Jing et al. [JHG*15] extend this approach by proposing a content-driven layout system that generates layouts dynamically rather than relying on pre-defined templates.

In contrast to all these methods, our work:

- allows as input any set of images (i.e., it is not limited to videos),
- does not require any additional information such as scripts,
- eschews summarization in favor of exploration, and
- focuses on interactivity and responsiveness.

The only significant overlap is the cartoon stylization step where we present a novel interactive tool for designing styles.

2.2. Video Summarization

Summarization techniques focus on selecting good individual frames or shots from a video sequence to create a summary.

Keyframe-based approaches attempt to extract the different interesting parts of the video using low-level features (e.g., image difference [ZWZS97]) or higher-order descriptions of the video (e.g., creating hierarchical visual models [LHC10]). Saliency detection [LYS*11] methods rely on finding interesting objects and using them to identify important frames of the video.

More recent work uses supervised deep learning to perform video summarization. Gygli et al. [GGVG15] present an approach to learn the importance of global characteristics optimized for multiple summarization objectives. Zhang et al. [ZCSG16] use training videos to transfer summary structures non-parametrically.

All these techniques look for a uniquely representative subsample of shots to summarize the video. While our work might benefit from these techniques, we require a simple and quick process because our objective is different: given an input, produce a wide range of options for a user to select among.

2.3. Stylization

Stylization of images and videos represents a broad research area. We review three main groups: filtering, video stylization, and style transfer.

• Filtering and Abstraction.

Winnemöller et al. [WKO12] explore Difference-of-Gaussians (XDoG) to create interesting sketch and hatching effects. Kang et al. [KLC07] extend XDoG by adding an edge tangent flow block to create smooth edges.

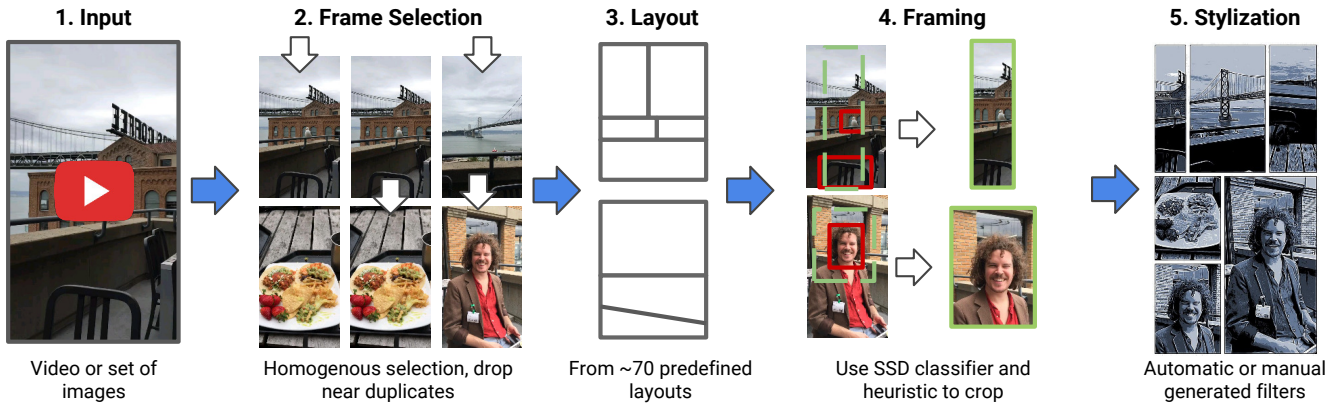


Figure 2: Pipeline.

Kyprianidis and Döllner [KD08] use oriented separable filters and XDoG to achieve a high level of image abstraction. Kang et al. [KLC09] improve the level of abstraction from this approach by adding a flow-based step. Other more complex algorithms simplify images using an advanced multi-scale detail image decomposition [TM16].

For a more comprehensive survey of artistic stylizations we refer the reader to Kyprianidis et al. [KCWI13].

We use some of these filters as building blocks for our stylization blocks as we explain in Section 7.

- **Video Stylization.**

As an extension of image filtering, some works have focused on speeding up the stylization process to run at interactive rates. Winnemöller et al. [WOG06] use a bilateral filter iteratively to abstract the input, then quantize the background color and overlay XDoG to produce strong outlines. Our system can achieve similar results using a different set of filters. Barnes et al. [BZL*15] precompute a multidimensional hash table to accelerate the process of finding replacement patches. This structure enables them to stylize a video in real time using a large collection of patch examples.

- **Style Transfer.**

As an alternative to explicit filter creation, a wide range of works have developed a technique called *style-transfer*. Style transfer is the process of migrating a style from a given image (reference) to the content of another (target), i.e., synthesizing a new image which is an aesthetic mixture of the two. Recent work on this problem uses Convolutional Neural Networks (CNN). Gatys et al. [GEB16a] pose the style-transfer problem as an energy minimization task, seeking an image close to the target using a CNN as score. Elad and Milanfar [EM17] present an approach based on texture synthesis. Their approach copies patches from the reference image to the target while maintaining the main features of the content image using a hierarchical structure.

Despite the visual appeal of these approaches, their complexity is a major drawback. The method described by Gatys et al. [GEB*16b] can take up to an hour to compute a single image. More recent works have focused on addressing this issue. Johnson et al. [JAFF16], achieve real-time style transfer with sim-

plified networks running on a high-end desktop GPU. However, achieving similar results on a full HD image on a mobile device would require tens of seconds. Elad and Milanfar’s [EM17] approach takes multiple minutes to run on device.

Style transfer has several other drawbacks. First, as we will show in the results, output quality depends directly on the reference image. While often considered an advantage, having a specific template can generate inconsistent and undesirable results for different inputs (e.g., very bright/dark images). Second, current style transfer approaches do not provide sufficient aesthetic control for creating storyboards. Gatys et al. [GEB*16b] extended their own work to introduce control over spatial location, color information, and scale. However, this does not allow the fine tuning required to design new stylizations. In this context, Barnes et al. [BZL*15] presented a method to efficiently query patches within a large dataset and replace each patch of the target image with one from the reference image. This approach does not allow control over color, line weight, or other aspects that characterize the hand drawn aesthetic of storyboards and comics.

Note that our method does not compete against these approaches, each of the techniques above can be incorporated into our system as a new block that further enriches expressiveness.

3. Overview

Our system takes as input a set of images such as an album or a video sequence and converts them into one or more stylized storyboards. The aim is to create a visually pleasing graphic narrative that shows in a glance the input. Rather than creating a single unique summarization, our approach makes it easy to explore countless alternatives and quick to choose preferred variations in images, layout, and stylization. Figure 2 describes the system pipeline:

1. **Input:** The input of our system can be any video (e.g., a personal video sequence that last few tens of seconds) or any set of images (e.g., a personal photo album of a trip).
2. **Frame Selection:** From the given input frames, we select a subset of frames that we will use to create the storyboard (Sec. 4).

3. **Layout:** We choose a layout to hold the selected frames (Sec. 5).
4. **Framing:** We crop and zoom the selected frames to fit the layout (Sec. 6).
5. **Stylization:** We stylize the resulting frame to give it a storyboard-like appearance (Sec. 7).

4. Image Selection

In this section, we describe how we select a subset of individual video frames or photos that will be used to create the final storyboard. Our approach needs to work for both videos and images and still run at interactive rates. Keyframe approaches and deep learning summarization proved too computationally expensive, so we explored simpler approaches.

For videos, we found that uniform frame sampling yields a representative set of input images. For image albums, we use all the available images as initial input set. In both cases, we remove duplicates and near-duplicates.

4.1. Duplicates

Albums and videos contain visually similar scenes that might result in repetitive and redundant storyboards. Thus, we require a fast and reliable approach to detect near-duplicate images. In videos, most duplicates appear as successive frames; however, in photo albums it might not be the case, since it is common to collect images from different sources where the time stamp and viewpoint often differs.

One approach to finding duplicate files is to use hash functions, such as SHA-1 or MD5 [Riv92]. However, this limits the detection to exact duplicates. Because we are interested in near duplication, we use a modified version of perceptual hashing [ND16]. The process is as follows:

1. The image is converted to grayscale.
2. The image is downsampled to 9×8 pixels. We downscale it recursively by 2 until it reaches the closest scale factor and then apply a final spatial downscaling filter.
3. Gradient is computed vertically on odd rows and horizontally on even rows.
4. We create a 64-bit fingerprint according to the sign of the gradient, vertically in the odd rows and horizontally in the even rows.

This `uint_64` can be used as a fingerprint of each image and compared with other images using the Hamming distance, i.e., to compare two images we just count the number of different bits of their fingerprints. This approach, since it aggregates many pixels into one, is very robust to noise, color change, and small changes in camera position. Moreover, the distance can give an estimation of change between images.

4.2. Detect Blurry Images

Given a set of near-duplicate images, we select the sharpest image.

There are many methods in the literature for detecting blurry images. Aydin et al. [ASG15] use the residuals of iterative bilateral filters to compute a sharpness metric of the image. Talebi and Milanfar [TM17] evaluate image quality using a CNN trained on images with different levels of distortion. However, these methods are too computationally expensive.

After analyzing different alternatives, we found the gradient of the image is the best indicator of sharpness. We compute a subsampled version of the gradient of the image as:

$$hist(i) = \#\{(x,y) \in \Omega : \lfloor \nabla(x,y) \rfloor = i\} \quad (1)$$

where $\#$ is the cardinality, Ω is the subsample domain of the image (in our case, we subsample by a factor of 3 in each direction), and ∇ is the gradient (in our case, computed using central differences).

We then compute the *sharpness metric* as the inverse of the *coefficient of variation* of the histogram:

$$sharpness = \frac{\overline{hist}}{\sqrt{hist_2}} \quad (2)$$

where $\overline{hist} = \frac{1}{\#\Omega} \sum_i hist(i)$, i.e., the mean value of the histogram and $\sqrt{hist_2} = \sqrt{\frac{\sum_i (hist(i) - \overline{hist})^2}{\#\Omega}}$, i.e., its standard deviation. The \overline{hist} is already a good metric for sharpness between near-duplicates; we divide by the standard deviation to normalize the *sharpness* measure so it is more robust to different content.

5. Layout

Because we aim to create many alternative storyboards (in contrast to a unique optimal solution), we initially believed automatic layout generation would provide the best approach. We based our initial efforts on the idea of parcel subdivision in urban planning [VGDA*12], i.e., given a city block (in this case the whole storyboard), how it can be divided into parcels (in this case panels).

This process yielded feasible but aesthetically unsatisfying layouts, so we opted instead to have designers manually create a set of over 70 predefined layouts. Some examples can be seen in Fig. 3.

Having a predefined set of layouts allows us to ensure consistently visually pleasing and non-repetitive results. Moreover, since we know ahead of time the position of all the panels, we can optionally merge content across adjacent panels while preserving their frames, a technique commonly employed by comic artists for dramatic effect.

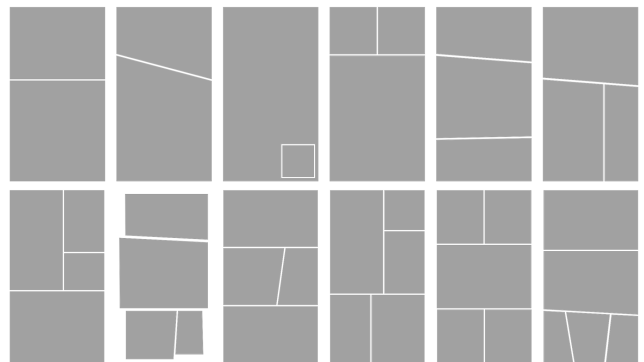


Figure 3: Example set of layout templates.

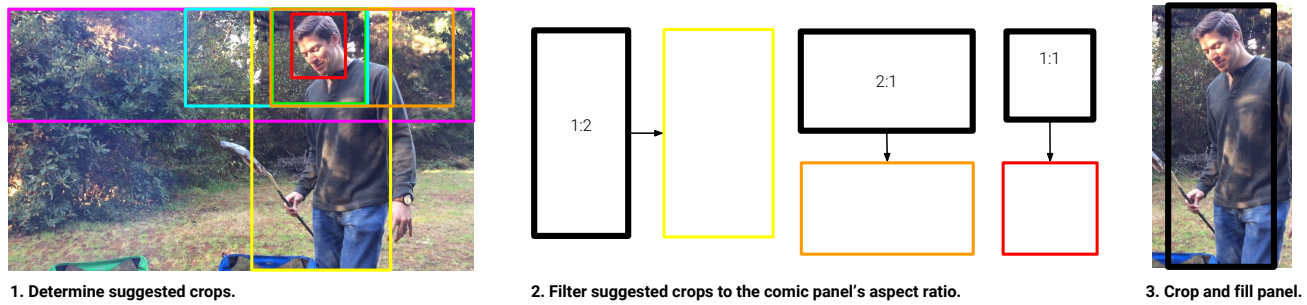


Figure 4: Framing: Cropping and Zoom.

6. Framing: Cropping and Zoom

In this section, we explain how we frame images into layout panels. We use cropping and zooming to conform input images to fit specific panels. This allows us to explore more compositions than optimizing the layout to the selected frames – especially and find more visually interesting crops when the image lacks recognizable objects (e.g., when the image is just a landscape).

Most summarization comic approaches rely on detecting faces to select the frames. We noticed that for our goal, despite faces being an essential part of videos and albums, just detecting faces would produce repetitive panels, especially in shorter videos. Basing frame selection on faces does not generalize to many other common video subjects such as places, vehicles, or pets. Thus, we use two CNNs: one trained to detect faces and another trained to detect objects. Since we wanted to be very fast and to be able to run on device, we use two *single shot multibox detector* (SSD) presented by Liu et al. [LAE*16].

The SSD returns an approximate bounding box for detected objects, and from this, we generate a set of suggested crops (Fig. 4.1). These crops are expansions of the detected bounding box (in red) with the following rules: extend 20 pixels, expand the crop bounding box’s height and/or width by 10% or 50%, or extend the bounding box to the image’s full height or width. Once we have the expanded bounding boxes, we find the best match for a given target panel’s aspect ratio (e.g., 1:2, 2:1, or 1:1) (Fig. 4.2). Finally, we crop and zoom the input image to fill the target panel (Fig. 4.3). We repeat this process until all the panels are populated.

7. Stylization

Stylization is the process of altering an image’s appearance to make it look like another image (i.e., style transfer) or to make it more interesting or artistic, through filtering. As we show in the results (Sec. 8), we discarded style transfer given its complexity and since it does not allow the required control (efforts such as [GEB*16b] just allows control within the same style). Therefore, we focus on stylization through filtering.

In this section, we describe the interactive tool we created for designing styles by combining filter blocks.

7.1. Interactive Style Design

While most stylization and abstraction works (Sec. 2.3) focus on creating one filter or a fixed set of filters to achieve a stylization, our goal is to create a flexible tool that allows anyone to design stylization filters, regardless of technical ability.

The final result of our framework is as shown in Fig. 5. The main components are:

1. A wide range of filter blocks. (See Sec. 7.2)
2. Parameters for each filter that can be controlled by sliders that appear on the right side of the screen.
3. Two layers: a black and white foreground layer used as alpha channel to display lines and contours (e.g., using XDoG or Sobel filtering); and a background layer for color stylization. This separation provides added design flexibility.
4. A visual flow diagram of filter blocks which allow any block filter to be added, moved, reordered, removed, and tuned at any time. This is the essence of the system’s interaction design.

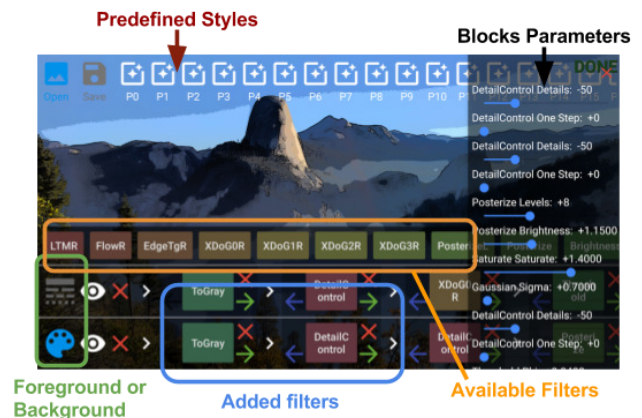


Figure 5: Style Design App.

7.2. Filter blocks

We implemented three kinds of blocks: pixel operations, advanced filters, and histogram modification filters. Fig. 6 shows the effect of each individual block filter.

We included the following list of pixel operations:

- **To Grayscale:** Converts the image into grayscale/luma and saves the chromatic channels (UV). This block is useful for applying any other block filter to just the luma channel (for aesthetics or performance reasons). (Fig. 6.b).
- **To Color:** Uses the current luma and converts the image back to RGB using the UV previously saved by the *To Grayscale* block.
- **Posterization:** Discretizes continuous image colors (e.g., 255 levels) to regions of fewer tones (e.g., $levels = 10$). (Fig. 6.c)
- **Luma Posterization:** Posterizes the image in the luma channel by converting the image to grayscale, applying *Posterization*, and converting it back to color. (Fig. 6.d).
- **Brightness:** Multiplies the luma channel by a user-selected *brightness* constant, clipping the output values. (Fig. 6.e).
- **Soft Threshold [WKO12]:** Performs the following expression for each pixel $out(x,y) = 255 * (1 + \tanh(\min(0, \phi * (in(x,y) - \epsilon))))$ where ϕ determines the slope and ϵ the cut-off. For grayscale images, this block behaves like a binary cut-off that preserves smooth transitions. For color images, it simplifies each RGB channel into two levels. (Fig. 6.f).
- **Saturation:** Makes the colors more vivid (saturates) or more muted (desaturates) by adding or subtracting in RGB the grayscale image tuned by a parameter. (Fig. 6.g).
- **Hue:** Performs a color rotation in UV space and adds a bias in RGB . This block is useful for changing the image's tint. (Fig. 6.h).
- **Colorize:** Convert to monochrome using an HSL palette transformation. (Fig. 6.i).

We also included more advance filters:

- **Gaussian Smoothing:** Blurs and removes details and noise from the image. The standard deviation (σ) is controlled by a parameter. (Fig. 6.j).
- **Edge Tangent Flow (ETF) [KLC07]:** Creates an impressionistic oil painting effect. This method uses a kernel-based nonlinear smoothing of vector field inspired by bilateral filtering. It is implemented as described at cite2017arXiv171110700G. (Fig. 6.k).
- **Total Variation Flow (TVF) [LM11]:** Makes image piecewise constant with an anisotropic diffusion filter. It is implemented as described at [GGI*17]. (Fig. 6.l)
- **Sobel Filter [Sob90]:** Fast edge-detection filter. (Fig. 6.m).
- **XDoG [WKO12]:** Uses *Difference-of-Gaussians* to find the edges of the image. In our implementation, the user can control the variance of the main Gaussian (σ) and the multiplier (p). (Fig. 6.n)
- **Size:** Upscales or downscales the image using a user-selected scale parameter. This block can help to speed up computation (computing other blocks in a lower resolution) and alter the behavior of other scale-dependent filters.
- **Pattern Filling Filter:** Uses *Luma Posterization* to discretize the image into a set of levels, then, each pixel is replaced by a texture depending on its level. This block is useful for creating cross-hatching patterns. (Fig. 6.o).

- **Halftone:** Replaces colors with a set of dots that vary in size and color. This style mimics the behavior of the four-color printing process traditionally used to produce comics. (Fig. 6.p)
- **Image Detail Control:** Inspired by [TM16], controls the details of the image by adding the residual of the image to its bilateral version multiplied by a δ . $\delta < 0$ smooths the image while $\delta > 0$ adds details. (Fig. 6.q).

Finally, we included two blocks to modify the histogram:

- **Linear Histogram Equalization:** Equalizes the luma channel expanding the p_l to zero and the p_h to 255. Normally we choose $l = 5$ and $h = 95$ such that the 5% percentile is moved to zero and the 95% is moved to 255, thereby increasing the image's dynamic range. (Fig. 6.r).
- **Histogram Minimum Dynamic Range:** Computes the percentile 5% and 95% on the luma histogram and expands (if necessary) the dynamic range to match the user parameter range DR .

These two filters are usually placed as the first filter in the pipeline to force the image to have a proper dynamic range to obtain a satisfactory result (e.g., XDoG on a hazy or too bright or dark image would result in an almost entire white output image).

8. Results

In this section, we present the following results: a set of styles produced by designers, a set of styles generated procedurally, an example of video converted to storyboards, and finally a few examples of style transfer.

Designed Styles and Timings:

We conducted several style design sessions with designers who generated dozens of styles that were suitable for storyboard-like stylization. Fig. 7 shows four examples of these styles and Table 1 lists the time for computing the filter blocks for each style. The first filter uses XDoG and Soft Threshold to create the main effect. The second filter, which produces the outlines, uses XDoG with Threshold (strong lines) and, for the background, uses Pattern. The third filter uses a smoothed version of the image to compute the outlines, and the color is achieved by smoothing the image and applying posterization. The final filter uses ETF, Posterize, and Colorize. The average time to compute a Full HD (1920×1080) image on desktop (Xeon E5-1650v3) is 72.45 ms and on a Nexus 6P is 291.625 ms. The code is optimized using the Halide programming language [RKAP*12]. Halide code decouples algorithm implementation from scheduling. This allows us to create optimized code that utilizes native vector instructions while parallelizing over multiple CPU cores.

Video to Storyboard:

Fig. 9 shows one example of our system converting a previously unedited video (Fig. 8) into storyboards. Our system randomly selects the layout, stylization, and creates the variations. Note that the pictured examples were not cherry-picked; they are representative of the system's normal behavior.

Procedural Styles:



Figure 6: Filters: a) Input, b) To GrayScale, c) Posterization, d) Luma Posterization, e) Brightness, f) Soft Threshold, g) Saturation, h) Hue, i) Colorize, j) Gaussian Smoothing, k) ETF, l) TVF, m) Sobel, n) XDoG, o) Pattern Filling, p) Halftone, q) Image Detail Control, and r) Linear Equalization. See text for details.

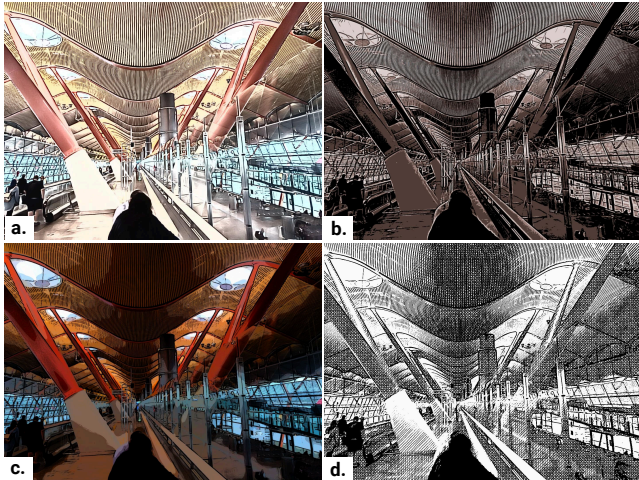


Figure 7: Four styles generated by our design tool.

Procedural modeling is a set of techniques that is able to generate hundreds or thousands of examples using a limited set of rules. In the computer graphics community, the technique has been used to generate buildings, cities, trees, and more complex models (e.g., [MWH*06, NGDA*16]). Based on this idea, we call *procedural styles* those styles created randomly using a set of simple rules. This allow us to create thousands of styles and explore Given that our system can assemble any sequence of filters (just restricted

by the number of input channels for some filters), we decided to randomly generate several thousands of styles and choose the best ones.

The simple set of rules for this experiment was to add a random number of filters between 4 and 9 chosen randomly with random input parameters: XDoG ($\sigma \in [0.5, 8.0]$ and $p \in [1, 40]$), TVF, Soft Threshold ($\phi \in [0.013, 0.059]$ and $\epsilon \in [50, 110]$), Detail Control ($\delta \in [-100, 60]$), Luma Posterization (*level* $\in [5, 12]$), Saturation (*saturation* $\in [1.5, 2.2]$), Size (*size* $\in [100, 300]$), and To GrayScale (20% probability). We also enforce that XDoG and TVF are the only filters that can be added more than once (duplicating the rest of filters has the same effect as selecting different parameters).

We developed a web-based visualization tool to quickly review the randomly generated alternatives (Fig. 10). In order to automate this process, we used a CNN approach [TM17] that evaluates the aesthetic quality of an image. Using this approach, we found the styles in Fig. 11.

Style Transfer:

One interesting alternative to using stylization through filtering, as described before, it is to use style transfer. We conducted a thorough analysis and could not find a template image that would work for all test inputs. Fig. 12 shows several examples we explored. We found the result was strongly dependent on the context, i.e., consistent results would occur only when the template and the input had similar properties (e.g., similar content with different styles). The high complexity of this algorithm is also significant. Gatys et



Figure 8: All the keyframes extracted from a 16-second input video. This is the input for the results pictured in Fig. 9.

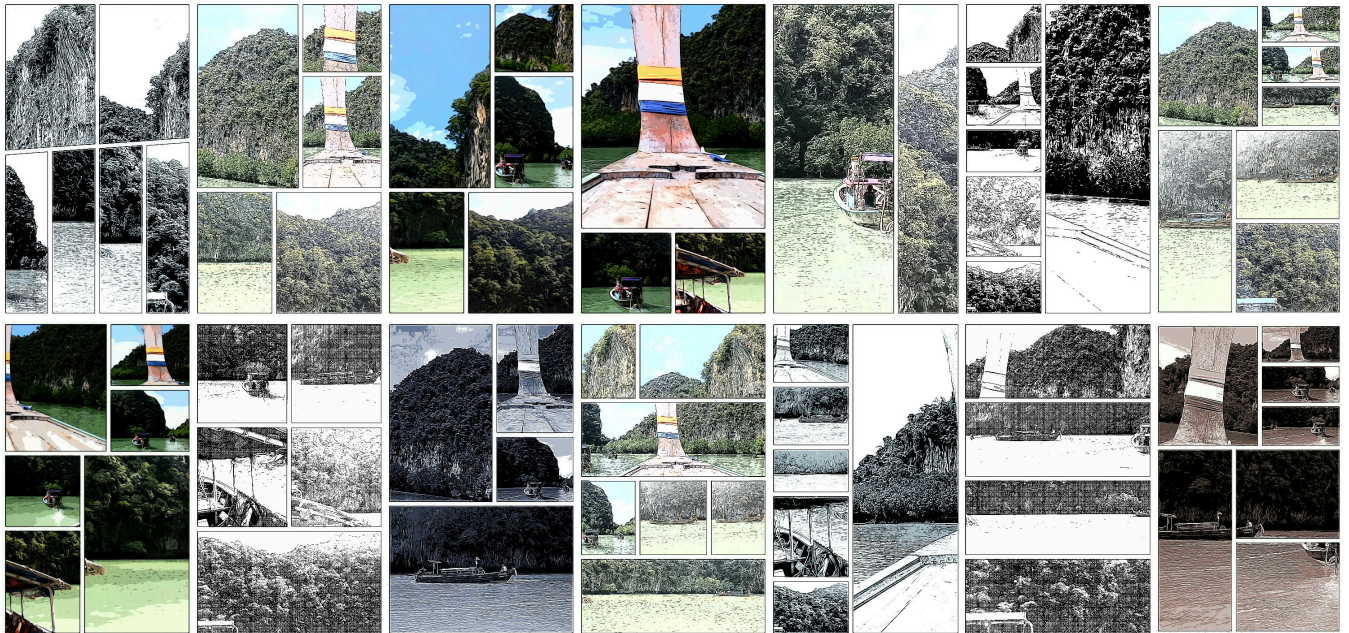


Figure 9: Example output of our system for the input images shown in Fig. 8. The images have not been cherry-picked; they are the first 14 layouts generated by our system for this input.

al. [GEB16a] report it takes up to an hour on a graphics card to compute a 512×512 image. Extensions of this work that focus on speed, e.g., [JAFF16], achieve real-time performance using a high-end desktop GPU, however, it would require dozens of seconds to run on a Full HD image on a mobile device, which makes these approaches impractical for our purposes.

9. Conclusions and Future Work

This paper described a system for automatically converting an album or a video into hundreds of storyboards. We described the process through its main steps: image selection, framing, and stylization. The end result takes a set of images and converts them into an interesting, visually pleasing graphic representation.

In this work, we also presented the first interactive framework

for designing filter-based stylization. This approach boosts the creativity by allowing the designer to tune, modify, and play with the filters directly. In parallel to this manual design, we presented a procedural stylization that follows a list of simple rules to create hundreds of styles automatically. These styles can be selected through manual visualization or using a previously trained aesthetics quality assessment neural network.

As future work, we may pursue several items. We would like to explore using the audio from the input video, possibly by detecting and labelling the audio to add sound effects to the resulting storyboard [GEF⁺17]. As an example, applause in the video might be converted into a stylized sticker with the words "clap clap." Also, we plan to use *scene descriptor* [KFF15] to caption the storyboard. As an alternative to our stylization, we would also like to explore how we could extend CNN-based approaches to fit our context.

Table 1: Style computation time of Fig. 7 (in ms) for ‘Desktop’ (Xeon E5-1650v3) and ‘Device’ (Nexus 6P). For each example, we list the background filters first and then the foreground filters (if any) separated by a line.

Fig. 7.a	Desktop	Device
ToGray	2.0	15.8
XDoG	27.9	107.6
Thresh.	0.4	1.9
ETF	18.3	72.4
ToRgb	4.3	15.1
Total	52.9ms	212.8ms

Fig. 7.c	Desktop	Device
Detail C.	27.0	125.4
ToGray	1.2	4.8
Posteriz.	0.3	2.0
ToRgb	3.8	15.3
Saturate	3.4	11.0
ToGray	1.6	4.7
Detail C.	17.8	92.2
XDoG	22.5	93.9
Threshold	0.5	2.3
TVF	20.6	85.3
Total	98.7ms	436.9ms

Fig. 7.b	Desktop	Device
ToGray	1.3	5.7
ETF	17.6	77.1
XDoG	23.5	92.1
Posteriz.	0.5	1.9
Colorize	7.5	12.0
Detail C.	38.1	132.8
Total	88.5ms	321.6ms

Fig. 7.d	Desktop	Device
ToGray	1.4	4.8
Threh.	0.3	2.0
Pattern	3.7	10.3
TVF	20.7	77.2
XDog	23.1	89.9
Thresh.	0.5	2.0
Total	49.7ms	186.2ms



Figure 10: Visualization tool to explore alternatives. Most alternatives (as seen in the figure) are not very interesting. However, the visualization tool makes it fast to explore and identify promising alternatives.

References

- [ASG15] AYDIN T. O., SMOLIC A., GROSS M.: Automated aesthetic analysis of photographic images. *IEEE transactions on visualization and computer graphics* 21, 1 (2015), 31–42. 4
- [BZL*15] BARNES C., ZHANG F.-L., LOU L., WU X., HU S.-M.: Patchtable: Efficient patch queries for large datasets and applications. *ACM Transactions on Graphics (TOG)* 34, 4 (2015), 97. 3
- [CYW15] CHU W.-T., YU C.-H., WANG H.-H.: Optimized comics-based storytelling for temporal image sequences. *IEEE Transactions on Multimedia* 17, 2 (2015), 201–215. 2
- [EM17] ELAD M., MILANFAR P.: Style transfer via texture synthesis. *IEEE Transactions on Image Processing* 26, 5 (2017), 2338–2351. 3
- [GEB16a] GATYS L. A., ECKER A. S., BETHGE M.: Image style transfer using convolutional neural networks. In *Proceedings of the*



Figure 11: Top-scored randomly selected styles.



Figure 12: Examples of style transfer done using [GEB16a].

IEEE Conference on Computer Vision and Pattern Recognition (2016), pp. 2414–2423. 2, 3, 7, 9

- [GEB*16b] GATYS L. A., ECKER A. S., BETHGE M., HERTZMANN A., SHECHTMAN E.: Controlling perceptual factors in neural style transfer. *arXiv preprint arXiv:1611.07865* (2016). 2, 3, 5
- [GEM*17] GEMMEKE J. F., ELLIS D. P., FREEDMAN D., JANSEN A., LAWRENCE W., MOORE R. C., PLAKAL M., RITTER M.: Audio set: An ontology and human-labeled dataset for audio events. In *IEEE ICASSP* (2017). 8
- [GGI*17] GETREUER P., GARCIA-DORADO I., ISIDORO J., CHOI S., ONG F., MILANFAR P.: BLADE: Filter Learning for General Purpose Image Processing. *ArXiv e-prints* (Nov. 2017). 6
- [GGVG15] GYGLI M., GRABNER H., VAN GOOL L.: Video summarization by learning submodular mixtures of objectives. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2015), pp. 3090–3098. 2
- [JAFF16] JOHNSON J., ALAHI A., FEI-FEI L.: Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision* (2016), Springer, pp. 694–711. 3, 8
- [JHG*15] JING G., HU Y., GUO Y., YU Y., WANG W.: Content-aware video2comics with manga-style layout. *IEEE Transactions on Multimedia* 17, 12 (2015), 2122–2133. 2
- [KCWI13] KYPRIANIDIS J. E., COLLOMOSSE J., WANG T., ISENBERG T.: State of the ‘art’ art: A taxonomy of artistic stylization techniques for images and video. *IEEE transactions on visualization and computer graphics* 19, 5 (2013), 866–885. 3
- [KD08] KYPRIANIDIS J. E., DÖLLNER J.: Image abstraction by structure adaptive filtering. In *TPCG* (2008), pp. 51–58. 3
- [KFF15] KARPATHY A., FEI-FEI L.: Deep visual-semantic alignments

- for generating image descriptions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2015), pp. 3128–3137. 8
- [KLC07] KANG H., LEE S., CHUI C. K.: Coherent line drawing. In *Proceedings of the 5th international symposium on Non-photorealistic animation and rendering* (2007), ACM, pp. 43–50. 2, 6
- [KLC09] KANG H., LEE S., CHUI C. K.: Flow-based image abstraction. *IEEE transactions on visualization and computer graphics* 15, 1 (2009), 62–76. 3
- [LAE*16] LIU W., ANGUELOV D., ERHAN D., SZEGEDY C., REED S., FU C.-Y., BERG A. C.: Ssd: Single shot multibox detector. In *European conference on computer vision* (2016), Springer, pp. 21–37. 5
- [LHC10] LIU D., HUA G., CHEN T.: A hierarchical visual model for video object summarization. *IEEE transactions on pattern analysis and machine intelligence* 32, 12 (2010), 2178–2190. 2
- [LM11] LOUCHET C., MOISAN L.: Total variation as a local filter. *SIAM Journal on Imaging Sciences* 4, 2 (2011), 651–694. 6
- [LYS*11] LIU T., YUAN Z., SUN J., WANG J., ZHENG N., TANG X., SHUM H.-Y.: Learning to detect a salient object. *IEEE Transactions on Pattern analysis and machine intelligence* 33, 2 (2011), 353–367. 2
- [MWH*06] MÜLLER P., WONKA P., HAEGLER S., ULMER A., VAN GOOL L.: Procedural modeling of buildings. In *Acm Transactions On Graphics (Tog)* (2006), vol. 25, ACM, pp. 614–623. 2, 7
- [ND16] NANDISHA M., DESHPANDE S.: Piracy detection app of android applications. *International Journal of Computer Applications* 146, 6 (2016). 4
- [NGDA*16] NISHIDA G., GARCIA-DORADO I., ALIAGA D. G., BENES B., BOUSSEAU A.: Interactive sketching of urban procedural models. *ACM Transactions on Graphics (TOG)* 35, 4 (2016), 130. 7
- [Riv92] RIVEST R.: The md5 message-digest algorithm. 4
- [RKAP*12] RAGAN-KELLEY J., ADAMS A., PARIS S., LEVOY M., AMARASINGHE S., DURAND F.: Decoupling algorithms from schedules for easy optimization of image processing pipelines. *ACM Trans. Graph.* 31, 4 (July 2012), 32:1–32:12. 6
- [RPL*08] RYU D.-S., PARK S.-H., LEE J.-W., LEE D.-H., CHO H.-G.: Cinetoon: A semi-automated system for rendering black/white comic books from video streams. In *Computer and Information Technology Workshops, 2008. CIT Workshops 2008. IEEE 8th International Conference on* (2008), IEEE, pp. 336–341. 2
- [Sob90] SOBEL I.: An isotropic 3×3 image gradient operator. *Machine vision for three-dimensional scenes* (1990), 376–379. 6
- [sto] Storyboard app. <https://play.google.com/store/apps/details?id=com.google.android.apps.photolab.storyboard>. Accessed: 2017-12-12. 2
- [TM16] TALEBI H., MILANFAR P.: Fast multilayer laplacian enhancement. *IEEE Transactions on Computational Imaging* 2, 4 (2016), 496–509. 3, 6
- [TM17] TALEBI H., MILANFAR P.: NIMA: Neural Image Assessment. *ArXiv e-prints* (Sept. 2017). 4, 7
- [VGDA*12] VANEGAS C. A., GARCIA-DORADO I., ALIAGA D. G., BENES B., WADDELL P.: Inverse design of urban procedural models. *ACM Transactions on Graphics (TOG)* 31, 6 (2012), 168. 4
- [Whi04] WHITEHEAD M.: *Animation*. Old Castle, 2004. 1
- [WHY*12] WANG M., HONG R., YUAN X.-T., YAN S., CHUA T.-S.: Movie2comics: Towards a lively video content presentation. *IEEE Transactions on Multimedia* 14, 3 (2012), 858–870. 2
- [WKO12] WINNEMÖLLER H., KYPRIANIDIS J. E., OLSEN S. C.: Xdog: an extended difference-of-gaussians compendium including advanced image stylization. *Computers & Graphics* 36, 6 (2012), 740–753. 2, 6
- [WOG06] WINNEMÖLLER H., OLSEN S. C., GOOCH B.: Real-time video abstraction. In *ACM Transactions On Graphics (TOG)* (2006), vol. 25, ACM, pp. 1221–1226. 3
- [ZCSG16] ZHANG K., CHAO W.-L., SHA F., GRAUMAN K.: Summary transfer: Exemplar-based subset selection for video summarization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2016), pp. 1059–1067. 2
- [ZWZS97] ZHANG H. J., WU J., ZHONG D., SMOLIAR S. W.: An integrated system for content-based video retrieval and browsing. *Pattern recognition* 30, 4 (1997), 643–658. 2