

# A review of image and video colorization: From analogies to deep learning

Shu-Yu Chen<sup>a</sup>, Jia-Qi Zhang<sup>b</sup>, You-You Zhao<sup>c</sup>, Paul L. Rosin<sup>d</sup>, Yu-Kun Lai<sup>d</sup>, Lin Gao<sup>a,e,\*</sup>

<sup>a</sup> Beijing Key Laboratory of Mobile Computing and Pervasive Device, Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China

<sup>b</sup> Beihang University, Beijing, China

<sup>c</sup> University of California, Santa Cruz, CA, US

<sup>d</sup> Cardiff University, Wales, UK

<sup>e</sup> University of Chinese Academy of Sciences, Beijing, China

## ARTICLE INFO

### Article history:

Received 6 December 2021

Received in revised form 11 April 2022

Accepted 13 May 2022

Available online 8 June 2022

### Keywords:

Image colorization

Sketch colorization

Manga colorization

## ABSTRACT

Image colorization is a classic and important topic in computer graphics, where the aim is to add color to a monochromatic input image to produce a colorful result. In this survey, we present the history of colorization research in chronological order and summarize popular algorithms in this field. Early work on colorization mostly focused on developing techniques to improve the colorization quality. In the last few years, researchers have considered more possibilities such as combining colorization with NLP (natural language processing) and focused more on industrial applications. To better control the color, various types of color control are designed, such as providing reference images or color-scribbles. We have created a taxonomy of the colorization methods according to the input type, divided into grayscale, sketch-based and hybrid. The pros and cons are discussed for each algorithm, and they are compared according to their main characteristics. Finally, we discuss how deep learning, and in particular Generative Adversarial Networks (GANs), has changed this field.

© 2022 The Authors. Published by Elsevier B.V. on behalf of Zhejiang University and Zhejiang University Press Co. Ltd. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

## 1. Introduction

Color plays a very important role in the process of human cognition of the world, and rich colors can not only express more information, but also enhance the human visual experience. Image colorization has been a very active research topic in the field of digital image processing, and is an inter-disciplinary area involving disciplines such as Computer Vision, Computer Graphics, Pattern Recognition and Human Computer Interaction. Colorization has been widely used in many fields, such as grayscale photo colorization, old film color restoration, cartoon automatic colorization, etc.

Based on different types of input to be colorized, colorization methods can be divided into two broad categories, one is colorization of grayscale images and black-and-white videos which are ordinary photos without color, and the other is colorization of monochrome art forms, including sketch images (or line art images), manga (or comics) and black-and-white cartoons (or line-art video). See Fig. 1 for some examples. We can see

that grayscale images contain rich intensity details, while sketch images (or other art forms) only contain relatively sparse details.

Therefore, when processing input images of different categories, researchers usually use different processing methods. For colorization of grayscale images, most methods convert the image in YUV or Lab color space (Cheng et al., 2015; Zhang et al., 2016), and restore the value of the chrominance channels of the image to be colored based on the similarity of the luminance channel (Levin et al., 2004). For black-and-white videos, most models use unsupervised or self-supervised learning from the visual tracking process to track the location of an object in different frames, and link corresponding pixels together, to colorize them based on a user-provided reference photo or based on data-driven deep learning technologies. The colorization of sketch images often involves segmenting it into different regions (Sato et al., 2014), and based on a learning model a color is assigned to each segment where the color information can come from reference images, users' color scribbles or input text hints. Although grayscale-based colorization methods can be directly used to predict the color value of each pixel in sketch images, they usually do not have good performance due to the lack of texture information. Therefore sketch-based colorization methods are required to propose new solutions for line feature extraction and region boundary determination, such as studying the temporal

\* Corresponding author at: Beijing Key Laboratory of Mobile Computing and Pervasive Device, Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China.

E-mail address: [gaolin@ict.ac.cn](mailto:gaolin@ict.ac.cn) (L. Gao).

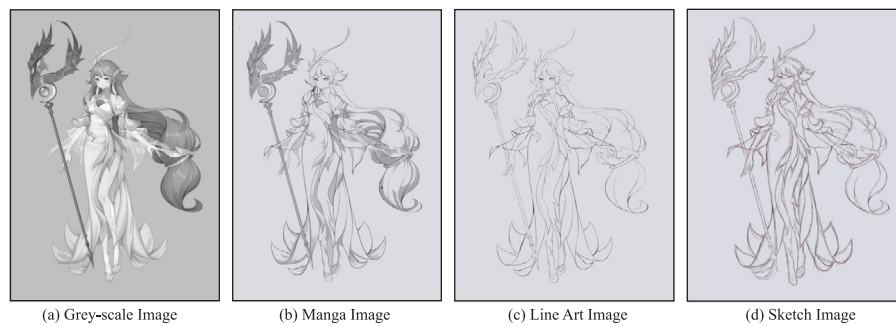


Fig. 1. Typical categories of images suitable for colorization.

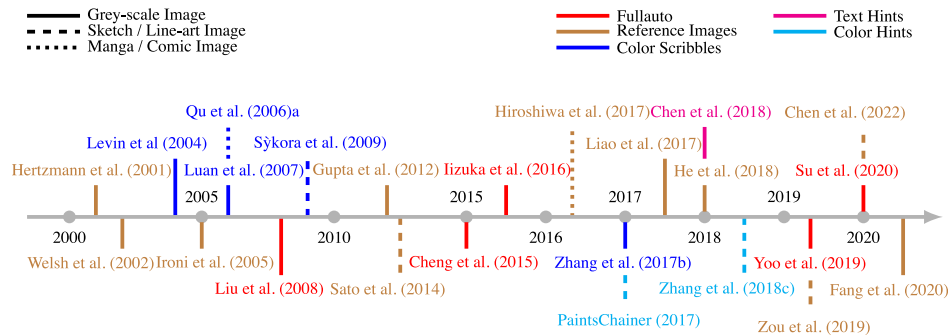


Fig. 2. The timeline of image colorization methods. Different line colors represent methods with different types of control, as shown in the upper-right. Different line types represent the different types of input, as indicated in the upper-left. From this timeline, it can be seen that there can be a wide variety in the type of input data, from gray-scale image to manga to sketch, which shows the difficulty of data processing. Compared with gray images, sketch images are sparse and the information available from just lines can be ambiguous. Early approaches were based on reference images, then user interaction was introduced, and finally fully automatic colorization.

and semantic relationships between lines (Zhang et al., 2018a; Yang et al., 2021; Ge et al., 2022).

Generally speaking, traditional colorization methods require a lot of manual interaction and are often sensitive to the parameter settings of the methods. As a result, adding manual interaction and parameter optimization will take a lot of time and effort. Especially in grayscale videos or cartoon film colorization, even a short film usually has thousands of images to process. To improve efficiency, using a DCNN (Deep Convolutional Neural Network) (Krizhevsky et al., 2012) to build up the model or a GAN (Generative Adversarial Network) (Goodfellow et al., 2014) for training is the most common approach used in recent methods. Both the image colorization effect and efficiency have been greatly improved. There exist both opportunities and challenges, and the development of deep learning technology has brought new directions to the work of image colorization.

Several major challenges remain. Some methods can only be used under certain restrictions, and moreover have some defects. For example, the colorization method can only handle gray-scale images, or the model needs to provide suitable reference color images. Some models need to identify different objects in the image, and then work out appropriate colors, but in particular for sketch image colorization, it is very difficult for the model to understand the sketch image and learn different artistic styles. The existing survey (Anwar et al., 2020) mainly summarizes works performing colorization of grayscale images and the datasets for colorization. However, the task of colorization is not restricted to grayscale images, but also includes manga and sketches. In this paper, we will summarize and discuss different colorization methods from three categories, including their advantages and drawbacks, to give an overview that should be useful for researchers and practitioners.

## 2. Overview

This survey paper divides existing colorization research work into the following three sections based on the different types of input images to be colored. In Section 3, we mainly introduce the colorization methods for grayscale images, which are further grouped into three subcategories: fully automatic colorization methods, semi-automatic colorization methods based on color strokes or reference images, and text-driven image colorization methods. In Section 4, we focus on methods related to colorization of line-art or sketch images, which are further classified into four subcategories: colorization methods based on color strokes, colorization methods based on reference images, text-driven colorization methods, and synthesis methods from line-art images to real images. In Section 5, we discuss the colorization work of comic or manga images. Finally, we summarize the colorization methods and discuss a possible area for future colorization work in Section 7. Fig. 2 shows a timeline of representative methods for image colorization.

## 3. Grayscale image colorization

The color value of each pixel on a grayscale image is between black and white. The grayscale channel can be extracted from color images, but images such as photos taken in the past and much comic art only have grayscale information, and can benefit from colorization. Colorization methods for grayscale images can be divided into two groups according to whether interaction is used, namely automatic colorization methods and semi-automatic colorization methods. In the former group, researchers have used data-driven deep learning technology to automatically colorize grayscale images based on training data (Cheng et al., 2015; Iizuka et al., 2016; Messaoud et al., 2018). For instance,

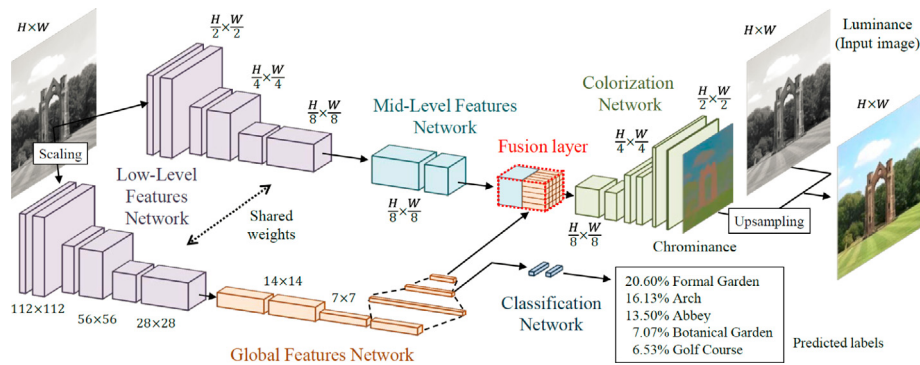


Fig. 3. An overview of an end-to-end network for grayscale images colorization (Iizuka et al., 2016).

there is an open source automatic model DeOldify (Antic, 0000) which is free to use. For the second group, colorization methods often take some guidance information from users, e.g., by drawing color strokes (Levin et al., 2004), providing reference color images (Liu et al., 2008) or giving an specific color theme (Wang et al., 2012). Incorporating user guidance usually increases the efficiency and correctness of the colorization model. It also helps resolve inherent ambiguities for the ill-posed colorization problem (such as tree leaves could be green in spring and yellow in autumn). In addition, in some recent studies, researchers have also studied the use of semantic information to guide image colorization, such as grayscale colorization based on text scripts (Bahng et al., 2018), which is introduced in the last subsection.

### 3.1. Automatic grayscale image colorization

Different from colorization methods based on guidance, in automatic image colorization methods, researchers can design the model to provide multiple colors for the same pixel to solve the problem of multi-modal colorization of monochrome images. For example, colorization models will generate green or yellow leaf images. In this subsection, we divide the automatic coloring method into two categories based on the diversity of the generated results; one is unimodal colorization in which methods can only generate one result, and the other is multi-modal colorization where methods can generate multiple diverse results, and introduce them respectively.

#### 3.1.1. Unimodal colorization

To reduce user interaction, Cheng et al. (2015) first propose a fully-automatic colorization method using deep learning with the SUN dataset (Patterson and Hays, 2012). Instead of direct taking the grayscale image as input, they take a combination of multi-level features to predicts the U and V channels. However the performance drops when similar reference images are not included from the training data-set. Concurrently, Deshpande et al. (2015) improved the learning model for image colorization and learned from examples. This learning model is built upon the LEARCH (Learning to Search) framework (Ratliff et al., 2009), and is able to minimize the quadratic objective function defined on the chromaticity maps, comparable to a Gaussian random field.

Larsson et al. (2016) proposed an automatic colorization method based on a self-supervised visual representation learning process. The network is built upon the fully convolutional network of VGG-16 (Simonyan and Zisserman, 2014) with the classification layer removed and a filter layer added. In addition, the model uses skip-layer connections to concatenate the features of different convolutional layers, to provide input to the classification layer which predicts the color histogram of each pixel. Iizuka et al. (2016) proposed to learn the global and local features separately

from an image and then combine them together for the final colorization process (see Fig. 3). However, for objects with multiple different colors, the result will most likely produce dominant colors which is learnt in training, like the leaves in green.

Previous colorization methods usually learn to colorize an entire image, and thus they are often unsuitable for colorizing instances in the image. Su et al. (2020) proposed a method to realize instance-aware colorization. They first detect the location of the target objects, and then colorize the object and the overall image respectively. When there is miss detection or overlap detection in the same class instance, the result will be affected, as shown in Fig. 4. Considering the constraints of training data, Yoo et al. (2019) attempt to colorize with little data. Their whole network consist of memory network and colorization networks. The memory network is trained in unsupervised way to help get the most similar color features that match the input image, and then give it as the condition to colorization networks.

Unimodal automatic colorization methods can generate reasonably colored images when given arbitrary grayscale images, but only one corresponding color image can be generated for a grayscale input. A single generated result may not meet user expectations and the user cannot specify local or object-specific colors. Various results could be generated by introducing random noise features or random color condition vectors into the network module.

#### 3.1.2. Multi-modal colorization

Aiming at solving the problem that colorization requires a lot of user interaction and that the color saturation of colorized images tends to be low, Zhang et al. (2016) proposed a fully automatic colorization method that can generate rich and realistic colorization images. This method transforms the colorization task into a self-supervised expression learning task by learning the semantic and texture mapping between the grayscale image and the color image. At the same time, the colorization problem is transformed in a novel manner into a classification task, and a color distribution is predicted for each pixel to solve the multi-modal colorization problem of the image, which maintains the diversity of the colorization results. Zhang et al. (2016) were inspired by the simulated annealing method (Kirkpatrick et al., 1983) and proposed the operation that works out the annealed mean of a distribution, to estimate the color value of the  $ab$  space from the color distribution of each pixel. The value of each pixel in the grayscale image colorization task is not fixed, and the same object in the real world can be colored in different ways.

Unlike (Zhang et al., 2016), Deshpande et al. (2017) not only considers the estimation of the color value of each pixel, but also considers the overall spatial continuity of the colorization results. This method uses a variational autoencoder (VAE) to learn the low-dimensional latent variable embedding of the color field, and

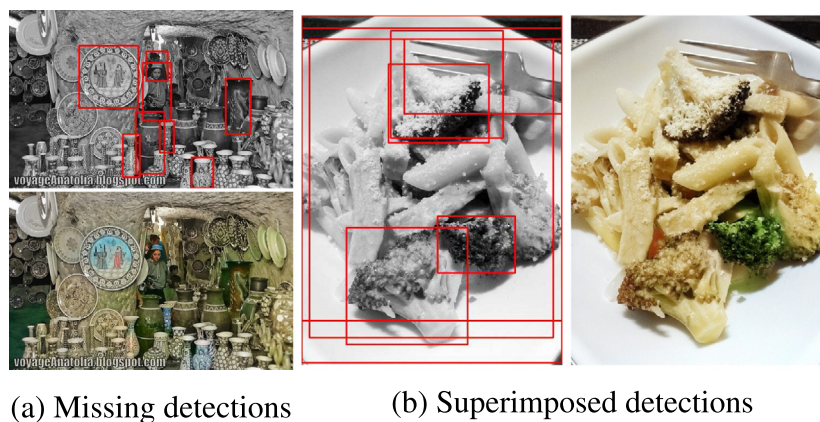


Fig. 4. Failure cases of Instance-Aware Image Colorization (Su et al., 2020).

uses a Mixed Density Network (MDN) to learn a multi-modal model conditioned on the grayscale image. Finally, multiple samples are taken from MDN, and combined with the VAE decoder to obtain multiple colorization results for each sample, so as to provide a rich set of colorization results.

Although the classification model based on color distribution and the generative model based on variational autoencoders can obtain a variety of colorization schemes, the colorization results lack the consistency of the spatial structure and the user controllability of color. Sometimes in the same semantic area, spots of different colors appear in the colorization result. In order to ensure global colorization consistency and user controllability, Messaoud et al. (2018) proposed a conditional random field based on VAE and use a Gaussian Conditional Markov Random Field (G-CRF) to capture global image statistics, modeling the output space of the VAE decoder and the encoding of user editing information.

When an image colorization method is directly applied to video colorization, discontinuity will appear. Lei and Chen (2019) proposed an automatic colorization model for black-and-white video without any user interaction or reference image. This method designs a self-regularization and diversity loss function in order to achieve the consistency and diversity of the grayscale video colorization. The self-regularization loss is mainly composed of a bilateral regularization term and a temporal regularization term, which adds color consistency constraints in the bilateral space of adjacent pixels and corresponding pixels of adjacent frames. Diversity Loss to constrain the multiple generated results to be consistent with real color images. Although the method achieves the generation of multiple colorization results, there are no rich colorization results between different results.

With the rapid development of Transformer (Vaswani et al., 2017) in the field of computer vision, Kumar et al. (2020) proposed a grayscale colorization network architecture (Colorization Transformer, ColTran) based on Transformer blocks. ColTran is mainly composed of an autoregressive Colorizer, a color upsampler and spatial upsampler. Autoregressive Colorizer enabled color information to be matched to input grayscale images at low resolution, and then the color upsampler and spatial upsampler sampled low resolution color images into high resolution images in a completely parallel way. Based on transformer's better matching ability, this method can provide a variety of colored gray images according to different reference color images.

Compared with the unimodal colorization, multi-modal colorization methods can generate multiple color results for a given grayscale input. Although those automatic methods do not require user interaction, the generated results rely on pretrained network models. The user cannot adjust the generated results, such as the overall colorization style or detail colors, making it difficult to generate the results the user expects.

### 3.2. Color strokes based colorization

In order to solve the problem that automatic methods cannot control the color of the details, some work attempt to take user color strokes and provide an intuitive approach for user control.

#### 3.2.1. Optimization colorization

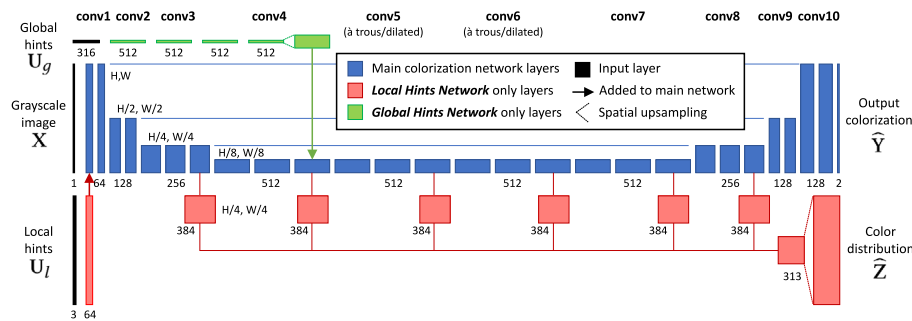
Levin et al. (2004) were one of the most important pioneers in the colorization area. In this method, the user needs to mark a grayscale image with color strokes to colorize the image in YUV color space. Then, based on the rule that adjacent pixels have similar intensities and their colors are similar, the method spreads the color of the strokes to the entire image. But when different object colors are diffused and mixed together, there are color bleeding problems in Levin et al. (2004). To solve it, Huang et al. (2005) modified the weighting function and proposed an adaptive edge detection algorithm to improve the accuracy of the edges. They use Sobel filters and iterative optimization to improve the edge detection. Further, the colorization method will be more accurate, while at the same time reducing color bleeding issues, and making the image color effect more realistic.

Previous colorization methods based on color strokes, such as Levin's method (Levin et al., 2004), usually require a lot of manual interaction for a complex scene. To reduce it, Luan et al. (2007) proposed a new interactive system that can quickly and easily color grayscale images. This method consists of two stages, the color labeling stage and the color mapping stage. The color labeling stage spreads the marked colors to similar areas by constraining the intensity smoothness and texture similarity of all pixels. The color mapping stage establishes a piece-wise linear mapping in luminance ( $Y$ ) space according to the scribble's luminance ( $Y$ ) and chroma ( $UV$ ) values, and finally the chromaticity values of other pixels are obtained by interpolation. The comparison between Levin et al. (2004) and Luan et al. (2007) is shown in Fig. 5. It can be seen that Luan et al. (2007) is faster and more effective.

For monochrome movies, Yatziv and Sapiro (2006) uses intrinsic, gradient weight, and the relationship between colorized point and the nearby spot to colorize in the YCbCr color space. This method first defines the intrinsic (geodesic) distance of any two points to calculate the smoothness between the luminance channel of the two points. Then for any point on the image, the method calculates the shortest intrinsic distance from the point to the known chrominance point, finds multiple chrominance values corresponding to it, and obtains the final chrominance value by blending different chrominance values. But this process involves a large amount of calculation and a complicated process to ensure color quality.



**Fig. 5.** Comparison of stroke based colorization. (a–c) Color strokes, strokes alone, colorization result by Levin et al. (2004), (d–f) Color strokes, two pixels in each region are labeled, colorization result by Luan et al. (2007).



**Fig. 6.** An overall user-guided image colorization network architecture (Zhang et al., 2017b).

In view of the value of each pixel in the grayscale image, the information contained in the grayscale image is rich, and the optimization method can determine the propagation range of the color according to the relationship between grayscale information or adjacent pixels. However, the computation cost of those methods is high. Therefore, most of them struggle to perform in real time. The emergence of deep neural networks has effectively reduced the time required for the process, and has subsequently inspired more colorization methods based on deep learning.

### 3.2.2. Deep neural network based colorization

In the recent work on colorization of grayscale images based on user guidance, researchers began to use Convolutional Neural Networks (CNNs) to learn the mapping relationship between grayscale images and color images, with user constraints. Based on deep learning, Zhang et al. (2017b) proposed a network model that takes user input as guidance for real-time image colorization. The model is mainly divided into three parts, the main colorization network learns to colorize, the local hints network predicts a color probability distribution and the global hints network learn to encode the input global histogram statistics and average image saturation into the middle part of the main colorization network. By combining the colorization network structure of global and local information, Zhang et al. (2017b) mapped the input image and user input to the output color image, and finally realized real-time user-guided image colorization. The specific network structure is shown in Fig. 6. Although the network can complete the colorization in real time with only a single forward propagation, it still requires the user to specify a large amount of color hints.

Neural network-based methods effectively improve the speed and quality of image colorization. However, the network training adopts an end-to-end training strategy. As a result, when the input and output are given, the output result cannot be controlled, and the result can only be optimized by editing the input. At the same time, this approach needs to specify the color of the area for each image, which cannot achieve batch colorization work. Therefore, some researchers study the colorization based on reference images.

### 3.3. Reference color image based colorization

Another approach that balances controllability with user effort is reference based colorization where the user provides reference images with desired color distribution to guide the colorization process. The reference images may be specified directly by the user, retrieved from the Internet or obtained from a large dataset (He et al., 2018). By referring to the reference image, the colorization results can better satisfy the user's expectations. Although there are substantial overall differences between the images, similarities between the images can still be found in local areas. For example, areas with similar color or texture often also have similarities in structure or lines. Therefore, we can guide the generation of images by finding similarities between the grayscale image and the reference image.

#### 3.3.1. Similarity with luminance features

To colorize a grayscale image, these methods need to have one or more reference images, and then use luminance channel mapping with the input image. Hertzmann et al. (2001) transferred color information into the input image from analogous regions of the reference. In the work by Welsh et al. (2002), the grayscale image only contains one dimensional information, and for a color reference image, its luminance channel can be used to match the grayscale input. So the algorithm converts the reference image into  $\alpha\beta$  color space, and selects a small subset of pixels as a sample. Then the pixels in the grayscale image are scanned in raster order and the best matching part is selected using neighborhood statistics. Welsh et al. (2002) described how their model could be applied to a single frame in a video sequence. They used the same colorized target swatch that was used in the first frame to colorize the remainder of the video. This process can effectively solve the problem of color inconsistency. After finding the corresponding pixel, they use the swatch model to produce a vivid colorization effect. In the equation, the error distance  $E(N_g, N_s)$  uses the  $L_2$  distance metric between neighborhood  $N_g$  in the grayscale image and neighborhood  $N_s$  in the colorized image (see Fig. 8).

Gupta et al. (2012) incorporated SIFT features (Liu et al., 2011) into Welsh's method (Welsh et al., 2002) and created a new

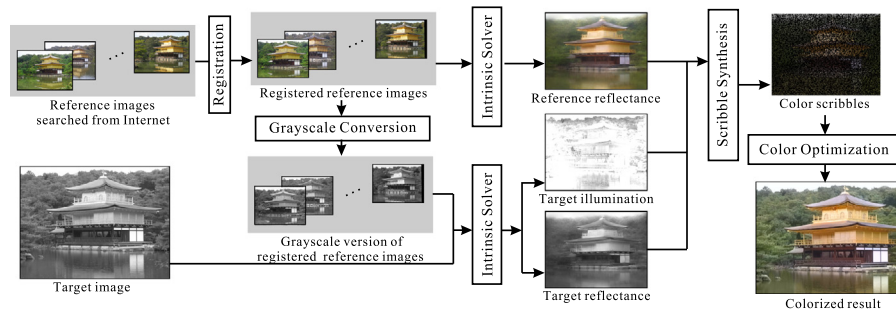


Fig. 7. An overview of intrinsic grayscale image colorization with reference images retrieved from the Internet (Liu et al., 2008).

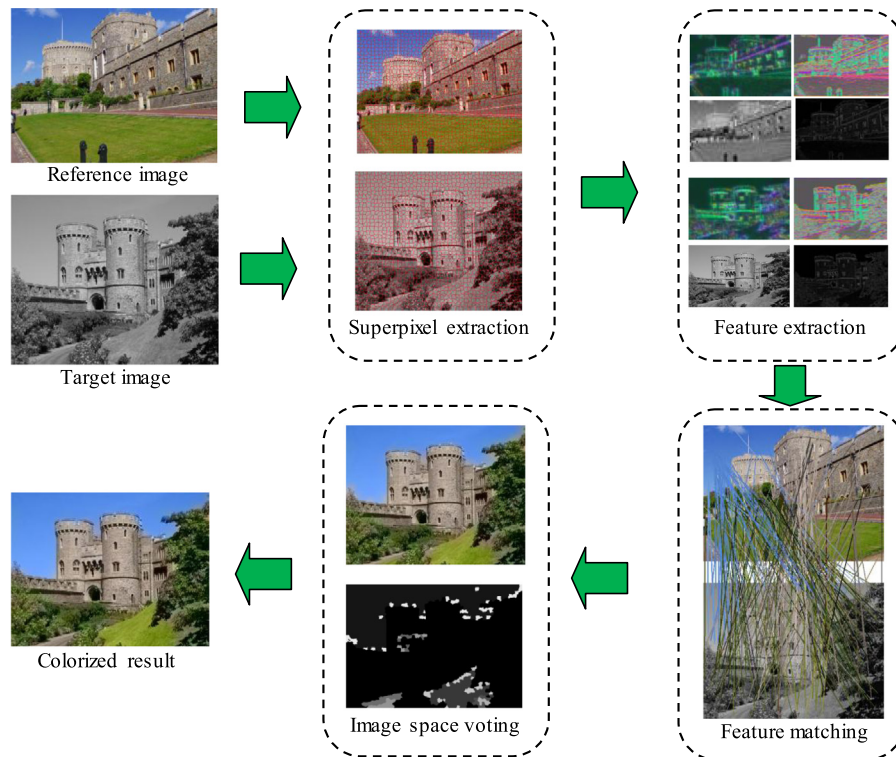


Fig. 8. Overview of the colorization method based on matching of similar images (Gupta et al., 2012).

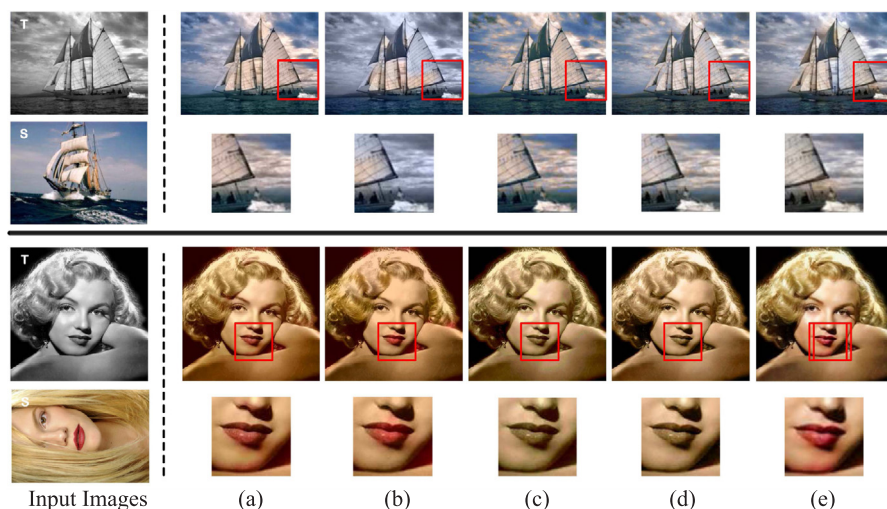
method that allows the transfer of color information from a reference image using multiple image features. Gupta et al. (2012) use SIFT features to obtain the correspondence between the reference image and the target image for color transfer. So the reference image needs to be a related image in order to obtain the best results. They also proposed image voting for color correction, which checks neighboring superpixels to identify and then correct invalid color assignment, in order to keep the color consistency.

Unlike other methods based on matching between pixels (Welsh et al., 2002; Gupta et al., 2012), the colorization method proposed by Ironi et al. (2005) tries to color images at a higher semantic level. Based on the method of Levin et al. (2004), Ironi et al. (2005) expect that the algorithm can automatically place color scribbles, and then use color optimization in Levin et al. (2004). Specifically, the method is mainly divided into four stages. They first train a supervised learning algorithm to build a low-dimensional feature space to discriminate which label the pixel belongs to. Then, they reliably determine the reference color value of each pixel by voting for the nearest neighbors in the feature space. Finally, the color is transferred to neighboring pixels in other spaces and the method of Levin et al. (2004) is used

for global optimization. Compared with scribbles, this method saves time, and adopts the spatial voting scheme to strengthen the spatial consistency, and has more robust color results than Welsh's method (Welsh et al., 2002).

Li et al. (2019) proposed a new location-aware cross-scale texture matching method to achieve grayscale colorization based on reference images. This method first uses the multi-label graph-cut algorithm to minimize global matching errors and spatial scale variations, and then uses the statistics of up-down relationships in the reference image to correct unreasonable color matches, and finally applies an optimization framework to propagate the high-confidence micro-scribbles to entire image. In the grayscale colorization method based on the reference image, it is very common that the provided reference image and the target image scale are inconsistent, and this method can handle this situation well, and performs well among methods based on texture matching.

The total variation (TV) minimizing denoising model proposed by Rudin et al. (1992) is used for image colorization. Kang and March (2007) proposed to use the total variation minimizing colorization model to deal with the problem of image color restoration. This method first minimizes the total variation, and then implements image colorization through weighted harmonic maps.



**Fig. 9.** Comparison of the results of colorizing grayscale images. The “T” and “S” of the first column input images are the target and source images. (a) is the result of Fang et al. (2020), (b) is the result of Gupta et al. (2012), (c) is the result of Welsh et al. (2002), (d) is the result of Pierre et al. (2015), and (e) is the result of He et al. (2018).

However, this method requires a large number of color scribbles to process complex images. Further, Bugeau et al. (2014) proposed a minimization variational formulation modeling which could colorize using a reference image. At the same time, a specific energy function is designed for modeling color selection and spatial consistency constraints. However, this method will produce a halo effect on edges with obvious texture contrast.

Fang et al. (2020) proposed a grayscale colorization method based on a reference image, which novelly takes the result of image superpixel segmentation as the target to be processed. The method first uses the Vcells (Wang and Wang, 2012) algorithm to segment an image, extracts the features of the segmented blocks, and then uses the method proposed by Gupta et al. (2012) to match the reference segmented feature and the target segmented feature. Different from Gupta’s method (Gupta et al., 2012), Fang et al. (2020) do not use the matched colors as micro-scribbles for color propagation, but instead select a set of candidate colors for each target superpixel. Finally, they used the TV based spatial consistency regularization and non-local self-similarity regularization to determine the most suitable color for each target superpixel from the color candidates. As the comparison shows in Fig. 9, with the same reference image, Welsh et al. (2002)(c) and Pierre et al. (2015)(d) are limited by the set of color candidates and cannot match enough correct colors. He et al. (2018)(e) and Gupta et al. (2012)(b) obtained more reliable color assignment results, but the results of (b) contain color inconsistency artifacts and those of (e) contain color blurring and color bleeding that appear in tiny objects. Although (a) achieves better colorization results than other methods, there are still incorrect color matching results, such as the hair edges of the characters in the second row.

Instead of specifying reference images by users, these methods achieve automatic grayscale colorization where the colorization algorithm automatically searched for similar color images on the Internet. Those images are often captured with different poses and illuminations. Given a specified reference image, Ironi et al. (2005) colorized images by a robust supervised classification scheme, but they cannot handle multiple reference images. Based on Ironi et al. (2005), Liu et al. (2008) proposed an automatic colorization method which could directly search the Internet for multiple color images similar to the target image. It ignored the illumination difference between the grayscale target image and the color reference image. As shown in Fig. 7, after getting the

images from Internet, the color reference images are registered by using the SIFT (Lowe, 2004) matching algorithm. And then the target illumination, target reflectance, and intrinsic reflectance images of the target scene are extracted from these reference images. Finally, the color is transferred from the above images, and the final result is obtained by combining the illumination component of the target image. Morimoto et al. (2009) also used web search for images with similar scene structure, and used the 20 most similar images for colorization. They then used color transfer based on the luminance value to colorize the input image.

The method of using similarity on luminance features has great dependency on the value of each pixel, and can solve the situation that the value of objects between different images is similar. However, when the lighting of the object or the structure changes, it is easy to get the wrong color. However, there is a correlation between the same objects or similar details. By introducing object category analysis or image feature analysis, etc., the color matching between the grayscale image and the reference image can be effectively improved.

### 3.3.2. Similarity with CNN features

He et al. (2018) proposed for the first time a fully automatic colorization method based on reference images, allowing users to use different reference images to achieve different colorization styles. Its network structure is mainly divided into a similarity sub-network and colorization sub-network, which is shown in Fig. 10. The similarity sub-network helps find the similarity maps between the reference image and the target image. Then the colorization sub-network mainly aligns pixels in the luminance channel based on the similarity sub-network, and then uses the big data learning ability of the colorization sub-network to refine misaligned pixel colors.

Inspired by He et al. (2018), Zhang et al. (2019) applied the method of deep exemplar-based colorization to grayscale video colorization. Similar to He et al. (2018), Zhang et al. (2019) obtained the dense correspondence between the target image feature and the reference image feature through computing a correlation matrix, and then feed it to a colorization sub-network. In this sub-network, the colorization result of the previous frame will be used as the condition for the current frame colorization. To reduce accumulated propagation errors, reference images are added. Through this recurrent framework, they achieved temporal consistency of video colorization. They further introduced a

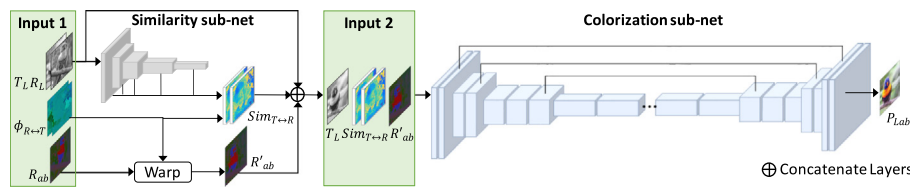


Fig. 10. An overview of the architecture of colorization by Deep Exemplar (He et al., 2018).



Fig. 11. An example of image analogy (Hertzmann et al., 2001).

temporal consistency loss (Chen et al., 2017) to reduce the color change along the flow trajectory during the video colorization process. Further, Wu et al. (2021) propose a method that can generate results with vivid colors by retrieving the matched features. Different from colorization based on reference images, they designed a GAN encoder to generate the color prior in colorization, which allows smooth interpolation between different colors and generates more diverse results.

Vondrick et al. (2018) proposed video colorization with self-supervision learning of visual tracking. This method copies colors from a reference frame, and the model needs to use the appropriate region in order to obtain the correct color. And it can also be applied to track people's movement through the video. This model uses a pointing mechanism to solve the problem of inconsistencies in video colorization to a large extent, and maintain the color stability of the frame. However, the pointing mechanism is still not precise enough, and the color edges of the colorization result are sometimes unclear.

Iizuka and Simo-Serra (2019) proposed a single end-to-end framework to tackle black-and-white vintage film remastering. They introduce source-reference attention to guide the color of the reference image into the target image. The source-reference attention mainly uses the extracted reference image and target image features to perform matrix operations to obtain non-local similarities. The method of calculating non-local similarities using the extracted features has been frequently used in recent years (He et al., 2018; Zhang et al., 2019; Shi et al., 2022; Lee et al., 2020; Siyao et al., 2021).

### 3.3.3. Colorization using image analogies

In addition, researchers attempt to use image analogy methods to achieve image style transformation and grayscale image colorization. The image analogy problem is illustrated in Fig. 11, in which we have image pair  $A$  and  $A'$ , and given a new image  $B$ , the image analogy is to find the image  $B'$ , so that the image pair  $B$  and  $B'$  has the same mapping relationship as the image pair  $A$  and  $A'$ .

Hertzmann et al. (2001) earlier proposed an image analogy method. The method mainly uses two matching processes, Best Approximate Match and Best Coherence Match. The Best Approximate Match process first uses a Gaussian pyramid to extract the feature information of pixels at different scales, and then uses an approximate-nearest-neighbor (ANN) search to search for the pixel  $p$  in the original image  $A$  that best matches each pixel  $q$  in the target image  $B$ . The Best Coherence Match guarantees

the spatial consistency of the matching results. This calculation method is mainly derived from the method proposed by Ashikhmin (2001). The specific calculation is as follows:

$$r^* = \arg \min_{r \in N(q)} \|F_\ell(s(r) + (q - r)) - F_\ell(q)\|^2 \quad (1)$$

where  $r$  represents a pixel that has been synthesized in the neighborhood of pixel  $q$  in  $B'$ ,  $s(r)$  represents the pixel corresponding to  $r$  in  $A'$ ,  $N(q)$  represents a pixel synthesized in the neighborhood of pixel  $q$ , and  $F_\ell(\cdot)$  represents a neighborhood feature vector of pixel in layer  $l$ .

Although the algorithm proposed by Hertzmann et al. (2001) can get good results, the algorithm requires pixel-by-pixel matching, which is particularly slow. Later, Liao et al. (2017) proposed a new image analogy method using deep learning technology, which has greatly improved the matching speed and effect. The method uses the pre-trained image feature extraction network VGG-19 (Simonyan and Zisserman, 2014) to extract the 5-layer high-dimensional features of images  $A$  and  $B'$ . Then the method uses Nearest-Neighbor Field Search (NNFs) to find dense correspondences with bidirectional constraints in each feature layer. Finally, the image features are gradually reconstructed from the roughly corresponding fifth layer to the finely corresponding first layer, and the final generated images  $A'$  and  $B$  are obtained. The specific pipeline of the method system is shown in Fig. 12.

There are other researchers who use the idea of image analogy of animation by extending the image analogy method to create a time-continuous animation sequence. Using the image analogy method, Jamriška et al. (2019) used image color, foreground object binary mask, position of SIFT Flow (Scale-invariant feature transform Flow) (Liu et al., 2011), and foreground object edge information as guidance information, and achieved video stylization.

Affected by the structure of the image and the color range of the image, when the pose and appearance of objects in the image change greatly, it is impossible to obtain good results by using the image analogy method, even if the two images to be compared are images from the same video at different times. In addition, if the color of a part of the new image does not appear in the original image, the algorithm cannot automatically fill this part of the color. At the same time, when the method is applied to sketches and color images, it is difficult to generate reasonable color results corresponding to sketches.



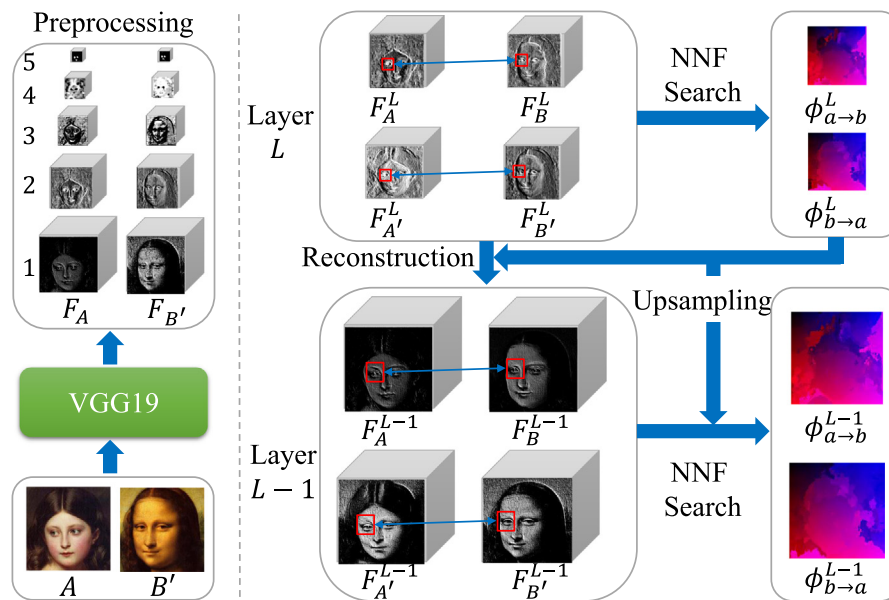


Fig. 12. System pipeline of deep image analogy (Liao et al., 2017).

### 3.4. Text hints based colorization

Text can be used to describe the color of objects, the color palette of images, etc. By analyzing the semantics information of color contained in the text, it can be used as a constraint to implement text-based coloring tasks. Using color semantics can effectively improve the color selection in the colorization process. The study (Hu et al., 2016) establishes upon the image segmentation from natural language expression, which correctly identifies and divides the specified content out of the whole image. Hu's method (Hu et al., 2016) is based on CNN and Long Short-Term Memory (LSTM) network to process the image and linguistic information, and can be trained end-to-end. Later, Chen et al. (2018) were the first to present a language-based colorization method with the help of natural language on Oxford-102 Flowers dataset. They also use the LSTM network to transform the word into a vector through a word matrix and produces a contextual vector. And then the recurrent attentive model combines the image and language features. Finally, get the colorized results from this feature.

Unlike Chen et al. (2018) who directly combine the extracted language features and visual features to control the colorization results, Manjunatha et al. (2018) apply the feature-wise linear modulation (FiLM) (Perez et al., 2018) structure to language-based colorization with fewer parameters. Since FiLM performs feature affine transformation on the output of each convolutional block, only two additional weight matrix parameters are required for each feature map. Instead of putting a specific object into the color form, they use semantic text input to generate the color palette to achieve color palette based on user's text input.

Bahng et al. (2018) have a similar focus on this area, instead of putting a specific object into the color form, they use semantic text input to generate the color palette to achieve color palette based on user's text input. This is based on previous related work on color palette design and image editing such as research by Heer and Stone (2012). Based on Bahng's model, the user can use both single and multi-word descriptions to create a color palette and colorize the grayscale image. Like Hu's method (Hu et al., 2016), the network is based on conditional generative adversarial networks (cGANs). Bahng used the palette-and-text (PAT) data set to train the model for predicting color palette parts. The data set contains 10,183 text and five-color palette pairs. The data set

is refined by harvesting user custom-made color palettes from community websites. To process raw data from the data set, they use four annotators to vote whether the semantic word matches the color palette.

Text2Color can be divided into two parts, a Text-to-Palette Generation Network (TPN) and Palette-based Colorization Network (PCN). TPN generates a reasonable color palette based on the text input. The objective for the first cGAN is expressed as:

$$L_{D_0} = \mathbb{E}_{y \sim P_{data}} [\log D_0(\bar{c}, y)] + \mathbb{E}_{x \sim P_{data}} [\log (1 - D_0(\bar{c}, \hat{y}))] \quad (2)$$

$$L_{G_0} = \mathbb{E}_{x \sim P_{data}} [\log (1 - D_0(\bar{c}, \hat{y}))] \quad (3)$$

For the discriminator  $D_0$  wants to maximize  $L_{D_0}$  in opposition to  $G_0$  which wants to minimize  $L_{G_0}$ . Vector  $x$  and real color palette  $y$  are from the data distribution  $P_{data}$ .

Bahng found that the Huber loss is the most effective way to increase color diversity in the color palette. They decide to use Huber loss to make the generated image closer to ground truth, and they added a Kullback–Leibler divergence regulation term. In this part, they use a conditioning augmentation technique.

$$\hat{y}_i = f(s_i) \text{ where } s_i = g(\hat{y}_{i-1}, c_i, s_{i-1}) \quad (4)$$

where  $s_i$  is a GRU (Gated Recurrent Unit) hidden state, and  $i$  is a time vector, previous generated colors are stored in  $\hat{y}_{i-1}$  and the content vector  $c_i$  and the previous state stored in  $s_{i-1}$  are provided as input. This state is used as input to a fully-connected layer to output the  $i$ th color into the palette, and results are the combination of five colors to form a single palette output  $\hat{y}$ .

The text-based method only needs to use the text description to realize the image colorization work, which can be used not only for the colorization of a single image, but also for the colorization of multiple images or videos. However, the text description has certain limitations on the specification of details and the selection of the color range, so it is more suitable for controlling the color palette of the whole image and the color of single object. The enhancement of detail colors can increase the control of detail by introducing color strokes, and researchers can conduct research on multi-modal colorization methods to integrate the advantages of different models.

#### 4. Colorization based on sketch images

Sketch images consist of sparse lines, and the information in the images is sparse compared to grayscale images. Colorization of grayscale images tends to use grayscale information from the L channel in Lab space, and it is easier to judge the same semantic area by pixel values, so it is difficult to directly apply to sketch images. Usually, sketch colorization methods are mostly sample-based or need users to provide guidance information, and contain both an automatic colorization option and an interactive user mode, since the input images do not carry texture information. In early research, color hints provided by the user are spread to the entire image, but these colorization methods are limited by the quality of sketch images, the richness of color information, and values of method parameters. In the last five years, most popular models are based on neural networks, such as CNN (Krizhevsky et al., 2012), GANs (Goodfellow et al., 2014) and U-net architecture (Ronneberger et al., 2015), which can replace the manual effort of carrying out colorization and can make monochrome images more attractive. In addition, we introduce a special sketch colorization method in the last subsection, which directly generates pictures based on the input sketch images.

##### 4.1. Color hints or strokes based colorization

In an early paper in 2009, Šykora et al. (2009) described the LazyBrush colorization model which needs users to carefully make color scribbles. They transform the coloring problem of sketch images into an optimization problem, and design an energy function which consists of two main terms: smoothness and data. The smoothness term mainly ensures to hide color discontinuity, and the data term mainly considers the color prompt information added by the user. Although the LazyBrush algorithm performs well in interactive colorization, a long time is needed to calculate the iteration. In order to reduce the algorithmic complexity, Fourey et al. (2018) focus on line art images and do not process black-and-white cartoon or manga, and perform fine analysis on the local geometry of stroke contours. Compared with LazyBrush, their CPU calculation time is reduced by more than 70% in different sizes of test images. Although the algorithm based on optimization is robust and has a high success rate for colorization, methods based on deep networks have natural advantages in reducing user interaction and colorization speed.

In the method based on deep networks, the GAN usually is used as a generative model to colorize sketch images, and U-net (Ronneberger et al., 2015) are instead of the traditional encoder-decoder structure as the GAN's generator. Liu et al. (2017b) used conditional generative adversarial networks (cGANs) to train the automatic painter model to produce compatible colors for a sketch. Moreover, their architecture allows users to control the color of generated images, which is based on Sangkloy et al. (2017) to add color strokes to the input sketch image. Ci et al. (2018) proposed a novel conditional adversarial synthesis architecture, which combined with a local feature network. The main branch of the generator is developed based on U-net (Ronneberger et al., 2015) and four sub-networks, each containing a convolution block to fuse features from skip connection. And the local feature extracted from local feature network is as conditional input of generator and discriminator, to avoid overfitting characteristic of the line art, to help the generator to produce vivid colorful color outputs.

In addition to the end-to-end colorization network, Zhang et al. (2018c) proposed a two-stage line art colorization network based on color hints (known as style to paint version V3). Each stage of the network consists of a generator based on the U-Net structure and a discriminator, where the color hints marked by

the user in the two stages are directly inputted into the network together with the line art image. It is worth noting that the second-stage network uses Inception V1 (Szegedy et al., 2015) to extract the features of the simulated color draft generated by the first-stage network and merge them into the intermediate features of the generator. This method can get a good colorization result by adding color hints, but each image requires a lot of user input as color hints.

Moreover, a commercial website PaintsChainer (Yonetsuji, 2017) also colorizes sketches based on color hints, which is an automatic colorization model for sketch images. The model provides three different painting styles for the user to choose from, and will produce different results based on the same sketch image. This product is user-friendly for the non-programmer artist and can process sketch images through their online web page with no need to download any software.

Color hints and strokes provide a more convenient interactive tool to support the user to specify the color of the local details of the image. For the colorization of a single image, such methods allow the user to iteratively optimize and produce the desired result. If used on multiple images and videos, there is still a lot of interaction. However, when artists color comics, they often design the clothing and color matching of the characters in advance, and color the content according to the preset image, which has inspired some colorization work based on reference images.

##### 4.2. Reference color image based colorization

Different from color hints based colorization methods, reference based colorization methods match segment shapes or even semantic similarity between reference images and input images to colorize different positions. We divide reference based colorization methods into deep networks and graph correspondence. In the deep network based colorization methods, the color of result only takes the reference image as conditional information, and does not completely depend on it. However, the method based on graph correspondence usually copies the color value of reference images, and color value of result images completely depends on the reference image.

###### 4.2.1. Colorization by deep networks

A series of research work of Style to Paint is an important work of applying deep network to color line art images. In style to paint version V1 (Zhang et al., 2017a), the model can color sketch images based on the input reference images. The main idea is to use U-net and AC-GAN (Odena et al., 2017) in a generator for the style image. They redesigned the new network residual U-net (Ronneberger et al., 2015), and initialized the network with a Gaussian random number. The network shows a stable gradient in the training process and to solve potential noise he adds two additional losses to avoid the vanishing gradient problem. In style to paint version V3 (Zhang et al., 2018c), as introduced in Section 4.1, they proposed a two-stage CNN-based colorization framework. Instead of using the reference image, in this model users use the color hint mark to provide colors for the sketch image. In the same year, Sun et al. (2019) designed a dual conditional generative adversarial network for the colorization problem in icon design. The model can get a good colorization result on simple strokes such as icons, but color bleeding still appears sometimes, and it cannot handle complex sketch images well.

In the work of colorizing line art videos, the videos are usually divided into multiple video sequences. Researchers usually provide one or more reference color images for each video sequence to color the remaining line art frames. Thasarathan et al. (2019) proposed a line art video colorization model called automatic

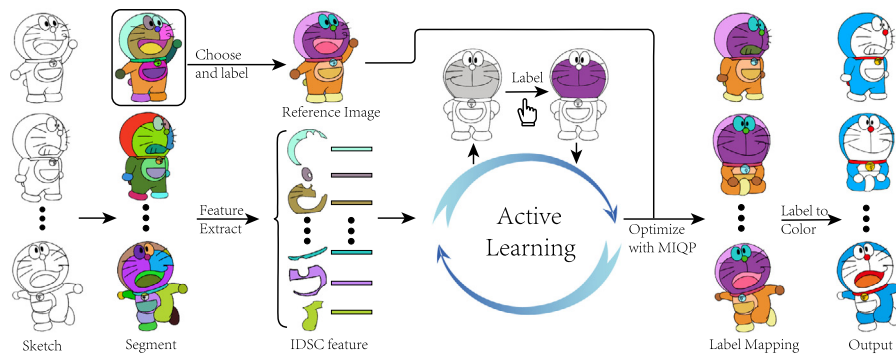


Fig. 13. An overview architecture of active colorization (Chen et al., 2022).

Temporally Coherent Video Colorization (TCVC), which extends the image-to-image translation model based on the conditional GAN (Isola et al., 2017). They input the line art image and the color image of the previous frame into the generator network for providing color prior information. The discriminator network uses the patch GAN structure proposed by Isola et al. (2017), and inputs the line art image and the corresponding color image at the same time to ensure the temporal color consistency.

TCVC (Thasarathan et al., 2019) regards the previously generated color image as the color condition of the current line art image. Although the temporal constraint is maintained, serious color error accumulation occurs. Shi et al. (2022) proposed a line art video colorization model based on a small number of reference images. The model is mainly divided into color transform network and temporal constraint network. The color transform network is mainly based on the conditional GAN, using an Embedder module to extract the style embedding vector from multiple reference images. The similarity based color transform layer is essentially an attention mechanism, but they designed the network to adaptively select the reference color input and the features after the attention calculation. The AdaIN (Huang and Belongie, 2017) structure is used to align the feature mean and variance learned from style embedding vector to the output feature's mean and variance of the similarity based color transform layer. The Temporal Constraint Network inputs multiple reference image pairs and target image pairs into a generation network composed of 3D gated convolutions (Chang et al., 2019) to obtain a final colorization result with a temporal constraint.

The method proposed by Shi et al. (2022) restricts the reference image to be in the same video clip. This method is similar to the video frame interpolation method. Subsequently, Siyao et al. (2021) propose a deep animation video interpolation network in the wild. The network also uses high-dimensional features extracted from input images to calculate a similarity matrix and learn the optical flow correspondence between adjacent frames. Compared with Shi et al. (2022) directly extracting line art image features, Siyao et al. (2021) first extract contours which are filled with the “trapped-ball” algorithm (Zhang et al., 2009) to generate color patches, and then are fed into the VGG-19 model to generate color patches features. In addition, after obtaining the bidirectional optical flow correspondence, a recurrent flow refinement network are proposed to optimize the final optical flow through multiple iterations of learning.

Given an arbitrary reference image, such methods can generate color images. Obvious artifacts such as color matching errors and color mixing are prone to occur, especially when the fill is a pure color, as shown in Fig. 14(b). Researchers can consider adding region segmentation information to optimize the colorization results.

#### 4.2.2. Colorization by graph correspondence

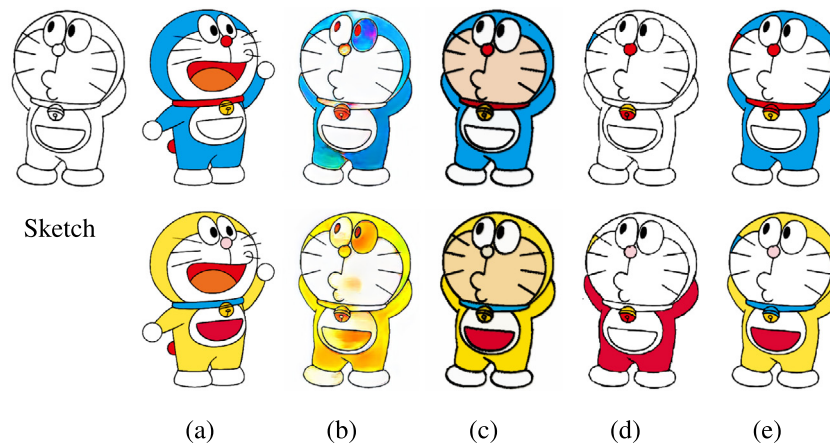
Since a line art image can be divided into different regions according to lines, the researcher can quickly transfer color from the reference image using the image matching method. Sato et al. (2014) proposed a colorization method for line art images based on image matching. First, the method uses the super-pixel method (Shi and Malik, 2000) to divide the reference image and the target image into different regions, and then the different regions are used as nodes to compose the graph structure. Then the method uses the area of the node and the length/angle of the vector composed of the centroid of different nodes to calculate the similarity between the two nodes. Finally, according to the reference image and the target image represented by the graph structure, they obtained the final matching result by solving a Quadratic Programming problem, and transferred the color of the reference image to the target image.

The method proposed by Sato et al. (2014) needs to specify the number of segmentation regions, and when the lines of the line art image are very complicated, the user needs to manually adjust the segmentation results. In addition, the method calculates the pairwise similarity of nodes, and does not make full use of the shape features of the segmented regions themselves. Chen et al. (2022) proposed a colorization method combined with active learning (Zhu et al., 2003). The overall pipeline of the method is shown in Fig. 13. First, the trapped-ball segmentation (Zhang et al., 2009) was used to automatically segment the line art image, and inner-distance shape context (IDSC) (Ling and Jacobs, 2007) is used to extract the features of the segmented regions. And then, the features of the segmented regions of the reference image and the adjacency relationship between the regions are expressed as a graph structure, and finally the mixed-integer quadratic programming method (MIQP) is used to solve the graph matching problem. As the comparison shows in Figs. 14(d)(e), with the same reference image, Sato et al. (2014)(d) failed when the number and structure of regions changed. Chen et al. (2022)(e) produce better results, however, there are still colorization errors in smaller segmented regions (such as the left side of the face).

The colorization method based on the graph correspondence needs to segment the image into sub-regions in advance, and construct the graph structure according to the adjacency relationship between sub-regions. Some methods introduce constraints such as shape similarity during optimization. But such methods suffer from segmentation mistakes and are only suitable for some simple-structured images. When the number of image areas is too large and the structure changes, the effect is slightly worse, such as the case where there are too many lines in the hair area.

#### 4.3. Text hints based method

Text based colorization methods not only learn the correspondence between the color and the semantic regions of images, but



**Fig. 14.** The comparison of reference based colorization method. The leftmost image is the sketch image. (a) is the reference image in different color styles. (b) is the result of Zhang et al. (2018c). (c) is the result of Hensman and Aizawa (2017). (d) is the result of Sato et al. (2014) and (e) is the result of Chen et al. (2022).

also the correspondence between the text and the color needs to be determined. Kim et al. (2019) purposed using text tags to colorize sketch images. Users can easily color line art images based on their input color variant tag (CVT) from which the generator produces a color result. The CVT module first extracts the features of input text information, and then merges the features to the output features through the SECat (Squeeze and Excitation with Concatenation) module.

For the first time, Zou et al. (2019) proposed a language-based interactive colorization system for scene sketches. Users can color the foreground objects and background sequentially through text scripts. They proposed a new instance matching model, which uses the DeepLab-v2 (Chen et al., 2018) network to extract sketch features, and the interactive model mLSTM (multimodal Long Short-Term Memory) (Liu et al., 2017a) is also added to the generator network, which can achieve the joint modeling of text script and images. Zou et al. (2019) also designed a foreground colorization network and a background colorization network to facilitate the processing of foreground objects and background areas with different image characteristics. It is worth noting that since the current research on semantic understanding in natural language processing is still in the preliminary stage, the input text of text based colorization methods is more similar to a text control instruction, and more accurate text understanding and color matching are still Need further research.

#### 4.4. Sketch to image synthesis

The synthesis of sketch to image is a kind of research related to the colorization of line art images. Different from the line art colorization, the sketch to image synthesis method does not strictly provide color values at different positions, but extracts the semantic features of the lines in sketch image and compares them with the existing images in the data set. After that, a new image is synthesized by fusing different matching results. Therefore, the method of sketch to image usually requires an amount of data to be matched.

##### 4.4.1. Image synthesis by internet search

Chen et al. (2009) proposed a system Sketch2Photo that can automatically synthesize realistic images from sketches with text labels. Sketches will be divided into background and multiple scene items, and then match the initial set of candidate images on the Internet based on text labels and lines. For the background image, the candidate images with inconsistent content and cluttered regions are removed. Inspired by Ben-Haim et al. (2006),

they used a clustering algorithm to filter the content consistency of the image, and count the number of regions covered by the convex hull of all scene items to determine whether the image is uncluttered or not. And then, after discarding candidate background images with salient areas and complex backgrounds, they use the grab-cut algorithm (Rother et al., 2004) to segment the expanded areas of the target object from the image. According to the shape feature and the clustering algorithm, the images with inconsistent shapes and inconsistent content are eliminated. Finally, they used the proposed hybrid method to fuse the background and scene item candidate images.

Later, Chen et al. (2013) proposed a method that can quickly build a large-scale human image database. Inspired by Sketch2-Photo Chen et al. (2009), they use a human detection algorithm (Felzenszwalb et al., 2008) to extract images containing humans from the Internet. Then they filter out the algorithm-unfriendly images, and segment the foreground and background of the image. The final human image is organized by action semantics and clothes attributes, and the user can retrieve images of the corresponding posture through the outline. In addition, Chen et al. (2013) demonstrated the use of this data to generate multi-frame personalized content image synthesis programs.

By extracting and assembling content from existing data, reasonable image results can be obtained. However, since the generated results depend on the existing data content, the image content outside the dataset cannot be generated. The emergence of generative networks solves this problem, training on existing datasets to generate new results in addition to existing ones. Therefore, the methods based on generative networks for image synthesis have become more popular.

##### 4.4.2. Image synthesis by generative networks

Sangkloy et al. (2017) proposed a sketch-to-photo architecture, similar to an image-to-image translation network (Isola et al., 2017). They are the first to use a feed-forward architecture model that can generate realistic images based on imperfect sketches. The architecture presented colorization in three different domains: faces, bedrooms and cars. Chen and Hays (2018) proposed a fully automatic model to synthesize realistic images from human-drawn sketches. The model takes the GAN model as the basic structure and proposes the Mask Residual Unit (MRU). The convolutional features of the previous layer and an additional image are fed into the MRU, then the MRU dynamically determines the final output features of the network by calculating an internal mask. This is conducive to generating results that are similar in content information to the input image and generating higher quality results.

Different from the previous image synthesis methods based on sketches and colors (Chen et al., 2009; Isola et al., 2017), Xian et al. (2018) proposed an image synthesis method based on sketches and textures. This method realizes the fine-grained control of the synthesized image based on the texture image. For this purpose, they mainly design two GAN network structures to produce preliminary synthesis results and fine-tuned textures of synthesized images respectively. In addition, a novel local texture loss  $\mathcal{L}_t$  was proposed. The loss function first randomly samples  $n$  patches of  $s \times s$  size from the generated result  $G(x)$  and the input texture image  $I_t$ , and then calculates its Local Adversarial Loss  $\mathcal{L}_{adv}$ , Local Style Loss  $\mathcal{L}_s$  and Pixel Loss  $\mathcal{L}_p$ . The specific calculation is as follows:

$$\mathcal{L}_t = \mathcal{L}_s + w_p \mathcal{L}_p + w_{adv} \mathcal{L}_{adv} \quad (5)$$

where  $\mathcal{L}_s$  use Gram matrix-based style loss and  $\mathcal{L}_p$  use L2 pixel loss.  $\mathcal{L}_{adv}$  is defined as follows:

$$\mathcal{L}_{adv} = - \sum_i (D_{txt}(h(G(x_i), R_i), h(I_t, R_i)) - 1)^2 \quad (6)$$

where  $D_{txt}$  is a local texture discriminator,  $h(x, R)$  represents cropping a patch from the segmentation mask  $R$  of image  $x$ .

Inspired by Gao et al. (2019), Chen et al. (2020) proposed a local–global network for sketch-based image synthesis. They take the facial structure into consideration and design a manifold projection to deal with rough/in-complete sketch. For easier use by ordinary users with little drawing skill, they design an interface with shadow-guided on the drawing board like *ShadowDraw* (Lee et al., 2011). They further proposed a disentanglement framework (Chen et al., 2021) which could disentangle the geometry and the appearance features from facial images. With this framework, it is possible to edit the appearance by changing the appearance image, and edit the geometry by using a sketch. For example, adding wrinkles or changing the color of hair.

The image synthesis method based on generative networks provides a new solution for sketch-to-image synthesis. However, due to the variety of styles of hand-drawn sketches by users, how to map sketches to reasonable results is a challenging research problem. Analysis of the semantic information of the lines in the sketch can give guidance when generating the image, thereby improving the correspondence between the image and the sketch.

## 5. Hybrid method—cartoon and manga image colorization

In this Section, we will introduce some hybrid methods, for which the input is a black-and-white cartoon or manga. Those images contain a clear sketch-based line art, and also carry some black-and-white shading information. In manga, the colorization process has an additional pattern colorization process. Colorizing black and white manga can make them more attractive, and using semi-automatic methods as a helper can increase the speed of drawing.

Early work by Sỳkora et al. (2005) carried out research on the colorization of black-and-white cartoons. They divided each frame into two parts: the static background and dynamic foreground, which were colorized separately. To separate frames into two parts, they used the Laplacian-of-Gaussian mask to detect outlines. When realizing continuous frame colorization, they used a probabilistic reasoning scheme to calculate the similar region and neighborhood relationship to process more video frames. Finally, to improve the quality they added color modulation, composition, and a technique to remove dust spots in order to improve the appearance of the final image. In this way, the image edges are sharp and can adapt to complicated drawing. The whole process is semi-automatic.

For manga colorization, Qu et al. (2006) proposed a manga colorization method based on user scribbles. In this model, they use a novel texture-based level set method for segmentation. The model provides two modes of propagation for segmentation, pattern-continuous and intensity continuous propagation. Users can easily alter those two methods in different steps. The model can identify the hatching and screening effects that are used in traditional paper comics. After confirming the segment region, they use three different colorization methods for various conditions, such as color replacement, stroke-preserving colorization, and pattern to shading.

Hensman and Aizawa (2017) used cGAN and post-processing to colorize manga images, reducing the degree of user interaction and producing better results. This method can be divided into the following steps: screentone removal, segmentation, color selection, saturation increase, color quantization, and generation of shading. Model training is based on the corresponding grayscale image and colorized reference image as a single image pair. Based on the model parameters obtained from the training of a single image pair, the model can colorize manga images similar to the reference image. However, when the character's clothing becomes complex, the model cannot achieve the correct correspondences between different frames. The results will not only contain artifacts, but also are not as colorful as the reference image, as shown in Fig. 14(c).

Furusawa et al. (2017) connected the common manga colorization problem and was able to color the same character at once. The overall pipeline of the semi-automatic manga colorization method is shown in Fig. 15. Users need to provide the reference image and the corresponding color image together as input. A segmentation step is based on Ishii et al. (2009), and semi-automatic colorization is based on the CNN architecture, using an encoder–decoder network with some refinement to improve the performance. After the palette model, the model produces the draft colorization result, and the user needs to carry out interactive revision of the details to ensure the correctness of the output result. During revision users can choose either color dots or a histogram to adjust the image.

Sketch images contain sparse lines, and grayscale images represent image information such as shading and texture. The manga image may contain one or more styles, such as region boundary, narrow structure, specific region representations, etc. This leads to a conflict in the understanding of image content. Distinguishing the texture of an object or the specific representation in the image is the key to improving the colorization of the manga.

## 6. Assessment of colorization

Researchers typically evaluate different colorization methods based on quantitative and qualitative aspects. In quantitative evaluation, researchers assume the existence of unique ground truth and provide unique ground truth, which facilitates simple analysis of colorization results. Then researchers typically use root mean square error (RMSE) (Deshpande et al., 2015; Larsson et al., 2016), peak signal to noise ratio (PSNR) (Larsson et al., 2016; Cheng et al., 2015; Messaoud et al., 2018), or structural similarity index measurement (SSIM) (Messaoud et al., 2018; Chen et al., 2022; Shi et al., 2022) image evaluation indicators, comparing the results of different colorization methods.

However, the widely used indicators such as PSNR and SSIM do not fully match human perception. In recent years, researchers have generally extracted the deep features of images to compare the perception similarity between images. A common measurement is to use the Inception Score (IS) (Salimans et al., 2016) to evaluate the quality of the generated images. Given a desired image set, we could use the Frechet Inception Distance score (FID)

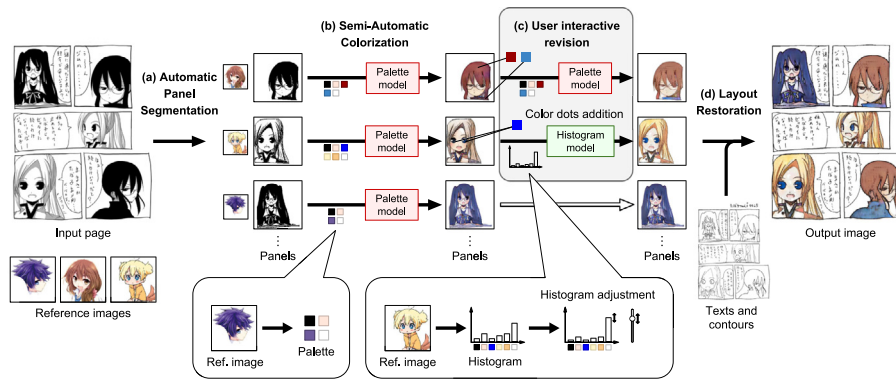


Fig. 15. The overall pipeline of the semi-automatic manga colorization (Furusawa et al., 2017).

(Heusel et al., 2017; Kumar et al., 2020) to measure the difference between the distribution of the real images and the distribution of the generated images. Zhang et al. (2018b) systematically analyzed the unreasonable effectiveness of deep features as a perceptual metric, and proposed the Learned Perceptual Image Patch Similarity (LPIPS) metric which better captures image perception similarity. LPIPS has been used to evaluate colorization by Yoo et al. (2019), and this has been shown to be an effective colorization evaluation measure by Žeger et al. (2021). In addition, researchers will also compare the colorization time of different methods for different resolution images (Liu et al., 2008; Iizuka et al., 2016) and the computing resources, such as network parameter size (Zhang et al., 2016; Chen and Hays, 2018).

In the qualitative comparison of different colorization methods, researchers usually use user study to evaluate, and then use box plots, radar plots, or bar graphs to display the results in the paper. Specifically, in the user study design scheme, the researcher will design multiple questions for different colorization results for users to score and evaluate. These questions mainly include: 1. The authenticity of the colored image (Gupta et al., 2012): They will mix colored images with real (natural) color images, allowing users to score the authenticity of the image; 2. Whether the colored result obeys the hint color or reference image color (Bahng et al., 2018; Li et al., 2019; Yoo et al., 2019): participants will be asked to evaluate the color consistency between the coloring result and the reference image or palette; 3. Color segmentation (Zhang et al., 2018c; Kim et al., 2019): Whether the coloring method can accurately identify the colors that should be used in different areas to prevent color bleeding and fusing problems; 4. Coloring time (Zhang et al., 2018c): In interactive colorization methods the researcher will compare the interactive coloring time needed for users to complete the colorization tasks; 5. Quality of the coloring results: Whether the coloring method can produce reasonable results (Li et al., 2019; Kim et al., 2019). In addition, researchers not only display user evaluation results, but also use analysis of variance to determine whether there are significant differences between different evaluation indicators, and perform statistical analysis on user survey results (Sun et al., 2019; Li et al., 2019; Chen et al., 2020).

Moreover, existing indicators for quantitative evaluation of colorization results are inaccurate, including PSNR, SSIM, LPIPS, and qualitative comparison can better evaluate the pros and cons of different methods. However, due to the lack of a unified public dataset, especially in the field of line art colorization task, qualitative evaluation cannot fairly compare the capabilities of different colorization methods. Researchers should be committed to open datasets of different coloring tasks, and then use multiple quantitative indicators and qualitative comparisons to evaluate different methods under the same dataset, so as to evaluate the pros and cons of different methods more reasonably.

## 7. Conclusions

In this survey, we summarize many different methods in image colorization and related areas. Based on the user interaction, we divided grayscale colorization methods into three categories: fully-automatic colorization, semi-automatic colorization and language-based colorization. For sketch images, most methods combine the fully-automatic colorization models with user guidance together, as simple editing can improve the result and make it closer to their expectations. For each area, they are pursuing different goals to meet the different demands, and with deep learning, huge progress has been made in recent years.

Exploring the connection between colorization for grayscale images and for sketch images, we found that there are many similarities, and the color processing techniques are the same. The common goal of grayscale colorization is to achieve the most realistic color which should be the same as the ground truth, so there are many methods that focus on how to fill in one or more right colors for each pixel in the image. Usually some object colors are not fixed to one color in the real-world such as balloons and the color of a dress. And for sketch image colorization, the methods require accuracy which is achieved by image segmentation, and prevents the color bleeding issue. Most anime and manga character colors are unknown, so using the ‘wrong’ color will not be a big issue in sketch colorization. We summarize the models in Table 1, which lists the grayscale image and sketch image colorization methods based on the publication year respectively. It can be seen from Table 1 that before 2016 most research was carried out on grayscale images, and thereafter research has increased on sketch images. Moreover, with the development of neural networks, the latest methods tend to be more automatic and less interactive. The two main problems of the current research are how to further improve the quality of generation while also reducing interactions, and how to accurately control the boundaries of colorization areas in a convenient way, such as text or sketch.

Developments in this area are closely connected with segmentation, semantic and style transfer study. Regardless of whether it is to improve the quality of image colorization or interactively control image colorization, researchers need to propose more accurate image region segmentation and semantic matching methods to determine the boundaries of coloring regions, so as to avoid color overflow and cross-coloring in colorized images. The current research trend is to use machine learning to solve problems instead of determining everything by hand. However the parameters of the network are hard to determine based on a lack of information, and they still do not provide a good solution for all cases. At the same time, the current deep learning-based methods cannot precisely control the colorization result regardless of whether it is processing grayscale images or line

**Table 1**

Summary for different categories of image colorization method. **A** marks fully-automatic colorization. **R**, **S**, **H** and **T** marks various types of color control, which are reference images, color-scribbles, color hints and text hints. Methods with a light blue background are neural network-based approaches.

Input category	Paper	Year	Condition	Paper	Year	Condition	
Gray-scale image	Hertzmann et al. (2001)	2001	R	Welsh et al. (2002)	2002	R	
	Levin et al. (2004)	2004	S	Huang et al. (2005)	2005	R S H	
	Ironi et al. (2005)	2005	R	Yatziv and Sapiro (2006)	2006	S	
	Luan et al. (2007)	2007	S	Liu et al. (2008)	2008	A R	
	Morimoto et al. (2009)	2009	A R	Gupta et al. (2012)	2012	R	
	Cheng et al. (2015)	2015	A	Deshpande et al. (2015)	2015	A	
	Iizuka et al. (2016)	2016	A	Larsson et al. (2016)	2016	A	
	Zhang et al. (2016)	2016	A	Liao et al. (2017)	2017	R	
	Zhang et al. (2017b)	2017	H	Deshpande et al. (2017)	2017	A	
	Chen et al. (2018)	2018	H	Manjunatha et al. (2018)	2018	H	
	Bahng et al. (2018)	2018	H	Messaoud et al. (2018)	2018	A	
	Vondrick et al. (2018)	2018	R	He et al. (2018)	2018	R	
	Jamriška et al. (2019)	2019	R	Zhang et al. (2019)	2019	R	
	Lei and Chen (2019)	2019	A	Yoo et al. (2019)	2019	A	
	Li et al. (2019)	2019	R	Iizuka and Simo-Serra (2019)	2019	R	
	Fang et al. (2020)	2020	R	Su et al. (2020)	2020	A R	
	Sketch image	Sýkora et al. (2009)	2009	S	Chen et al. (2009)	2009	R
		Chen et al. (2013)	2013	R	Sato et al. (2014)	2014	R
		Yonetsuji (2017)	2017	A S	Zhang et al. (2017a)	2017	R
		Liu et al. (2017b)	2017	A S	Sangkloy et al. (2017)	2017	S T
Xian et al. (2018)		2018	R	Zhang et al. (2018c)	2018	S H	
Ci et al. (2018)		2018	A S	Kim et al. (2019)	2019	H T	
Zou et al. (2019)		2019	T	Sun et al. (2019)	2019	R	
Thasarathan et al. (2019)		2019	R	Shi et al. (2022)	2020	R	
Chen et al. (2022)		2020	R	Chen et al. (2020)	2020	A	
Chen et al. (2021)		2021	A R	Siyao et al. (2021)	2021	A R	
Manga image	Sýkora et al. (2005)	2005	R	Qu et al. (2006)	2006	S	
	Furusawa et al. (2017)	2017	R H	Hensman and Aizawa (2017)	2017	A	

art images, and lacks an effective and precise method to control interaction. This is mainly reflected in the fact that the interactive colorization methods based on hints, sketch or text can only provide color prior information for the region, but cannot control the region boundary or the propagation of color prior information.

In the existing methods, interactive control methods include simple color feature fusion, intermediate network output feature normalization, and color palette control after quantization of the reference image. These feature control methods need to be further improved. In the future, researchers may consider using the feature fusion method of the transformer model (Vaswani et al., 2017) or the CLIP (Contrastive Language-Image Pre-training) model (Radford et al., 2021). In addition, the information contained in the sketch is sparse and ambiguous. The semantic analysis of the lines in the sketch can effectively reduce the ambiguity of the lines and further improve the quality of the image. To reach it, we need to refocus on the interpretability of the neural network structure, and have a deeper understanding of how color is formed in the colorization process. There are many commercial applications of colorization, and these can be further extended by exploring new formulations and solutions.

**Ethical approval**

This study does not contain any studies with human or animal subjects performed by any of the authors. All data used in the study are taken from public databases that were published in the past.

**CRedit authorship contribution statement**

**Shu-Yu Chen:** Writing – original draft. **Jia-Qi Zhang:** Writing – original draft. **You-You Zhao:** Writing – original draft. **Paul L. Rosin:** Writing – review & editing. **Yu-Kun Lai:** Writing – original draft. **Lin Gao:** Writing – review & editing, Project administration.

**Declaration of competing interest**

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: Shu-Yu Chen reports financial support was provided by Foundation for Innovative Research Groups of the National Natural Science Foundation of China. Lin Gao reports financial support was provided by Royal Society Newton Advanced Fellowship.

**Acknowledgments**

We thank the anonymous reviewers for the constructive comments. This work was supported by grants from the National Natural Science Foundation of China (No. 61872440, No. 62061136007 and No. 62102403), the Beijing Municipal Natural Science Foundation for Distinguished Young Scholars (No. JQ21013), the Youth Innovation Promotion Association CAS, Royal Society Newton Advanced Fellowship (No. NAF\R2\192151) and the Open Project Program of State Key Laboratory of Virtual Reality Technology and Systems, Beihang University (No. VRLAB2022C07).

## References

- Antic, J., 0000. Deoldify, <https://deoldify.ai/>.
- Anwar, S., Tahir, M., Li, C., Mian, A., Khan, F.S., Muzaffar, A.W., 2020. Image colorization: A survey and dataset. *arXiv:2008.10774*.
- Ashikhmin, M., 2001. Synthesizing natural textures. In: Proceedings of the 2001 Symposium on Interactive 3D Graphics. In: I3D '01, Association for Computing Machinery, New York, NY, USA, pp. 217–226. <http://dx.doi.org/10.1145/364338.364405>.
- Bahng, H., Yoo, S., Cho, W., Keetae Park, D., Wu, Z., Ma, X., Choo, J., 2018. Coloring with words: Guiding image colorization through text-based palette generation. In: Proceedings of the European Conference on Computer Vision (ECCV), pp. 431–447.
- Ben-Haim, N., Babenko, B., Belongie, S., 2006. Improving web-based image search via content based clustering. In: 2006 Conference on Computer Vision and Pattern Recognition Workshop (CVPRW'06). IEEE, p. 106.
- Bénard, P., Cole, F., Kass, M., Mordatch, I., Hegarty, J., Senn, M.S., Fleischer, K., Pesare, D., Breeden, K., 2013. Stylizing animation by example. *ACM Trans. Graph.* 32 (4), <http://dx.doi.org/10.1145/2461912.2461929>.
- Bugeau, A., Ta, V., Papadakis, N., 2014. Variational exemplar-based image colorization. *IEEE Trans. Image Process.* 23 (1), 298–307.
- Chang, Y.-L., Liu, Z.Y., Lee, K.-Y., Hsu, W., 2019. Free-form video inpainting with 3D gated convolution and temporal PatchGAN. In: Proceedings of the International Conference on Computer Vision (ICCV).
- Chen, T., Cheng, M.-M., Tan, P., Shamir, A., Hu, S.-M., 2009. Sketch2photo: Internet image montage. *ACM Trans. Graph.* 28 (5), 1–10. <http://dx.doi.org/10.1145/1618452.1618470>.
- Chen, W., Hays, J., 2018. SketchyGAN: Towards diverse and realistic sketch to image synthesis. In: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 9416–9425.
- Chen, D., Liao, J., Yuan, L., Yu, N., Hua, G., 2017. Coherent online video style transfer. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 1105–1114.
- Chen, S.-Y., Liu, F.-L., Lai, Y.-K., Rosin, P.L., Li, C., Fu, H., Gao, L., 2021. Deep-faceediting: Deep face generation and editing with disentangled geometry and appearance control. *ACM Trans. Graph.* 40 (4), <http://dx.doi.org/10.1145/3450626.3459760>.
- Chen, L., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L., 2018. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs. *IEEE Trans. Pattern Anal. Mach. Intell.* 40 (4), 834–848.
- Chen, J., Shen, Y., Gao, J., Liu, J., Liu, X., 2018. Language-based image editing with recurrent attentive models. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 8721–8729.
- Chen, S.-Y., Su, W., Gao, L., Xia, S., Fu, H., 2020. DeepFaceDrawing: Deep generation of face images from sketches. *ACM Trans. Graph.* 39 (4), <http://dx.doi.org/10.1145/3386569.3392386>.
- Chen, T., Tan, P., Ma, L.-Q., Cheng, M.-M., Shamir, A., Hu, S.-M., 2013. Poseshop: Human image database construction and personalized content synthesis. *IEEE Trans. Vis. Comput. Graphics* 19 (5), 824–837. <http://dx.doi.org/10.1109/TVCG.2012.148>.
- Chen, S.-Y., Zhang, J.-Q., Gao, L., He, Y., Xia, S., Shi, M., Zhang, F.-L., 2022. Active colorization for cartoon line drawings. *IEEE Trans. Vis. Comput. Graphics* 28 (2), 1198–1208. <http://dx.doi.org/10.1109/TVCG.2020.3009949>.
- Cheng, Z., Yang, Q., Sheng, B., 2015. Deep colorization. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 415–423.
- Ci, Y., Ma, X., Wang, Z., Li, H., Luo, Z., 2018. User-guided deep anime line art colorization with conditional adversarial networks. In: 2018 ACM Multimedia Conference on Multimedia Conference. ACM, pp. 1536–1544.
- Deshpande, A., Lu, J., Yeh, M.-C., Chong, M.J., Forsyth, D., 2017. Learning diverse image colorization. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE, <http://dx.doi.org/10.1109/cvpr.2017.307>.
- Deshpande, A., Rock, J., Forsyth, D., 2015. Learning large-scale automatic image colorization. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 567–575.
- Fang, F., Wang, T., Zeng, T., Zhang, G., 2020. A superpixel-based variational model for image colorization. *IEEE Trans. Vis. Comput. Graphics* 26 (10), 2931–2943.
- Felzenszwalb, P., McAllester, D., Ramanan, D., 2008. A discriminatively trained, multiscale, deformable part model. In: 2008 IEEE Conference on Computer Vision and Pattern Recognition, pp. 1–8.
- Fourey, S., Tschumperlé, D., Revoy, D., 2018. A fast and efficient semi-guided algorithm for flat coloring line-arts. In: Proceedings of the Conference on Vision, Modeling, and Visualization. In: EG VMV '18, Eurographics Association, Goslar, DEU, pp. 1–9. <http://dx.doi.org/10.2312/vmv.20181247>.
- Furusawa, C., Hiroshiba, K., Ogaki, K., Odagiri, Y., 2017. Comicolorization: Semi-automatic manga colorization. In: SIGGRAPH Asia 2017 Technical Briefs. In: SA '17, Association for Computing Machinery, New York, NY, USA, <http://dx.doi.org/10.1145/3145749.3149430>.
- Gao, L., Yang, J., Wu, T., Yuan, Y.-J., Fu, H., Lai, Y.-K., Zhang, H., 2019. SDM-net: Deep generative network for structured deformable mesh. *ACM Trans. Graph.* 38 (6), 243:1–243:15, Proceedings of ACM SIGGRAPH Asia 2019.
- Ge, C., Sun, H., Song, Y.-Z., Ma, Z., Liao, J., 2022. Exploring local detail perception for scene sketch semantic segmentation. *IEEE Trans. Image Process.*
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y., 2014. Generative adversarial nets. In: Advances in Neural Information Processing Systems. pp. 2672–2680.
- Gupta, R.K., Chia, A.Y.-S., Rajan, D., Ng, E.S., Zhiyong, H., 2012. Image colorization using similar images. In: Proceedings of the 20th ACM International Conference on Multimedia. ACM, pp. 369–378.
- He, M., Chen, D., Liao, J., Sander, P.V., Yuan, L., 2018. Deep exemplar-based colorization. *ACM Trans. Graph.* 37 (4), <http://dx.doi.org/10.1145/3197517.3201365>.
- Heer, J., Stone, M., 2012. Color naming models for color selection, image editing and palette design. In: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems. ACM, pp. 1007–1016.
- Hensman, P., Aizawa, K., 2017. cGAN-based manga colorization using a single training image. In: 2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR). Vol. 3, IEEE, pp. 72–77.
- Hertzmann, A., Jacobs, C.E., Oliver, N., Curless, B., Salesin, D.H., 2001. Image analogies. In: Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques. ACM, pp. 327–340.
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S., 2017. GANs trained by a two time-scale update rule converge to a local Nash equilibrium. *Adv. Neural Inf. Process. Syst.* 30.
- Hu, R., Rohrbach, M., Darrell, T., 2016. Segmentation from natural language expressions. In: European Conference on Computer Vision. Springer, pp. 108–124.
- Huang, X., Belongie, S., 2017. Arbitrary style transfer in real-time with adaptive instance normalization. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 1501–1510.
- Huang, Y.-C., Tung, Y.-S., Chen, J.-C., Wang, S.-W., Wu, J.-L., 2005. An adaptive edge detection based colorization algorithm and its applications. In: Proceedings of the 13th Annual ACM International Conference on Multimedia. ACM, pp. 351–354.
- Iizuka, S., Simo-Serra, E., 2019. Deepremaster: Temporal source-reference attention networks for comprehensive video enhancement. *ACM Trans. Graph.* 38 (6), <http://dx.doi.org/10.1145/3355089.3356570>.
- Iizuka, S., Simo-Serra, E., Ishikawa, H., 2016. Let there be color! joint end-to-end learning of global and local image priors for automatic image colorization with simultaneous classification. *ACM Trans. Graph.* 35 (4), <http://dx.doi.org/10.1145/2897824.2925974>.
- Ironi, R., Cohen-Or, D., Lischinski, D., 2005. Colorization by example. In: Rendering Techniques. Citeseer, pp. 201–210.
- Ishii, D., Kawamura, K., Watanabe, H., 2009. A study on control parameters of frame separation method for comic images. Tech. rep., IEICE technical report, The Institute of Electronics, Information and Communication Engineers.
- Isola, P., Zhu, J., Zhou, T., Efros, A.A., 2017. Image-to-image translation with conditional adversarial networks. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 5967–5976.
- Isola, P., Zhu, J.-Y., Zhou, T., Efros, A.A., 2017. Image-to-image translation with conditional adversarial networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1125–1134.
- Jamriška, O., Sochorová, v., Texler, O., Lukáč, M., Fišer, J., Lu, J., Shechtman, E., Sýkora, D., 2019. Stylizing video by example. *ACM Trans. Graph.* 38 (4), <http://dx.doi.org/10.1145/3306346.3323006>.
- Kang, S.H., March, R., 2007. Variational models for image colorization via chromaticity and brightness decomposition. *IEEE Trans. Image Process.* 16 (9), 2251–2261.
- Kim, H., Jhoo, H.Y., Park, E., Yoo, S., 2019. Tag2pix: Line art colorization using text tag with secant and changing loss. *arXiv preprint arXiv:1908.05840*.
- Kirkpatrick, S., Gelatt, C.D., Vecchi, M.P., 1983. Optimization by simulated annealing. *Science* 220 (4598), 671–680.
- Krizhevsky, A., Sutskever, I., Hinton, G.E., 2012. Imagenet classification with deep convolutional neural networks. In: Advances in Neural Information Processing Systems. pp. 1097–1105.
- Kumar, M., Weissenborn, D., Kalchbrenner, N., 2020. Colorization transformer. In: International Conference on Learning Representations.
- Larsson, G., Maire, M., Shakhnarovich, G., 2016. Learning representations for automatic colorization. In: European Conference on Computer Vision. Springer, pp. 577–593.
- Lee, J., Kim, E., Lee, Y., Kim, D., Chang, J., Choo, J., 2020. Reference-based sketch image colorization using augmented-self reference and dense semantic correspondence. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5801–5810.
- Lee, Y.J., Zitnick, C.L., Cohen, M.F., 2011. Shadowdraw: Real-time user guidance for freehand drawing. In: ACM SIGGRAPH 2011 Papers. In: SIGGRAPH '11, Association for Computing Machinery, New York, NY, USA, <http://dx.doi.org/10.1145/1964921.1964922>.
- Lei, C., Chen, Q., 2019. Fully automatic video colorization with self-regularization and diversity. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3753–3761.



- Levin, A., Lischinski, D., Weiss, Y., 2004. Colorization using optimization. *ACM Trans. Graph.* 23 (3), 689–694. <http://dx.doi.org/10.1145/1015706.1015780>.
- Li, B., Lai, Y., John, M., Rosin, P.L., 2019. Automatic example-based image colorization using location-aware cross-scale matching. *IEEE Trans. Image Process.* 28 (9), 4606–4619.
- Liao, J., Yao, Y., Yuan, L., Hua, G., Kang, S.B., 2017. Visual attribute transfer through deep image analogy. *ACM Trans. Graph.* 36 (4), <http://dx.doi.org/10.1145/3072959.3073683>.
- Ling, H., Jacobs, D.W., 2007. Shape classification using the inner-distance. *IEEE Trans. Pattern Anal. Mach. Intell.* 29 (2), 286–299. <http://dx.doi.org/10.1109/TPAMI.2007.41>.
- Liu, C., Lin, Z., Shen, X., Yang, J., Lu, X., Yuille, A., 2017a. Recurrent multimodal interaction for referring image segmentation. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV).
- Liu, Y., Qin, Z., Luo, Z., Wang, H., 2017b. Auto-painter: Cartoon image generation from sketch by using conditional generative adversarial networks. *arXiv preprint arXiv:1705.01908*.
- Liu, X., Wan, L., Qu, Y., Wong, T.-T., Lin, S., Leung, C.-S., Heng, P.-A., 2008. Intrinsic colorization. In: ACM SIGGRAPH Asia 2008 Papers. In: SIGGRAPH Asia '08, Association for Computing Machinery, New York, NY, USA, <http://dx.doi.org/10.1145/1457515.1409105>.
- Liu, C., Yuen, J., Torralba, A., 2011. SIFT flow: Dense correspondence across scenes and its applications. *IEEE Trans. Pattern Anal. Mach. Intell.* 33 (5), 978–994.
- Lowe, D.G., 2004. Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vision* 60 (2), 91–110. <http://dx.doi.org/10.1023/B:VISI.0000029664.99615.94>.
- Luan, Q., Wen, F., Cohen-Or, D., Liang, L., Xu, Y.-Q., Shum, H.-Y., 2007. Natural image colorization. In: EGSR'07, Eurographics Association, Goslar, DEU, pp. 309–320.
- Manjunatha, V., Iyyer, M., Boyd-Graber, J., Davis, L., 2018. Learning to color from language. In: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers). Association for Computational Linguistics, New Orleans, Louisiana, pp. 764–769. <http://dx.doi.org/10.18653/v1/N18-2120>, URL <https://www.aclweb.org/anthology/N18-2120>.
- Messaoud, S., Forsyth, D., Schwing, A.G., 2018. Structural consistency and controllability for diverse colorization. In: Proceedings of the European Conference on Computer Vision (ECCV), pp. 596–612.
- Morimoto, Y., Taguchi, Y., Naemura, T., 2009. Automatic colorization of grayscale images using multiple images on the web. In: SIGGRAPH '09: Posters. In: SIGGRAPH '09, Association for Computing Machinery, New York, NY, USA, <http://dx.doi.org/10.1145/1599301.1599333>.
- Odena, A., Olah, C., Shlens, J., 2017. Conditional image synthesis with auxiliary classifier GANs. In: Proceedings of the 34th International Conference on Machine Learning-Volume 70. *JMLR.org*, pp. 2642–2651.
- Patterson, G., Hays, J., 2012. Sun attribute database: Discovering, annotating, and recognizing scene attributes. In: 2012 IEEE Conference on Computer Vision and Pattern Recognition. pp. 2751–2758. <http://dx.doi.org/10.1109/CVPR.2012.6247998>.
- Perez, E., Strub, F., de Vries, H., Dumoulin, V., Courville, A.C., 2018. Film: Visual reasoning with a general conditioning layer. In: AAAI.
- Pierre, F., Aujol, J.-F., Bugeau, A., Papadakis, N., Ta, V.-T., 2015. Luminance-chrominance model for image colorization. *SIAM J. Imaging Sci.* 8 (1), 536–563.
- Qu, Y., Wong, T.-T., Heng, P.-A., 2006. Manga colorization. In: ACM SIGGRAPH 2006 Papers. In: SIGGRAPH '06, Association for Computing Machinery, New York, NY, USA, pp. 1214–1220. <http://dx.doi.org/10.1145/1179352.1142017>.
- Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al., 2021. Learning transferable visual models from natural language supervision. In: International Conference on Machine Learning. PMLR, pp. 8748–8763.
- Ratliff, N.D., Silver, D., Bagnell, J.A., 2009. Learning to search: Functional gradient techniques for imitation learning. *Auton. Robots* 27 (1), 25–53.
- Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: Convolutional networks for biomedical image segmentation. In: Navab, N., Hornegger, J., Wells, W.M., Frangi, A.F. (Eds.), *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*. Springer International Publishing, Cham, pp. 234–241.
- Rother, C., Kolmogorov, V., Blake, A., 2004. Grabcut: Interactive foreground extraction using iterated graph cuts. *ACM Trans. Graph.* 23 (3), 309–314. <http://dx.doi.org/10.1145/1015706.1015720>.
- Rudin, L.I., Osher, S., Fatemi, E., 1992. Nonlinear total variation based noise removal algorithms. *Physica D* 60 (1–4), 259–268.
- Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., Chen, X., 2016. Improved techniques for training GANs. *Adv. Neural Inf. Process. Syst.* 29, 2234–2242.
- Sangkloy, P., Lu, J., Fang, C., Yu, F., Hays, J., 2017. Scribbler: Controlling deep image synthesis with sketch and color. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 5400–5409.
- Sato, K., Matsui, Y., Yamasaki, T., Aizawa, K., 2014. Reference-based manga colorization by graph correspondence using quadratic programming. In: SIGGRAPH Asia 2014 Technical Briefs. In: SA '14, Association for Computing Machinery, New York, NY, USA, <http://dx.doi.org/10.1145/2669024.2669037>.
- Shi, J., Malik, J., 2000. Normalized cuts and image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.* 22, 888–905.
- Shi, M., Zhang, J.-Q., Chen, S.-Y., Gao, L., Lai, Y., Zhang, F.-L., 2022. Reference-based deep line art video colorization. *IEEE Trans. Vis. Comput. Graphics*.
- Simonyan, K., Zisserman, A., 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
- Siyao, L., Zhao, S., Yu, W., Sun, W., Metaxas, D., Loy, C.C., Liu, Z., 2021. Deep animation video interpolation in the wild. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 6587–6595.
- Su, J.-W., Chu, H.-K., Huang, J.-B., 2020. Instance-aware image colorization. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 7968–7977.
- Sun, T.-H., Lai, C.-H., Wang, S.-K., Wang, Y.-S., 2019. Adversarial colorization of icons based on contour and color conditions. In: Proceedings of the 27th ACM International Conference on Multimedia. In: MM '19, Association for Computing Machinery, New York, NY, USA, pp. 683–691. <http://dx.doi.org/10.1145/3343031.3351041>.
- Sýkora, D., Buriánek, J., Žára, J., 2005. Colorization of black-and-white cartoons. *Image Vis. Comput.* 23 (9), 767–782.
- Sýkora, D., Dingliana, J., Collins, S., 2009. Lazybrush: Flexible painting tool for hand-drawn cartoons. In: *Computer Graphics Forum*. Vol. 28, (2), Wiley Online Library, pp. 599–608.
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A., 2015. Going deeper with convolutions. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1–9.
- Thasarathan, H., Nazeri, K., Ebrahimi, M., 2019. Automatic temporally coherent video colorization. *arXiv:1904.09527*.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I., 2017. Attention is all you need. In: *Advances in Neural Information Processing Systems*. pp. 5998–6008.
- Vondrick, C., Shrivastava, A., Fathi, A., Guadarrama, S., Murphy, K., 2018. Tracking emerges by colorizing videos. In: Proceedings of the European Conference on Computer Vision (ECCV), pp. 391–408.
- Wang, X.-H., Jia, J., Liao, H.-Y., Cai, L.-H., 2012. Affective image colorization. *J. Comput. Sci. Tech.* 27 (6), 1119–1128.
- Wang, J., Wang, X., 2012. Vcells: Simple and efficient superpixels using edge-weighted centroidal voronoi tessellations. *IEEE Trans. Pattern Anal. Mach. Intell.* 34 (6), 1241–1247.
- Welsh, T., Ashikhmin, M., Mueller, K., 2002. Transferring color to greyscale images. *ACM Trans. Graph.* 21 (3), 277–280. <http://dx.doi.org/10.1145/566654.566576>.
- Wu, Y., Wang, X., Li, Y., Zhang, H., Zhao, X., Shan, Y., 2021. Towards vivid and diverse image colorization with generative color prior. In: International Conference on Computer Vision (ICCV).
- Xian, W., Sangkloy, P., Agrawal, V., Raj, A., Lu, J., Fang, C., Yu, F., Hays, J., 2018. TextureGAN: Controlling deep image synthesis with texture patches. In: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8456–8465.
- Yang, L., Zhuang, J., Fu, H., Wei, X., Zhou, K., Zheng, Y., 2021. Sketchgnn: Semantic sketch segmentation with graph neural networks. *ACM Trans. Graph.* 40 (3), 1–13.
- Yatziv, L., Sapiro, G., 2006. Fast image and video colorization using chrominance blending. *IEEE Trans. Image Process.* 15 (5), 1120–1129.
- Yonetsuji, T., 2017. Paintschainer. [github.com/pfnnet/Paintschainer](https://github.com/pfnnet/Paintschainer) 1, 2.
- Yoo, S., Bahng, H., Chung, S., Lee, J., Chang, J., Choo, J., 2019. Coloring with limited data: Few-shot colorization via memory augmented networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).
- Žeger, I., Grgic, S., Vuković, J., Šišul, G., 2021. Grayscale image colorization methods: Overview and evaluation. *IEEE Access* 9, 113326–113346. <http://dx.doi.org/10.1109/ACCESS.2021.3104515>.
- Zhang, J., Chen, Y., Li, L., Fu, H., Tai, C.-L., 2018a. Context-based sketch classification. In: Proceedings of the Joint Symposium on Computational Aesthetics and Sketch-Based Interfaces and Modeling and Non-Photorealistic Animation and Rendering, pp. 1–10.
- Zhang, S.-H., Chen, T., Zhang, Y.-F., Hu, S.-M., Martin, R.R., 2009. Vectorizing cartoon animations. *IEEE Trans. Vis. Comput. Graphics* 15 (4), 618–629.

- Zhang, B., He, M., Liao, J., Sander, P.V., Yuan, L., Bermak, A., Chen, D., 2019. Deep exemplar-based video colorization. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 8052–8061.
- Zhang, R., Isola, P., Efros, A.A., 2016. Colorful image colorization. In: *European Conference on Computer Vision*. Springer, pp. 649–666.
- Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O., 2018b. The unreasonable effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 586–595.
- Zhang, L., Ji, Y., Lin, X., Liu, C., 2017a. Style transfer for anime sketches with enhanced residual U-net and auxiliary classifier GAN. In: *2017 4th IAPR Asian Conference on Pattern Recognition (ACPR)*. IEEE, pp. 506–511.
- Zhang, L., Li, C., Wong, T.-T., Ji, Y., Liu, C., 2018c. Two-stage sketch colorization. In: *SIGGRAPH Asia 2018 Technical Papers*. ACM, p. 261.
- Zhang, R., Zhu, J.-Y., Isola, P., Geng, X., Lin, A.S., Yu, T., Efros, A.A., 2017b. Real-time user-guided image colorization with learned deep priors. arXiv preprint [arXiv:1705.02999](https://arxiv.org/abs/1705.02999).
- Zhu, X., Lafferty, J., Ghahramani, Z., 2003. Combining active learning and semi-supervised learning using gaussian fields and harmonic functions. In: *ICML 2003 Workshop on the Continuum from Labeled to Unlabeled Data in Machine Learning and Data Mining*, Vol. 3.
- Zou, C., Mo, H., Gao, C., Du, R., Fu, H., 2019. Language-based colorization of scene sketches. *ACM Trans. Graph.* 38 (6), [http://dx.doi.org/10.1145/3355089.3356561](https://doi.org/10.1145/3355089.3356561).