

國立臺灣大學電機資訊學院資訊工程學系

碩士論文

Department of Computer Science and Information Engineering

College of Electrical Engineering and Computer Science

National Taiwan University

Master Thesis

互動式視覺多媒體拼貼

Dynamic Media Assemblage: Summarizing and Presenting Visual  
Media Interactively



指導教授：陳炳宇 博士

陳維超 博士

Advisor: Chen Bing-Yu, Ph.D., Chen Wei-Chao, Ph.D.

中華民國 100 年 7 月

July, 2011

## **ACKNOWLEDGEMENT**

**First and foremost, I owe my deepest gratitude to my advisor Prof. Robin Bing-Yu Chen for the continuous support of my master study and research. He allowed me to join his team since I was a senior student in NTU and also supported me to attend SIGGRAPH2010 conference last summer. Without those experiences, it is impossible for me to develop the original idea of this thesis.**

**I especially want to thank my advisor Prof. Wei-Chao Chen's helping. Without his enthusiasm, his inspiration, his great efforts to explain things clearly and his helping in paper-writing period, this thesis could not be successfully accomplished.**

**I would also like to thank all teachers in CML with their good-teaching, accompany and encourage. Especially prof. Yung-Yu Chuang, he not only provided a lot of learning resources in his class but also give suggestion to this thesis.**

**Last but not the least, my deepest gratitude goes to my family for their unflagging love and support throughout my life. Besides, I am indebted to many of my classmates in CML with their encouragement, helping and discussion. They make my master career to be meaningful and interesting.**

**Chun-Yu Tsai**



## 中文摘要

隨著人們普遍的使用照相機與攝影機以及網路的普及化，人們自製和分享多媒體影音也更加的容易。由於每個人所擁有以及觀看的多媒體影音內容數量越來越多，種類也越來越複雜，如何有效率地觀看這些多媒體內容也成了一個需要解決的課題。在本論文中，我們基於兩個目的提出了一種以拼貼和互動的方式呈現多媒體影音內容的方法。第一個目的是，由於希望多個視訊多媒體能在同一個畫面呈現，我們希望能更有效的利用畫面的空間，此即每個視訊多媒體不重要的區域可被遮蓋住。第二，我們希望在視訊多媒體播放的同時，可以提供一些互動功能讓使用者在觀看的時候更有參與感。此拼貼和互動的方法分為三個步驟：

一、我們首先對單個影片做分鏡偵測 (shot detection)：由於我們希望畫面上每一個影片都以一個分鏡為單位，以利拼貼的空間利用和視覺效果，此部份我們採用顏色統計差異 (color histogram difference) 的方法來把每個影片分成多個分鏡。

二、找出每個分鏡中的顯著區域 (saliency region)：為了拼貼空間的有效利用，我們對每個影片做一些事前處理，希望保留影片中物體的移動並使重要內容盡量在分鏡的中心播放以利拼貼空間的使用效率。

三、設計互動式隨機拼貼的演算法，使得使用者在觀看多個影片拼貼結果的時候，能夠任意添加、刪除和移動這些影片分鏡。

關鍵字： 圖片、影片、拼貼、瀏覽、互動





## Abstract

In this thesis, we present a dynamic media assemblage method for summarizing and presenting visual media interactively. We propose an iterative packing algorithm to arrange the cropped visual media while considering their temporal-spatial salient regions to efficiently utilize the 2D canvas. We first divide longer videos into individual shots by detecting their shot boundaries. We then compute the temporal-spatial salient regions within each shot and use them to remove the apparent camera or object motion for more efficient packing. Our packing algorithm respects the salient regions and screen aspect ratio. Its interactive and iterative nature means that our method is particularly well-suited for interactive manipulations such as moving, insertion, and deletion of media files while composing media assemblages in real-time for summarization and presentation purposes.

**Terms:** Algorithms, Design, Performance

**Keywords:** video collage, video saliency, media assemblage, media browser, iterative packing



# Contents

<b>Chapter 1</b>	<b>Introduction</b>	<b>1</b>
<b>Chapter 2</b>	<b>Related Work</b>	<b>5</b>
2.1	Automatic Image Collage . . . . .	5
2.2	Video Summarization . . . . .	7
2.3	Our Contributions . . . . .	10
<b>Chapter 3</b>	<b>Algorithm</b>	<b>13</b>
3.1	Video Shot Detection . . . . .	15
3.2	Saliency Analysis . . . . .	15
3.3	Salient Volume Extraction . . . . .	17
3.4	Packing . . . . .	18
3.4.1	Packing Initialization . . . . .	20
3.4.2	Packing Optimization . . . . .	20
3.5	Media Assemblage . . . . .	23

<b>Chapter 4</b>	<b>Implementation and Results</b>	<b>25</b>
4.1	Comparisons . . . . .	31
4.2	Video Results . . . . .	33
<b>Chapter 5</b>	<b>Conclusion and Future Work</b>	<b>35</b>
<b>Bibliography</b>		<b>38</b>
<b>Appendix A</b>		<b>43</b>
A.1	Convex Hull Region Calculation . . . . .	43
A.2	Voronoi Calculation . . . . .	44



# List of Figures

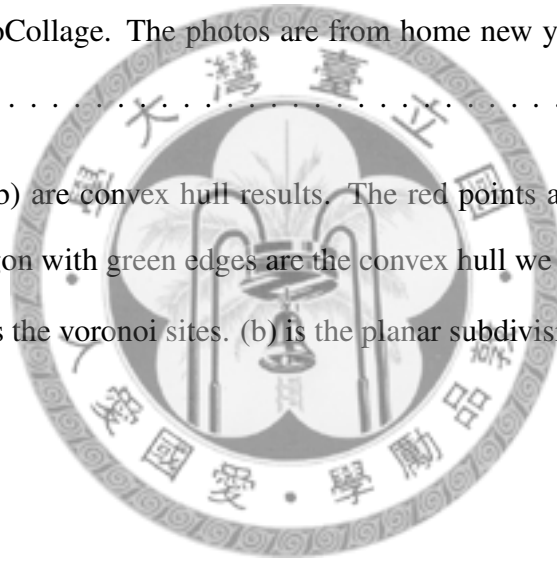
1.1	Example assemblages generated by our system, including a set of cartoon movie trailers (upper) and photos on cats (lower). These are generated automatically and can be interactively manipulated by the users. The cartoon movie trailers are from films <i>Ratatouille</i> (Pixar animation studios), <i>Shark Tale</i> (DreamWorks Animation), <i>Shrek</i> (DreamWorks pictures), <i>How to Train Your Dragon</i> (DreamWorks nimation), <i>Ice Age</i> (Blue Sky Studios) and <i>Toy Story</i> (Walt Disney Pictures). The cats colleition are home photos. . . . .	4
2.1	Some Automatic Image Collage examples listed according to the publish dates. This four works are all about image summaries. The first style, page layout, aims at maximize page coverage without having photo overlap. The second style, Picture Collage, imitates the collage style created by human and later Smart Photo Sticking adds semantic information to make more plausible results. The final one, AutoCollage, presents a seamless and visually appealing collage style. . . . .	6
2.2	Previous works for video presentation(image summaries methods) listed according to the publish dates. . . . .	9

3.1	The overview of our assemblage system. (a) The media collections with shots are detected. (b) The saliency analysis is performed on every frame. (c) The salient volumes are extracted in every shot. (d) These media are packed into the canvas. (e) The final assemblage is generated. . . . .	13
3.2	The examples of saliency analysis on the video frame. (a) The input video frame $i$ . (b) The input video frame $i + 1$ . (c) The original motion magnitude of frame $i$ . (d) The motion contrast magnitude of frame $i$ . (The dog video clip is licensed as Creative Commons.) . . . . .	16
3.3	Examples of our saliency analysis process. (a) The input video frames. (b) The motion contrast saliency $S_M$ . (c) The image saliency $S_I$ . (d) The face saliency $S_F$ . (e) The combined saliency $S$ . . . . .	16
3.4	An initial packing example. The $I$ -Region are added to canvas at the position that minimizes the empty space while respecting the aspect ratio of the canvas. If the placement orders are different, it will produce different initial packing results. . . . .	19
3.5	The results of packing initialization with different preferred aspect ratio: (a) 4:3; (b) 3:1. . . . .	21
3.6	Comparison of assemblage methods. The red polygons are $I$ -Regions. (a) A straightforward Voronoi segmentation. Notice some $I$ -Regions are eroded in this example. (b) Our approximated region-based Voronoi approach. . . . .	22

3.7	Different sample rates of the sites. (a)Only sample the I-Region's vertex. (b)Sample per 40 pixels length. (c)Sample per 20 pixels length. (d)Sample per 5 pixels length, the resulting edges are much more smooth. The photos are from New York trips photo collection. . . . .	24
4.1	Interacting with the assemblage. The user drags the blue region (a), causing several salient regions to become occluded (b). Our system iteratively refines the assemblage and resolves the problem in just a few iterations (c). . . . .	25
4.2	This figure shows the insertion processing of four cartoon movie trailers. In (a), there is just a video. (b),(c) and (d) are the reoptimization results after insertion. The cartoon movie trailers are from films Ratatouille(Pixar animation studios), Shark Tale(DreamWorks Animation), How to Train Your Dragon(DreamWorks nimation) and Ice Age(Blue Sky Studios). . . . .	26
4.3	An assemblage of animal videos. (These animal video clips are licensed as Creative Commons.) . . . . .	27
4.4	Assembly of a collection of movie trailers, including (a) one frame from a movie trailers, (b) its saliency map, and (c) the extracted volume and the <i>I-Region</i> (red polygon). (d) shows the final assemblage of the collection. The cartoon movie trailers are from films Ratatouille(Pixar animation studios), Shark Tale(DreamWorks Animation), Shrek(DreamWorks pictures), How to Train Your Dragon(DreamWorks nimation), Ice Age(Blue Sky Studios) and Toy Story(Walt Disney Pictures). . . . .	27



4.5	Summarization of a movie trailer(Harry Potter). Regions of videos in this assemblage are different video shots come from the same trailer.	30
4.6	A media file browser application. (a) Preview of a folder with travel photos and videos. After clicking on the bridge photo, the system brings up an assemblage of media files related to the trip to New York (b). The photos are from home foreign trips photo collections. . . . .	30
4.7	Compare media compilation to AutoCollage. Top row: Central Park. Bottom row: New Year's Eve 2011. Column (a): our results. Column (b): AutoCollage. The photos are from home new year photo collection. . . . .	32
A.1	(a) and (b) are convex hull results. The red points are point sets and the polygon with green edges are the convex hull we found. . . . .	45
A.2	(a) shows the voronoi sites. (b) is the planar subdivision result. . . . .	46



# List of Tables

4.1 Processing Time . . . . .	28
-------------------------------	----





# Chapter 1

## Introduction

The popularity of digital cameras and video recorders put the power of media creation at the hands of millions of people. On the one hand, we can now record our surroundings and share them almost instantly through the Internet. On the other hand, with the sheer amount of the contents we are exposed to nowadays, we now face an increasing burden to categorize, browse, or look at our own personal video and photo collections.

To reduce the effort of managing the ever-growing media collections, many commercial applications started to incorporate search interfaces to retrieve pre-tagged media files. Keyword tags, locations, or face identities are popular methods particularly for Internet applications such as Flickr where people wish to have their uploaded media files noticed and discovered [2]. However, the majority of personal media files remain untagged, and they can become forgotten in the presence of those that are heavily tagged. What we find lacking, therefore, is a new breed of image browser that allows

us to quickly rediscover what we already have within our own piles of untagged media store.

There have been various prior research that are useful for our purposes. For example, image collage research such as [21] allow us to generate compact and aesthetically pleasing views of image collections. Image retargeting techniques such as seam carving [22] focus on reducing the size of images without throwing away important contents. Similarly, video summarization methods focus on producing the gist of videos by either reducing their length or by summarizing each of them into one, or a collection, of images [4, 25].

Our goals are related to the research stated above, with several important distinctions. First, ours goal is to create a tool that allows a user to simultaneously display several media files at the same time. As such, the tool should not dictate what is displayed where. Instead, it ought to respond to user requests such as insertion, deletion, and rearrangement while using the screen real-estate effectively. Second, we wish to reduce the effort and time for a user to browse through video and image collections, and therefore we are not limited to a specific category of video summarization techniques – we can choose to break up a single video into multiple images or even video segments, and they can all be displayed on the same screen canvas simultaneously. Fig. 1.1 shows a few sample results generated by our system.

Our techniques provides a new style to present several videos and images on the same screen canvas. After developing the techniques, it not only can be a browser to help quickly rediscover what we already have within our own piles of untagged media store but also can be used at several purposes. For example, we can divide a single

video into several shots or just select several key frames and use this technique to do video summarization. Moreover, it provide interaction for user to edit the final result. Those application and comparison will show in detail later.





Figure 1.1: Example assemblages generated by our system, including a set of cartoon movie trailers (upper) and photos on cats (lower). These are generated automatically and can be interactively manipulated by the users. The cartoon movie trailers are from films *Ratatouille*(Pixar animation studios), *Shark Tale*(DreamWorks Animation), *Shrek*(DreamWorks pictures), *How to Train Your Dragon*(DreamWorks nimation), *Ice Age*(Blue Sky Studios) and *Toy Story*(Walt Disney Pictures). The cats collection are home photos.

## Chapter 2

### Related Work

#### 2.1 Automatic Image Collage

An image collage refers to an image created from an assemblage of a collection of images. A variety of automatic image collage techniques have been developed both for research and commercial purposes. Fig. 2.1 shows some examples listed according to the publish dates.

Google's Picasa<sup>1</sup>, for example, incorporates a feature that generates collages of complete input images. It also provides different composition styles to the users. Atkins [3] proposed an efficient method of organizing images in a page. The method attempts to maximize page coverage without having photo overlap and provides explicit control over the aspect ratios and relative areas of the photos. Wang *et al.* [27]

---

<sup>1</sup><http://picasa.google.com/>





Figure 2.1: Some Automatic Image Collage examples listed according to the publish dates. This four works are all about image summaries. The first style, page layout, aims at maximize page coverage without having photo overlap. The second style, Picture Collage, imitates the collage style created by human and later Smart Photo Sticking adds semantic information to make more plausible results. The final one, AutoCollage, presents a seamless and visually appealing collage style.

presented picture collage as a kind of visual image summary, which optimizes the layout of rectangular images to maximize the portion of visual visible information (salient regions) in the result. Battiato *et al.* [5] improved the result of picture collage both considering a self-adaptive image cropping algorithm, exploiting visual and semantic information, and introducing an optimization process based on a genetic algorithm. Rother *et al.* [21] presented AutoCollage which constructed a visually appealing collage from a collection of input images. The aim is that the resulting collage should be representative of the collection, summarising its main themes. It is also assembled largely seamlessly, using graph-cut, Poisson blending of alpha-masks, to hide the joins between input images. Goferman *et al.* [11] presented Puzzle-like Collage which is based on assembling regions of interest of arbitrary shape in a puzzle-like manner. They also show that this approach produces collages that are informative, compact, and eye-pleasing.

## 2.2 Video Summarization

Truong and Venkatesh [25] provided an excellent review of video summarization techniques, which are divided into two classes, namely video skims and still image summaries. Video skims generate a shorter summary video to summarize the whole video, while still image summaries extract a number of keyframes from a video to pack the summary image. For video skims, Christel *et al.* [8] presented studies that measure effectiveness of video skim techniques. Divakaran *et al.* [10] devised a method to adjust video framerates by analyzing temporal motion activity, and speed up parts of

the video with less activity. Peker and Divakaran [20] used motion activity as well as various semantic cues such as face, skin color, or speech to control the video playback rate. Simakov *et al.* [23] propose an optimization method based on bi-directional similarity measure to retarget (or summarization) of image/video data into smaller sizes. Bennett *et al.* [6] present methods for generating novel time-lapse videos that address the inherent sampling issues that arise with traditional photographic techniques.

The second class of techniques, image summaries, is related to image collage and similar to our applications in several respects. Some previous works are showed in Fig. 2.2.

Zhu *et al.* [30] proposed the video booklet system, which extracts a number of thumbnails from a video, and then reshaped by a set of predefined shape templates. Wang *et al.* [28] presented video collage, which first select the representative keyframes from video, and then blends the selected keyframes to produce seamless video summary. Unlike Video Collage in which both the shapes of ROI and final collage are fixed as rectangle, Yang *et al.* [29] and Mei *et al.* [19] both extended the seamless blending to generate arbitrary shape collages. They also design three ROI arrangement schemes (i.e., book, diagonal, and spiral ) for satisfying different video genres.

Correa and Ma [9] presented a method to interactively generate seamless video summaries and their result is the static 2D panorama background with moving object on it. Barnes *et al.* [4] proposed a method to automatically generate video tapestries in different level of details that allow for continuous panning and zooming. Chiu *et al.* [7] ,for small displays on mobile devices, presented a method for creating highly condensed video summaries called Stained-Glass visualizations and generated non-



Stained-Glass (2004)



Figure 3: Sample shape templates



Figure 4: More flexible shape templates

Video Booklet (2005)



(a) Video keyframes



(b) BS-Video Collage



(c) FS-Video Collage

Video Collage (2009)



Dynamic Video Narrative(2010)



Video Tapestries (2010)

Figure 2.2: Previous works for video presentation(image summaries methods) listed according to the publish dates.

triangular layout to effectively summarize the salient regions. Kang *et al.* [15] proposed the space-time video montage. The method simultaneously analyzes both the spatial and temporal information distribution in a video sequence, and extracts the visually informative space-time portions of the input videos and then incorporates several extracted video parts in one video volume and try to maximize total salience.

## 2.3 Our Contributions

Our goals are similar to automatic image collage in that we wish to preserve salient regions on the canvas. Additionally, we have to incorporate both images and videos in our input, and automatically generate a collage that we can efficiently edit and manipulate afterwards. This means we have to forgo some of the more time-consuming techniques such as graph-cut, Poisson blending and Markov chain optimizations used by several image collage research. In the end, ours is an incremental algorithm that iteratively computes locally optimal collage configurations. We also choose to play the videos as-is without skimming, and we do not consider this a drawback because we can readily transform our input videos with any suitable skimming techniques. Specifically, our contributions are as follows.

- We propose a novel method for analyzing temporal-spatial salient regions of a video,
- a method for extracting the temporal-spatial salient regions while removing apparent camera or object motions,

- an efficient, greedy technique for packing a collection of irregularly-shaped visual media, and
- a scheme to iteratively optimize the packing when it is disturbed.





## Chapter 3

### Algorithm

We design our media assemblage techniques based on a few visual guidelines. First, we wish to reduce visual complexities while navigating these media. This means the non-essential parts of a video can be covered up or eliminated in the assemblage. We also plan to support interactive operations such as addition, deletion, and rearrangement, and while the assemblage would change during these operations, its layout should stabilize and stay static soon after these operations are complete so as to minimize disturbances while playing individual videos within the assemblage. For

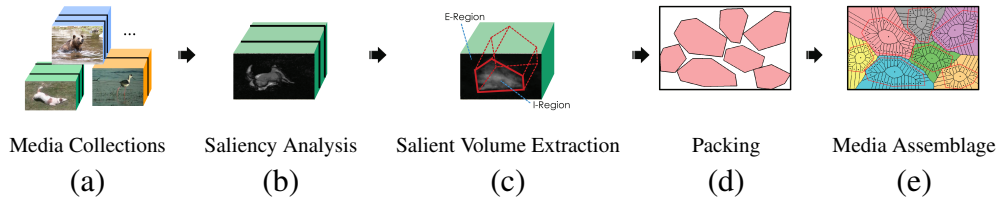


Figure 3.1: The overview of our assemblage system. (a) The media collections with shots are detected. (b) The saliency analysis is performed on every frame. (c) The salient volumes are extracted in every shot. (d) These media are packed into the canvas. (e) The final assemblage is generated.



this purpose, the essential, or salient, region of a video needs to have a static outline throughout its timeline. To ensure efficient extraction of the salient regions, we need to be aware of the camera motions for videos where the subject matters are moving.

Starting with a set of videos and photos  $\mathbf{M}$ , our system computes a configuration  $\mathbf{X} = \{p_i, R_i\}_{i=1}^{|\mathbf{M}|}$ , where  $p_i$  is the position and  $R_i$  is the high salient region of the  $i$ -th element of  $\mathbf{M}$ . Our goal is to find the optimal configuration  $\mathbf{X}^*$  subject to a packing energy  $E$ ,

$$\mathbf{X}^* = \arg \min_{\mathbf{X}} E(\mathbf{X}). \quad (3.1)$$

Fig. 3.1 illustrates the steps of our algorithm. We treat photos as static videos, and perform the process in two stages, namely *media analysis* (Fig. 3.1(a-c)) and *media packing* (Fig. 3.1(d-e)). The *media analysis* stage is a preprocessing step designed to discover informative regions within each individual videos. Given an input video, we use color histogram to split it into a number of video shots. Then, for each individual shot  $i$ , we compute temporal-spatial saliency of its frames (Fig. 3.1(b)) and extract the region of interest  $R_i$  by considering the saliency distribution within the shot (Fig. 3.1(c)). The *media packing* stage combines all the input media and efficiently packs them together. We first compute an initial packing by greedily minimizing the blank region on the canvas, followed by an iterative process that adjusts the packing configuration  $\mathbf{X}$  according to the packing energy  $E$  (Fig. 3.1(d)). Finally, the system decides what regions are visible for every media and generates a final assemblage (Fig. 3.1(e)). In the rest of this section, each step of the algorithm will be described in more details.

### 3.1 Video Shot Detection

As stated before, we would like the salient boundaries of each individual element in the final assemblage to stay fixed while playing back videos. For this purpose, we want each element in  $\mathbf{M}$  to be as coherent as possible, and ideally each element should consist of only one single shot. Since how to detect perfect shots is not our main concern, we find solution from previous shot boundary works. There are many possible methods for shot boundary detection and they shows that different kinds of video will need different kinds of algorithms to produce suitable result. After experimenting with several possible methods, we found the method by Lienhart [16] to be effective for our video sample sets. This technique measures color histogram differences between two adjacent frames and declares a shot boundary when a large color discontinuity occurs.

### 3.2 Saliency Analysis

Psychological studies show that visual signals contrast such as motion and color are likely to attract people's visual attentions [24]. We adopt similar visual attention formulation by Liu *et al.* [17] to compute a saliency map per frame per video, and emphasize the salient regions in the final assemblage in order to utilize the 2D canvas more efficiently. In this method, the saliency of each pixel  $p$  is calculated as a weighted sum of the motion contrast saliency ( $S_M$ ), the image saliency ( $S_I$ ) and the face saliency ( $S_F$ ), as follows

$$S(p) = w_M S_M(p) + w_I S_I(p) + w_F S_F(p). \quad (3.2)$$

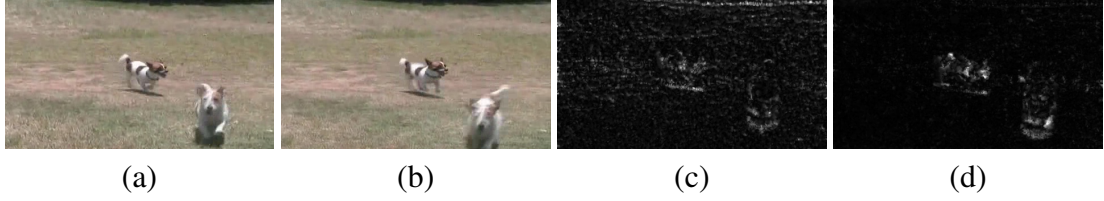


Figure 3.2: The examples of saliency analysis on the video frame. (a) The input video frame  $i$ . (b) The input video frame  $i + 1$ . (c) The original motion magnitude of frame  $i$ . (d) The motion contrast magnitude of frame  $i$ . (The dog video clip is licensed as Creative Commons.)

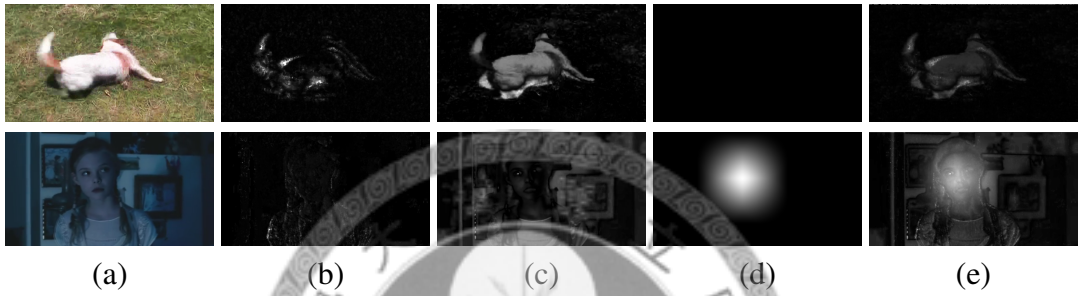


Figure 3.3: Examples of our saliency analysis process. (a) The input video frames. (b) The motion contrast saliency  $S_M$ . (c) The image saliency  $S_I$ . (d) The face saliency  $S_F$ . (e) The combined saliency  $S$ .

We use  $w_M = w_I = w_F = 1/3$  in our implementation. Our method differs in that we use a simple panning motion model to preserve the styles of the original shot, and we adopt a different image saliency measure. An example of this process is shown in Fig. 3.3.

**Motion Contrast Saliency.** Moving objects should be assigned higher saliency values because humans are particularly good at perceiving them. Therefore, given the observed motion (the optical flow that assigns a motion vector to each pixel) with an image, we would expect that regions with motion are likely to be salient. However, the simple approach of only considering the amount of motion at each pixel is insufficient because people are good at factoring out global motion, such as that induced by head

or camera movement. This behavior is encoded as motion contrast saliency as shown in Fig. 3.2. First, we use Lucas-Kanade method [18] to analyze the relative motion between two adjacent frames, and then approximate a global camera motion by using a voting scheme where the motion vectors are used to vote both on a consensus motion direction and magnitude. The motion contrast is then obtained by subtracting the original motion vector with the global camera motion and normalized to  $0 \sim 1$  into motion contrast saliency.

**Image Saliency.** There exist various methods to measure image saliency based on low-level feature contrast [14, 13, 1, 12]. We choose the approach by Achanta *et al.* [1] which calculates the saliency of each pixel based on its color and luminance differences with respect to its neighbors and outputs full resolution saliency maps with well-dened boundaries of salient objects. Fig. 3.3(c) shows the image saliency results.

**Face Saliency.** To emphasize the importance of human faces, we detect the presence of faces with the methods proposed by Viola [26]. The face saliency value is then calculated by applying a Gaussian attenuation function surrounding the area of the detected faces. Fig. 3.3(d) shows the image saliency results.

### 3.3 Salient Volume Extraction

Given the saliency map of each individual frame, our next goal is to extract volumes of the video for the following packing stages. For a video shot where a single salient object moves at a constant speed, we may compute the cumulated saliency  $S_c$  along a

temporal skew  $(x, y)$ ,

$$S_c(i, j) = \sum_f S_f(i + fx, j + fy), \quad (3.3)$$

where  $S_f(i, j)$  is the saliency value of the  $f$ -th frame at pixel index  $(i, j)$ . An optimal direction  $(x, y)$  is where the cumulated saliency values are concentrated in a region as small as possible. To determine this direction, for each frame we iteratively select the pixels with highest saliency values until the sum of these values exceeds half of the total saliency values within this frame, and construct a bounding box using the selected pixels. Then, we fit a least-square line over the centers of the bounding boxes over time, and use the line direction as the optimal direction to align the video volume, accumulate the saliency values, and determine the region that should be preserved in the final assemblage.

To calculate this region, we again select those pixels with highest cumulated saliency value until the sum of these values reaches a predefined threshold of say, in our case, 50% of the sum of cumulated saliency from all the pixels. We then construct this *I-Region*  $R_i$  of the selected pixels. The region outside the *I-Region* is defined as the external region *E-Region*, whose pixels can be discarded when a tighter packing is desirable.

### 3.4 Packing

After extracting the salient volumes, we pack the media set  $\mathbf{M}$  by following a few criteria. First, we wish to use the canvas space efficiently. Second, salient regions of the media should never be occluded. Third, the canvas should observe the aspect

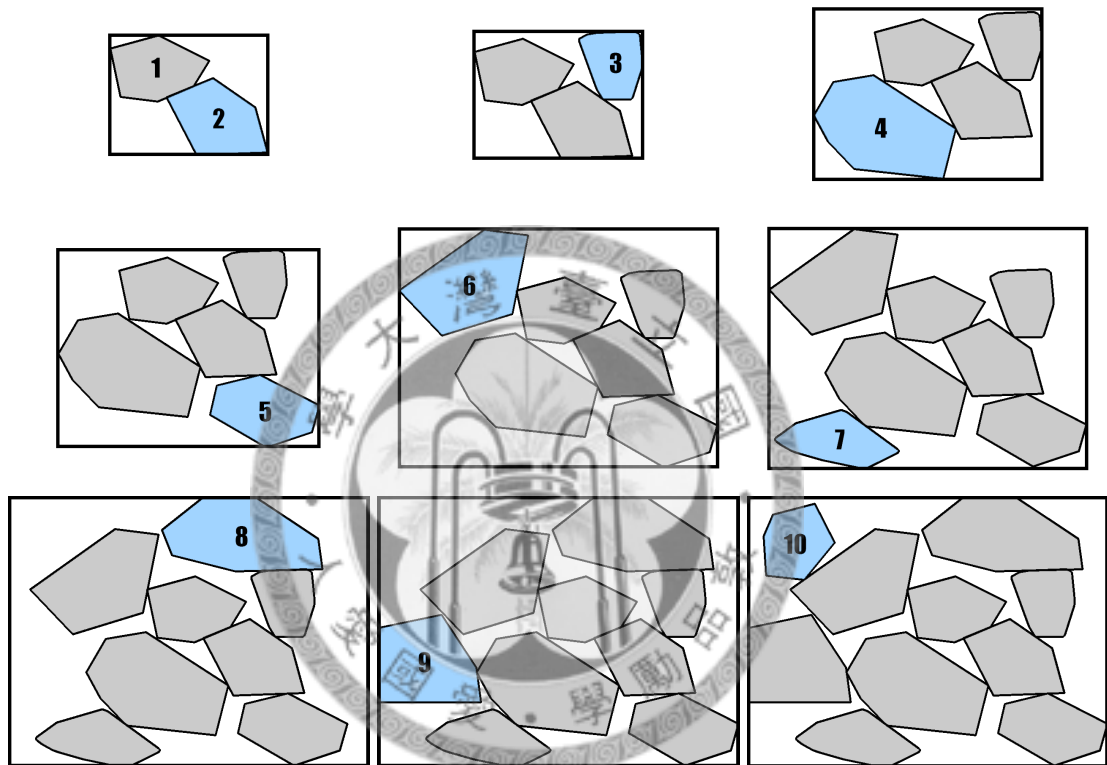


Figure 3.4: An initial packing example. The *I-Region* are added to canvas at the position that minimizes the empty space while respecting the aspect ratio of the canvas. If the placement orders are different, it will produce different initial packing results.

ratios of the display devices. With these goals in mind, our objective is to find an optimal configuration  $\mathbf{X}^*$  without occluding each of the *I-Regions* in  $\mathbf{M}$ . This packing problem, unfortunately, is a NP-complete problem, and we propose to approximate the optimal solution by a two-stage heuristic. First, we use a greedy algorithm to initialize a layout configuration. Then, this configuration is iteratively optimized to reach a local minimum.

### 3.4.1 Packing Initialization

As shown in Fig. 3.4, the packing initialization process is as follows. The first *I-Region* is first placed at the center of the canvas. Then, we place each remaining *I-Region*  $R_i$  (in arbitrary order) radially around the canvas center while ensuring no overlap between  $R_i$  and all other *I-Regions* already on the canvas. We pick an optimal direction that minimizes the empty space while respecting the aspect ratio of the canvas. Because of our packing criteria, the initial packing result can have predefined aspect ratio. Fig. 3.5 shows two examples of the packing initialization under two different pre-selected aspect ratios of 4:3 and 3:1, respectively.

### 3.4.2 Packing Optimization

After initialization, we iteratively optimize for the configuration  $\mathbf{X}^*$  by randomly selecting an *I-Region* and moving it toward a unit direction that reduces the packing energy  $E$  by the greatest amount. Each step of this process is guaranteed to reduce the packing energy, and we repeat the process until it stabilizes to a local minima. Based

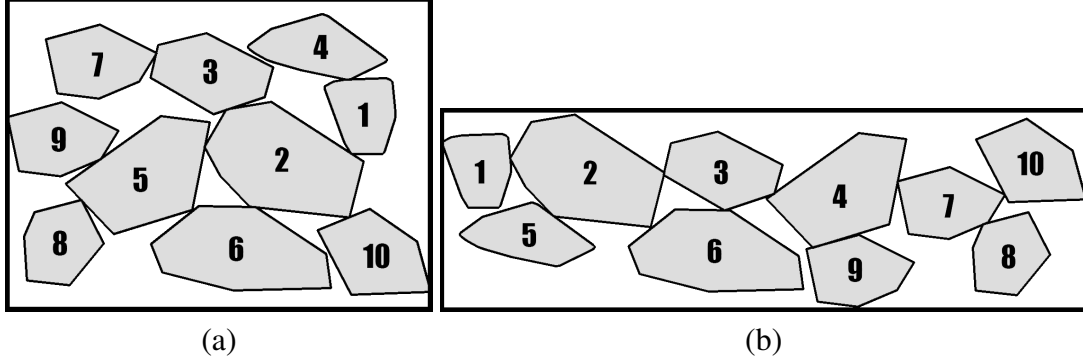


Figure 3.5: The results of packing initialization with different preferred aspect ratio: (a) 4:3; (b) 3:1.

on the packing criteria described before, we design our energy function based on a combination of penalty measures on empty space, *I-Region* occlusion, and aspect ratio deviation, as follows

$$E = E_{es}^{\alpha} + kE_{occ} \quad (3.4)$$

where the occlusion weight  $k$  is set to  $1 \times 10^5$  in our implementation to ensure that the *I-Regions* never get occluded. We now describe each of the energy terms ( $E_{es}$ ,  $\alpha$ , and  $E_{occ}$ ) in the following paragraphs.

**Empty Space Penalty.** We define the empty spaces on the canvas as regions that are not covered by any *I-Region*. This energy term represents the percentage of empty spaces on the canvas, as follows

$$E_{es} = \frac{Area(R_B - \cup_i R_i)}{Area(R_B)}, \quad (3.5)$$

where  $R_B$  is the bounding box formed by all *I-Regions*  $R_i$ .

***I-Region* Occlusion Penalty.** This energy term penalizes coverage of salient video



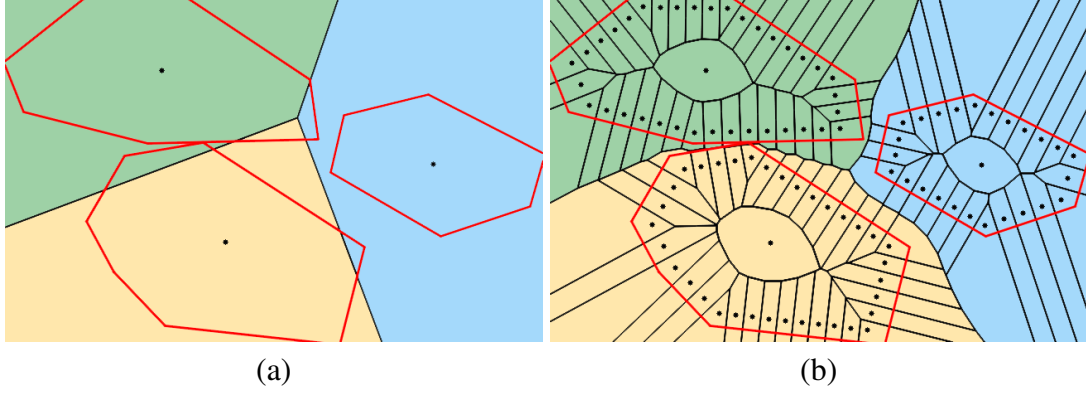


Figure 3.6: Comparison of assemblage methods. The red polygons are *I-Regions*. (a) A straightforward Voronoi segmentation. Notice some *I-Regions* are eroded in this example. (b) Our approximated region-based Voronoi approach.

regions, and is simply defined as the total areas of covered *I-Regions*,

$$E_{occ} = \sum_i Area(R_i) - Area(\cup_i R_i). \quad (3.6)$$

**Aspect Ratio Deviation Penalty.** The optimal packing should respect an aspect ratio specified by the user. This can be described by the following term

$$\alpha = \frac{1}{(q_c - q_p)^2 + \varepsilon} \quad (3.7)$$

where  $q_c$  is the aspect ratio of the bounding box  $R_B$ ,  $q_p$  is the desired aspect ratio, and we use a small number  $\varepsilon = 10^{-6}$  to set an upper bound for  $\alpha$ . Since the magnitude of the empty space penalty  $E_{es}$  is always less than 1, the first term in Eq. 3.4 becomes very small once we approach the desired aspect ratio.

### 3.5 Media Assemblage

Now that we have a layout configuration of the media set, we may begin generating an assemblage on the canvas. The choice of rendering styles is a rather artistic one. For example, AutoCollage adopts a seamless blending style between adjacent images [21]. For our purposes, we need to clearly distinct media boundaries for playback and interaction purposes. Simply rendering all *I-Regions* may suffice, but we would like to equally allocate the empty spaces to its surrounding videos and fill up as much empty spaces as possible by rendering non-essential *E-Regions*.

A straightforward approach is to segment the canvas using complete *I-Regions* as Voronoi sites. As this method proved to be too slow to run at interactive rate, we approximate this algorithm by sampling a number of Voronoi sites along the boundaries of the *I-Regions*. This approach, in addition to its speed advantages, has an additional benefit where we can control the smoothness of Voronoi region boundaries by changing the sampling rate of the sites. Fig. 3.6 shows the assemblage using our approach (b) compared to a simple Voronoi segmentation (a). Notice that the sites are sampled a short distance away from the *I-Region* boundaries to prevent nearby regions from eroding into each other. Fig. 3.7 shows some results of different sample rates.

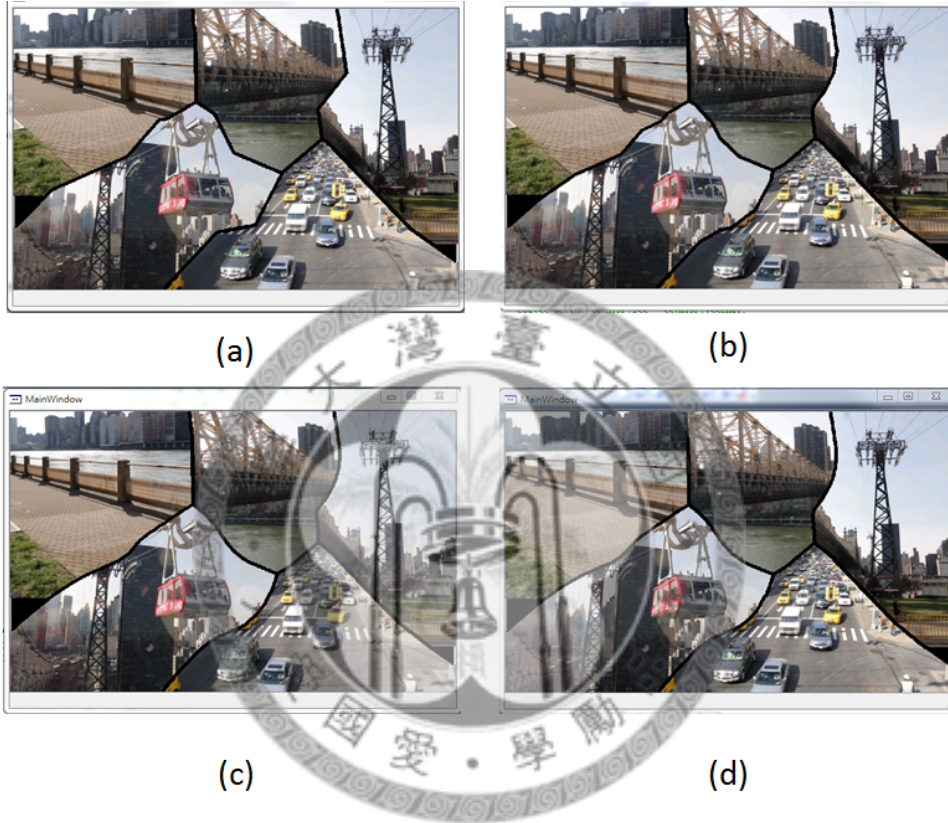


Figure 3.7: Different sample rates of the sites. (a)Only sample the I-Region's vertex. (b)Sample per 40 pixels length. (c)Sample per 20 pixels length. (d)Sample per 5 pixels length, the resulting edges are much more smooth. The photos are from New York trips photo collection.

## Chapter 4

### Implementation and Results

In this section we present our results as well as several applications. Fig. 4.4 shows the intermediate results of summarizing a collection of movie trailers, including saliency maps (b), *I-Region* (c), and final media assemblage (d).

Our implementation runs on a 3.20GHz desktop PC with 4.00GB of memory and the running time and space used by our algorithm varies according to the amounts of

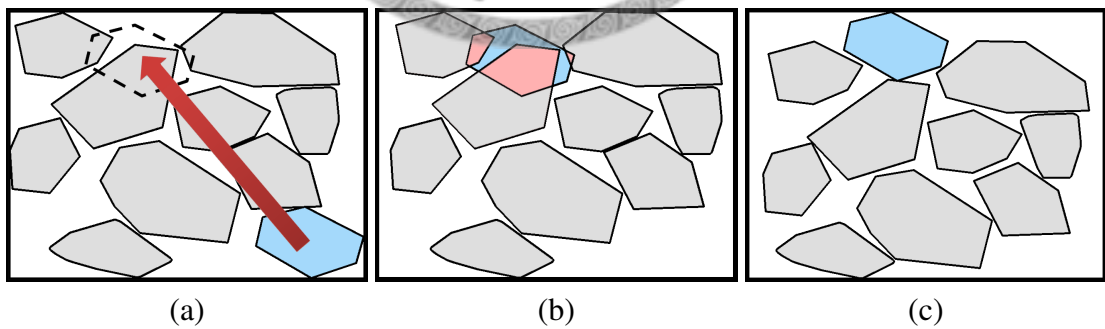


Figure 4.1: Interacting with the assemblage. The user drags the blue region (a), causing several salient regions to become occluded (b). Our system iteratively refines the assemblage and resolves the problem in just a few iterations (c).

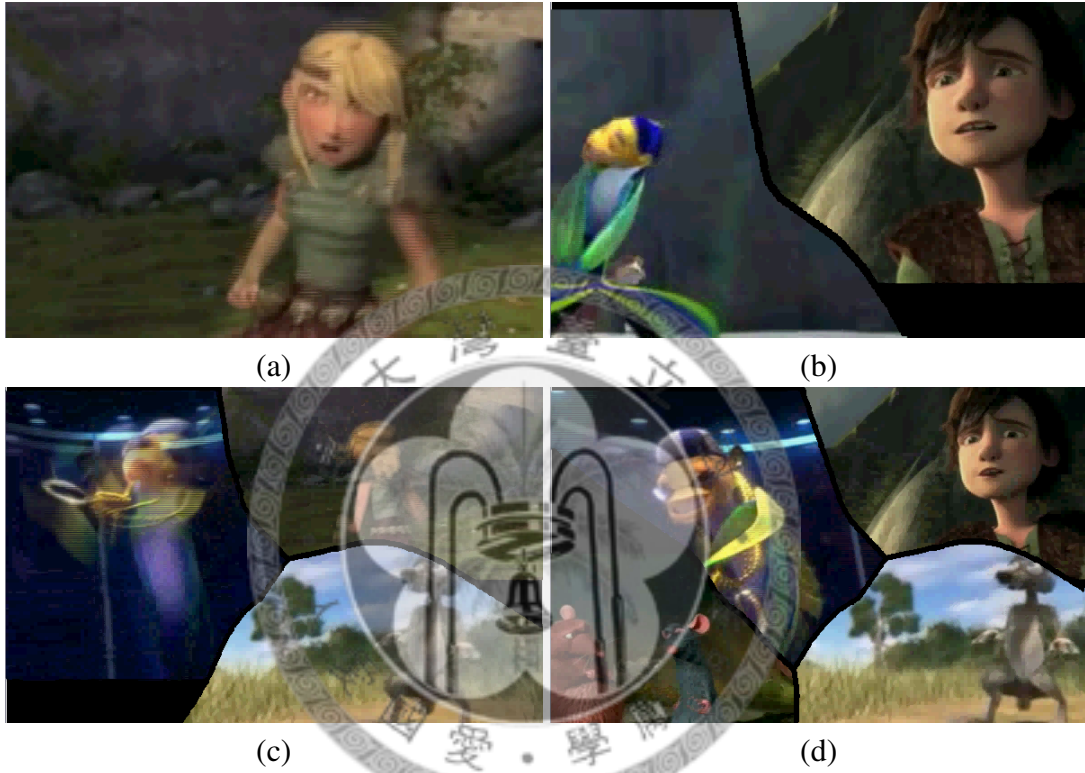


Figure 4.2: This figure shows the insertion processing of four cartoon movie trailers. In (a), there is just a video. (b),(c) and (d) are the reoptimization results after insertion. The cartoon movie trailers are from films Ratatouille(Pixar animation studios), Shark Tale(DreamWorks Animation), How to Train Your Dragon(DreamWorks animation) and Ice Age(Blue Sky Studios).





Figure 4.3: An assemblage of animal videos. (These animal video clips are licensed as Creative Commons.)

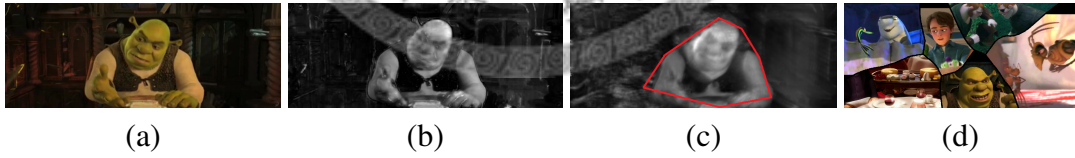


Figure 4.4: Assembly of a collection of movie trailers, including (a) one frame from a movie trailers, (b) its saliency map, and (c) the extracted volume and the *I-Region* (red polygon). (d) shows the final assemblage of the collection. The cartoon movie trailers are from films *Ratatouille*(Pixar animation studios), *Shark Tale*(DreamWorks Animation), *Shrek*(DreamWorks pictures), *How to Train Your Dragon*(DreamWorks nimation), *Ice Age*(Blue Sky Studios) and *Toy Story*(Walt Disney Pictures).

images/videos. The preprocessing, saliency map computing and video volume extraction, for a short shot(about 300 frames) takes roughly 1.5 minutes to compute. After the preprocessings finished offline, the initial packing needs another 2 seconds to compute the initial result showed on the system UI.

Our system creates interactive assemblages that users can interact with via insertion, deletion, and dragging operations, and our system automatically adjust the layout to the next optimal configuration after these operations. Fig. 4.2 and Fig. 4.1 shows examples of insertion, dragging and dropping of a media file, and our system's response to refine the configuration and Table 4.1 shows time consumed by our algorithm.

In addition to these manipulation operations, our system can be used for a variety of applications such as video collection summarization, single video summarization, personal media folder visualization, and interactive video board. We now discuss a few of the possibilities as follows.

**Media Collection Presentation** Previous works like image collages and slideshow can only present static image collections. For video collection presentation, video-sharing websites usually use a single frame to represent a whole video clip. Our system

Table 4.1: Processing Time

Step	Cartoon	New Year
Num of Medias	6	8
Initial Packing(secs)	0.071	0.134
Optimization util Stabilized(iterations/secs)	61/3.605	107/3.689

Note: The table shows time consumed by our algorithm and these experiments are generated on a desktop PC with an Intel Core 2 Duo 3.2GHz CPU and 4GB RAM. According to the table, we can know Initial Packing time is positive correlative to num of medias. From the optimization time, we also can show our algorithm is real-time.

can provide a presentation of a set of visual media collections, including both videos and images, and dynamically play the videos. A user can preview all these videos simultaneously on a single screen and then modify the presentation as she sees fit. If the video shots are of un-equal duration, they are played in loop in our current design. However, this is just one of many possible presentations of the assemblage and is not the major concern of our core algorithm. A possible alternative is, for example, for the system (or the users) to remove finished shots from the canvas and insert new shots whenever desired. This can be realized without modification to our core real-time packing algorithm. Fig. 4.4(d) and Fig. 4.3 are a few example of video collection summarization.

**Single Video Summarization.** Our system can summarize a single video by assembling each individual shot onto a canvas. Unlike traditional video summarizations, which select some key frames and summarize them with still image collage, our approach can play all shots simultaneously, or sequentially with time overlaps, and the users can get a quick temporal review through this dynamic summarization. In current implementation, we fix the shots amount in the canvas. When shot change, we simply substitute the old shot with the new shot in the same position. After substitution, the system will real-time reoptimize the shot layout. Fig. 4.5 shows an example of summarizing a video.

**Personal Media File Browser.** Fig. 4.6 shows an example of visualizing a collection of media related to traveling. As a user browses through this collection in the root folder, she can choose to preview representative photos and videos from different trips





Figure 4.5: Summarization of a movie trailer(Harry Potter). Regions of videos in this assemblage are different video shots come from the same trailer.

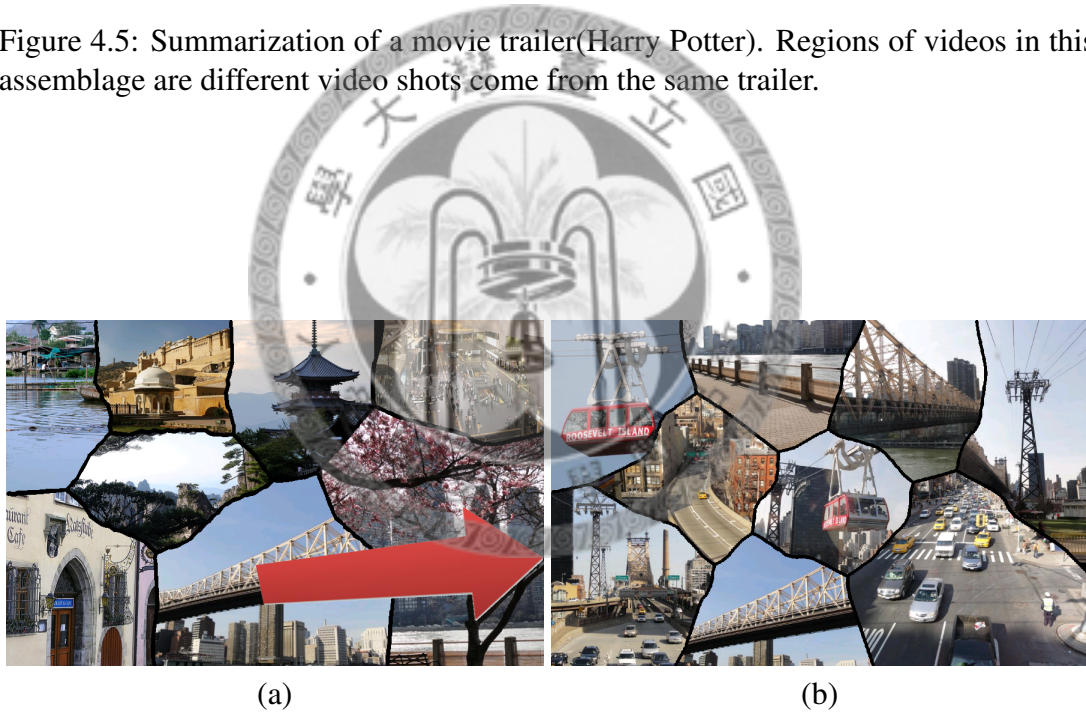


Figure 4.6: A media file browser application. (a) Preview of a folder with travel photos and videos. After clicking on the bridge photo, the system brings up an assemblage of media files related to the trip to New York (b). The photos are from home foreign trips photo collections.

(a). She can click on one of the interest region on the assemblage which brings up media files from the sub-folder, presented as another assemblage (b).

## 4.1 Comparisons

The main difference between our approach and the prior work is that all of the previous methods are all offline solutions, whereas ours is the first real-time interactive collage system. Moreover, our packing algorithm is designed to support real-time media insertion, deletion, or relocation operations, whereas previously proposed methods do not have related mechanisms to real-time optimize for the media layout while respecting user intentions.

Here we present a qualitative comparison between AutoCollage [21] and our approach. AutoCollage operates on images, and we generate both of the results using the same photo sets. Fig. 4.7 shows the comparisons on two different photo sets. The most obvious difference is an aesthetic choice where AutoCollage generates seamless collages while ours have clear boundaries around the photos. Both AutoCollage and our system are completely automatic. However, users can interactively adjust and refine the assemblage generated by our system, and this is not possible with AutoCollage.

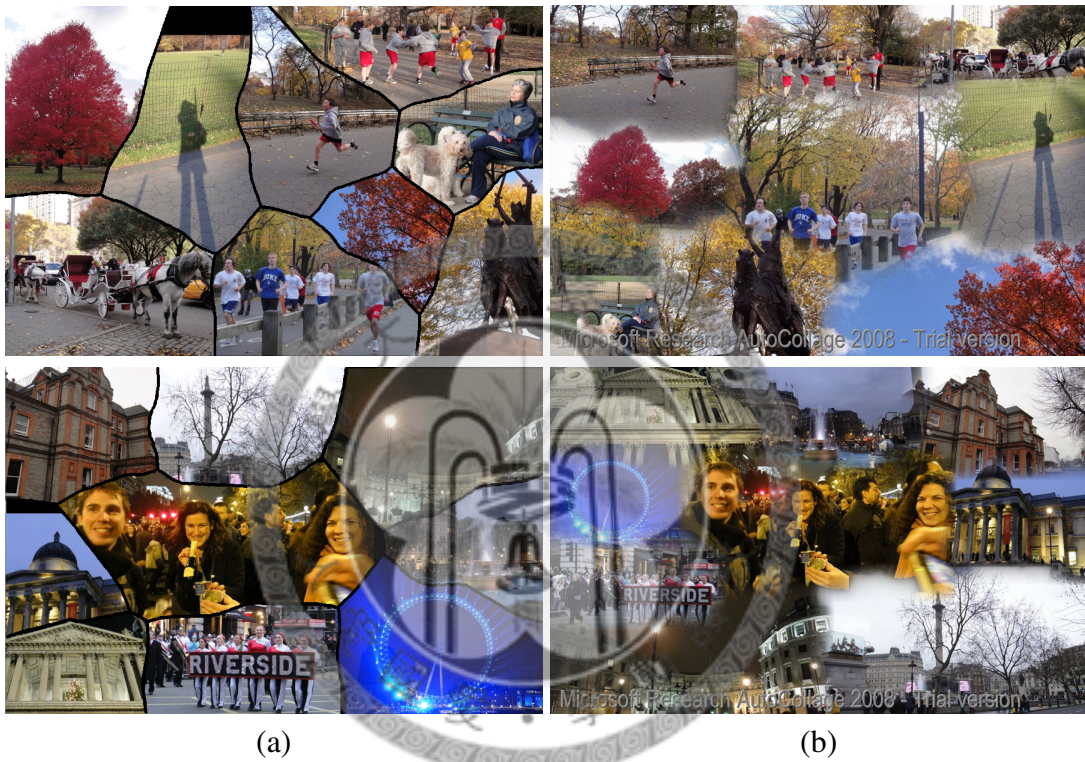


Figure 4.7: Compare media compilation to AutoCollage. Top row: Central Park. Bottom row: New Year's Eve 2011. Column (a): our results. Column (b): AutoCollage. The photos are from home new year photo collection.

## 4.2 Video Results

Readers are strongly recommended to watch our video results via the following anonymous link:


<http://www.youtube.com/watch?v=9tzazErILCE>





## Chapter 5

### Conclusion and Future Work



In this paper, we have presented a dynamic media assemblage method for summarizing and presenting visual media interactively. We analyze the temporal-spatial salient regions within each shot for more efficient packing. Our energy function and iterative optimization process guarantees occlusion-free packing of salient media regions while ensuring appropriate canvas aspect ratio. We also showed that our method can be applied to many applications, such as image and video collection presentation, single video summarization, and hierarchical media browser.

In our current implementation, the packing algorithm does not respect any user-specified order, and we are working on packing algorithms that take the order into consideration. Furthermore, our algorithm, while being fairly interactive, may get stuck on local minima, and we would like to see a better re-initialization scheme when this happens. Finally, unlike a traditional file browser, a media assemblage is inherently limited by the size of its canvas, and therefore we plan to introduce methods that allow

panning and scrolling as well as smart hierarchical layout within the assemblage. With this, it becomes possible to jump seamlessly from file browsers, media previewer and media assemblage browsers, and thus endowing users with more choices of managing their ever-growing media collections.



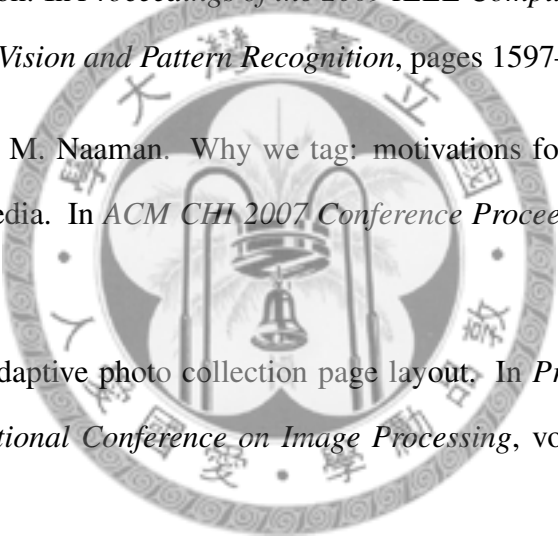
## Acknowledgements

The cartoon movie trailers are films Ratatouille (Pixar animation studios), Shark Tale (DreamWorks Animation), Shrek (DreamWorks pictures), How to Train Your Dragon (DreamWorks animation), Ice Age (Blue Sky Studios) and Toy Story (Walt Disney Pictures). The all animal video clips are licensed as Creative Commons. The photos, new year, new york and cats, are homemade collections.





# Bibliography

- 
- [1] R. Achanta, S. Hemami, F. Estrada, and S. Susstrunk. Frequency-tuned salient region detection. In *Proceedings of the 2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 1597–1604, 2009.
- [2] M. Ames and M. Naaman. Why we tag: motivations for annotation in mobile and online media. In *ACM CHI 2007 Conference Proceedings*, pages 971–980, 2007.
- [3] B. Atkins. Adaptive photo collection page layout. In *Proceedings of the 2004 IEEE International Conference on Image Processing*, volume 5, pages 2897 – 2900, 2004.
- [4] C. Barnes, D. B. Goldman, E. Shechtman, and A. Finkelstein. Video tapestries with continuous temporal zoom. *ACM Transactions on Graphics*, 29(4):89:1–89:9, 2010. (SIGGRAPH 2010 Conference Proceedings).
- [5] S. Battiato, G. Ciocca, F. Gasparini, G. Puglisi, and R. Schettini. Smart photo sticking. In *Proceedings of the 5th International Workshop on Adaptive Multimedia Retrieval*, pages 211–223, 2007.

- [6] E. P. Bennett and L. McMillan. Computational time-lapse video. In *ACM SIGGRAPH 2007 papers*, SIGGRAPH '07, New York, NY, USA, 2007. ACM.
- [7] P. Chiu, A. Girgensohn, and Q. Liu. Stained-glass visualization for highly condensed video summaries. In *Proceedings of the 2004 IEEE International Conference on Multimedia and Expo*, volume 3, pages 2059 – 2062, 2004.
- [8] M. G. Christel, M. A. Smith, C. R. Taylor, and D. B. Winkler. Evolving video skims into useful multimedia abstractions. In *ACM CHI 1998 Conference Proceedings*, pages 171–178, 1998.
- [9] C. D. Correa and K.-L. Ma. Dynamic video narratives. *ACM Transactions on Graphics*, 29(4):88:1–88:9, 2010. (SIGGRAPH 2010 Conference Proceedings).
- [10] A. Divakaran, K. A. Peker, and H. Sun. Constant pace skimming and temporal sub-sampling of video using motion activity. In *Proceedings of the 2001 IEEE International Conference on Image Processing*, volume 3, pages 414–417, 2001.
- [11] S. Goferman, A. Tal, and L. Zelnik-Manor. Puzzle-like collage. *Computer Graphics Forum*, 29:459–468, 2010.
- [12] S. Goferman, L. Zelnik-Manor, and A. Tal. Context-aware saliency detection. In *Proceedings of the 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 2376–2383, 2010.
- [13] X. Hou and L. Zhang. Saliency detection: A spectral residual approach. In *Proceedings of the 2007 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 1–8, 2007.

- [14] L. Itti, C. Koch, and E. Niebur. A model of saliency-based visual attention for rapid scene analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20:1254–1259, 1998.
- [15] H.-W. Kang, X.-Q. Chen, Y. Matsushita, and X. Tang. Space-time video montage. In *Proceedings of the 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 2, pages 1331–1338, 2006.
- [16] R. Lienhart. Comparison of automatic shot boundary detection algorithms. In *Proceedings of SPIE Storage and Retrieval for Image and Video Databases VII*, volume 3656, pages 290–301, 1999.
- [17] F. Liu and M. Gleicher. Video retargeting: automating pan and scan. In *ACM Multimedia 2006 Conference Proceedings*, pages 241–250, 2006.
- [18] B. D. Lucas and T. Kanade. An iterative image registration technique with an application to stereo vision. In *Proceedings of the 7th international joint conference on Artificial intelligence*, volume 2, pages 674–679, 1981.
- [19] T. Mei, B. Yang, S.-Q. Yang, and X.-S. Hua. Video collage: presenting a video sequence using a single image. *The Visual Computer*, 25:39–51, 2009.
- [20] K. A. Peker and A. Divakaran. An extended framework for adaptive playback-based video summarization. In *Proceedings of SPIE Internet Multimedia Management Systems IV*, volume 5242, pages 26–33, 2003.
- [21] C. Rother, L. Bordeaux, Y. Hamadi, and A. Blake. AutoCollage. *ACM Trans-*

- actions on Graphics*, 25(3):847–852, 2006. (SIGGRAPH 2006 Conference Proceedings).
- [22] A. Shamir and S. Avidan. Seam carving for media retargeting. *Communications of the ACM*, 52:77–85, 2009.
- [23] D. Simakov, Y. Caspi, E. Shechtman, and M. Irani. Summarizing visual data using bidirectional similarity. In *CVPR*. IEEE Computer Society, 2008.
- [24] A. M. Treisman and G. Gelade. A feature-integration theory of attention. *Cognitive Psychology*, 12(1):97–136, 1980.
- [25] B. T. Truong and S. Venkatesh. Video abstraction: A systematic review and classification. *ACM Transactions on Multimedia Computing, Communications, and Applications*, 3:3:1–3:37, 2007.
- [26] P. Viola and M. Jones. Rapid object detection using a boosted cascade of simple features. In *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 1, pages 511–518, 2001.
- [27] J. Wang, L. Quan, J. Sun, X. Tang, and H.-Y. Shum. Picture collage. In *Proceedings of the 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 1, pages 347–354, 2006.
- [28] T. Wang, T. Mei, X.-S. Hua, X.-L. Liu, and H.-Q. Zhou. Video collage: A novel presentation of video sequence. In *Proceedings of the 2007 IEEE International Conference on Multimedia and Expo*, pages 1479–1482, 2007.

- [29] B. Yang, T. Mei, L.-F. Sun, S.-Q. Yang, and X.-S. Hua. Free-shaped video collage. In *Proceedings of the 14th International Multimedia Modeling Conference*, pages 175–185, 2008.
- [30] C.-Z. Zhu, X.-S. Hua, T. Mei, and X.-Q. Wu. Video booklet: a natural video searching and browsing interface. In *Proceedings of the 7th ACM SIGMM International Workshop on Multimedia Information Retrieval*, pages 113–120, 2005.



# Appendix A

## A.1 Convex Hull Region Calculation

In order to construct the *I-Region* mentioned in 3.3, we need to calculate the convex hull region. There are five common algorithms<sup>1</sup>:

**Incremental.** The incremental algorithm is an algorithm for computing the convex hull of a set of points in two or more dimensions. The basic idea is to add points one at a time updating the hull as we proceed.

**Gift wrap(Jarvis' march).** Start at some extreme point, which is guaranteed to be on the hull. At each step, test each of the points, and find the one which makes the largest right-hand turn. That point has to be the next one on the hull.

**Graham's Scan.** First, we need to find an extreme point with the largest y coordinate as pivot on hull. Sort the points in order of increasing angle about the pivot. We end

---

<sup>1</sup><http://sls.weco.net/blog/josh68/07-jun-2007/2060>

up with a star-shaped polygon. To build the hull, by marching around the star-shaped poly, adding edges when we make a left turn, and back-tracking when we make a right turn.

**Quickhull.** We need to first find good chord to start the algorithm goes from the leftmost to the rightmost point in the set and then do partition iteratively until find the convex hull.

**Divide and Conquer.** Recursively divide the points into two equal sized sets and find convex hull of each smaller sets, and merge the final result.

The method here we use is proposed by Sklansky<sup>2</sup>. It's the first  $O(n)$  algorithm. First, find an convex vertex and label it  $p_0$ . Second, Label the remaining  $n-1$  vertices in a clockwise order, starting at  $p_0$ . Third, Place three coins on vertices  $p_0$ ,  $p_1$ ,  $p_2$  and label them "back", "center", and "front" respectively. Finally, through the remove and relabel process to decide the convex hull region. Fig. A.1 shows the calculation results.

## A.2 Voronoi Calculation

In 3.6, we use voronoi algorithm to equally allocate the empty spaces to its surrounding videos. The problem states as follow, given a set of points  $S$ (Voronoi sites) in the plane, do the planar subdivision<sup>3</sup>Planar subdivision is the subdivision of a plane

---

<sup>2</sup><http://cgm.cs.mcgill.ca/beezer/cs507/3coins.html>

<sup>3</sup>(

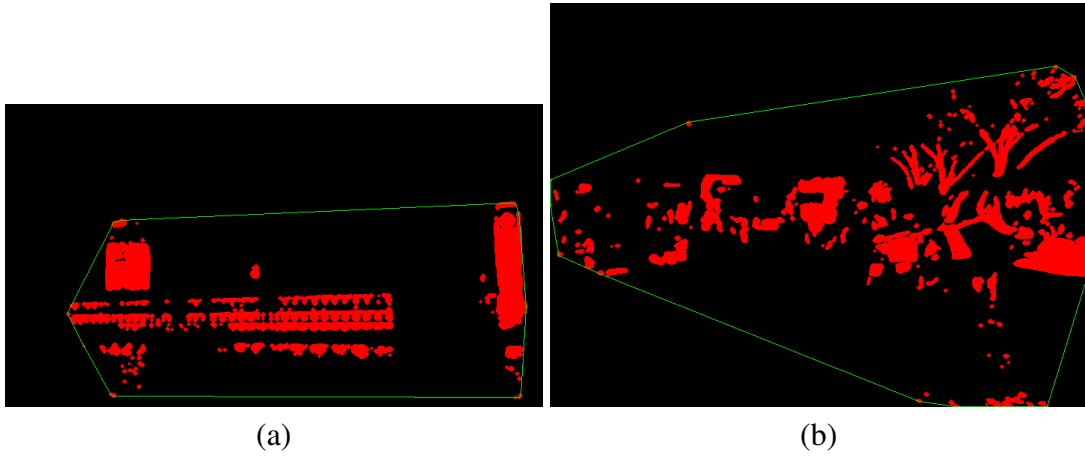


Figure A.1: (a) and (b) are convex hull results. The red points are point sets and the polygon with green edges are the convex hull we found.

into a set of non-overlapped regions (facets) that cover the whole plane.) that edges of each subdivision facets are all the points in the plane that are equidistant to the two nearest sites.

The algorithm we implement is based on delaunay triangulation which . The idea is that for every subdivision there exists a dual subdivision in which facets and points (subdivision vertices) swap their roles, so if we can do delaunay triangulation for those voronoi sites, the planar subdivision is result from the dual answer. The voronoi segmentation is calculated by the following steps:

1. With voronoi sites, subdivide a plane into triangles using Delaunays algorithm.
2. The dual subdivision is a Voronoi diagram of the input 2d point set.

Fig. A.2 shows the calculation results.



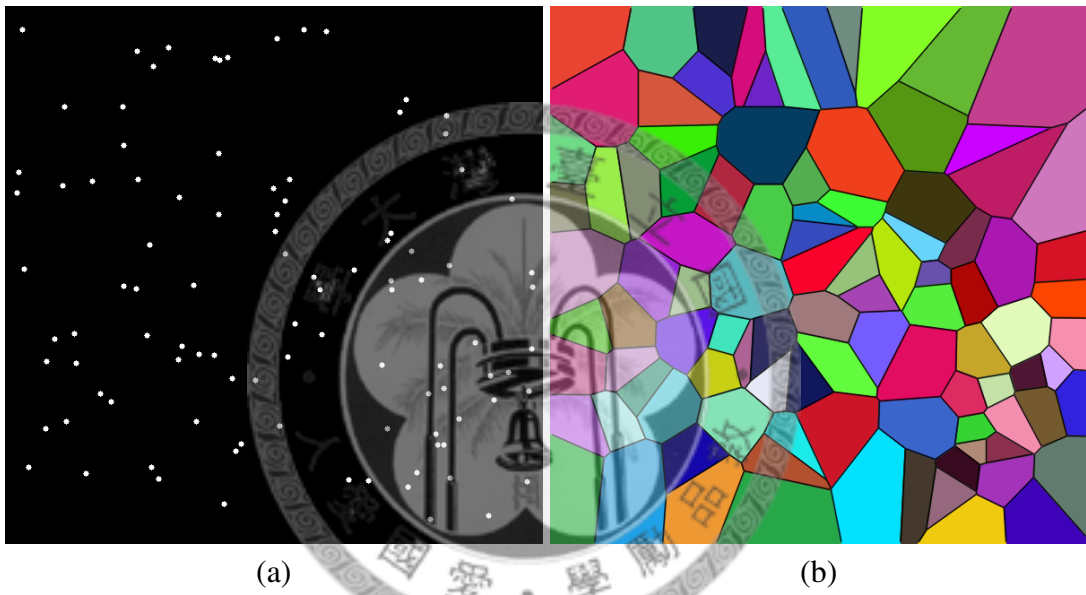


Figure A.2: (a) shows the voronoi sites. (b) is the planar subdivision result.