

# Content-based Video Retrieval for Pattern Matching Video Clips

Adam Jaamour

Bachelor of Science in Computer Science with Honours  
The University of Bath  
May 2019

This dissertation may be made available for consultation within the University Library and may be photocopied or lent to other libraries for the purposes of consultation.

Signed:

# **Content-based Video Retrieval for Pattern Matching Video Clips**

Submitted by: Adam Jaamour

## **COPYRIGHT**

Attention is drawn to the fact that copyright of this dissertation rests with its author. The Intellectual Property Rights of the products produced as part of the project belong to the author unless otherwise specified below, in accordance with the University of Bath's policy on intellectual property (see <http://www.bath.ac.uk/ordinances/22.pdf>).

This copy of the dissertation has been supplied on condition that anyone who consults it is understood to recognise that its copyright rests with its author and that no quotation from the dissertation and no information derived from it may be published without the prior written consent of the author.

## **Declaration**

This dissertation is submitted to the University of Bath in accordance with the requirements of the degree of Bachelor of Science in the Department of Computer Science. No portion of the work in this dissertation has been submitted in support of an application for any other degree or qualification of this or any other university or institution of learning. Except where specifically acknowledged, it is the work of the author.

Signed:

## **Abstract**

Your abstract should appear here. An abstract is a short paragraph describing the aims of the project, what was achieved and what contributions it has made.



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Problem Description . . . . .	1
1.2	Related CBR Systems and Their Applications . . . . .	3
1.3	Project Aim . . . . .	4
<b>2</b>	<b>Literature Survey</b>	<b>7</b>
2.1	Content-Based Retrieval . . . . .	7
2.1.1	Video Retrieval Methods . . . . .	7
2.1.2	Temporal Aspects of Videos . . . . .	9
	Temporal Structure of a Video . . . . .	9
	Challenges of Temporality . . . . .	10
2.1.3	CBVR for Mobile Devices . . . . .	12
	Query Video Quality . . . . .	12
	User Experience . . . . .	12
2.2	Visual Content Extraction for Pattern Matching . . . . .	13
2.2.1	Static Features . . . . .	13
	Colour-based Features . . . . .	13
	Texture-based Features . . . . .	16
	Shape-based Features . . . . .	16
2.2.2	Dynamic Features . . . . .	17
	Points of Interest . . . . .	17
	Object Features . . . . .	18
	Motion Features . . . . .	18
2.2.3	Learning Models for Pattern Matching . . . . .	20
	Bag-of-Visual-Words . . . . .	20
	Constellation Models . . . . .	21
	Deep Learning Models . . . . .	21
2.3	Structural Movie Pre-Processing . . . . .	22
2.3.1	Temporal Movie Segmentation . . . . .	22
	Feature Extraction . . . . .	23
	Similarity Measurements . . . . .	24
	Detection . . . . .	25
2.3.2	Key Frame and Thumbnail Extraction . . . . .	27

Key Frames . . . . .	27
Thumbnails for Initial Shortlisting . . . . .	28
2.3.3 Feature-Length Movie Optimisation . . . . .	29
2.4 Chapter Summary . . . . .	29

# List of Figures

1.1	Wireframe showing the basic high-level concept that the system will try to achieve. . . . .	5
2.1	Illustrations of a text-based content-retrieval system ( <i>a</i> ), a content-based video retrieval system ( <i>b</i> ), and the project’s desired CBVR system ( <i>c</i> ). . . . .	8
2.2	Temporal structure of videos, including the different terms used to describe video temporal units. Figure courtesy of A. Araujo et al. . . . .	10
2.3	Frames from the famous running scene in Forrest Gump extracted at intervals of 10 seconds. Video frames courtesy of “ <i>Forrest Gump long run scene</i> ” YouTube video available online: <a href="https://youtu.be/QgnJ8GpsBG8?t=325">https://youtu.be/QgnJ8GpsBG8?t=325</a> . . . . .	11
2.4	Example of static features, including colour-based features (with an RGB colour histograms of 256 bins for colour-based features), texture-based features (with an co-occurrence matrix showing the frequency of pairs of pixels) and shape-based features (with a shape matrix descriptor). . . . .	14
2.5	Example of a histogram counting the location of points relative to a vertical grid. Image courtesy of Bradski and Kaehler. . . . .	15
2.6	If the range of the bins is too large, then the distribution is coarse (left). If the range of the bins is too small, then the distribution is not accurately represented and spikes cells appear (right). Image courtesy of Bradski and Kaehler. . . . .	16
2.7	An image of the Palace of Monaco with coloured windows representing poor features in blue (flat), edges in green, and good features in red (corners). . . . .	18
2.8	Visualisation of sparse and dense optical flow algorithms on different shots. . . . .	20
2.9	A 3D and a 2D representation of dense optical flow. Each arrow represents the direction and magnitude of a specific pixel location, which can variate in direction and length. . . . .	21

2.10	Shot boundary detection example of a video scene made up of two shots with a gradual transition between the two shots. Figure courtesy of Michael Gygli available online at: <a href="https://medium.com/gifs-ai/_ridiculously-fast-shot-boundary-detection-with-fully-convolutional-neural-networks-da9d8c73e86c">https://medium.com/gifs-ai/_ridiculously-fast-shot-boundary-detection-with-fully-convolutional-neural-networks-da9d8c73e86c</a> . . . . .	22
2.11	. Visual examples of a quick cut, a dissolve cut and a fading cut. Frames courtesy of Koprinska & Carrato (2001) “Temporal video segmentation: A survey” available online at: <a href="https://www.sciencedirect.com/science/article/pii/S0923596500000114">https://www.sciencedirect.com/science/article/pii/S0923596500000114</a> . . . . .	23
2.12	Threshold-based approach for shot boundary detection. In this example, the shot boundary is detected at <i>A</i> , it is missed at <i>B</i> , and it is misinterpreted at <i>C</i> as two cuts rather than a single one. . . . .	25
2.13	Example of frames to sample for low-visual content analysis. The first frame for each second (one frame every thirty seconds) is retrieved for a 30 fps 3-second video of a ball rolling from the left-hand side of the screen to the right-hand side. Video frames courtesy of “ <i>How to Animate a Rolling Ball</i> ” YouTube video available online: <a href="https://youtu.be/cgbLAreE1NI?t=130">https://youtu.be/cgbLAreE1NI?t=130</a> . . . . .	28

# List of Tables

1.1 The four types of visual search involving images and videos, classified by type of query used and by type of database being queried. . . . .	2
--	---



# Acknowledgements

Add any acknowledgements here.

# Chapter 1

## Introduction

During the past decade, the amount of data in the world has grown at exponential rates and is currently showing no signs of slowing down. According to infographics released by IBM, already 2.7 ZB<sup>1</sup> of data existed in the world in 2012 (Karr, 2012), enough to fill up almost 40 billion 64 Gb iPhones. This number has since then risen to 8 Zb in 2015 and is expected to reach a yearly production of 35 Zb by 2020 (Deutscher, 2012) (Karr, 2012). 85% of all the world data is considered to be unstructured data (Blumberg and Atre, 2003), which is mainly made up of multimedia in the form of images and videos. An important factor in this growth is the rise of mobile devices usage and social media. Indeed, the amount of mobile-dependant users has grown at impressive rates, with now 95% of the United States citizens owning a mobile phone (Fanning, Mullen and McAuley, 2012). People tend to use their mobile devices for most day-to-day activities, with a majority of this time spent on social media services such as YouTube or Facebook. The social networks are one of the primary causes for the exponential growth of unstructured data mentioned earlier, with over 300 hours of new video content uploaded to YouTube every minute of the day. The combination of mobile devices usage and social media data are the main contributors in today's flood of unstructured data. This project will look at creating a system combining the problem of unstructured data in the form of video data on mobile devices.

### 1.1 Problem Description

Content-based video retrieval is a computer vision solution to the problem of searching large databases of videos, where "content" corresponds to information from the videos such as colours or shapes that can be used to efficiently retrieve a desired video from the database. The main difficulty with content-based video retrieval lies with the reference database of videos itself. Videos in

---

<sup>1</sup>Zettabytes. 1 Zb = 1 trillion Gb.

a database correspond to unstructured data, meaning only the video data<sup>2</sup> and their metadata<sup>3</sup> are stored in the database. Computations involved in querying a database of videos using only the metadata available without the ability to analyse the video files would be too expensive and highly inefficient (Patel and Meshram, 2012). Therefore, the inspection of the video files is a necessity to find solutions to the database complication in content-based video retrieval. It is important to note that the metadata can still be used as a complement to the video files analysis to fine-tune the pattern matching. Another complication involves the large database size. Indeed, more complications arise from databases populated with a large number of videos, especially if their duration is lengthy e.g. feature-length films. Solutions to those problems hence require intelligent and efficient pattern matching algorithms to overcome the specified difficulties. This is where visual search and the variety of fields it includes such as artificial intelligence, machine learning and database management come in.

Several visual search techniques exist for querying databases of unstructured data such as video or image files. These techniques can be classified based on the type of query and on the type of database used, as shown in table 1.1. The most common forms of visual search consist in querying a database of images either with an image (I2I) or with a video (V2I), as depicted in section 1.2 where similar existing systems are mentioned. Another less frequent variant consists of querying a database of images with a video query (V2I) (Araujo and Girod, 2018). However, this dissertation will solely focus on querying a database of videos with a video query (V2V). Because algorithms from other variants of visual search can be relevant to V2V, existing solutions for these variants will also be explored to find ways potential techniques that could be implemented with this project.

Image Query	I2I	I2V
Video Query	V2I	V2V
Database of Images		Database of Videos

Table 1.1: The four types of visual search involving images and videos, classified by type of query used and by type of database being queried.

To pattern match the query video to a video in the database, key elements from the query video, called *"features"*, are extracted and compared to the same features from videos in the database to find similarities. These features can include elements such as colour histograms or shapes of interest and their on-screen movements (Patel and Meshram, 2012), as well colour percentages,

---

<sup>2</sup>Video data includes the video frames and the audio.

<sup>3</sup>Examples of metadata related to a video file includes captions, file name, file type, video length, file size, etc.

colour layouts or textures (Petković, 2000), or other unique features such as on-screen text or audio components e.g. soundtracks, dialogues or sound effects. A wide spectrum of existing solutions exists to efficiently retrieve videos using video queries, which are discussed in the section below.

## 1.2 Related CBR Systems and Their Applications

Visual search technology has an endless amount of applications in fields such as education, navigation, utilities and security. Some of the systems using this technology have already found their way to commercial applications. For example, YouTube’s Content ID system uses visual search, and more precisely V2V, to detect copyright infringements on user-uploaded videos by automatically comparing an uploaded video with a database of protected videos provided by their partners. If a match is found, YouTube allows these partners to take action on the uploaded video that infringed their copyrights (Lars, 2012). Another example is Google Lens (*Google Lens*, n.d.), a mobile application that recognises the input context through a mobile device’s camera in order to relay information about objects of interest. For instance, if pointing the camera at a landmark, information such as historical facts or opening hours about that location will be displayed, or if pointing it at a WiFi router’s label, the mobile device will automatically connect to the network (Villas-Boas, 2017). Other profitable applications include A9’s Amazon Flow (*A9: Visual Search*, n.d.) that uses deep-learning based computer vision for Amazon’s search services, as well as Shazam (*Shazam*, n.d.), a mobile application that matches short user-recorded sounds with a piece of music.

Aside from the commercialised systems stated in the previous paragraph, many visual search applications are yet to be implemented for large-scale use. For example, a system allowing companies to detect all appearances of their logos during television broadcasts, or a system enabling students to find a section of a recorded lecture by using a lecture slide as a query (Araujo and Girod, 2018) are possibilities that need to be explored. Content-based video retrieval academic papers have also been published in recent years. For instance, Liu et al. (Liu, Mei and Zhang, 2014) discuss a concept similar to this project’s idea consisting of a mobile visual search system allowing users to discover videos by pointing their phone at a screen. Recently, research in visual search has gained momentum, notably with the annual TRECVID<sup>4</sup> conference (Awad, Butt, Curtis, Lee, Fiscus, Godil, Joy, Delgado, Smeaton, Graham, Kraaij, Quénnot, Magalhaes, Semedo and Blasi, 2018) hosted by the NIST<sup>5</sup>. This conference’s goal is to host workshops that focus on information retrieval research, with an emphasis on content-based retrieval of digital videos

---

<sup>4</sup>Text REtrieval Conference Video Retrieval Evaluation

<sup>5</sup>National Institute of Standards and Technology

(Smeaton, Over and Kraaij, 2006). The combination of academic research and specialised conferences will hopefully help bring more content-based retrieval solutions to real-life problems.

### **1.3 Project Aim**

This project aims to find solutions to address the previously mentioned problems from Section 1.1, namely:

- the exponential rise of unstructured data in the form of media, with a focus on videos,
- the over-reliance people have on their mobile devices for basic day-to-day tasks,
- the sparse collection of existing solutions that could be combined to create more efficient systems,
- the lack of focus on videos, especially feature-length movies.

This project will attempt to translate solutions to the mentioned problems by creating a system that will combine existing solutions into a single one, with a focus on recognising a feature-length movie from a short video query recorded from a mobile device, as depicted in Figure 1.1.

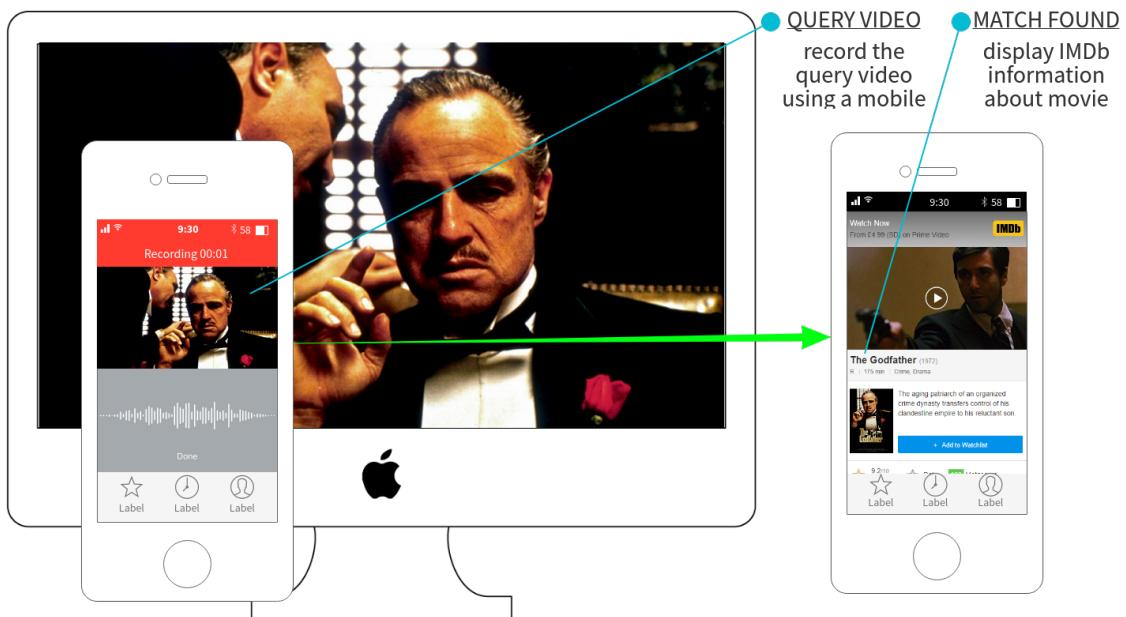


Figure 1.1: Wireframe showing the basic high-level concept that the system will try to achieve.



# Chapter 2

## Literature Survey

### 2.1 Content-Based Retrieval

Content-based retrieval is a type of visual search technique where large databases of either images or videos are queried to find the closest match to a query image or video. Although this project focuses on content-based video retrieval, also referred to as CBVR<sup>1</sup>, image retrieval techniques (CBIR<sup>2</sup>) will be discussed as well due to their relevance in video retrieval.

This section will first review the different video retrieval methods and their evolution, starting from text-based retrieval to content-based retrieval, before addressing the various challenges that exist in video retrieval, such as the additional difficulty caused by the temporal aspect of videos compared to images, and the complications of targeting mobile devices for such a system.

#### 2.1.1 Video Retrieval Methods

Drastic advances in video capturing technology have caused important amounts of unstructured data in the form of videos to be produced in recent years. This has lead to a high-demand to develop new efficient solutions for processing this data, with video retrieval being the answer.

Throughout the years, video retrieval has improved in parallel with the breakthroughs in video recording devices. Early video retrieval techniques used a text-based approach where the system accepted text input to search the database of videos (Lai and Yang, 2015), as seen in Figure 2.1.a. For example, the user would input the string query “*De Niro*”, which would return all movies in which Robert De Niro starred or “*Coppola*” to find all movies directed by Francis Ford Coppola. Unique aspects of the video clip such as

---

<sup>1</sup>Content-Based Video Retrieval

<sup>2</sup>Content-Based Image Retrieval

movie credits or sports scores were often analysed using OCR<sup>3</sup> technology (Li and Doermann, 2002). The query text was then compared to a video file's content, such as colours, shapes, texture, luminance or objects, or to the file's metadata<sup>4</sup>, such as the video title, author, date, content description, commentaries, captions or keywords (Li and Doermann, 2002) (Feng, Cao, Bao, Bao, Zhang, Lin and Yun, 2011) (Patel and Meshram, 2012). However, these techniques were highly inefficient compared to content-based techniques as they often relied on manually noted annotations and textual descriptions to find similarities for matching the query video to a video in the database and did not make use of the actual visual content that describes a video.

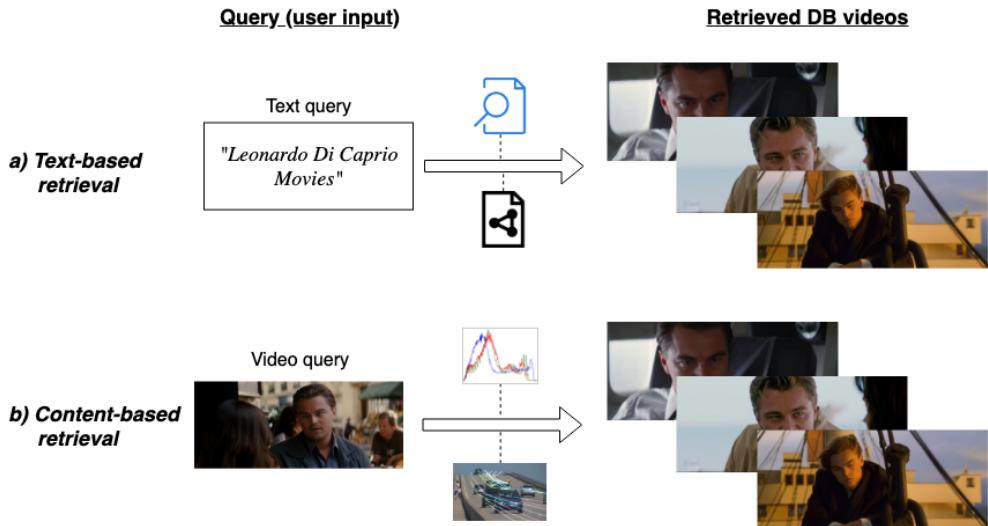


Figure 2.1: Illustrations of a text-based content-retrieval system (a), a content-based video retrieval system (b), and the project's desired CBVR system (c).

Content-based retrieval techniques quickly replaced text-based retrieval techniques towards the end of the 20th century by making use of the visual content to compute similarities between videos (Lai and Yang, 2015). Instead of accepting a text string, the system takes a video as input to extract its visual contents, as shown in Figure 2.1.b. According to Petković et al. (Petković, 2000), this visual content can be broken down into three different categories:

- *Raw data*, which corresponds to individual raw video frames and the video file's attributes such as the frame rate per seconds, the number of bits per pixel or the colour model used.
- *Low-level visual content* consists of the visual features that describe a

<sup>3</sup>Optical Character Recognition

<sup>4</sup>The data associated to a video file

video. This content includes colours, shapes, textures and motion. Low-level visual content can be extracted into static features such as histograms (see Section 2.2.1) or dynamic features such as objects or motion (see Section 2.2.2) using a wide variety of existing techniques. Once the content has been extracted, it can be used to first compute similarities between videos (Lai and Yang, 2015) and later pattern match them.

- *Semantic content* contains the high-level concepts that are present in a video. These high-level concepts can be described as objects or events using the features. To extract semantic content from a video, a grammar of rules for objects must be provided. An example of an object rule could be "if the shape is round, the colour is orange and the object is moving, then that object is a basketball".

In comparison to raw data, low-level visual content provides the most relevant visual information that can be extracted from a video for the purpose of a CBVR system. Semantic content extraction adds an additional layer of complexity compared to low-level visual content extraction as it requires domain knowledge and user interaction (Petković, 2000). Therefore, this project will focus on using low-level visual content to extract information about the video and compute the similarities between the query video and database videos. It is important to note that this project's goal differs from classic CBVR systems where a list of videos is returned (see in Figure 2.1.b), as it must return a specific video that matches the most the query video. To improve the pattern matching accuracy phase, raw data (e.g. audio) and metadata (e.g. captions) may be used to improve the pattern matching accuracy (Patel and Meshram, 2012).

### 2.1.2 Temporal Aspects of Videos

#### Temporal Structure of a Video

The most important difference between content-based image retrieval and video retrieval lies within the temporal aspect of the video. Naturally, the temporal aspect of a video clip stores supplementary information about the content, including dynamic low-level visual content e.g. an object's motion, and semantic content e.g. actions and events. According to A. Araujo et al. (Araujo and Girod, 2018), a video's temporal structure can be subdivided into three units, as shown in Figure 2.2:

- *Frames* correspond to the smallest temporal unit of a video file. A single segment of a video is referred to as a frame. Frames are also used to describe the frame rate (the frequency at which consecutive stills appear on a screen every second) e.g. "24 fps" corresponds to a video made up of 24 stills per second.

- *Shots* are grouped sequences of visually similar frames. They are usually described in seconds.
- *Scenes* are collections of shots which are related based on the action and objects present in the shot, thus giving them a semantic aspect. The length of a scene is generally calculated in minutes rather than seconds.

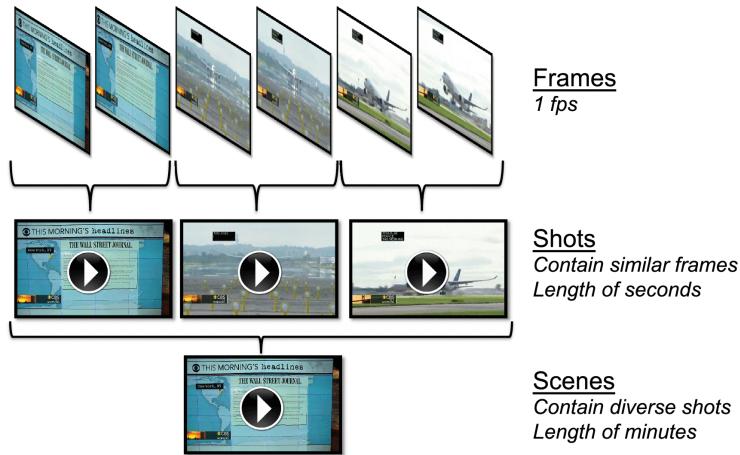


Figure 2.2: Temporal structure of videos, including the different terms used to describe video temporal units. Figure courtesy of A. Araujo et al.

Because this project will explore possible solutions to create a CBVR system targeting databases of feature-length movies, a fourth video temporal structure category relevant to this project can be added to Araujo et al.'s initial list:

- *Movies* can be described as a large group of scenes that are used to tell a story. Movie durations commonly range from one to three hours.

### Challenges of Temporality

Multiple challenges arise when dealing with CBVR systems in contrast to CBIR systems as the videos' temporal aspect adds a new dimension of complexity when extracting visual information. While the low-level visual content describing images remains mostly the same for videos, some new information that did not exist in images can be extracted, such motion. As mentioned previously, videos are made up of frames, which make up shots when a combination of similar frames are played in succession. This means that videos carry information about motion, such as the trajectory of objects. This introduces unique challenges to the algorithms used to extract motion.

Because videos are made up of numerous stills, usually around 24 frames per second (Brownlow, 1980), two consecutive frames are near-identical. The pixels describing an object in one frame will remain the same in the next frame, except for the edge pixels perpendicular to the motion's trajectory (Bradski and Kaehler, 2008).

Figure 2.3 shows six frames from Forrest Gump's famous running shot, which lasts 44 seconds, making up a total of 1056 frames. Each frame in the figure was captured with 10-second intervals, meaning 240 frames separate each still. In the first three frames, the group of people running in the background barely moves in the space of 20 seconds. The pixels that describe the group in the shot remain mostly unchanged for all of the frames between the 3 samples, equivalent to 480 frames, with a few additional pixels describing the group as it advances towards the camera. The same can be said about the red cap in the last two frames. Most of the pixels making up the cap in the fifth frame remain the same in the sixth frame. This example perfectly betrays the reason why analysing a video frame by frame would be extremely inefficient when it comes to CBVR. Due to the similarities between consecutive frames, these should be aggregated (Araujo and Girod, 2018) to describe a shot by using a selection of frames, such as taking the six frames in Figure 2.3 to describe the entire 44 seconds of video, rather than keeping the original 1056 frames to describe it.

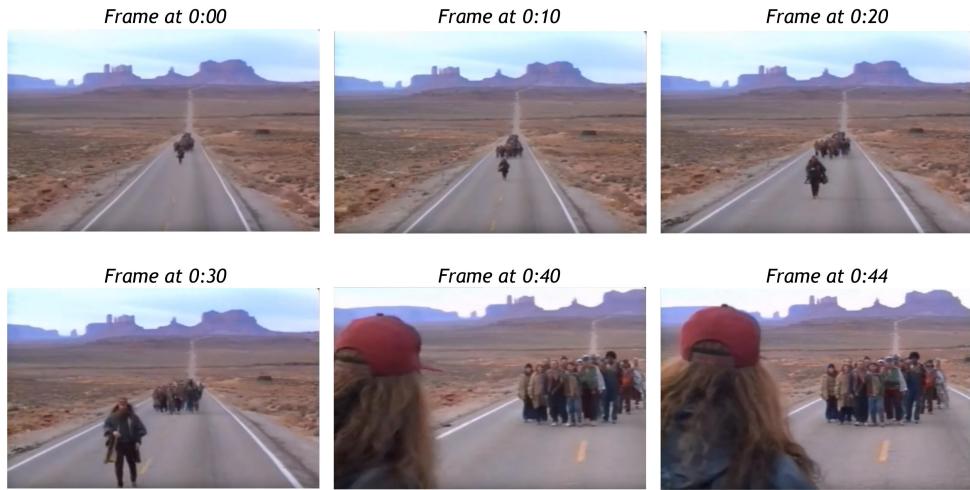


Figure 2.3: Frames from the famous running scene in Forrest Gump extracted at intervals of 10 seconds. Video frames courtesy of “*Forrest Gump long run scene*” YouTube video available online: <https://youtu.be/QgnJ8GpsBG8?t=325>.

### 2.1.3 CBVR for Mobile Devices

The aforementioned project aim is for the system to work on mobile devices for two reasons. The first reason is to allow users to directly use their mobile phone to record the query video by pointing their camera to a screen displaying a movie, which will, in turn, tell them which movie is being played (see the wireframe in Figure 1.1). Additional information retrieved from IMDb<sup>5</sup>, such as cast, crew, ratings, runtime and synopsis could also be displayed. The second reason is the popularity of mobile devices, which may be due to the improvements made on mobile phones' processing power, allowing more tasks to be carried out through this medium. However, such a system on a mobile device causes many problems regarding the query video recording method and the computational power available on mobile devices.

### Query Video Quality

Large visual differences are caused between the query clip and the actual clip stored in the database due to the capture conditions (Liu et al., 2014) (Wang, Oneata, Verbeek and Schmid, 2016) such as:

- Undesired camera movements due to unstable recording e.g. unstable recording, hand shaking.
- Low-quality recording due to poor user recording e.g. scaling and rotation, and due to the environmental conditions e.g. lighting, reflections, blurring.
- Video noise because of the camera sensor.
- Decoding artefacts caused by various file compression.

These low-quality conditions add difficulty to the pattern matching phase where the similarities between the query video and database videos have to be computed. Indeed, if the query video is very different from the actual video, then the noisy elements of the video query must be filtered out. For example, if the recorded video is shaky, then this shaking motion has to be pruned before analysing the recorded clip's motion. However, processing power must be used from the actual visual content extraction and pattern matching phases to be used for video noise filtering.

### User Experience

According to Liu et al. (Liu et al., 2014), the majority of mobile device users expect a polished product with quick video query and instant or progressive

---

<sup>5</sup>Internet Movie Database

results, meaning that the searching algorithms must be efficient. However, one of the downsides of mobile devices is the computation power constraints. Despite the improvements in mobile processors, desktop devices still remain more powerful than mobile ones. A solution that Liu et al. suggest is to retrieve the low-level visual content locally on the mobile device and send the query to a server where the pattern matching will take place (Liu et al., 2014). This allows heavy computations to be off-loaded from the mobile device. Once a match is found, the result is returned to the user on his mobile device. A downside to this approach is the new constraint on network bandwidth rather than computational power.

## 2.2 Visual Content Extraction for Pattern Matching

Extracting the visual content from a video allows this content to be used to describe videos and compute similarities between them. This visual content is extracted from the aforementioned low-level visual content (see Section 2.1.1) (Petković, 2000) and stored in the form features, also referred to as visual descriptors. The term “features” is very broad and can be used to describe many different visual aspects in an image or in a video, ranging from colours, shapes and textures to points, edges, objects and motion.

These features can be divided into two categories: static features and dynamic features (Petković, 2000). This section will first survey examples of static visual descriptors and methods to extract them from videos, and will then focus on examples and methods of extracting dynamic visual descriptors from videos.

### 2.2.1 Static Features

Methods to extract static features operate on stills, which can correspond to individual video frames, thumbnails or key frames (see Section 2.3). This means that traditional image techniques can be applied to those stills (Hu, Xie, Li, Zeng and Maybank, 2011). They are organised in three different categories: colour-based features, texture-based features and shape-based features, which can be visualised in Figure 2.4.

#### Colour-based Features

The main colour-based features model are colour histograms. In general, a histogram consists of counts of some underlying data that is organised into predetermined bins to a statistical representation of the distribution of that data. Figure 2.5 depicts an example of a histogram where a collection of points

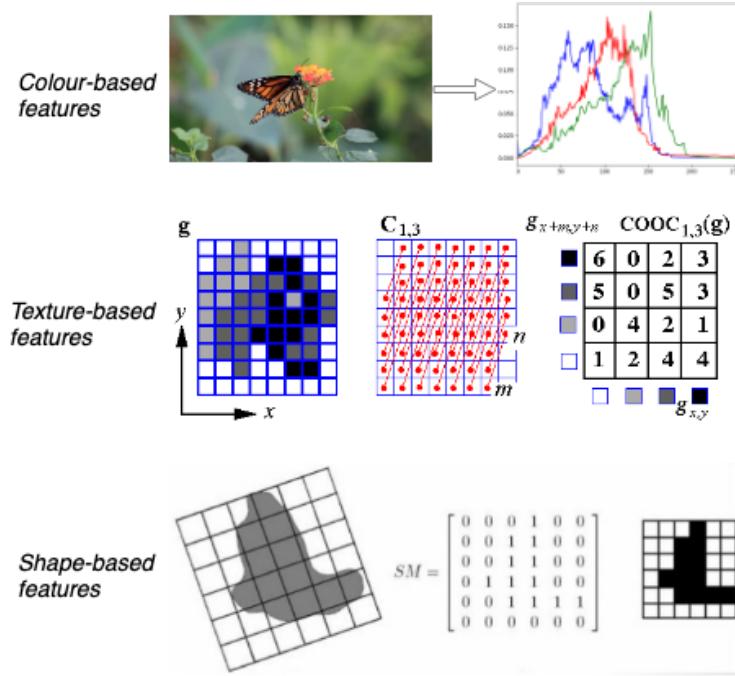


Figure 2.4: Example of static features, including colour-based features (with an RGB colour histograms of 256 bins for colour-based features), texture-based features (with an co-occurrence matrix showing the frequency of pairs of pixels) and shape-based features (with a shape matrix descriptor).

is organised into specific pre-defined bins based on their location relative to a vertical grid.

In the case of a colour histogram, the underlying data that the histogram is trying to represent is the distribution of the colour pixels throughout an image or a video frame. Different kinds of colour histograms exist as they depend on the chosen colour space, which includes RGB<sup>6</sup>, HSV<sup>7</sup>, HSL<sup>8</sup> or YPbPr colour spaces to name a few. These may vary based on the applications of the colour histograms.

Typically, for an RGB colour histogram, 256 bins are used to accurately represent all the possible values that the pixels can take (ranging from 0 to 255) for each of the three RGB channels, which are then plotted as three individual graphs. Choosing the right range for the histogram's bins is crucial to represent the distribution efficiently. If the range of pixels that defined the bins is wider, meaning there are less overall bins, then the histogram's distribution would

<sup>6</sup>Red Green Blue

<sup>7</sup>Hue Saturation Value

<sup>8</sup>Hue Saturation Light

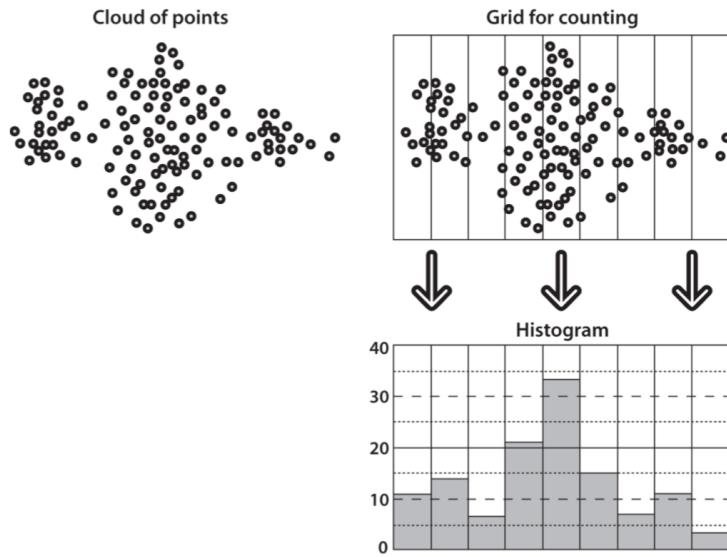


Figure 2.5: Example of a histogram counting the location of points relative to a vertical grid. Image courtesy of Bradski and Kaehler.

be too coarse-grained and the general structure of the histogram would be lost, as pointed out by the left half of Figure 2.6. On the other hand, if the range of the bins is too narrow, meaning there are more overall bins, then the histogram's distribution would not be represented accurately and there would be many spiky cells, as betrayed in the right half of Figure 2.6 (Bradski and Kaehler, 2008).

One of the inaccuracies with colour histograms lies within the scope of the distribution. If the histogram represents the global distribution of all the pixels in the still, then two images might have very similar histograms (Petković, 2000). For example, a histogram containing 60% white pixels and 40% blue pixels could either describe both a blue sky with white clouds or a snowy landscape with a blue sky. Despite both histograms being good colour-based features, the actual result is still poor when it will be used for matching the histograms. A solution consists in segmenting the still into multiple local images, and extract the local colour-based features for each segment. For example, the still could be partitioned into a 5x5 grid, and a colour histogram could then be computed for each grid (Yan and Hauptmann, 2007). This would enable colour-based features to represent specific regions of the still rather than globally describing an image. However, the same problem mentioned earlier could occur if the still is segmented into too many regions, causing the overall histogram to be coarse.

Other types of colour-based features can be extracted from images and videos such as colour moments, colour correlograms (Huang, Kumar, Mitra,

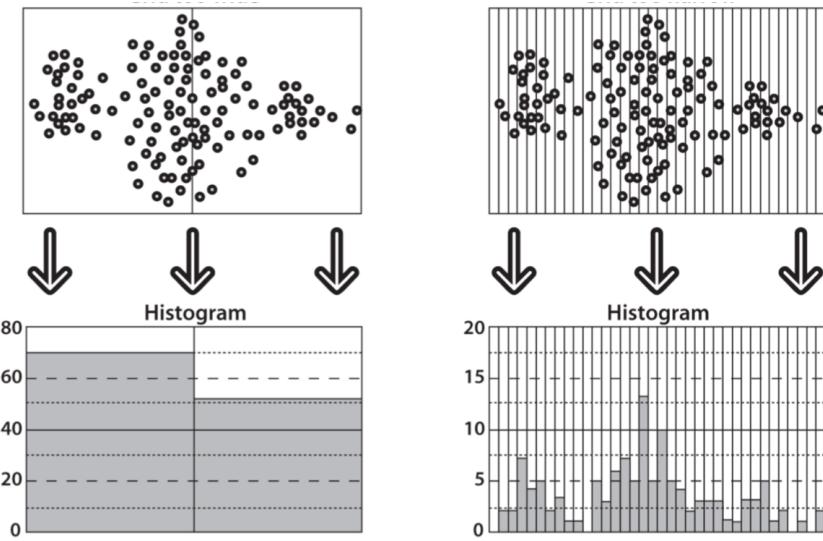


Figure 2.6: If the range of the bins is too large, then the distribution is coarse (left). If the range of the bins is too small, then the distribution is not accurately represented and spikes cells appear (right). Image courtesy of Bradski and Kaehler.

Zhu and Zabih, 1997) and Gaussian models. However, colour-based features have their limitations as they cannot describe textures and shapes, rendering them inefficient in certain applications (Hu et al., 2011).

### Texture-based Features

Texture-based features are often used in parallel with colour-based features. The aforesaid problem where two different objects might share a similar histogram can be solved by using texture-based features to differentiate them e.g. green tree leaves and green grass. These can be discerned by using a variety of features such as Tamura features, which extract information including coarseness, contrast, and directionality of the objects (Amir, Berg, Chang, Hsu, Iyengar, Lin, Naphade, Natsev, Neti, Nock et al., 2003).

### Shape-based Features

Shape-based features are used to describe the overall shape of objects present in the image. The most common approach consists in computing an EHD<sup>9</sup>, which consists in detecting the edges present in the image (see Section edges) and then plot their spatial distribution in a histogram by counting the number of pixels for each edge (Hauptmann, Baron, Chen, Christel, Duygulu, Huang,

---

<sup>9</sup>Edge Histogram Descriptor

Jin, Lin, Ng and Moraveji, 2004). The same image segmentation technique can be used to localise the EHD. These shape-based features have many applications but are harder to extract and require more computing power than colour-based features and texture-based features.

### 2.2.2 Dynamic Features

In contrast to static features, which can be extracted from individual video frames, dynamic features require the continuity between consecutive frames to extract relevant visual descriptors, making use of the temporal aspect of the video mentioned in section 2.1.2). These features can be divided into two subcategories: object features and motion features. However, before reviewing the techniques used to extract these features, it is important to specify what defines a good visual feature.

#### Points of Interest

Figure 2.7 shows an image with coloured windows used to make the difference between poor and good potential features that could be used for object and motion features:

- *Flat surfaces* are portrayed in blue in Figure 2.7. These blue windows are spread over large areas of the image, meaning it is difficult to find their specific location. Moving the blue window along the image in any direction will result in the same visual content being represented in the window. Therefore, these flat regions are the worst structures as they do not contain any useful information.
- *Edges* are characterised in green in Figure 2.7. These are more informative than flat surfaces as they can be more accurately localised, but pinpointing an exact location is still hard as the patch can be moved in the direction parallel to the edge. Moving the green window along the edge will again result in the same visual content being represented in the window. Edges are efficient to detect object boundaries, but not for tracking specific points.
- *Corner* are characterised in red in Figure 2.7. These are the most descriptive points as they are often unique and can be precisely located in an image. Moving the red window in any direction will cause it to look different. Corners are therefore the ideal candidate for features used in object matching and tracking.

Once expressive and unique descriptors like corners are detected, they can be used to extract object features and motion features, and to compute similarities between the query video and the videos in the database.



Figure 2.7: An image of the Palace of Monaco with coloured windows representing poor features in blue (flat), edges in green, and good features in red (corners).

Many different algorithms exist to find robust features and interest points. For example, Sobel operators can be used to detect edges, and Shi and Tomasi algorithms and Harris operators can be used to detect corners (Bradski and Kaehler, 2008). Combinations of both can be used, by detecting edges to facilitate corner detection.

### Object Features

Object features correspond to objects that are detected using the colour, texture and size of image regions. Some of the most common objects usually detected in videos are faces, as many CBVR systems use them to compute similarities between videos (Sivic, Everingham and Zisserman, 2005). However, extracting object features is time-consuming and expensive in terms of required processing power, which is why CBVR algorithms either focus on detecting specific sets of objects rather than general objects that may be present in a scene, or on static features.

### Motion Features

Motion features are a unique characteristic of videos that are absent from images. All of the above-mentioned features can be extracted from images, apart from motion features, which are unique to videos. These can originate

from two sources: either background camera movement or foreground object movement. Motion features can therefore be classified into two categories: camera-based motion features and object-based motion features. *Camera-based motion features* include movement such as camera zooms, pannings and tiltings. *Object-based motion features* are more interesting than the former since they can describe key objects in the shot, which can be further classified into (Hu et al., 2011):

- *Statistics-based* motion features are extracted to model the local and global distributions of motion points in the video. For example, a causal Gibbs model can be used to represent the spatiotemporal distribution of local motion measurements in a shot once the trajectory-based motion features have been used to prune undesired camera motions (Fablet, Bouthemy and Perez, 2002). This type of motion features is very cheap to extract but lacks depth as they cannot represent object actions and object relationships accurately.
- *Trajectory-based* motion features are extracted by representing object trajectories in a shot. A long list of algorithms exist to obtain these features such as the, but this section will focus on differentiating sparse and dense techniques with one of the most popular methods: optical flow. Optical flow can be defined as the relative motion between the visual content in a shot and the camera (Bradski and Kaehler, 2008):

In *sparse optical flow*, only a few pixels from the frames in a shot are used. Features describing objects, such as corners, are used to track an object's motion across the shot. However, these only give information about where certain parts of the object are going, which are then used to estimate where the overall object is moving towards. Sparse optical flows. For example, Figure 2.8 indicates the overall motion of a car by tracking only a few pixels from it. On the other hand, *dense optical flow* calculates the optical flow for every pixel in the frame. At each pixel of the frame, the direction and magnitude of that specific location are represented by an arrow, which can variate in direction and length (see Figure 2.9).

Object-based motion features have many applications outside of CBVR systems. One of these applications is calculating the optical flow of a shot to prune camera movements caused by unstable recording e.g. hand shaking movement (Wang et al., 2016). Compared to statistics-based features, trajectory-based features can be used to describe object actions, but they rely on multiple challenging tasks to function efficiently such as correct object tracking and automatic trajectory recording. Finally, object-based motion features can also be classified in a third category: objects' relationship-based motion features. However, these are not de-

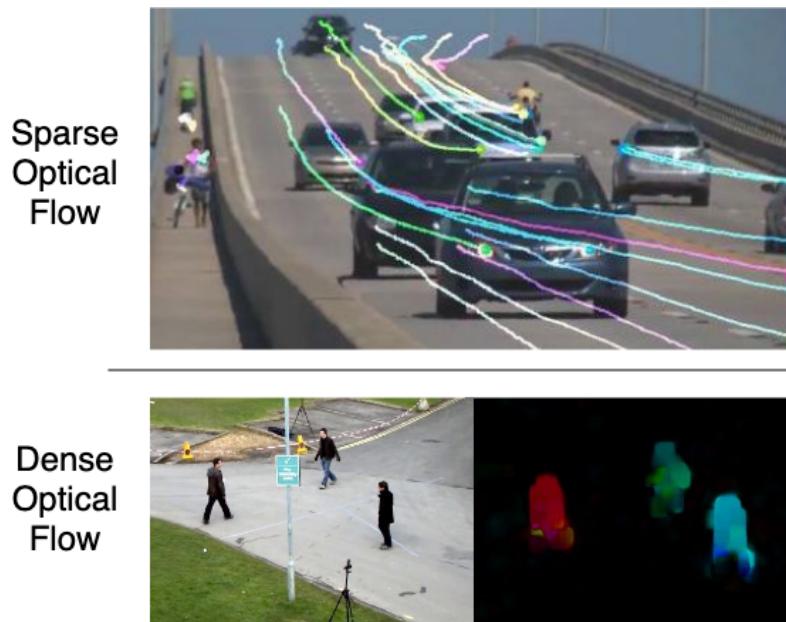


Figure 2.8: Visualisation of sparse and dense optical flow algorithms on different shots.

scribed in this literature review as they make use of the objects' symbolic representations, which is not relevant to this topic.

### 2.2.3 Learning Models for Pattern Matching

This section retraces the evolution of pattern matching models, ranging from Bags-of-Visual-Words and Constellation Models to Deep Learning models such as Neural Networks.

#### Bag-of-Visual-Words

The idea of BoVW<sup>10</sup> originates from document analysis, where the frequency words appeared in the text depended on the document type, and the document type could predict the word frequency.

Visual words are created by extracting distinctive features such as SIFT, SURF or HoG to train videos from a training set. These features are then clustered using classic models such as Mean-shift or GMMs. The visual words correspond to the middle of each cluster.

---

<sup>10</sup>Bag-of-Visual-Words

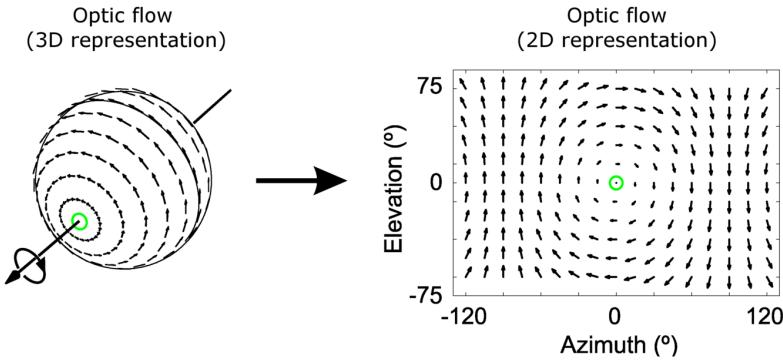


Figure 2.9: A 3D and a 2D representation of dense optical flow. Each arrow represents the direction and magnitude of a specific pixel location, which can variate in direction and length.

Histograms can be used by inputting a collection of visual words and a labelled training set of videos. For each video in the dataset, features are extracted and incremented in the related histogram. Once the training is accomplished, a new histogram can be generated for a query video (in a CBVR system) which will then be compared to the previously generated histograms for each dataset video. Example: H. Wang's implementation of BoW histogram (Wang et al., 2016). Fisher Vectors are an improvement on BoW (Araujo and Girod, 2018).

### Constellation Models

Constellation models are mainly used for detection. Features are organised in space such as a graph model where nodes are labelled with appearance and a Gaussian distribution is present on each node location. There is one model per class,  $s$  in a CBVR system, one model per video. The model is moved over image (could be video keyframes), which fires responses at each location. The biggest responses that are larger than a threshold locate the model.

### Deep Learning Models

- Classical Neural Networks
- Convolutional Neural Networks
- State of the art: AlexNet, VGG, Inception

## 2.3 Structural Movie Pre-Processing

The database of videos can be pre-processed to optimise the visual content extraction and pattern matching phases. As stated in Section 2.1.2 on the temporal aspects of videos, movies can be defined as a logically ordered collection of scenes, which contain multiple shots, all made up of individual frames. Shots are therefore the logical fundamental unit to organise and segment a video into (Hu et al., 2011). Pre-processing the database of movies by segmenting it into a list of shots and representing each shot with a single key frame is a profitable solution that will exponentially improve a CBVR system's efficiency. However, a limit must be set on the amount of data that is segmented to balance the efficiency and the accuracy of the CBVR system.

### 2.3.1 Temporal Movie Segmentation

The frames that make up a shot usually show strong content correlation. This means that features extracted from one frame will be extremely similar in another frame from the same shot. Detecting shot boundaries would allow the movie to be segmented and organised by shots, as depicted in Figure 2.10. These boundaries are defined by the type of transition between two different shots, which can either be defined as a quick cut when the transition is direct, or as gradual when it is a dissolve or a fade in/out transition that runs over multiple frames (Yuan, Wang, Xiao, Zheng, Li, Lin and Zhang, 2007), as shown in Figure 2.11.

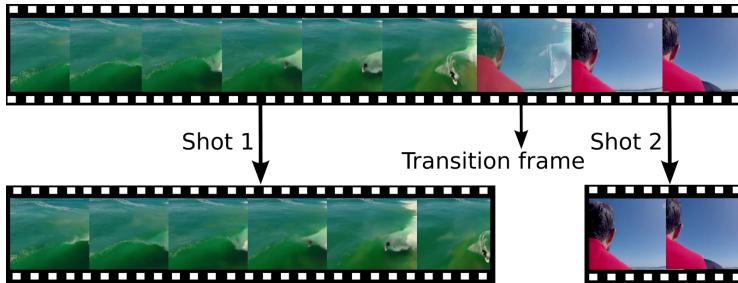


Figure 2.10: Shot boundary detection example of a video scene made up of two shots with a gradual transition between the two shots. Figure courtesy of Michael Gygli available online at: [https://medium.com/gifs-ai/\\_ridiculously-fast-shot-boundary-detection-with-fully-convolutional-neural-networks-da9d8c73e86c](https://medium.com/gifs-ai/_ridiculously-fast-shot-boundary-detection-with-fully-convolutional-neural-networks-da9d8c73e86c)

The goal of SBD<sup>11</sup> is to detect the discontinuity between consecutive frames. Also known as temporal video segmentation, this process requires three steps:

---

<sup>11</sup>Shot Boundary Detection



Figure 2.11: . Visual examples of a quick cut, a dissolve cut and a fading cut. Frames courtesy of Koprinska & Carrato (2001) “Temporal video segmentation: A survey” available online at: <https://www.sciencedirect.com/science/article/pii/S0923596500000114>

1. feature extraction
2. similarity measurements
3. shot boundary detection

### Feature Extraction

The first step in video segmentation consists in extracting features. These features include all the diverse types mentioned in Section 2.2, ranging from static features such as colour-based features (e.g. colour histograms) (Hoi, Wong and Lyu, 2006) to dynamic features such as corner points and SIFT<sup>12</sup>. Although colour histograms are simple to compute and work well with most shots as long as there is at most minor camera movement, they are inefficient when the shots contain major camera movements (Hu et al., 2011). For example, if in the shot the camera moves from the inside a house towards the outside by going through the window, the frames retrieved from the beginning of the shot will be extremely different to frames from the end of the shot.

Dynamic features are therefore more efficient for shot boundary detection than histograms due to their robustness. On the one hand, edge features are more vigorous than histograms when dealing with major camera movement and can handle changes in luminosity. On the other hand, corner and motion features can additionally handle camera motion and the impact of objects in the shot such rapid motion e.g. people walking in front of the camera.

---

<sup>12</sup>Scale Invariant Feature Transform

However, these dynamic features are more complicated to extract and do not always outperform simple features like colour histograms. Due to their simplicity, colour histograms remain the most common feature extraction technique for shot boundary detection (Yuan et al., 2007).

### Similarity Measurements

The second step in video segmentation is to use these extracted features to compute the similarities between frames. Many metrics exist to compute the similarities between two the extracted features. The most basic method consists in computing the Euclidean distance or the absolute distance  $g(n, n + k)$  between a pair of frames (Janwe and Bhoyar, 2013, p.476), as shown in Equation 2.1, where  $n$  and  $n + k$  represent the two frames being compared, and  $I_n(x, y)$  the intensity level of frame  $n$  at pixels locations  $(x, y)$ :

$$g(n, n + k) = \sum_{x,y} |I_n(x, y) - I_{n+k}(x, y)| \quad (2.1)$$

More advanced distance metrics can be used based on the type of feature that was extracted. For instance, intersection and chi-square distances can be used to compute the distance between the histograms of two frames (Camara-Chavez, Precioso, Cord, Phillip-Foliguet and Araujo, 2007, p.197). During the 2006 TRECVID conference, *Hoi et al.*, used grey scale histograms for their extracted features and calculated the colour differences between frames using the EMD<sup>13</sup>. EMD is used to calculate the distance between two probability distributions, which are represented by the two histograms that represent each frame in (*Hoi et al.*, 2006, p.2)'s case. These similarities can be measured using two techniques: pair-wise similarities and window similarities.

The first technique rests on using pairs of frames to compute the similarities between them. The similarities between two consecutive frames  $I_1$  and  $I_2$  are compared, before comparing the pair of frames  $I_2$  and  $I_3$ , and so on and so forth until frames  $I_{n-1}$  and  $I_n$  are compared. This straightforward technique precisely detects rough changes between the visual content in a shot (e.g. quick cuts between two shots) but is more sensitive to intensity changes caused by noise and disturbances such as camera motion and objects, leading to a high level of wrong detections (Janwe and Bhoyar, 2013).

The second technique is the window-based similarity measure, that measures the similarities between multiple frames specified within a range (Cernekova, Pitas and Nikou, 2006). This helps counter the effects of noise and disturbances captured by the pair-based approach but is more expensive to compute,

---

<sup>13</sup>Earth Mover's Distance

therefore less used (Hu et al., 2011).

### Detection

The final step in segmenting a video is the detection of the frames to cut, which is achieved by using the measured similarities between the frames. Two methods exist for this step: a threshold-based approach and a statistical learning-based approach.

The *threshold-based approach* compares the measured similarities between a pair of frames with a threshold. Whenever the similarity measure is smaller than the threshold, a shot boundary is detected and the shot can be cut. Figure 2.12 betrays an example of a threshold-based approach for different types of transitions, with a correct detection (*A*), a missed detection (*B*) and a misinterpreted detection (*C*).

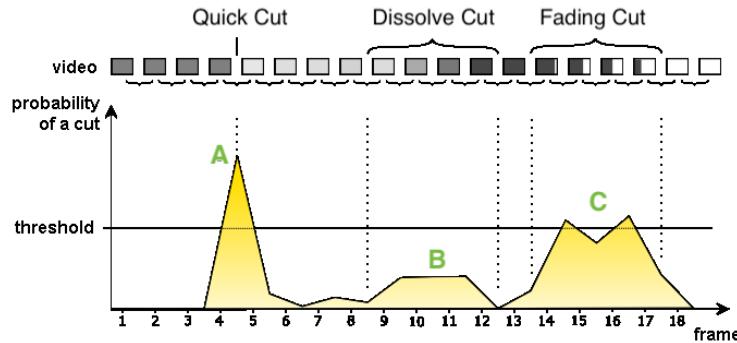


Figure 2.12: Threshold-based approach for shot boundary detection. In this example, the shot boundary is detected at *A*, it is missed at *B*, and it is misinterpreted at *C* as two cuts rather than a single one.

Either global thresholds, adaptive thresholds, or combinations of both can be used (Cernekova et al., 2006) (Hu et al., 2011):

- *Global thresholds* are set empirically and predefined for the entire video. They are therefore inefficient when trying to detect local variations as these are not integrated into the initial global threshold calculations.
- *Adaptive thresholds* are local thresholds that are continuously calculated for windows of frames. As the window slides across the video, the threshold is re-estimated accordingly to account for local variations. These are more accurate than global thresholds but are harder to estimate and require some knowledge about the video to set variables such as the window's size.

- *Combination of global and adaptive thresholds* are adaptive thresholds that take into account the values of global thresholds when estimating the frames within the window. Different thresholds can be determined for different types of transitions such as cut and dissolve transitions. However, these add increased complexity to the threshold estimation.

The *statistical learning-based approach* is a two-class classification task where each frame is classified either as a “shot change” or as a “no shot change” based on the extracted features and the similarity measurements between those features. In more complex systems, additional classes can be added by adding a classification of quick and gradual transitions for example. Two types of classifiers can be used for this approach:

- *Supervised learning-based classifiers*, such as the Support Vector Machine (Camara-Chavez et al., 2007) and Adaboost (Zhao and Cai, 2006, p.617). These supervised learning approaches have multiple advantages compared to threshold-based approaches as they do not require thresholds to be set and many different types of features can be combined in the feature extraction step to increase the classifier’s accuracy. However, this approach requires more expertise as the classifier relies on a well-trained data set to remain accurate (Hu et al., 2011).
- *Unsupervised learning-based classifiers* are divided into frame similarity-based and frame-based algorithms. *Frame similarity-based classifiers* groups the similarity measurements between a pair of frames into two clusters. The first cluster holds the similarity measurements with low values, which correspond to the “shot change” class, while the second cluster holds the similarity measurements with high values, which correspond to the “no shot change” class. K-Means and Fuzzy K-Means clustering algorithms can be used for frame similarity-based classifiers (Lo and Wang, 2001). Alternatively, *frame-based classifiers* consider each shot as a cluster of frames with similar features. Clustering ensembles can be used to classify frames into their respective shots (Chang, Lee, Hong and Archibald, 2007). This approach is advantageous over supervised learning-based classifiers as no prior training is required for, but the logical temporal order of the shots in scenes is not preserved.

All the aforementioned techniques in this section that are used in the three steps required in the process of segmenting a video can be applied to feature-length movies but require some improvements to work efficiently with these large videos. For instance, (Hanjalic, Lagendijk and Biemond, 1999) suggests considering the semantic aspects of movies, since people relate to the story, such as a marking dialogue, when they remember a movie.

### 2.3.2 Key Frame and Thumbnail Extraction

Frames from the same shot usually describe the same visual content, meaning there is a lot of redundancy amongst the frames that make up a shot. It is therefore more efficient to deal with a single frame that represents the entire shot or the entire video rather than dealing with an entire video.

#### Key Frames

Once shot boundaries have been detected and the video has been divided into a series of shots, a single frame, named the “key frame”, can be selected to represent the entire shot. (Heo, Lee and Jung, 2016) suggests only considering key frames in videos to operate on. These key frames would be used for all the sections mentioned in this Literature Survey, including the visual content extraction phases, the pattern matching phases, and database videos pre-processing phase. A wide variety of approaches exist to extract a key frame from a video.

The most simple approach is *Sequential Comparison Between Frames* where frames are sequentially compared until a frame very different to the previous key frame is found, at which point that frame is set as a new key frame (Heo et al., 2016). Frames are compared using the features mentioned in Section 2.3.1 e.g. histograms, and similarities are calculated using the same distance functions mentioned in Section 2.3.1 e.g. chi-square distance (Hu et al., 2011). More advanced approaches can be used:

- Global Comparison Between Frames
- Reference Frame
- Clustering
- Curve Simplification
- Objects/Events

To illustrate the advantage of using key frames, a 3-seconds long shot recorded at 30 fps<sup>14</sup> of a ball rolling on the ground can be used (90 frames in total). Analysing all the frames individually as stills would obviously be highly inefficient. However, selecting key frames to work on, as depicted in Figure 2.13 where a single frame is chosen for each second, would mean that 3 key frames can be used for feature extraction and pattern matching instead of using all of the 90 frames that make up the video.

---

<sup>14</sup>Frames Per Second

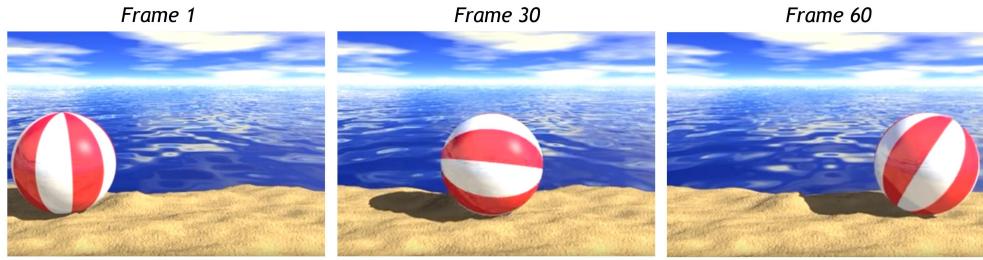


Figure 2.13: Example of frames to sample for low-visual content analysis. The first frame for each second (one frame every thirty seconds) is retrieved for a 30 fps 3-second video of a ball rolling from the left-hand side of the screen to the right-hand side. Video frames courtesy of “*How to Animate a Rolling Ball*” YouTube video available online: <https://youtu.be/cgbLAreElNI?t=130>.

### Thumbnails for Initial Shortlisting

Thumbnails can be generated for each video in a database, which are stored as additional data along with the original video file (Okabe, Dobashi and Anjyo, 2018). In a CBVR system, the thumbnails for the query video and for the database videos would be generated using the same algorithm in order to create similar outputs. This technique can be used in parallel to (Araujo and Girod, 2018), who states that an initial shortlist of potentially matching videos can be generated before the main pattern matching phase. This shortlist can be created by computing the similarities between the query video’s thumbnail and the thumbnails of the database videos.

Many advantages can be gained from this small initial step which could have an important impact on the system’s overall speed and efficiency performance. Indeed, database videos that share no similarities to the query video will not be considered at all during the main pattern matching phase e.g. if the query video corresponds to a colourful sunset, then database videos of cloudy environments will be immediately filtered out as it is unlikely that they will match with the query video in the main pattern matching phase. This step diminishes the number of videos to compare to the query.

This initial step of generating thumbnails for the query video and each video in the database is extremely speedy as it only uses a single frame that describes the entire video. Therefore the entire process will not be noticeably slowed down. However, this concept could only be applied to relatively small databases where the videos contain a single scene with multiple similar shots. Indeed, shots can be characterised by a single still as they are made up of frames with similar visual content, whereas scenes can be made up of shots that change a lot and describe different visual content. Using thumbnails to

shortlist a database of feature-length movies would greatly diminish the accuracy of the system due to the high number of scenes that make up a movie that cannot be resumed to a single thumbnail. The database videos' thumbnails will have already been generated during the database's pre-processing phase, which only occurs a single time. A possible solution could be to create a thumbnail for each shot in a scene and store each thumbnail in a list. For example, if a scene contains six different shots, then six thumbnails will be generated to describe that scene.

### 2.3.3 Feature-Length Movie Optimisation

For quicker pattern matching and feature extraction, 16 pixels separate two feature points, allowing for rough pattern matching rather than exact pattern matching. This could be used when using feature-length movies as database videos to pattern match (Okabe et al., 2018).

Fisher Vectors can be used to generate a compact signature for videos, which will then be used for computing similarities between the query video's FV and the database videos' FVs (Araujo and Girod, 2018).

## 2.4 Chapter Summary

Include a chapter summary here.

This is the chapter in which you review the major achievements in the light of your original objectives, critique the process, critique your own learning and identify possible future work.



# Bibliography

- A9: Visual Search (n.d.), <https://www.a9.com/what-we-do/visual-search.html>. [Online] Accessed: 2018-10-28.
- Amir, A., Berg, M., Chang, S.-F., Hsu, W., Iyengar, G., Lin, C.-Y., Naphade, M., Natsev, A., Neti, C., Nock, H. et al. (2003), ‘Ibm research trecvid-2003 video retrieval system’, *NIST TRECVID-2003* **7**(8), 36.
- Araujo, A. and Girod, B. (2018), ‘Large-scale video retrieval using image queries’, *IEEE Transactions on Circuits and Systems for Video Technology* **28**(6), 1406–1420.
- Awad, G., Butt, A., Curtis, K., Lee, Y., Fiscus, J., Godil, A., Joy, D., Delgado, A., Smeaton, A. F., Graham, Y., Kraaij, W., Quénot, G., Magalhaes, J., Semedo, D. and Blasi, S. (2018), Trecvid 2018: Benchmarking video activity detection, video captioning and matching, video storytelling linking and video search, in ‘Proceedings of TRECVID 2018’, NIST, USA.
- Blumberg, R. and Atre, S. (2003), ‘The problem with unstructured data’, *Dm Review* **13**(42-49), 62.
- Bradski, G. and Kaehler, A. (2008), *Learning OpenCV: Computer vision with the OpenCV library*, ” O’Reilly Media, Inc.”, pp. 493–636.  
**URL:** <http://shop.oreilly.com/product/0636920044765.do>
- Brownlow, K. (1980), ‘Silent films: What was the right speed?’, *Sight and Sound Summer*, 164–167.  
**URL:** [https://web.archive.org/web/20110708155615/http://www.cinemaweb.com/silentfilm/bookshelf/18\\_kb\\_2.htm](https://web.archive.org/web/20110708155615/http://www.cinemaweb.com/silentfilm/bookshelf/18_kb_2.htm)
- Camara-Chavez, G., Precioso, F., Cord, M., Phillip-Foliguet, S. and Araujo, A. d. A. (2007), Shot boundary detection by a hierarchical supervised approach, in ‘2007 14th International Workshop on Systems, Signals and Image Processing and 6th EURASIP Conference focused on Speech and Image Processing, Multimedia Communications and Services’, IEEE, pp. 197–200.

- Cernekova, Z., Pitas, I. and Nikou, C. (2006), ‘Information theory-based shot cut/fade detection and video summarization’, *IEEE Transactions on circuits and systems for video technology* **16**(1), 82–91.
- Chang, Y., Lee, D.-J., Hong, Y. and Archibald, J. (2007), ‘Unsupervised video shot detection using clustering ensemble with a color global scale-invariant feature transform descriptor’, *EURASIP Journal on Image and Video Processing* **2008**(1), 860743.
- Deutscher, M. (2012), ‘Big data is the new natural resource: New ibm tool powers decision making’, <https://siliconangle.com/2012/04/03/big-data-is-the-new-natural-resource-new-ibm-tool-powers-decision-making/>. [Online] Accessed: 2018-10-28.
- Fablet, R., Bouthemy, P. and Perez, P. (2002), ‘Nonparametric motion characterization using causal probabilistic models for video indexing and retrieval’, *IEEE Transactions on Image Processing* **11**(4), 393–407.
- Fanning, J., Mullen, S. P. and McAuley, E. (2012), ‘Increasing physical activity with mobile devices: a meta-analysis’, *Journal of medical Internet research* **14**(6).
- Feng, B., Cao, J., Bao, X., Bao, L., Zhang, Y., Lin, S. and Yun, X. (2011), ‘Graph-based multi-space semantic correlation propagation for video retrieval’, *The Visual Computer* **27**(1), 21–34.  
**URL:** <https://doi.org/10.1007/s00371-010-0510-6>
- Google Lens* (n.d.), <https://lens.google.com/>. [Online] Accessed: 2018-10-28.
- Hanjalic, A., Lagendijk, R. L. and Biemond, J. (1999), ‘Automated high-level movie segmentation for advanced video-retrieval systems’, *IEEE Transactions on Circuits and Systems for Video Technology* **9**(4), 580–588.
- Hauptmann, A., Baron, R. V., Chen, M.-y., Christel, M., Duygulu, P., Huang, C., Jin, R., Lin, W.-H., Ng, T. and Moraveji, N. (2004), Informedia at trecvid 2003: Analyzing and searching broadcast news video, Technical report.  
**URL:** <https://www-nlpir.nist.gov/projects/tvpubs/tvpapers03/cmu.final.paper.pdf>
- Heo, Y. S., Lee, S. and Jung, H. Y. (2016), ‘Consistent color and detail transfer from multiple source images for video and images’, *The Visual Computer* **32**(10), 1273–1289.  
**URL:** <https://doi.org/10.1007/s00371-015-1162-3>

- Hoi, S. C., Wong, L. L. and Lyu, A. (2006), Chinese university of hongkong at trecvid 2006: Shot boundary detection and video search, in ‘TRECVid 2006 Workshop’, Gaithersburg, Maryland, NIST, pp. 76–86.
- Hu, W., Xie, N., Li, L., Zeng, X. and Maybank, S. (2011), ‘A survey on visual content-based video indexing and retrieval’, *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)* **41**(6), 797–819.
- Huang, J., Kumar, S. R., Mitra, M., Zhu, W.-J. and Zabih, R. (1997), ‘Image indexing using color correlograms’, *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition* **1**, 762.
- Janwe, N. J. and Bhoyar, K. K. (2013), Video shot boundary detection based on jnd color histogram, in ‘2013 IEEE Second International Conference on Image Information Processing (ICIIP-2013)’, IEEE, pp. 476–480.
- Karr, D. (2012), ‘Big data brings marketing big numbers’, <https://martech.zone/ibm-big-data-marketing/>. [Online] Accessed: 2018-10-28.
- Lai, Y. and Yang, C. (2015), ‘Video object retrieval by trajectory and appearance’, *IEEE Transactions on Circuits and Systems for Video Technology* **25**(6), 1026–1037.
- Lars, K. (2012), ‘Improving rights management at youtube’, <https://youtube-creators.googleblog.com/2012/04/improving-rights-management-at-youtube.html>.
- Li, H. and Doermann, D. (2002), Video indexing and retrieval based on recognized text, in ‘IEEE Workshop on Multimedia Signal Processing.’, pp. 245–248.
- Liu, W., Mei, T. and Zhang, Y. (2014), ‘Instant mobile video search with layered audio-video indexing and progressive transmission’, *IEEE Transactions on Multimedia* **16**(8), 2242–2255.
- Lo, C.-C. and Wang, S.-J. (2001), ‘Video segmentation using a histogram-based fuzzy c-means clustering algorithm’, *Computer Standards & Interfaces* **23**(5), 429–438.
- Okabe, M., Dobashi, Y. and Anjyo, K. (2018), ‘Animating pictures of water scenes using video retrieval’, *The Visual Computer* **34**(3), 347–358.  
**URL:** <https://doi.org/10.1007/s00371-016-1337-6>
- Patel, B. V. and Meshram, B. B. (2012), ‘Content based video retrieval’, *The International Journal of Multimedia & Its Applications* **4**(5), 77–98.  
**URL:** <http://arxiv.org/abs/1211.4683>

- Petković, M. (2000), Content-based video retrieval, in ‘Extending DataBase Technology PhD Workshop’.
- URL:** <http://citeserx.ist.psu.edu/viewdoc/summary?doi=10.1.1.33.7317>
- Shazam* (n.d.), <https://www.shazam.com/gb/company>. [Online] Accessed: 2018-10-28.
- Sivic, J., Everingham, M. and Zisserman, A. (2005), Person spotting: video shot retrieval for face sets, in ‘International conference on image and video retrieval’, Springer, pp. 226–236.
- Smeaton, A. F., Over, P. and Kraaij, W. (2006), Evaluation campaigns and trecvid, in ‘MIR ’06: Proceedings of the 8th ACM International Workshop on Multimedia Information Retrieval’, ACM Press, New York, NY, USA, pp. 321–330.
- Villas-Boas, A. (2017), ‘Google’s new lens product uses your phone’s camera to do clever tricks, like connecting your phone to a wifi network’, <http://uk.businessinsider.com/googles-lens-feature-can-connect-your-phone-to-wifi-using-your-camera-2017-5?op=1&r=US&IR=T>.
- Wang, H., Oneata, D., Verbeek, J. and Schmid, C. (2016), ‘A robust and efficient video representation for action recognition’, *International Journal of Computer Vision* **119**(3), 219–238.
- URL:** <https://doi.org/10.1007/s11263-015-0846-5>
- Yan, R. and Hauptmann, A. G. (2007), ‘A review of text and image retrieval approaches for broadcast news video’, *Information Retrieval* **10**(4–5), 445–484.
- Yuan, J., Wang, H., Xiao, L., Zheng, W., Li, J., Lin, F. and Zhang, B. (2007), ‘A formal study of shot boundary detection’, *IEEE transactions on circuits and systems for video technology* **17**(2), 168–186.
- Zhao, Z.-C. and Cai, A.-N. (2006), Shot boundary detection algorithm in compressed domain based on adaboost and fuzzy theory, in ‘International Conference on Natural Computation’, Springer, pp. 617–626.