

Determining Computable Scenes in Films and their Structures using Audio-Visual Memory Models

Hari Sundaram Shih-Fu Chang
Dept. of Electrical Engineering,
Columbia University,
New York New York 10027.

{sundaram, sfchang}@ctr.columbia.edu

ABSTRACT

In this paper we present novel algorithms for computing scenes and within-scene structures in films. We begin by mapping insights from film-making rules and experimental results from the psychology of audition into a computational scene model. We define a computable scene to be a chunk of audio-visual data that exhibits long-term consistency with regard to three properties: (a) chromaticity (b) lighting (c) ambient sound. Central to the computational model is the notion of a causal, finite-memory viewer model. We segment the audio and video data separately. In each case we determine the degree of correlation of the most recent data in the memory with the past. The respective scene boundaries are determined using local minima and aligned using a nearest neighbor algorithm. We introduce a periodic analysis transform to automatically determine the structure within a scene. We then use statistical tests on the transform to determine the presence of a dialogue. The algorithms were tested on a difficult data set: five commercial films. We take the first hour of data from each of the five films. The best results: scene detection: 88% recall and 72% precision, dialogue detection: 91% recall and 100% precision.

Keywords

Computable scenes, scene detection, shot-level structure, films, periodic analysis transform, memory models.

1. INTRODUCTION

This paper deals with the problem of computing scenes within films using audio and visual data. We also derive algorithms for shot-level structures that exist within each scene. The problem is important for several reasons: (a) automatic scene segmentation is the first step towards greater semantic understanding of the film (b) breaking up the film into scenes will help in creating film summaries, thus enabling a non-linear navigation of the film. (c) the determination of visual structure within each scene (e.g. dialogues), will help in the process of visualizing each scene in the film summary.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ACM Multimedia 2000 Los Angeles CA USA

Copyright ACM 2000 1-58113-198-4/00/10...\$5.00

There has been prior work on video scene segmentation using image data alone [8], [19]. In [19], the authors derive scene transition graphs to determine scene boundaries. Their method assumes a presence of repetitive shot structure within a scene. While this structure is present in scenes such as interviews, it can be absent from many scenes in commercial films. This can happen, for example, when the director relies on fast succession of shots to heighten suspense or uses a series of shots to merely develop the plot of the film. In [8], the authors use a infinite, non-causal memory model to segment the video.

Prior work [14], [15], [20] concerning the problem of audio segmentation dealt with very short-term (100 ms) changes in a few features (e.g. energy, cepstra). This was done to classify the audio data into several predefined classes such as speech, music ambient sounds etc. They do not examine the possibility of using the long-term consistency found in the audio data for segmentation. Audio data has been used for identifying important regions [6] or detecting events such as explosions [9] in video skims. These skims do not segment the video data into scenes; the objective there is to obtain a compact representation.

There has been prior work on structure detection [19], [20]. There, the authors begin with time-constrained clusters of shots and assign labels to each shot. Then, by analyzing the label sequence, they determine the presence of dialogue. This method critically depends upon cluster threshold parameters that need to be manually tuned.

There are constraints on what we see and hear in films, due to rules governing camera placement, continuity in lighting as well as due to the psychology of audition. In this paper, we derive the notion of a computable scene by making use of these constraints. A computable scene exhibits long-term consistency with respect to three properties: (a) chromatic composition of the scene (b) lighting conditions and (c) ambient audio. We term such a scene as *computable*, since it can be reliably computed using low-level features alone. In this paper, we do not deal with the *semantics* of a scene. Instead, we focus on the idea of determining a computable scene, which we believe is the first step in deciphering the semantics of a scene.

We present algorithms for determining computable scenes and periodic structures that may exist within such scenes. We begin with a idea of a memory model [8]. Our memory model is causal and finite. The model has two parameters: (a) an analysis window that stores the most recent data (the attention span) (b) the total amount of data (memory).

In order to segment the data into audio scenes, we compute correlations amongst the envelopes of the audio features in the

attention-span with feature envelopes in the rest of the memory. The video data comprises shot key-frames. The key-frames in the attention span are compared to the rest of the data in the memory to determine a coherence value. This value is derived from a color-histogram dissimilarity. The comparison takes also into account the relative shot length and the time separation between the two shots. In both cases, we use a local minima for detecting a scene change and the audio and video scene boundaries are aligned using a simple time-constrained nearest neighbor approach.

We introduce the idea of a periodic analysis transform to determine visual structure within each computable scene. The transform computes the degree of periodicity amongst a time-ordered sequence of images (key-frames of shots). We use the Student's t-test in conjunction with a simple rule on this transform, to detect the presence of a dialogue. In contrast to [19], [20] which require manually tweaked cluster diameter threshold parameters, this algorithm is almost parameter free. Our experiments show that the scene change detector and the intra-scene structure detection algorithm show good results.

The rest of this paper is organized as follows. In the next section, we formalize the definition of a computable scene. In section 3, we present an algorithm for the detection of such computational scenes. In section 4, we derive algorithms for automatically determining periodic structures within a scene. In section 5, we present our experimental results. In section 6 we discuss shortcomings of our model and finally in the section 7, we present our conclusions.

2. WHAT IS A COMPUTABLE SCENE?

In this section we shall define the notion of a computable scene. We begin with a few insights obtained from understanding the process of film-making and from the psychology of audition. We shall use these insights in creating our computational model of the scene.

2.1 Insights from Film Making Techniques

The line of interest is an imaginary line drawn by the director in the physical setting of a scene [4]. During the filming of the scene, all the cameras are placed on one side of this line (also referred to as the 180 degree rule). This is because we desire successive shots to maintain the spatial arrangements between the characters and other objects in the location. As a consequence there is no confusion in the mind of the viewer about the spatial arrangements of the objects in the scene. He (or she) can instead concentrate on the

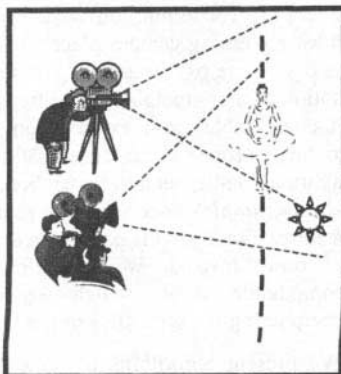


Figure 1: showing the line of interest (thick line) in a scene. We also see the fields-of-view of the two cameras intersecting.

dramatic aspect of the scene. It is interesting to note that directors willingly violate this rule only in very rare circumstances¹.

The 180 degree rule has interesting implications on the computational model of the scene. Since all the cameras in the scene remain on the same side of the line in all the shots, there is an overlap in the field of view of the cameras (see Figure 1). This implies that there will be a consistency to the chromatic composition and the lighting in all the shots.

Film-makers seek to maintain continuity in lighting amongst shots within the same physical location. This is done even when the shots are filmed over several days. This is because viewers perceive the change in lighting as indicative of the passage of time. For example, if two characters are shown talking in one shot, in daylight, the next shot cannot show them talking at the same location, at night.

2.2 The Psychology of Audition

The term *auditory scene analysis* was coined by Bregman in his seminal work on auditory organization [1]. In his psychological experiments on the process of audition, Bregman made many interesting observations, a few of which are reproduced below:

- Unrelated sounds seldom begin and end at the same time.
- A sequence of sounds from the same source seem to change its properties smoothly and gradually over a period of time. The auditory system will treat the sudden change in properties as the onset of a new sound.
- Changes that take place in an acoustic event will affect all components of the resulting sound in the same way and at the same time. For example, if we are walking away from the sound of a bell being struck repeatedly, the amplitude of all the harmonics will diminish gradually. At the same time, the harmonic relationships and common onset² are unchanged.

Bregman also noted that different auditory cues (i.e. harmonicity, common-onset etc.) compete for the user's attention and depending upon the context and the knowledge of the user, will result in different perceptions. However, the role played by higher forms of knowledge in grouping is yet to be ascertained.

Different computational models (e.g. [3]) have emerged in response to those experimental observations. While these models differ in their implementations and differ considerably in the physiological cues used, they focus on short-term grouping strategies of sound. Notably, Bregman's observations indicate that long-term grouping strategies are also used by human beings (e.g. it is easy for us to identify a series of footsteps as coming from one source) to group sound.

¹ This is so infrequent that directors who transgress the rule are noted in the film theory community. e.g. Alfred Hitchcock willingly violates this rule in his film *North by Northwest* thus adding suspense to the scene [4].

² Different sounds emerging from a single source begin at the same time.

2.3 The Computable Scene Model

The constraints imposed by production rules in film and the psychological process of hearing lead us to the following definition of a scene: It is a continuous segment of audio-visual data that shows *long-term*³ consistency with respect to three properties:

- Chromaticity
- Lighting conditions
- Ambient sound

Table 1: Several scenarios are examined using our c-scene definition. These would normally be viewed as a single “normal” scene.

Scenario	Shot-sequence	C-Scenes	Explanation
Alice goes home to read a book.	(a) She is shown entering the room. (b) she picks up the book. (c) we see her reading silently (d) while she is reading, we hear the sound of rain.	2	One consistent visual but two consistent chunks of audio.
Alice goes to sleep.	(a) She is shown reading on her bed. (b) she switches off the light and room is dark.	2	Two consistent visuals but the audio is consistent over both video segments.
Bob goes for a walk	He switches on his handy-cam inside the house and walks out of the house. Note, this is one single camera take.	2	There are two consistent visuals (inside/outside) as well as two consistent chunks of audio.

We denote this to be a *computable* scene since these properties can be reliably and automatically determined using low-level features present in the audio-visual data. We need to examine the relationship between a computable scene (abbreviated as c-scene) and normal notions of a shot and a scene. A shot is a segment of audio-visual data filmed in a single camera take. A scene is normally defined to be sequence of shots that share a common semantic thread. Table 1 examines the impact of the c-scene definition for several scenarios.

The semantics of a normal scene within a film, are often difficult to ascertain. While a collection of shots may have objects that are meaningful without context (e.g. a house, a man, a woman the colors of the dress etc.), the collection of shots are infused with meaning only with regard to the context.

The context is established due to two factors: the film-maker and the viewer. The film-maker infuses meaning to a collection of

³ Analysis of experimental data (one hour each, from five different films) indicates that the scenes in the same location (e.g. in a room, in the marketplace etc.) are typically 40~50 seconds long.

shots in three ways: (a) by *deciding* the action in the shots (b) the kind of shots that precede this scene and the shots that follow it (c) and finally by the manner in which he *visualizes*⁴ the scene. All three methods affect the viewer, whose *interpretation* of the scene depends on his world-knowledge. Hence, if the meaning in a scene is based on factors that cannot be measured directly, it is imperative that we begin with a scene definition in terms of those attributes that are measurable and which lead to a consistent interpretation. We believe that such a strategy will greatly help in deciphering the semantics of the c-scene at a later stage.

2.4 The C-Scene Definition

We wished to validate the computable scene definition, which appeared out of intuitive considerations, with actual film data. The data was diverse with one hour segments from three English language films and two foreign films⁵.

The definition for a scene works very well in many film segments. In most cases, the c-scenes are usually a collection of shots that are filmed in the same location and time and under similar lighting conditions. The definition does not work well for montage⁶ sequences. However, in such sequences, we observed a long-term consistency of the ambient audio. We need to define a c-scene in order to accommodate different production styles. We now make two distinctions:

1. **N-type:** These scenes (or normal scenes) fit our original definition of a scene: they are characterized by a long-term consistency of chromatic composition, lighting conditions and sound.
2. **M-type:** These scenes (or montage/Mtv scenes) are characterized by widely different visuals (differences in location, time of creation as well as lighting conditions) which create a unity of theme by manner in which they have been juxtaposed. However, M-type scenes will be assumed to be characterized by a long-term consistency in the audio track. Transient scenes are M-type scenes that are characterized by shots of long duration⁷.

In this paper, we narrow our focus to derive algorithms that detect two adjacent N-type scenes. We will not handle the two cases when we have either (a) two adjacent M-type scenes or (b) an N-type scene that borders an M-type scene. Analysis of the ground truth indicates that these two transitions constitute about 25% of all the transitions. Henceforth, for the sake of brevity, we shall use the term “scene” for our notion of a computable scene (c-scene).

⁴ In order to show a tense scene, one film-maker may have fast succession of close-ups of the characters in a scene. Others may indicate tension by showing both characters but changing the music.

⁵ The English films: *Sense and Sensibility*, *Pulp Fiction*, *Four Weddings and a Funeral*. The foreign films: *Farewell my Concubine (Chinese)*, *Bombay (Hindi)*.

⁶ In classic Russian montage, the sequence of shots are constructed from placing shots together that have no immediate similarity in meaning. For example, a shot of a couple may be followed by shots of two parrots kissing each other etc. The meaning is derived from the way the sequence is arranged.

⁷ Mtv videos are good examples of M-type scenes with shots of short duration. Transient scenes can occur when the director wants to show the passage of time e.g. a scene showing a journey.

3. DETECTING SCENES

We begin the process of scene detection by first detecting audio and video scene segments separately and then aligning the two by a simple nearest neighbor algorithm.

This section has four subsections. In section 3.1, we develop the idea of a memory model. In sections 3.2 and 3.3, we build upon some early techniques in [16], [17] for automatic audio and video scene detection. In section 3.4 we present a simple nearest-neighbor algorithm for aligning the two scene detector results.

3.1 A Memory Model

In order to segment data into scenes, we use a causal, first-in-first-out (FIFO) model of memory (figure 2). This model is derived in part from the idea of coherence [8].

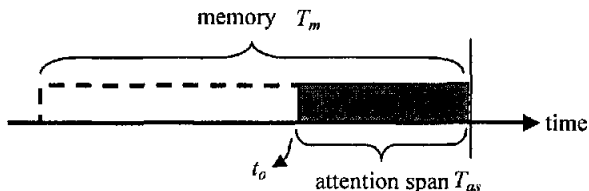


Figure 2: The attention span (T_{as}) is the most recent data in the buffer. The memory (T_m) is the size of the entire buffer. Clearly, $T_m \geq T_{as}$.

In our model of a listener, two parameters are of interest: (a) memory: this is the net amount of information (T_m) with the viewer and (b) attention span: it is the most recent data (T_{as}) in the memory of the listener. This data is used by the listener to compare against the contents of the memory in order to decide if a scene change has occurred.

The work in [8] dealt with a non-causal, infinite memory model based on psychophysical principles, for video scene change detection. We use the same psychophysical principles to come up with a causal and finite memory model. Intuitively, causality and a finite memory will more faithfully mimic the human memory-model than an infinite model. We shall use this model for *both* audio and video scene change detection.

3.2 Determining Audio Scenes

In this section we present our algorithm for audio-scene segmentation. We model the *audio-scene* as a collection of a few dominant sound sources. These sources are assumed to possess stationary properties that can be characterized using a few features. An audio-scene change is said to occur when the majority of the dominant sources in the scene change. A detailed description of audio scene segmentation can be found in [16].

3.2.1 Features and Envelope Models

We use ten different features [13], [14], [15], [16], [20] in our algorithm: (a) cepstral-flux (b) multi-channel cochlear decomposition (c) cepstral vectors (d) low energy fraction (e) zero crossing rate (f) spectral flux (g) energy (h) spectral roll off point. We also use the variance of the zero crossing rate and the variance of the energy as additional features. The cochlear decomposition was used because it was based on a psychophysical ear model. The cepstral features are known to be good discriminators [13]. All the other features were used for their ability to distinguish

between speech and music [14], [15], [20]. Features are extracted *per frame* (100ms. duration) for the duration of the analysis window.

Given a particular feature f and a finite time-sequence of values, we wish to determine the behavior of the envelope of the feature. The feature envelopes are force-fit into signals of the following types: constant, linear, quadratic, exponential, hyperbolic and sum of exponentials. All the envelope (save for the sum of exponentials case) fits are obtained using a robust curve fitting procedure [5]. We pick the fit that minimizes the least median error. The envelope model analysis is only used for the scalar variables. The vector variables (cepstra and the cochlear output) and the aggregate variables (variance of the zero-crossing rate and the spectral roll off point) are used in the raw form.

3.2.2 Detecting a Scene Change

Let us examine the case where a scene change occurs just to the left of the listeners attention span. First, for each feature, we do the following:

1. Place an analysis window of length T_{as} (the attention-span length) at t_o and generate a sequence by computing a feature value for each frame (100 ms duration) in the window.
2. Determine the optimal envelope fit for these feature values.
3. Shift the analysis window back by Δt and repeat steps 1. and 2. till we have covered all data in the memory.

We then define a local correlation function per feature, using the sequence of envelope fits. The correlation function C_f for each feature f is then defined as follows:

$$C_f(m\delta) = 1 - d(f(t_o, t_o + t_{as}), f(t_o + m\delta, t_o + m\delta + t_{as})) \quad (1)$$

where, $f(t_1, t_2)$ represents the envelope fit for feature f for the duration $[t_1, t_2]$. Now, $m \in [-N, 0]$, where $N \equiv (T_m - T_{as})/\delta$. δ is the duration by which the analysis window is shifted back and d is the Euclidean metric⁸ on the envelopes. For the vector and the aggregate data, we do not compute the distance between the windows using envelope fits but use a L^2 metric on the raw data. In our experiments we use $\delta = 1$ sec.

We model the correlation decay as a decaying exponential [16]: $C_i(t) = \exp(b_i t)$, $t < 0$ where C_i is the correlation function for feature i , and b_i is the exponential decay parameter. The audio-scene decision function $D(t_o)$ at any instant t_o is defined as

follows: $D(t_o) = \sum_{i=1}^N b_i$. Where, N is the number of features and b_i are the estimates at t_o .

We chose the exponential decay model based on empirical observations on the correlation data. More sophisticated modeling techniques could be easily employed. The parameter δ , effectively determines the number of samples of the correlation data. This in turn affects our estimate of the parameter b_i . For example, in a

⁸This metric is intuitive: it is a point-by-point comparison of the two envelopes. More sophisticated predictor based schemes are being investigated at present.

typical case of memory (T_m) size of 31 sec and attention span (T_{as}) of 16 sec. and $\delta=1$ sec., we have 16 data points.

The audio-scene change is detected using the local minima of the decision function. In order to do so, we use a sliding window of length $2w_a+1$ sec. to slide across the data. We then determine if the minima in the window coincides with the center of the window. If it does, the location is labeled as an audio scene change location. The result for a single film is shown figure 3. The figure shows that the results agree within an ambiguity window of w_a sec.

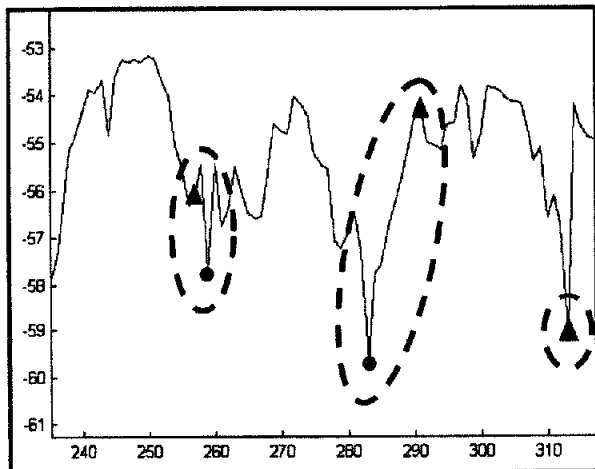


Figure 3: Audio detector results for a part of the audio track of the film *Sense and Sensibility*. The triangles show the ground truth label while the dots show the detector result. In the first two cases the result is very close while in the third case there is an exact match. The x -axis shows the time in sec. while y -axis shows the detector magnitude.

3.3 Determining Video Scenes

In this section, we shall describe the algorithm for video-scene segmentation. The algorithm is based on notions of *recall* and *coherence*. We model the video-scene as a contiguous segment of visual data that is chromatically coherent and also possesses similar lighting conditions. A video-scene is said to occur when there is a change in the long-term chromaticity and lighting properties in the video. This stems from the film-making constraints discussed in Section 2.1.

Ideally, we would like to work with raw frames and avoid having to detect shots. However, this would lead to an enormous increase in the computational complexity of the algorithm. Hence, the video stream is converted into a sequence of shots using a simple color and motion based shot boundary detection algorithm [10]. A frame at a fixed time after the shot boundary is extracted and denoted to be the key-frame.

3.3.1 Recall

In our visual memory model, the data is in the form of key-frames of shots (Figure 4) and each shot occupies a definite span of time. The model also allows for the most recent and the oldest shots to be partially present in the buffer. A point in time (t_o) is defined to be a scene transition boundary if the shots that come after that point in time, do not recall [8] the shots prior to that

point. The idea of recall between two shots a and b is formalized as follows:

$$R(a,b) = (1 - d(a,b)) \cdot f_a \cdot f_b \cdot (1 - \Delta t / T_m), \quad (2)$$

where, $R(a,b)$ is the recall between the two shots a, b . $d(a,b)$ is a L^1 color-histogram based distance between the key-frames corresponding to the two shots, f_i is the ratio of the length of shot i to the memory size (T_m). Δt is the time difference between the two shots.

The formula for recall indicates that recall is proportional to the length of each of the shots. This is intuitive since if a shot is in memory for a long period of time it will be recalled more easily. Again, the recall between the two shots should decrease if they are further apart in time.

We need to introduce the notion of a “shot-let.” A shot-let is a fraction of a shot, obtained by breaking individual shots into δ sec. long chunks but could be smaller due to shot boundary conditions. Each shot-let is associated with a single shot and its representative frame is the key-frame corresponding to the shot. In our experiments, we find that $\delta = 1$ sec. works well. Figure 4 shows how shot-lets are constructed. The formula for recall for shot-lets is identical to that for shots.

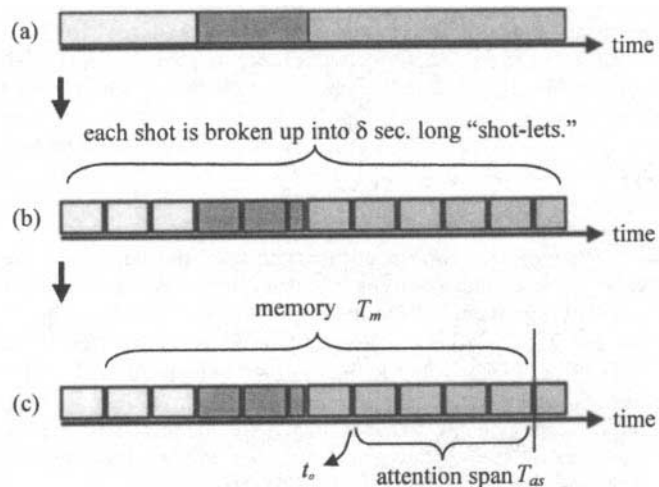


Figure 4: (a) Each solid colored block represents a single shot. (b) each shot is broken up into “shot-lets” each at most δ sec. long. (c) the bracketed shots are present in the memory and the attention span. Note that sometimes, only fractions of shots are present in the memory.

3.3.2 Computing Coherence

Coherence is easily defined using the definition of recall:

$$C(t_o) = \left(\sum_{a \in T_{as}} \sum_{b \in (T_m \setminus T_{as})} R(a,b) \right) / C_{\max}(t_o) \quad (3)$$

where, $C(t_o)$ is the coherence across the boundary at t_o and is just the sum of recall values between all pairs of shot-lets across the boundary at t_o . $C_{\max}(t_o)$ is obtained by setting $d(a,b)=0$ in the formula for recall (Equation (2)) and re-evaluating the numerator of Equation (3). This normalization compensates for the different number of shots in the buffer at different instants of time.

We compute coherence at the boundary between every adjacent pair of shot-lets. Then, similar to the procedure for audio scene detection, we determine the local minima. This we do by using a sliding window of length $2w_v+1$ sec. and determine if the minima in the window coincides with the center of the window. If it does, the location is labeled as a video scene change location.

3.3.3 The need for shot-lets

Shot-lets become necessary in films since they can contain c-scenes with shots that have a long duration. In [8], the authors evaluated their coherence function only at shot boundaries. Evaluating the coherence function ($C(t)$, Equation (3)) at shot boundaries has the effect of coarsely sampling the coherence function. This can cause problems in locating the scene change minima. For example, if in a segment of size 90 sec. if we have three shots, of 30 sec. each, where the first two belong to one scene and the third to another, then we cannot determine the minima as we will only have two values for coherence. A thought will indicate that interpolating the coherence values will not be of much use.

Shot-lets have two main advantages: (a) they preserve the location of existing shot boundaries and (b) they help us evaluate the coherence function at a fine time-scales. Simply uniformly segmenting the video stream into δ sec. chunks has two disadvantages: (a) shot boundaries are missed and (b) high computational complexity. Evaluating Equation (3) takes $O(k^2)$ operations where k is the number of shots in the buffer⁹. The idea of shot-lets can be shown to significantly improve the detection rate results in [17], [8].

3.4 Aligning the Detector Results

We generate correspondences between the audio and the video scene boundaries using a simple time-constrained nearest-neighbor algorithm. Let the list of video scene boundaries be V_i , where $i \in \{1..N_v\}$. Let the list of audio scene boundaries be A_i , where $i \in \{1..N_a\}$. The ambiguity window around each video scene is w_v sec. long. The ambiguity window width around each audio scene boundary is w_a sec long. Note that these sizes are the same size of the windows used for local minima location. For each video scene boundary, do the following:

- Determine a list of audio scene boundaries whose ambiguity windows intersect the ambiguity window of the current video scene boundary.
- If the intersection is non-null, pick the audio scene boundary closest to the current video scene boundary. Remove this audio scene boundary from the list containing audio scene boundaries.
- If the intersection is null, add the current video scene boundary to the list of singleton (i.e. non-alignable) video scene changes.

⁹ The primary contributor to the computational complexity of Equation (3) is the distance computation in Equation (2). A little thought indicates that with some book-keeping, this is avoided with shot-lets.

At the end of this procedure, if there are audio scene boundaries left, collect them and add them to the list of singleton audio scene changes. Figure 5 illustrates this scenario.

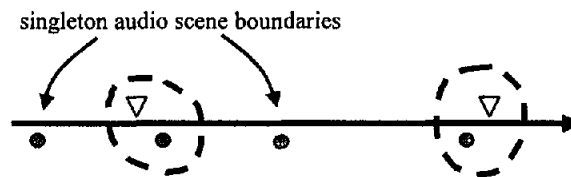


Figure 5: The figure shows video (triangles) and audio (solid circles) scene change locations. The dashed circles show audio/video scene boundaries which align.

Films exhibit interesting interactions between audio and video scene changes. Singleton audio and video scene boundaries can be caused due to the following reasons:

1. **Audio scene change but no video scene change:** this can happen for example : the director wants to indicate a change the mood of the scene, by using a sad / joyous sounding audio track. In *Sense and Sensibility*, one character was shown singing and once she was finished, we had conversation amongst the characters *in the same location*.
2. **Video scene changes within an audio scene:** This happens when a sequence of video scenes have the same underlying semantic . For example, we can have a series of video scenes showing a journey and these scenes will be accompanied by the same audio track.

Now that we have determined the scene boundaries, we will now present algorithms that determine structure within a scene.

4. THE SCENE LEVEL STRUCTURE

In this section we shall discuss possible structures that could exist within a scene and technique to detect and classify such structures. The detection of these structures will help in summarizing the scene. Here, we focus on detection of visual structures.

4.1 Postulating Scene-level Structures

We postulate the existence of two broad category of scenes: N-type (based on the initial definition) and the M-type scene. The N-type scenes are further subdivided into three types: (a) pure dialogue (b) progressive and (c) hybrid. We use an abstract graph representation for representing the shot structure within a scene. Each node in the graph represents one cluster of shots. Figure 6 shows a hybrid scene containing an embedded dialogue.

4.1.1 N-type Scenes

An N-type scene has unity of location, time and sound. We now look at three sub-categories:

Dialogue: A simple repetitive visual structure (amongst shots) can be present if the action in the scene is a dialogue. Note that sometimes, directors will *not* use an alternating sequence to represent a dialogue between two characters. He (or she) may use a single shot of long duration that shows both the characters talking. A repetitive structure is also present when the film-maker

shuttles back and forth between two shots to indicate an idea (e.g. man watching television). We denote this as a thematic dialogue.

Progressive: There can be a linear progression of visuals without any repetitive structure (the first part of figure 6 is progressive). For example, consider the following scene: Alice enters the room looking for a book. We see the following shots (a) she enters the room (b) she examines her book-shelf (c) looks under the bed (d) locates the book and the camera follows her as she sits on the sofa to read.

Hybrid: This is the most common case, when we have a dialogue embedded in an otherwise progressive scene. For example, in the scene mentioned above, assume Bob enters the room while Alice is searching for the book. They are shown having a brief dialogue that is visualized using an alternating sequence. Then Bob leaves the room and Mary continues her search.

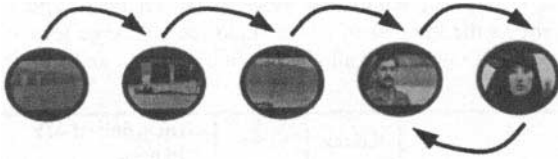


Figure 6: A hybrid scene with an embedded dialogue sequence.

4.1.2 M-type Scenes

In M-type scenes (in classic montage, commercials and MTV videos) we assume there to be no unity of visuals either in terms of location, time or lighting conditions¹⁰. However, we expect that the audio track will be consistent over the scene. This condition can be converted into a detection rule: A sequence of highly dissimilar shots with unity of sound will be labeled as a M-type scene. As noted earlier, transient scenes are M-type scenes with long shot durations.

4.2 Determining the Structure

In this section we shall describe techniques to identify structures within N-type scenes. We begin by first describing the periodic analysis transform. Then we show how to use this series in conjunction with statistical tests to determine the presence of a dialogue. Finally we show a simple algorithm that determines the exact location of the dialogue.

4.2.1 The Periodic Analysis Transform

The periodic analysis transform helps us estimate the periodicity in an time-ordered sequence of N key-frames. Let o_i , where $i \in \{0, N-1\}$ be a time ordered sequence of key-frames. Then:

$$\Delta(n) \triangleq 1 - \frac{1}{N} \sum_{i=0}^{N-1} d(o_i, o_{\text{mod}(i+n, N)}), \quad (4)$$

where, $\Delta(n)$ is the transform, d is the L^1 color-histogram based distance function, mod is the usual modulus function. The modulus function simply creates a periodic extension¹¹ of the

¹⁰ There will be a unity of theme which shall be brought about by how the director assembles the component shots of the scene.

¹¹ Defining the transform using a symmetric extension should improve our detector results.

original input sequence. Note that the transform definition will work on any time ordered sequence of arbitrary objects, provided we define a suitable metric on the objects.

4.2.2 Statistical Tests

We shall use two statistical tests: the students t-test for the means and the F-test for the variances [12]. The F-test is used to determine the appropriate¹² Student's t-test. These tests are used to compare two series of numbers and determine if the two means and the variance differ significantly.

4.2.3 Detecting Dialogues

We can easily detect dialogues using the transform. In a dialogue, every 2nd frame will be very similar while adjacent frames will differ. This is also to be observed in figure 7. Let us assume that we have a time-ordered sequence of N key-frames representing different shots in a scene. Then we do the following:

1. Compute the series $\Delta(n)$.
2. Check if $\Delta(2) > \Delta(1)$ and $\Delta(2) > \Delta(3)$.
3. A dialogue is postulated to exist if one of two conditions in step 2 is at least significant at $\alpha = 0.05$ and the other one is at least significant at $\alpha = 0.1$ ¹³. Note that $\Delta(n)$ for each n is the mean of N numbers. We use the Student's t-test to determine whether the two means are different in a statistically significant sense.

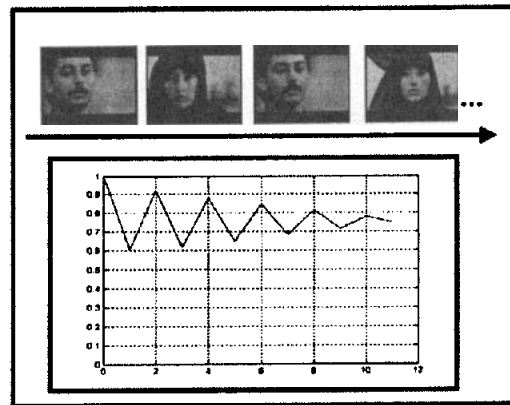


Figure 7: A dialogue scene and its corresponding periodic analysis transform. Note the distinct peaks at $n = 2, 4, \dots$

We use a simple technique to make a distinction between thematic and actual spoken dialogue. From test data we observe that the average shot length for a thematic dialogue is much shorter than for a spoken dialogue. The reason is that there is a minimum time required to utter a meaningful phrase. In [7], the authors assume that phrases last between 5~15 sec. An analysis of hand-labeled data reveals that dialogues with average shot length of less than 4 sec. are thematic.

¹² There are two Student's t-tests depending upon whether the variances differ significantly.

¹³ We are rejecting the null hypothesis that the two means are equal. We reject the hypothesis if we believe that the observed difference between the means occurred by chance with a probability less than α .

4.2.4 The Sliding Window Algorithm

We use a sliding window algorithm to detect the presence of a dialogue (thematic or spoken dialogue) within a hybrid N-type scene. Assuming that the total number of frames in the scene to be N, we set the size of the initial window to be k frames. Starting with the leftmost key-frame, the algorithm is as follows:

1. Run the dialogue detector on the current window.
2. If no dialogue is detected, keep shifting the window to the right by one key-frame to the immediate right until either a dialogue has been detected or we have reached the end of the scene.
3. If a dialogue has been detected, keep growing the window by adding the key-frame to the immediate right of the current window until either the end of the scene has been reached or the sequence of key-frames in the window is no longer a statistically significant dialogue. The dialogue is the longest statistically significant sequence.
4. Move the start of the new window to the immediate right of the detected dialogue. Go to step 1.

Setting the initial window size in terms of number of key-frames instead of time has the advantage of compensating for shot length variations across films.

5. EXPERIMENTAL RESULTS

In this section we shall discuss the experimental results of our algorithms. The data used to test our algorithms is complex: we have five one hour segments from five diverse films. There are three English films: (a) *Sense and Sensibility* (b) *Pulp Fiction* and (c) *Four Weddings and a Funeral*. We also have two foreign language films: (a) *Bombay* (Hindi) and (b) *Farewell my Concubine* (Chinese).

This section is organized as follows. We begin with a section that explains how the labeling of the ground truth data was done. The section following that section deal with the experimental results of the two detectors.

5.1 Labeling the Ground Truth

The audio and the video data were labeled separately (i.e. label audio without watching the video and label video without hearing the audio). This was because when we use *both* the audio and the video (i.e. normal viewing of the film) we tend to label scene boundaries based on the semantics of the scene. Only one person (the first author) labeled the data.

The video scene changes were labeled as follows. While watching the video if at any time there was a distinct change in lighting or color, this was labeled as a video scene change. This usually meant a change in location, but walking from a lit room to a dark one was also labeled as a scene change. Scene boundaries for M-type scenes (transient, montage) were difficult to locate. For example, if there was an M-type \rightarrow N-type transition, we need to see a few shots from the N-type scene, to determine that the first N-type shot was not part of the previous M-type scene.

For audio, we adopted the following policy: label a scene change if it was felt that the ambient audio properties had changed. For

example, if we heard the sounds of a marketplace immediately followed a conversation, this was labeled as a scene change. Correctly labeling the audio scene boundaries is challenging since we don't see the associated video. Often, with the beginning and the end of dialogues since there is silence, it becomes very hard to place the boundary accurately¹⁴. We need to wait to wait till the end of the silence and then work back to place the boundary.

Labeling the data became particularly challenging when labeling the Chinese film. Since the labeler (the first author) had no semantic understanding of the film, it became hard to determine if a conversation had ended if there was a pause after the last sentence or if the speakers had changed (if one dialogue sequence followed the other).

Table 2: The ground truth data derived from labeling the audio and video data of each film separately. Columns two and three show the number of audio and video scene changes. The last column shows the number of audio/video scene change locations that align. The table shows all changes including to and from M-type scenes.

Film	Audio	Video	Synchronized A/V Changes
Bombay	77	46	33
Farewell my Concubine	91	58	44
Four Weddings and a Funeral	76	57	37
Pulp Fiction	45	39	31
Sense and Sensibility	52	65	41

It is clear from the data in Table 2 that the audio and video scene changes in a film, are not random events i.e. there is a high degree of synergy between the two thus lending support to our joint audio-visual computable scene model.

5.2 Scene Change Detector Results

There are three parameters of interest in each scene change algorithm (i.e. audio and video). They are: (a) memory (T_m) (b) attention-span (T_{at}) and the (c) ambiguity-window size. For both audio and video scene change algorithms, the attention-span and the memory parameters follow intuition: results improve with a large attention-span and a large memory. For both scene change algorithms, large windows have the property of smoothing the audio decision function and the video coherence function. Larger windows decrease the number of false alarms but also increase the number of misses.

The audio and video ambiguity parameters¹⁵ are used in the location of local minima in both scene change algorithms. Hits and misses are determined by looking at whether the time

¹⁴ The accuracy in labeling that we refer to is with a comparison to the where the label would have been had we labeled the film with both audio and the video.

¹⁵ This is half the size of the windows used for location of the audio and video minima (w_a and w_v sec. respectively).

difference between the scene change location and the ground truth location is less than the ambiguity window size (i.e. less than w_a and w_v sec.). The memory buffer parameters for the entire data set was fixed as follows: audio: $T_m=31$ sec. $T_{as}=16$ sec., video: $T_m=16$ sec., $T_{as}=8$ sec.

We now present results for the five films in Table 3. These results are for two adjacent N-type transitions only since our algorithms cannot handle N-type \rightarrow M-type or M-type \rightarrow M-type transitions. Note that: recall = hits/(hits + misses) while precision = hits/(hits + false alarms).

Table 3: The table shows c-scene change detector results for the five films. We only deal with adjacent N-type scenes, the other transitions were manually excluded. The columns are: Hits, Misses, False Alarms, Recall and Precision.

Film	H	M	FA	Recall	Precision
Bombay	24	3	9	0.88	0.72
Farewell my Concubine	28	9	10	0.75	0.73
Four Weddings and a Funeral	17	11	4	0.60	0.80
Pulp Fiction	19	9	11	0.67	0.63
Sense and Sensibility	27	7	7	0.79	0.79

The results show that our detector works well, achieving a best result of recall of 0.88 and precision of 0.72 for the film *Bombay*. There are two types of errors that decrease our algorithm performance: (a) uncertainty in the location of the audio labels due to human uncertainty and (b) misses in the video shot boundary detection algorithm. Shot misses cause the wrong key-frame to be present in the buffer, thus causing an error in the minima location.

Prior work done in video scene segmentation used visual features alone [19], [8]. There, the authors focus on detecting scene boundaries for sitcoms (and other TV shows) and do not consider films. However, since we expect the c-scenes in sitcoms to be mostly long, coherent, N-type scenes, we expect our combined audio visual detector to perform very well.

5.3 Structure Detection Results

The statistical tests that are central to the dialogue detection algorithm make it almost parameter free. These test are used at the standard levels of significance ($\alpha = 0.05$). We do need to set two parameters: The initial sliding window size T_w (8 frames) and the threshold for the thematic dialogue test (4 sec.).

The results of the dialog detector (Table 4) show that it performs very well. The best result is a precision of 1.00 and recall of 0.91 for the film *Sense and Sensibility*. The misses are primarily due to misses by the shot-detection algorithm. Missed key-frames will cause a periodic sequence to appear less structured. The thematic/true dialog detector's performance is mixed: with a best detection result (precision) of 0.84 for the Indian film *Bombay* and a worst result of 0.50 for *Pulp Fiction*. Thematic dialogues seem to vary significantly with the film genre and the directorial style; hence instead of global time threshold, an algorithm that

adapts to the film, perhaps based on the average shot length will work better.

Table 4: The table shows the dialogue detector results for the five films. The table includes both thematic and spoken dialogues. The columns are: Hits, Misses, False Alarms, Recall and Precision.

Film	H	M	FA	Recall	Precision
Bombay	10	2	0	0.83	1.00
Farewell my Concubine	10	2	1	0.83	0.90
Four Weddings and a Funeral	16	4	1	0.80	0.94
Pulp Fiction	11	2	2	0.84	0.84
Sense and Sensibility	28	3	0	0.91	1.00

6. DISCUSSING MODEL BREAKDOWNS

In this section we shall discuss three situations that arise in different film-making situations. In each instance, the 180 degree rule is adhered to and yet our assumption of chromatic consistency across shots is no longer valid.

1. **Change of scale:** Rapid changes of scale cannot be accounted for in simple model as they show up as change in the chrominance of the shot. For example, the director might show two characters talking in a medium-shot¹⁶. Then he cuts to a close up. This causes a change in the dominant color of the shots.
2. **Widely differing backgrounds:** This results from the two opposing cameras having no overlap in their field-of-view causing an apparent change in the background. This can happen for example, when the film shows one character inside the house, talking through a widow to another character who is standing outside.
3. **Background changes with time:** This can happen for example if the film shows several characters talking in a party (or in a crowd) . Then the stream of people in motion can cause the dominant chrominance/lighting of the scene to change.

If these situations occur over long time scales, they will cause errors (misses, incorrect boundary placement) in the segmentation algorithm. However, short time scale chromatic (or lighting) changes will be handled by our algorithm. Clearly, our computational model makes simplifying assumptions on the possible scenarios even when film-makers adhere to the 180 degree rule.

7. CONCLUSIONS

In this paper we have presented a novel paradigm for film segmentation using audio and video data and an algorithm for visual structure detection within scenes. We developed the notion of computational scenes. The computational model for the c-scenes was derived from camera placement rules in film-making and from experimental observations on the psychology of

¹⁶ The size (long/medium/close-up/extreme close-up) refers to the size of the objects in the scene relative to the size of the image.

audition. A c-scene exhibits long-term consistency with regard to (a) lighting conditions (b) chromaticity of the scene (c) ambient audio. We believe that the c-scene formulation is the first step towards deciphering the semantics of a scene.

We showed how a causal, finite memory model formed the basis of our scene segmentation algorithm. In order to determine audio scene segments we first determine the correlations amongst the envelope fits for each feature extracted in the memory buffer. We then determine the decision function based on exponential fits to the envelope correlations. We use ideas of recall and coherence in our video segmentation algorithm. The algorithm works by determining the coherence amongst the shot-lets in the memory. A local minima criterion determines the scene change points and a nearest neighbor algorithm aligns the scenes.

We introduced the formulation of the periodic analysis transform to determine the periodic structure within a scene. We showed how one can use the Student's t-test to detect the presence of statistically significant dialogues.

The scene segmentation algorithms were tested on a difficult test data set: five hours from commercial films. They work well, giving a best scene detection result of 88% recall and 72% precision. The structure detection algorithm was tested on the same data set giving excellent results: 91% recall and 100% precision. We believe that the results are very good when we keep the following considerations in mind: (a) the data set is complex (b) the audio ground truth labeling was difficult and introduced errors (c) the shot cut detection algorithm had misses that introduced additional error.

There are some clear improvements possible to this work: (a) the computational model for the c-scene is limited, and needs to be tightened in view of the model breakdowns pointed out in section 6. (b) we need to come up with a technique that handles N-type scenes that abut M-type scenes and also the case when M-type scenes are in succession. A possible solution is to introduce a short-term self-coherence function followed by audio-scene based grouping.

8. ACKNOWLEDGEMENTS

The authors would like to thank Di Zhong for help with the shot boundary detection algorithm.

9. REFERENCES

- [1] A.S. Bregman *Auditory Scene Analysis: The Perceptual Organization of Sound*, MIT Press, 1990.
- [2] M. Christel et. al. *Evolving Video Skims into Useful Multimedia Abstractions* Proc. of the Conference on Human Factors in Computing System, CHI'98, pp 171-178, Los Angeles, CA, Apr. 1998.
- [3] D.P.W. Ellis *Prediction-Driven Computational Auditory Scene Analysis*, Ph.D. thesis, Dept. of EECS, MIT, 1996.
- [4] Bob Foss *Filmmaking: Narrative and Structural techniques* Silman James Press LA, 1992.
- [5] F. R. Hampel et. al. *Robust Statistics: The Approach Based on Influence Functions*, John Wiley and Sons, 1986.
- [6] A. Hauptmann M. Witbrock *Story Segmentation and Detection of Commercials in Broadcast News Video* Advances in Digital Libraries Conference, ADL-98, Santa Barbara, CA., Apr. 22-24, 1998.
- [7] Liwei He et. al. *Auto-Summarization of Audio-Video Presentations*, ACM MM '99, Orlando FL, Nov. 1999.
- [8] J.R. Kender B.L. Yeo, *Video Scene Segmentation Via Continuous Video Coherence*, CVPR '98, Santa Barbara CA. Jun. 1998.
- [9] R. Lienhart et. al. *Automatic Movie Abstracting*, Technical Report TR-97-003, Praktische Informatik IV, University of Mannheim, Jul. 1997.
- [10] J. Meng S.F. Chang, *CVEPS: A Compressed Video Editing and Parsing System*, Proc. ACM Multimedia 1996, Boston, MA, Nov. 1996
- [11] R. Patterson et. al. *Complex Sounds and Auditory Images*, in *Auditory Physiology and Perception* eds. Y Cazals et. al. pp. 429-46, Oxford, 1992.
- [12] W.H. Press et. al *Numerical recipes in C, 2nd ed.* Cambridge University Press, 1992.
- [13] L. R. Rabiner B.H. Huang *Fundamentals of Speech Recognition*, Prentice-Hall 1993.
- [14] Eric Scheirer Malcom Slaney *Construction and Evaluation of a Robust Multifeature Speech/Music Discriminator* Proc. ICASSP '97, Munich, Germany Apr. 1997.
- [15] S. Subramaniam et. al. *Towards Robust Features for Classifying Audio in the CueVideo System*, Proc. ACM Multimedia '99, pp. 393-400, Orlando FL, Nov. 1999.
- [16] H. Sundaram S.F. Chang *Audio Scene Segmentation Using Multiple Features, Models And Time Scales*, ICASSP 2000, International Conference in Acoustics, Speech and Signal Processing, Istanbul Turkey, Jun. 2000.
- [17] H. Sundaram S.F. Chang *Video Scene Segmentation Using Audio and Video Features*, to appear in IEEE International Conference on Multimedia and Expo, New York, NY, Aug. 2000.
- [18] S. Uchihashi et. al. *Video Manga: Generating Semantically Meaningful Video Summaries* Proc. ACM Multimedia '99, pp. 383-92, Orlando FL, Nov. 1999.
- [19] M. Yeung B.L. Yeo *Time-Constrained Clustering for Segmentation of Video into Story Units*, Proc. Int. Conf. on Pattern Recognition, ICPR '96, Vol. C pp. 375-380, Vienna Austria, Aug. 1996.
- [20] M. Yeung B.L. Yeo *Video Content Characterization and Compaction for Digital Library Applications*, Proc. SPIE '97, Storage and Retrieval of Image and Video Databases V, San Jose CA, Feb. 1997.
- [21] T. Zhang C.C Jay Kuo *Heuristic Approach for Generic Audio Segmentation and Annotation*, Proc. ACM Multimedia '99, pp. 67-76, Orlando FL, Nov. 1999.