

論文 / 著書情報  
Article / Book Information

論題	動画像インデクシングのためのシーン時系列の確率的言語モデル
著者	安藤 亮一, 篠田 浩一, 古井貞熙, 望月 貴裕
出典	第12回 画像センシングシンポジウム 予稿集, Vol. , No. , pp. 513-518
発行日 / Issue date	2006, 6
Note	第12回 画像センシングシンポジウム講演論文集より転載

# 動画画像インデクシングのためのシーン時系列の確率的言語モデル

## Probabilistic language model of scene sequence for automatic video indexing

安藤 亮一 †, 篠田 浩一 †, 古井 貞熙 †, 望月 貴裕 ‡

Ryoichi Ando† Koichi Shinoda† Sadaaki Furui† Takahiro Mochizuki‡

† 東京工業大学 情報理工学専攻 計算工学専攻, 東京都目黒区大岡山 2-12-1

‡ NHK 放送技術研究所 東京都世田谷区砧 1-10-11

† Department of Computer Science, Graduate School of Information Science and Engineering  
Tokyo Institute of Technology

‡ NHK Science & Technical Research Laboratories

E-mail: ando@ks.cs.titech.ac.jp

### Abstract

本論文では、野球放送のインデクシングにおいて、シーン間のコンテキストを表現するシーン時系列の確率的言語モデルを用いる手法を提案する。動画画像のモデルとして、隠れマルコフモデルを用い、確率的言語モデルとして N-gram モデルを用いた。25 試合分の野球放送データを用いシーン認識の評価実験を行った。確率的言語モデルを用いない場合と比較して、16 種類のシーン認識における F 値の平均は 1.3% 改善した。これらの結果により提案手法の有効性が確認された。

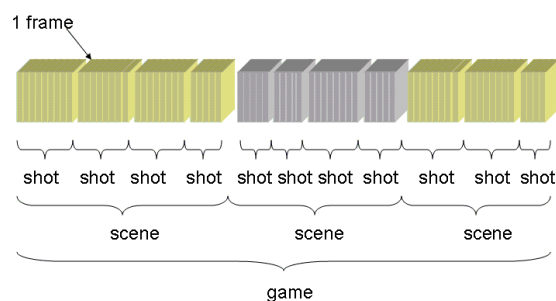


図 1 野球放送の構造

## 1 まえがき

近年、コンピューター技術、特に、ストレージ技術の進歩により、マルチメディアコンテンツが急増している。そのような背景の下、マルチメディアコンテンツを効率よく利用するためには、検索や要約が必要となる。そのため、インデックスを付与することが必要であるが、現状ではこの作業は人手によるものとなり、コストが大きい。そのような背景から、パターン認識技術を用いて自動でインデックスを付与する技術 (Contents Based Video Information Retrieval; CBVIR) が研究対象として注目されており、ニュース、スポーツ、映画など様々なコンテンツに対して研究が行われている [1, 2, 3, 4, 5].

この論文では野球放送のシーン認識を対象とする。野球放送のデータ構造の最小単位はフレームで、一枚の静止画像である。1 つの固定カメラで撮影された多数フレームからショットが構成される。更に、ショットのシーケンスからシーン (イベント) が構成される (図 1)。ショットの遷移の情報はシーンの特徴を表し、シーン認識において重要な情報となる (図 2)。ゲームの内容を理解するために重要となるシーンはハイライトと呼ば

れる。

各々のシーンに対するショットの遷移は多様であり、統計的なモデル化が必要となる。近年、統計モデルとして、隠れマルコフモデル (Hidden Markov Model; HMM) を用いる手法が、Chang らにより提案されている [3]。この手法ではまず動画画像データをショットに分割し、これらのショットから特徴量を抽出する。次に、この特徴量に基づきショットを分類し、ショットの遷移を HMM を用いてモデル化し、シーン認識を行う手法である。しかし、この手法では、ショットの境界検出やショットを分類する精度が低い場合にシーン認識率が低下するという問題がある。また、シーン HMM の状態間遷移 (トポロジー) が経験的に設計されているため、未知のデータに対する頑健性が不足している。それに対し、Nguyen ら [1] は各シーンに対し同一のトポロジーからなるマルチストリーム HMM を用いてシーンのモデル化を行い、さらにゲーム適応を用いる手法を提案し、その有効性を示した。しかし、この研究では、シーン認識に有効であるシーン間のコンテキストは考慮していない。

これまでスポーツのシーン認識において、ルールや放送データの構造に基づくシーンコンテキストの表現手



図2 ホームラン，内野ゴロ，四球のショットシーケンスの例

法はいくつか提案されている。野球のシーン認識では、例えば Liang ら [5] により提案されている。Liang らの手法では、画面上の試合のステータス表示を利用して野球ルールに基づいたシーン認識を行っている。しかしながら、一般にスポーツ番組では例外的なシーンの出現が多く、ルールも複雑である。そのため、文法を用いてスポーツのルールを完全に表現することは困難である。また、他のスポーツコンテンツに対して、シーン認識を行う場合、それに合わせたルールを人手で記述しなければならない。

本論文ではシーンコンテキストを確率的に表現した確率的言語モデルを用いる手法を提案する。この手法は、様々なコンテキストの違いに対し、頑健に動作する。シーン時系列において、より離れたシーンの関係を表現するために、出現頻度の高い連続シーンはチャンクとしてまとめ、N-gram モデルを用いてモデル化する。

本論文の構成は以下の通りである。2章ではシーン認識問題を定式化する。3章で特徴量を説明した後、4章で提案する手法を説明する。5章で評価実験結果を報告し、6章で全体をまとめ、今後の課題を述べる。

## 2 シーン認識問題の定式化

本研究では、入力された映像から自動的にシーン認識を行い、インデックスを作成することを目的とする。連続音声認識とのアナロジーから、シーン認識の問題を次のように定式化する。観測された特徴ベクトルのシーケンス  $O$  が与えられたとき、シーン時系列  $H$  の出現確率は次のようになる。

$$P(H|O) \propto P(O|H)P(H) \quad (1)$$

ここで  $P(O|H)$  はシーン時系列  $H$  から観測ベクトルの時系列  $O$  が出現する確率を表し、また、 $P(H)$  はシーン時系列  $H$  の出現確率を表す。ここでは、 $P(O|H)$  を表現するモデルをビデオモデルと呼び、 $P(H)$  を表現するモデルを言語モデルと呼ぶ。連続音声認識とのアナロジーにおいて、ショットは音素、シーンは単語にそれぞれ対応する。

### 2.1 ビデオモデル

ビデオモデルには HMM を用いる。HMM は時間的に変化するパターンのモデル化に広く用いられ、多くの

シーン認識の研究で用いられている [1, 2, 3]。HMM のパラメータは、学習データの画像から特徴ベクトル系列を抽出し、その特徴ベクトル系列を用い学習を行うことにより推定される。本研究では、全てのシーンに対し同じトポロジーをもつ HMM を用いる。これにより、シーン HMM の作成が容易になり、未知のデータに対しても頑健性をもつことができる。すなわち、新しいシーンを加えるときやデータが増加したとき、モデル設計をやり直す必要がない。このようなデータ駆動型のアプローチは、音声認識の分野でよく用いられており、高い認識性能と頑健性を得ている。

### 2.2 言語モデル

本研究では確率的言語モデルとして、音声認識の分野で広く使われている N-gram モデルを用いる。N-gram モデルにおいて、与えられた単語列  $w_1^n = w_1 w_2 \dots w_n$  の生成確率  $P(w_1^n)$  を計算する際に、各単語の生起確率は直前の  $(N - 1)$  単語にのみ依存すると仮定する。

$$P(w_n | w_1^{n-1}) = P(w_n | w_{n-N+1}^{n-1}) \quad (2)$$

なお、 $N = 1, 2, 3$  の場合をそれぞれユニグラム (1-gram)、バイグラム (2-gram)、トライグラム (3-gram) と呼ぶ。N-gram 確率は、学習データ中に出現する単語の  $N$  個組と  $(N - 1)$  個組の頻度から、最尤推定 (maximum likelihood estimation) により次のように推定することができる。

$$P(w_n | w_{n-N+1}^{n-1}) = \frac{C(w_{n-N+1}^n)}{C(w_{n-N+1}^{n-1})} \quad (3)$$

ここで、単語列  $w_1^n$  が学習データ中に出現する回数を  $C(w_1^n)$  で表す。なお、N-gram の推定に式 (3) をそのまま適用すると、学習データ中に出現しない単語組の確率値が 0 になる。また、学習データ中に出現しても出現頻度が小さな N-gram に対しては精度が劣化する。これらの問題に対処するため、N-gram の推定には、単語組の出現頻度をそのまま使うのではなく、出現頻度を補正した値を使う。本研究では学習データ中に出現しない N-gram の値を、低次の  $(N-1)$ -gram の値から推定する Back-off スムージング [6] を用いる。なお、このスムージングを用いる際に、学習データに出現しなかった N-gram に確率を割り当てるために、学習データに出現した N-gram の確率を割り引く必要がある。本研究では、その手法として good-turing 法 [7] を用いた。

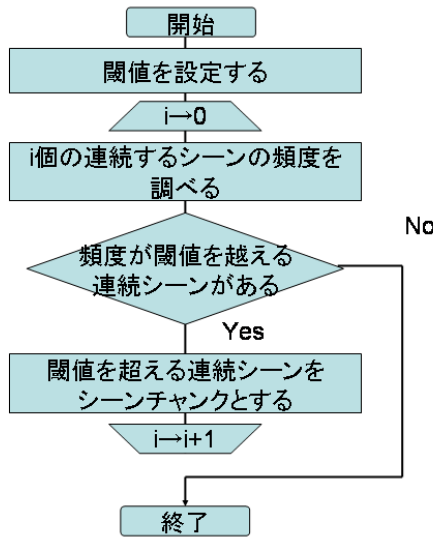


図3 シーンチャンクの生成手順

### 2.3 N-gram モデルによるシーンコンテキストのモデル化

N-gram モデルを用いて言語モデルを生成する際に、シーン時系列において、より離れたシーンの関係を表現するため、出現頻度の高い連続シーンを1つシーンとして扱う。以降、このシーンのまとまりをシーンチャンクと呼ぶ。シーンチャンクの生成手順は次の通りである。まず2つの連続するシーンの組み合わせの出現頻度を求め、ある閾値以上となる2つの連続するシーンをチャンクとしてまとめる。そして、3,4つと調べる幅を広げ、閾値を越える連続シーンが出現しなくなるまでその処理を繰り返す。シーンチャンク生成の手順を図3に示す。閾値以下のシーンはそのまま1つのシーンから構成されるシーンチャンクとなる。

## 3 特徴量

### 3.1 低周波数成分特徴量 (LF)

Nguyenらの手法[1]で用いられている主成分分析による特徴量(PF)は画像の全体的な特徴を表せるが、主成分の計算に用いる画像の選び方により、表せる特徴は変化するため、頑健な特徴量ではない。そこで本研究では、一枚の画像から画像の全体的な特徴を表すことが可能である、低周波数成分を特徴量として用いる。まず、計算量を減少させるために、720×480ピクセルの各フレーム画像を72×48サイズに圧縮する。次に、RGB画像からグレースケール画像を作る。その画像に対し2次元DCTを施し、低周波数成分30次元を特徴ベクトルとする。図4(b)は図4(a)の画像を圧縮した低周波数成分30次元の画像である。

### 3.2 低周波数成分差分特徴量 (DLF)

移動物体の情報として、Nguyenらの手法[1]で用いられる差分特徴量(DF)と同様に連続するフレームの差分情報に注目する。DFでは、動きの激しい部分の平均、分散と動き密度を特徴量として用いている。動き密度が小



図4 特徴量を表す画像

さい場合、ノイズが動きの激しい部分の平均、分散に影響を与える。そこで、LFと同様に低周波数成分を用い、差分画像の全体的な特徴を特徴量として用いる。まず、計算量を減少させるために、各フレームを72×48サイズに圧縮する。次に、2つの連続するフレーム画像の差分をとり、グレースケール画像(図4(c))を作る。その画像に対しLF特徴量と同様に2次元DCTを施し、低周波数成分30次元を特徴ベクトルとする。

### 3.3 オプティカルフロー特徴量 (OF)

野球放送データにおいて、移動物体の情報は重要であるが、カメラの動きもシーン認識に有用な情報となる。それを表現するためにオプティカルフローを用いる。まず、計算量を減少させるために、連続するフレームの画像をそれぞれ240×160のサイズに圧縮する。その画像に対して、20ピクセル間隔に追跡点を $N$ 個配置する。それぞれの追跡点の $x, y$ 座標を $(x_i, y_i) (i = 1, \dots, N)$ とし、L-K法[8]を用いて次のフレームでの対応点を求め、対応点の $x, y$ 座標を $(x'_i, y'_i) (i = 1, \dots, N)$ とする。オプティカルフローベクトルの $x, y$ 成分は、 $x''_i = x'_i - x_i$ 、 $y''_i = y'_i - y_i$ となる。次のように5次元特徴ベクトル $v = (\mu_x, \sigma_x^2, \mu_y, \sigma_y^2, z)^T$ を計算する。

$$\begin{aligned} \mu_x &= \frac{\sum_{i=1}^N x''_i}{N}, & \sigma_x^2 &= \frac{\sum_{i=1}^N (x''_i - \mu_x)^2}{N} \\ \mu_y &= \frac{\sum_{i=1}^N y''_i}{N}, & \sigma_y^2 &= \frac{\sum_{i=1}^N (y''_i - \mu_y)^2}{N} \\ z &= \frac{\sum_{i=1}^N \sqrt{v_i^2}}{N} \end{aligned} \quad (4)$$

ただし、画像の中心の $x, y$ 座標を $(x_c, y_c)$ とし、

$$v = (x''_i - \mu_x)(x'_i - x_c) + (y''_i - \mu_y)(y'_i - y_c)$$

とする。 $\mu_x, \mu_y$ はオプティカルフローベクトルの $x, y$ 成分それぞれの平均であり、 $\sigma_x^2, \sigma_y^2$ は分散となる。 $\mu_x, \mu_y$ はカメラの移動方向を表す。 $z$ は画像の中心位置に対して、オプティカルフローが全体的に内向きであるか、外向きであるかを表す。これはカメラのズームの度合いに対応する。この5次元ベクトルを特徴ベクトルとする。

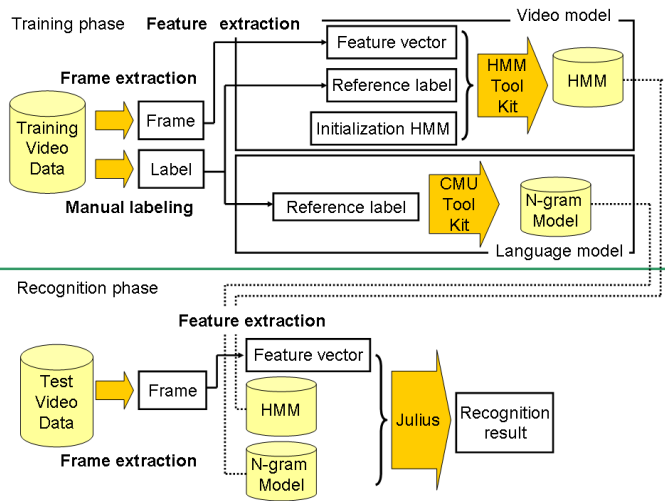


図 5 シーン認識フレームワーク

## 4 シーン認識手法

本章では、野球放送データからシーン認識を行う手法について説明する。本研究では図5に示すフレームワークを用い、シーンの学習及び認識を行う。

### 4.1 学習フェーズ

まず、ビデオモデルのモデル化においては、予め各々のシーンに対して初期 HMM を準備しておく。学習用の動画データから学習時間と冗長な情報を削減をするため、3 フレーム間隔でフレームを抽出する。抽出した各フレームから特徴ベクトルを計算する。計算した特徴ベクトル時系列と人手により付与された正解シーンラベルを用い、HMM のパラメータを学習する。

次に言語モデルのモデル化において、学習データより出現頻度の高い連続するシーンを 2.3 章の手法でまとめ、シーンチャンクとする。人手により付与された正解シーンラベルを用い、シーンチャンクから、N-gram モデルを生成する。

### 4.2 認識フェーズ

認識フェーズでは、学習フェーズと同様に、動画データからフレームを抽出し、同じ方法で特徴ベクトルを計算する。次に、学習フェーズでモデル化した HMM と N-gram モデルを用い、連続シーン認識を行う。出力はシーン時系列とシーン境界の時間情報である。

## 5 評価実験

### 5.1 実験条件

提案手法の評価データとして、NHK 放送技術研究所より提供された 25 試合分 (75 時間) の大リーグ野球放送データを用いた。評価データを 5 試合ずつの 5 つのグループに分割し、交差検定 (Cross Validation) を行った結果を平均したものを認識結果とした。表 1 に認識するシーンラベルを示す。5 つに分割したグループに含まれるシーンラベルの内訳を表 2 に示す。

表 1 シーンラベル

シーンラベル	内容
base hit(bh)	シングルヒット
extra-base hit(ebh)	長打
clutch hit(ch)	タイムリーヒット
home run(hr)	ホームラン
ground out(go)	内野ゴロ
fly out(fo)	フライアウト
strike out(so)	三振
strike(s)	ストライク
ball(b)	ボール
fall(f)	ファール
walk(wk)	四球 (デッドボールを含む)
pickoff(po)	牽制
steal(st)	盗塁 (盗塁失敗を含む)
out of play(op)	アウトオブプレー (試合とは関係のないシーン)
replay(rp)	リプレイ
effect(ef)	エフェクト (選手の成績表示などの CG)

### 5.2 評価方法

シーン認識の評価のために、F 値を用いた。F 値は Precision( $P$ ) と Recall( $R$ ) の調和平均である。本研究では各シーンに対する、Precision( $P$ ) と Recall( $R$ ) は以下のように計算する。

$$P = \frac{C}{S}, R = \frac{C}{T} \quad (5)$$

ここで、 $C$  は認識結果に含まれる正解のシーンの長さ (フレーム数)、 $S$  は認識結果におけるそのシーンの全体の長さ (フレーム数)、 $T$  はデータに含まれるそのシーンの全体の長さ (フレーム数) を示す。

### 5.3 シーン認識

ビデオモデルとして、各々のシーンに対し 1 つの HMM を準備した。ここで用いた HMM のトポロジーは自状態遷移と次の状態への遷移のみを許す単純な構造 (left-to-right HMM) を採用した。HMM の各状態の出力分布は単一ガウス分布である。HMM の状態数は予備実験の結果より 30 とした。学習及び認識に用いた特徴量は 3 章で説明した LF,DLF,OF の組合せである。HMM の学習には HTK[9] を用い、認識には Julius[10] を用いた。

2.3 章の手法を用い閾値を 70 とし、作成したシーンチャンクを表 3 に示す。言語モデルもビデオモデルと同様に、学習データからモデル化を行い、2-gram と 3-gram を生成した。N-gram モデルの学習には CMU-SLM-Toolkit[11] を用いた。言語モデルを評価するための基準としてしばしば用いられる、パープレキシティーの平均は 2-gram で 42.4、3-gram で 44.9 となった。音声

表2 5つに分割したグループに含まれるシーンラベル

	1	2	3	4	5	合計
b	271	371	353	308	398	1701
rp	235	351	270	356	366	1578
s	192	192	221	194	263	1062
op	185	197	181	199	175	937
f	160	166	199	200	170	895
go	78	67	73	86	76	380
fo	70	81	66	60	75	352
ef	58	72	52	52	38	272
so	41	36	48	49	49	223
bh	38	49	38	41	37	203
po	24	19	25	35	36	139
wk	25	38	23	23	27	136
ch	6	17	10	9	17	59
ebh	10	9	12	11	8	50
hr	9	9	4	10	6	38
st	4	4	6	9	1	24

認識におけるディクテ - ションの場合, 3-gram モデルのパープレキシティーは, 80 ~ 120 である。それと比較すると, 本研究のシーン認識は比較的単純なタスクであると言える。以上のようなビデオモデルと言語モデルを用いて, イニング境界は既知とし, イニングごとのシーン認識を行った。言語重みは, 事後的に最適な重みを用いた。表4は, baseline, 2-gram, 3-gram を用いた場合のシーン認識の結果である。baseline は, 任意のシーンの連続を許す単純な文法を用いた条件下での認識である。図6は, シーンチャンクを生成するときの閾値を変化させ, 3-gram を用いた場合のシーン認識の結果である。

表4のbaselineと確率的言語モデルの認識率を比較すると, 平均のF値は2-gramを用いた場合で0.6%, 3-gramでは1.3%改善した。この結果より, あるシーン時系列の後に出現しやすい(しにくい)シーンを考慮することで, 認識率が改善したと考えられる。特にN-gramを用いることでリプレイの認識率が高くなった。これは, ホームランや長打などの比較的重要なシーンの後にリプレイが出現しやすいというコンテキストが表現されたためと考えられる。また, 2-gramと3-gramの結果を比較すると, 3-gramの方がより離れたシーンの関係性を表現できるため, 認識率の改善は大きい。図6より, 70を境に閾値を大きくした場合と小さくした場合, とともに認識率が低下することがわかる。シーンチャンクの生成では, 閾値を小さくすればシーンチャンクの個数は増加する。シーンチャンクの個数が増加すれば, シーン時系列においてより離れたシーンの関係性を表現することができ, 認識率が高くなる。しかし, シーンチャンクが多過ぎる場

表3 連続するシーンより生成されたシーンチャンク

(閾値:70, 2つ組のシーン上位: 25個, 3つ組のシーン上位: 7個)

シーン列	出現頻度	シーン列	出現頻度	シーン列	出現頻度
b b	379	bh rp	151	go s	81
rp b	377	f f	149	go b	80
b s	302	s s	142	rp s	72
s b	272	go rp	136	rp b b	90
b f	262	b rp	134	rp s b	80
rp s	236	b fo	123	b b s	80
f b	200	b go	118	b s rp	75
s rp	164	so rp	106	b b b	74
s f	159	fo rp	103	s b b	73
rp fo	158	fo b	84	rp b s	73
f rp	157	go s	84		

合, 認識する単語数が増えてタスクが複雑になり, 認識率が低下する。

## 6 まとめと今後の課題

本論文では野球放送のインデクシングにおいて, シーンコンテキストを表現するために, シーン時系列の確率的言語モデルを用いる手法を提案した。野球放送のデータに対しシーン認識システムを構築し, 提案手法の評価実験を行った。言語モデルを用いない場合と比較して, 16種類のシーン認識におけるF値の平均は1.3%改善した。これらの結果により提案手法の有効性を確認した。本手法は, 学習によりシーンコンテキストを表現できるので, 様々なスポーツコンテンツのシーン認識に応用可能と考えられる。

今後は確率的言語モデルの効果をより明らかにするために, ルールに基づく文法を用いた場合との比較や融合を行う必要がある。また, 学習データを拡充し, 言語モデルの精度を上げ, 認識率の改善を目指す。現在は事後的に最適な言語重みを用いているが, その重みは未知のデータに対しては, 最適な重みではない。そのため, 自動的に最適な重みを決定する手法の検討が必要である。さらに, 現段階ではイニング境界を既知としてシーン認識を行っているが, 今後はイニング境界を検出する手法を検討し, それを用いてシーン認識を行う必要がある。

本論文では言語モデルに焦点を当て実験を行ってきた。そのため, ビデオモデルに関しては改善の余地が大きい。Nguyenら[1]が提案しているマルチストリームHMMやゲーム適応を用いれば, 認識率の全体的な底上げが期待できる。また, 音声情報と合わせて用いることで, さらなる認識率の改善が期待できる。

表 4 シーン認識の F 値 (%)

	baseline	2-gram	3-gram
bh	31.8	32.8	33.2
ebh	12.8	12.1	10.7
ch	26.1	24.5	26.8
hr	51.8	50.0	54.1
go	57.0	56.5	57.2
fo	46.6	48.0	48.2
so	57.8	57.4	58.1
s	38.2	35.9	36.3
b	41.3	46.2	47.9
f	43.1	43.2	43.0
wk	37.7	39.0	39.0
po	39.0	39.1	39.2
st	0.0	0.0	0.0
op	48.8	51.5	53.4
rp	48.7	54.2	55.2
ef	89.2	89.3	88.2
average	41.9	42.5	43.2

## 謝辞

この研究は 21 世紀 COE プログラム「大規模知識資源の体系化と活用基盤構築」の援助を受けた。

## 参考文献

- [1] H. B. Nguyen, K. Shinoda, and S. Furui, "Robust highlight extraction using multi-stream hidden markov models for baseball video," *Proc. International Conference on Image Processing 2005*, vol. 3, pp. 621–624, 2005.
- [2] T. Mochizuki, M. Tadenuma, and N. Yagi, "Baseball video indexing using patternization of scenes and hidden markov model," *Proc. International Conference on Image Processing 2005*, vol. 3, pp. 1212–1215, 2005.
- [3] P. Chang, M. Han, and Y. Gong, "Extract highlights from baseball game video with hidden markov models," *Proc. International Conference on Image Processing 2002*, vol. 1, pp. 609–612, 2002.
- [4] S. Eickeler and S. Müller, "Content-based video indexing of tv broadcast news using hidden markov models," *Proc of Int. Conf. on Acoustics, Speech, and Signal Processing (ICASSP)*, pp. 2997–3000, March 1999.
- [5] C.-H. Liang, W.-T. Chu, J.-H. Kuo, J.-L. Wu, and W.-H. Cheng, "Baseball event detection using game-

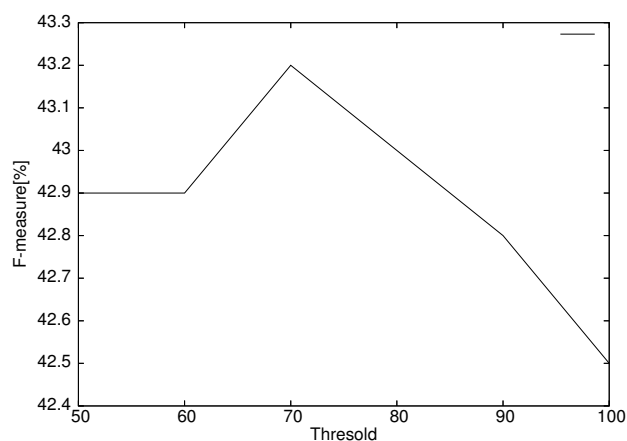


図 6 シーンチャンクの作成に用いた閾値と認識率の関係

specific feature sets and rules," *Proc. IEEE International Symposium on Circuits and Systems 2005*, pp. 3829–3832, 2005.

- [6] S. M. Katz, "Estimation of probabilities from sparse data for the language model component of a speech recognizer," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, pp. ASSP-35(3), 400–401, 1987.
- [7] I. J. Good, "The population frequencies of species and the estimation of population parameters," *Biometrika*, pp. 40(3), 237–264, 1953.
- [8] B. Lucas and T. Kanada, "An iterative image registration technique with an application to stereo vision," *Proc. of 7th International Joint Conference on Artificial Intelligence(IJCAI)*, pp. 674–679, 1981.
- [9] <http://htk.eng.cam.ac.uk>.
- [10] <http://julius.sourceforge.jp>.
- [11] <http://svr-www.eng.cam.ac.uk/~prc14/toolkit.html>.