



Structuration audiovisuelle par composantes primaires

Plate-forme AFIA

Atelier : Connaissance et Documents Temporels

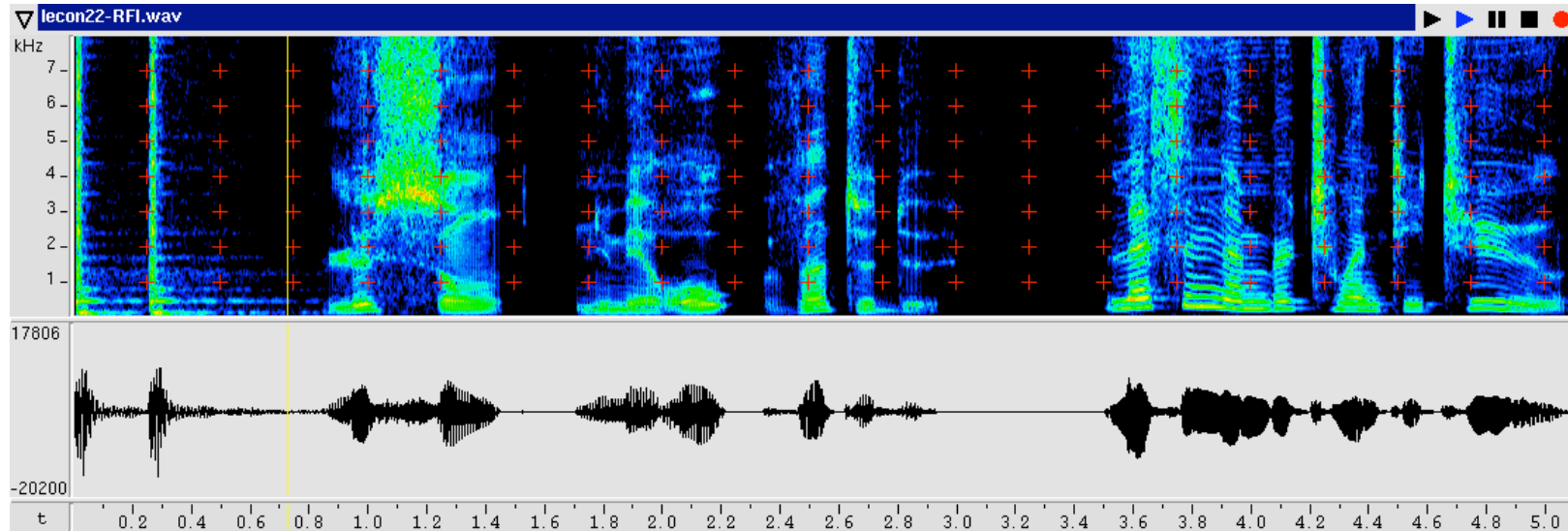
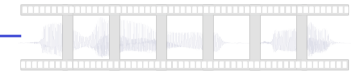
Julien Pinquier et Régine André-Obrecht

Équipe SAMoVA

(Structuration Analyse et Modélisation de la Vidéo et de l'Audio)



Indexation sonore : que faire ?



jingle 1

leçon

lesson



locuteur 1 (homme)

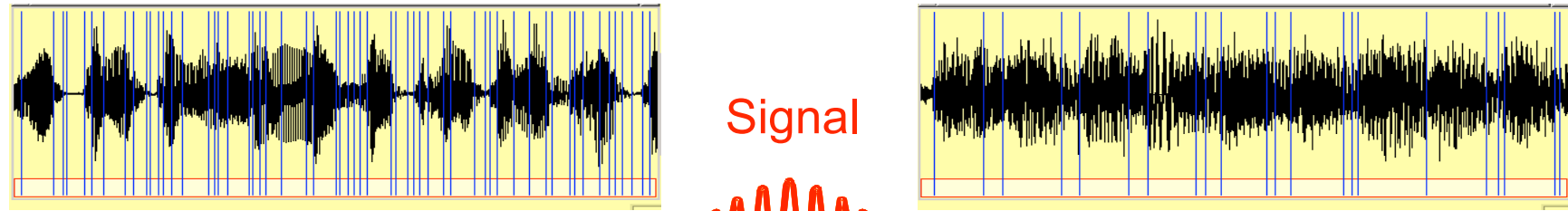
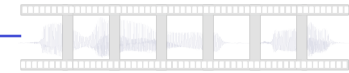
locuteur 2 (femme)

français

anglais

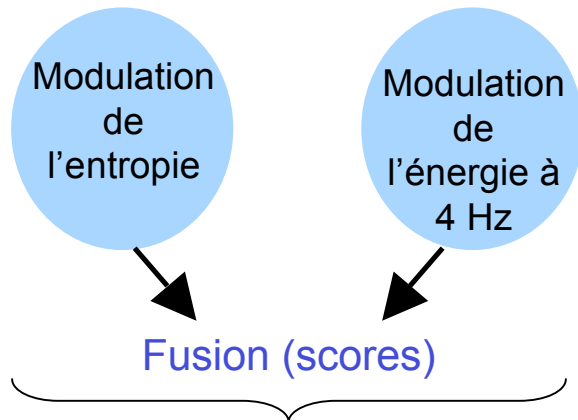


Détection PMB



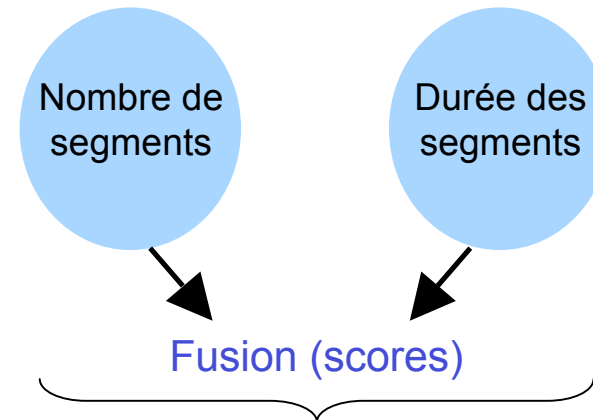
Détection de parole

Détection de musique



Classification Parole / NonParole

Segmentation



Classification Musique / NonMusique



Détection PMB

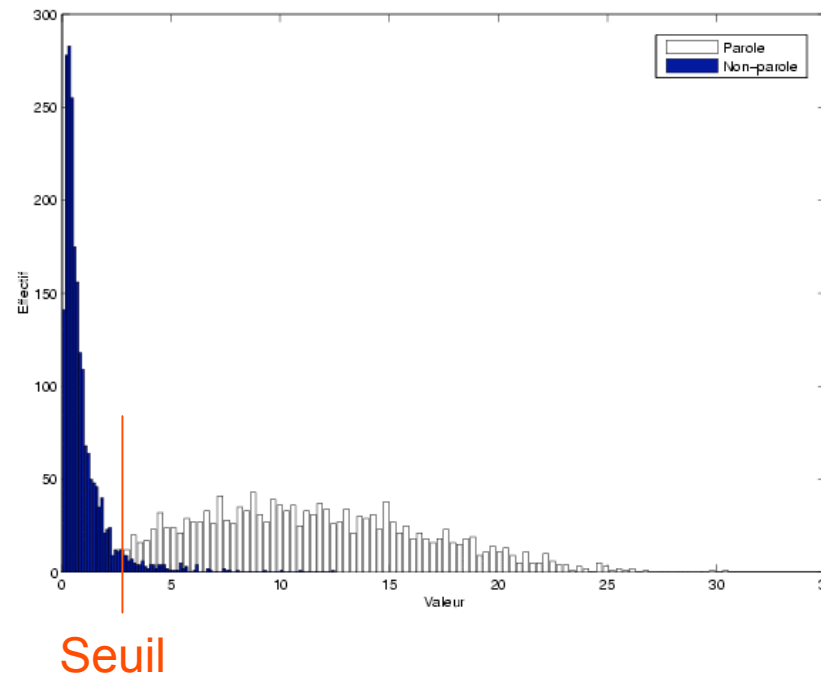


✧ Les seuils

- ◆ Parole : corpus MULTTEXT [Campione98]
- ◆ Musique : base personnelle

Exemple :

Modulation de l'énergie à 4 Hertz



Détection PMB : résultats



Paramètres	Score
------------	-------

P
A
R
O
L
E

Modulation de l'énergie à 4 Hertz	87,3 %
Modulation de l'entropie	87,5 %

Fusion (max)	90,5 %
--------------	--------

**CORPUS RFI
(6 heures)**

Campagne d'évaluation ESTER (2004-2005) :
1^{er} rang sur 7 participants (tâche SES)

M
U
S
I
Q
U
E

Nombre de segments	86,4 %
Durée des segments	78,1 %

Fusion (max)	89 %
--------------	------

**Sorties outils
: Format FDL**



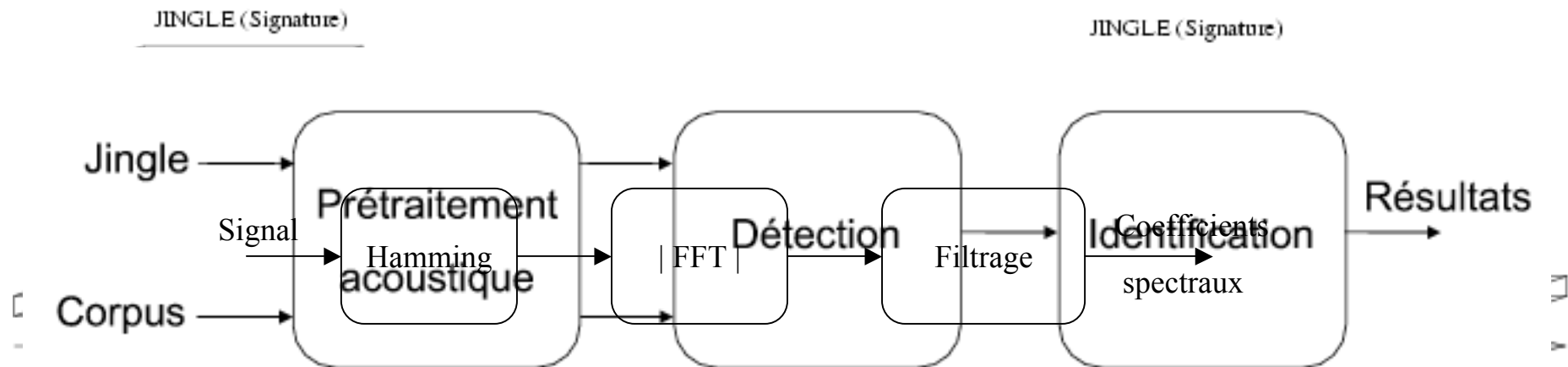
Détection de jingles



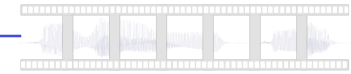
✧ Extrait sonore

✧ Système

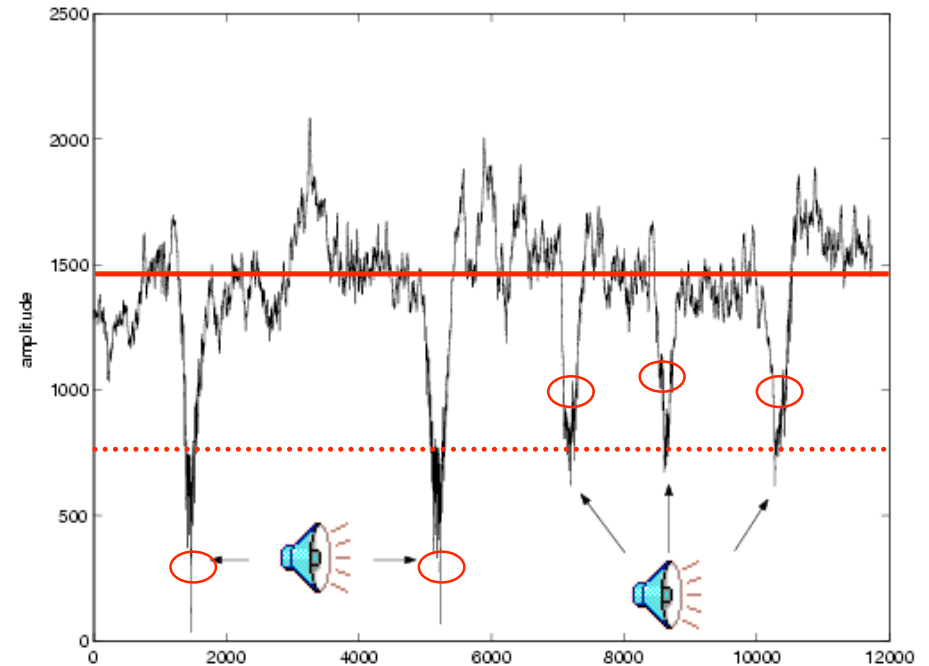
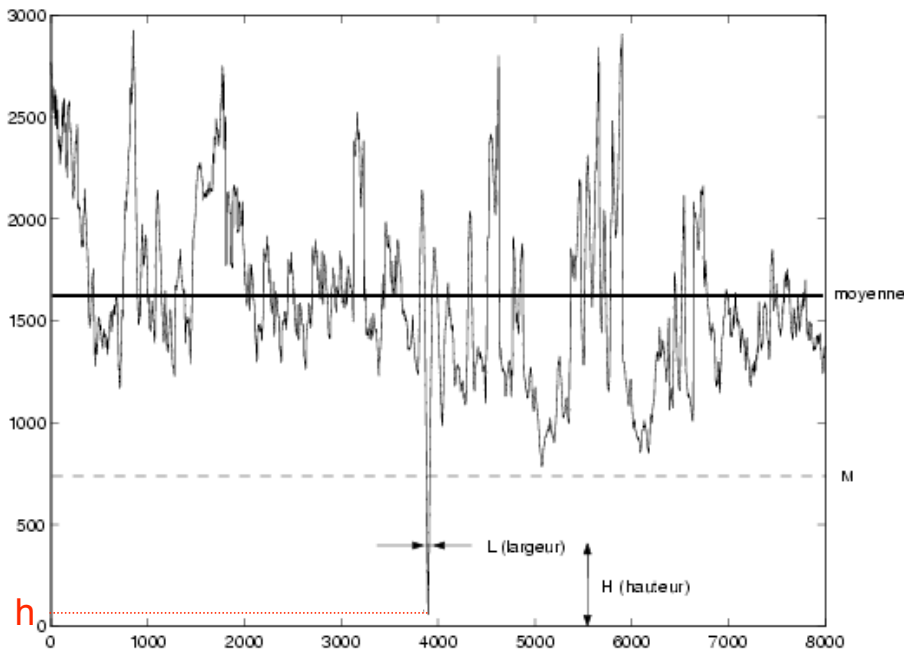
- ◆ Analyse spectrale (29 coefficients)
- ◆ Comparaison (distance Euclidienne)
- ◆ Analyse des « pics »



Détection de jingles



✧ Méthode d'analyse des pics

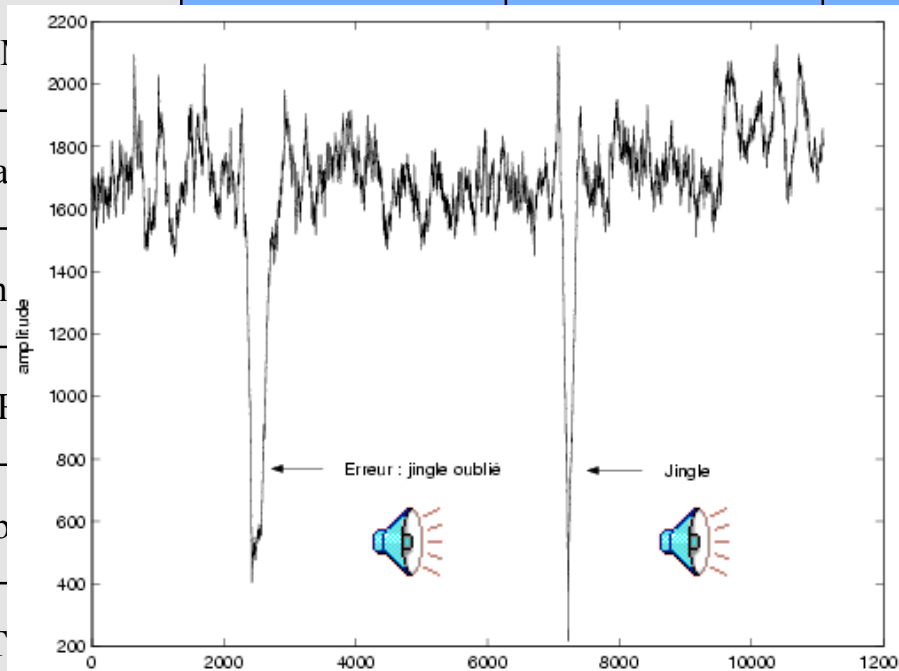


Détection de jingles



✧ Résultats

Corpus	Durée	Jingles	Détection manuelle
France 3	15 min	1	4
I			16
Ca			6
Fran			12
F			60
Pub			34
T			132



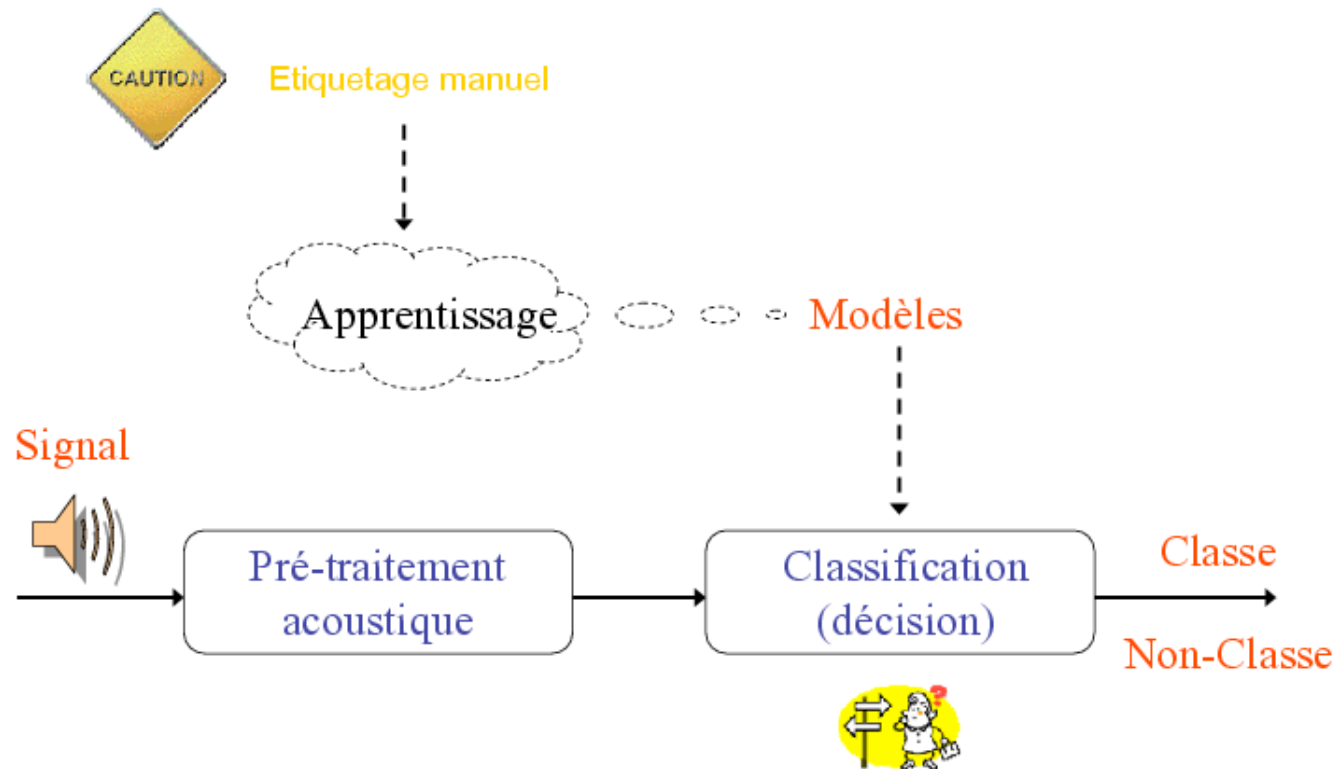
France Info



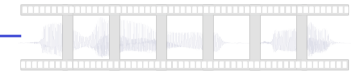
Détection des applaudissements



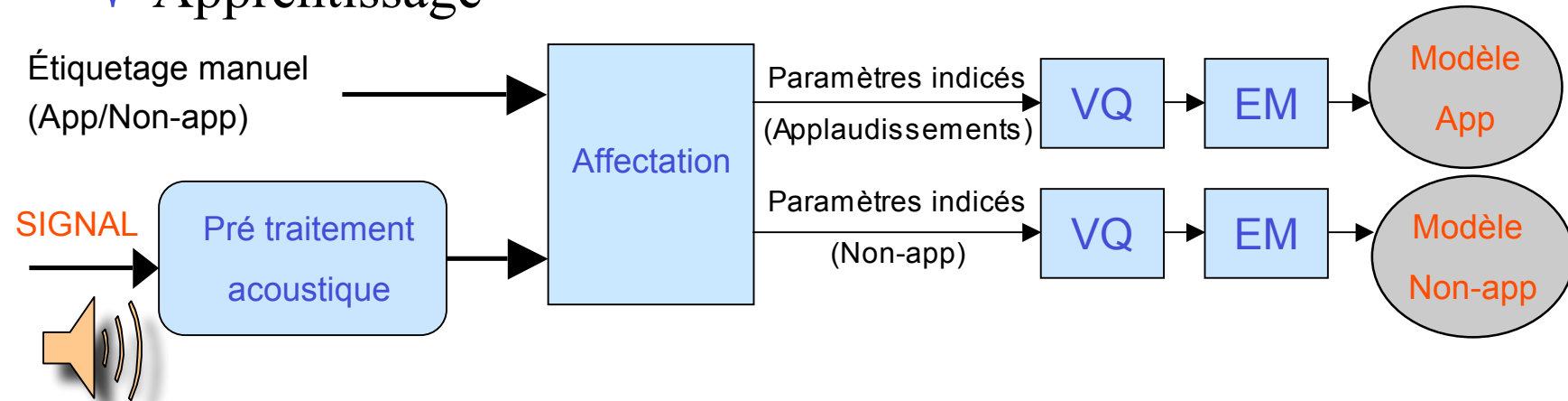
✧ Système



Détection des applaudissements



✧ Apprentissage



✧ Corpus : « Le Grand Échiquier », projet FERIA

- ◆ Apprentissage : 1 émission
- ◆ Reconnaissance : 1 émission

✧ Résultats

- ◆ Problème du critère évaluation : taux de reconnaissance > 98 %
- ◆ Segments significatifs : excellents






- ✧ Détection de motif dans une collection d'émissions
 - ◆ « Le grand Échiquier » 54 émissions de 3h
 - ◆ Motif :
 - présentateur / [APP] / spectacle / [APP/spectacle] / APP / présentateur
 - ◆ Détections automatiques :
 - détection de musique (chansons, spectacle)
 - détection de parole, puis du présentateur
 - détection des applaudissements
 - ◆ Apprentissage : 3 min (présentateur) et 6 min (applaudissements)
 - ◆ Résultats : une émission → détection de 10 motifs
 - ◆ Application aux autres émissions de la collection (vérité terrain)





✧ Macrosegmentation automatique (exemple du motif)

- ◆ Annotations automatiques
- ◆ Recherche de suites récurrentes [Haidar04]
- ◆ Inférence d'un motif
- ◆ Structuration

 Important : difficile manuellement



Conclusion



- ✧ Indexation sonore : étude de composantes primaires

- ✧ « Unités communes »
 - ◆ Parole et musique : → robustesse (plus d'apprentissage)

- ✧ « Unités exotiques »
 - ◆ Jingles : résultats excellents → 1 occurrence
 - ◆ Applaudissements : résultats très bons → universel

- ✧ Étude de structuration sonore
 - ◆ Détection d'un motif → très intéressante



Bonus : une autre étude de structuration

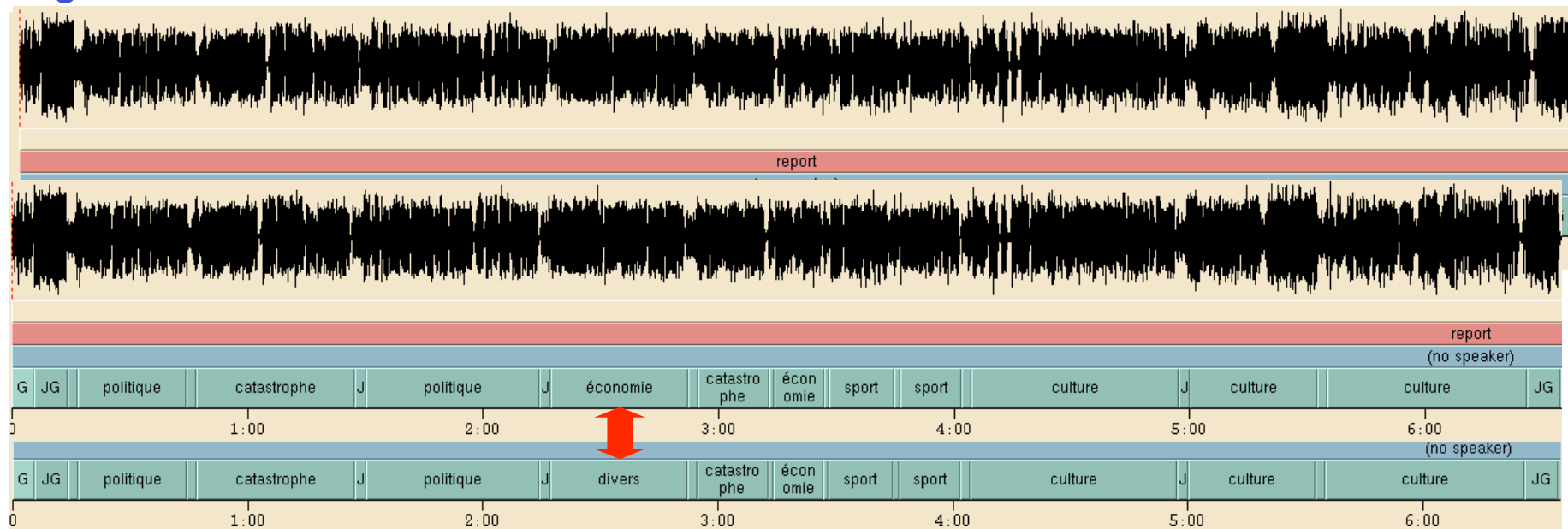


❖ Structuration d'un journal télévisé (« 6 minutes » de M6)

- ◆ Détection de jingles (J et JG)
- ◆ Détections de parole et de musique
- ◆ Détection de mots clés → thème du reportage

J

- ◆ 1 erreur

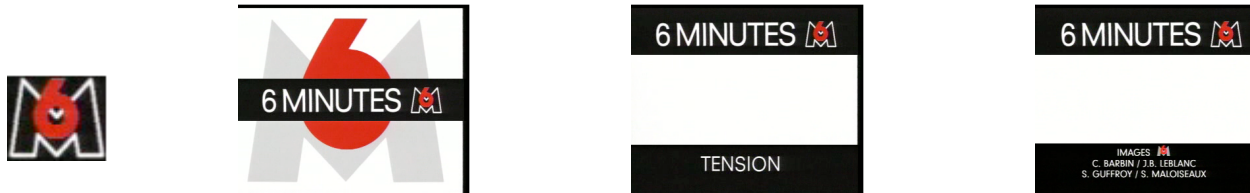


Perspectives

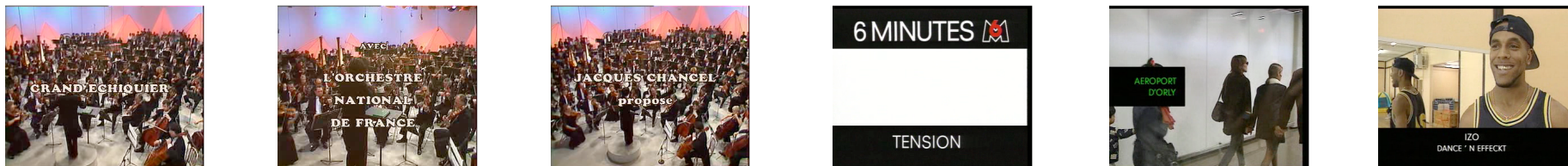


✧ Apport de la vidéo

◆ Détection de logos



◆ Extraction de texte



◆ Reconnaissance de l'intervenant



Perspectives



✧ Caractérisation de l'intervenant

- ◆ Voix IN, OUT, OFF
- ◆ Vidéo : mouvements
→ Corps, tête, lèvres
- ◆ Audio : conditions enregistrements

