

LIMSI
Séminaire TLP
5 octobre 2010

Hervé BREDIN

- bâtiment S
- poste 81 84
- bredin@limsi.fr
- www.limsi.fr/Individu/bredin

PARCOURS

- 2001-2004
 - Elève ingénieur @ Télécom Paris
 - *Traitement du signal / reconnaissance des formes*
- 2004-2007
 - Doctorat @ Télécom Paris
 - *Vérification audiovisuelle de l'identité*
- 2008
 - Post-doctorat @ Dublin City University
 - *Résumé automatique de séquences audiovisuelles*
- 2008-2010
 - CNRS @ Institut de Recherche en Informatique de Toulouse
 - *Indexation sémantique de documents audiovisuels*

Travaux de thèse
à Télécom Paris
(2004-2007)

VÉRIFICATION AUDIOVISUELLE DE L'IDENTITÉ

VÉRIFICATION AUDIOVISUELLE DE L'IDENTITÉ

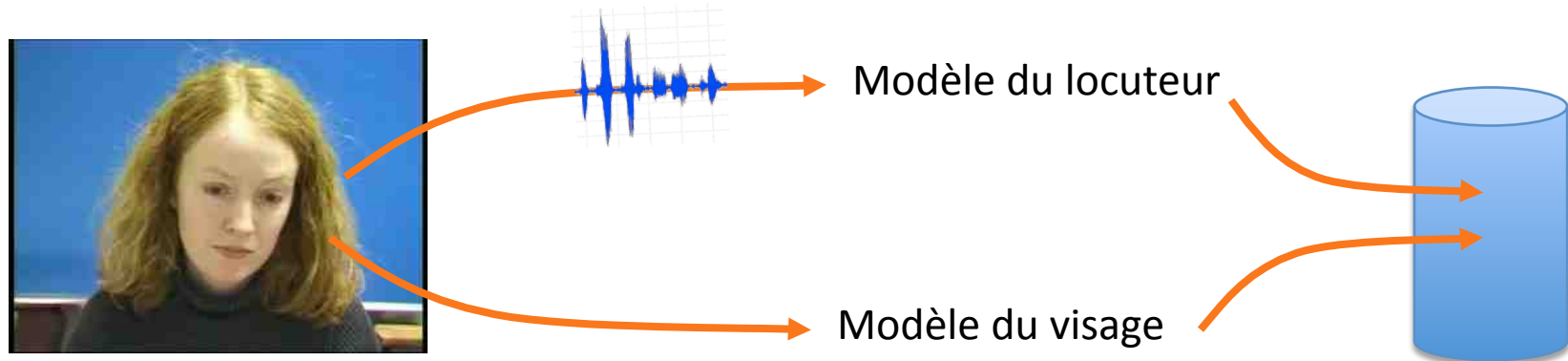
Contexte

- réseau d'excellence **BioSecure**
 - « ***Biometrics for Secure Authentication*** »
 - Empreinte digitale, main, iris, signature, **voix**, **visage**
 - **Fusion** multi-modale
- projet **SecurePhone**
 - « *enabling biometrically authenticated users to **deal m-contracts during a mobile phone call*** »
 - Empreinte digitale, **voix**, **visage**, signature

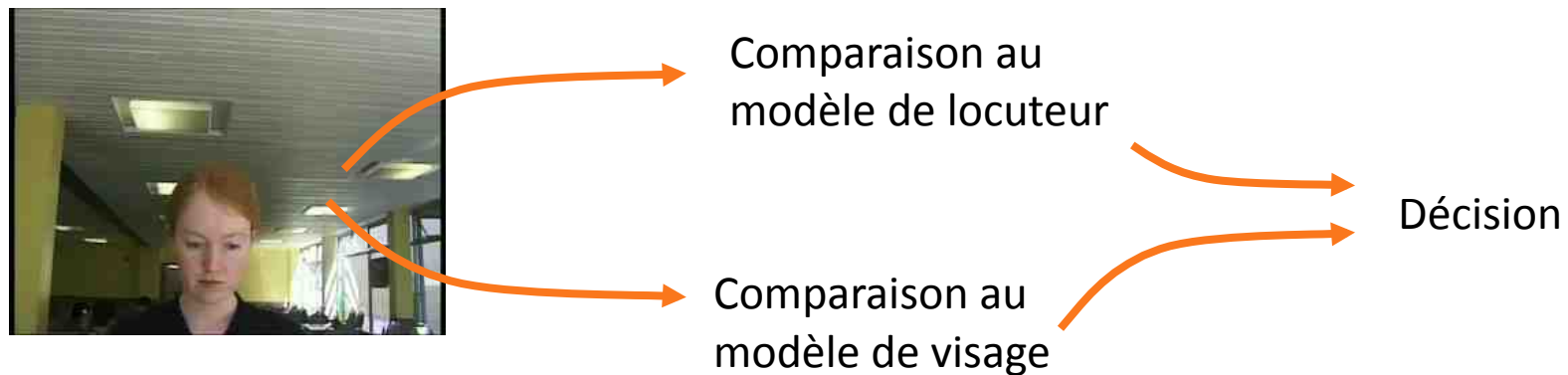
VÉRIFICATION AUDIOVISUELLE DE L'IDENTITÉ

Principe général

- Enrôlement



- Test



VÉRIFICATION AUDIOVISUELLE DE L'IDENTITÉ

Fusion multimodale (voix + visage)

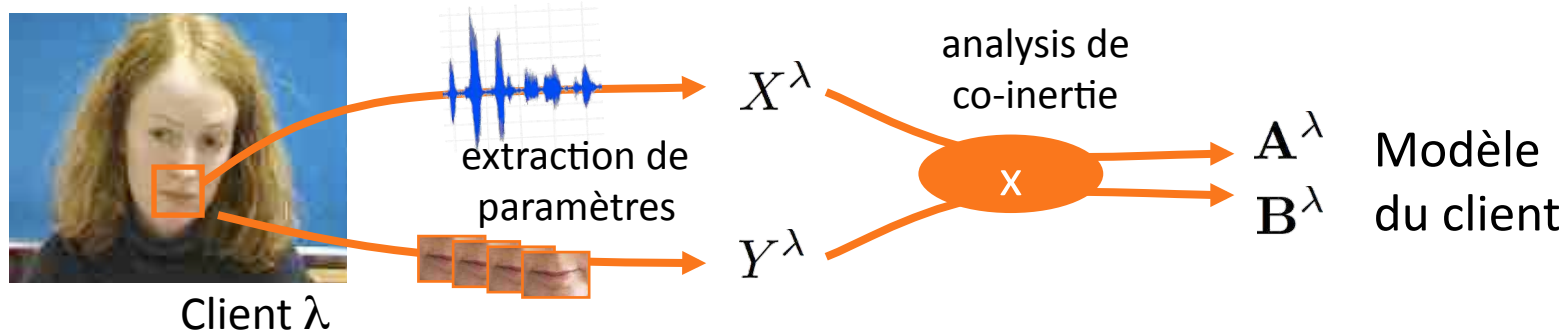
- Système de référence (baseline)
 - Vérification du locuteur
 - MFCC / GMM-UBM
 - Reconnaissance du visage
 - *Eigenfaces* / Distance euclidienne
 - Fusion tardive
 - Somme pondérée des scores normalisés

VÉRIFICATION AUDIOVISUELLE DE L'IDENTITÉ
Robustesse à l'imposture



VÉRIFICATION AUDIOVISUELLE DE L'IDENTITÉ

Modalité voix x lèvres (enrôlement)



extraction de
paramètres

X^λ Paramètres audio :
coefficients MFCC

Y^λ Paramètres visuels :
coefficients DCT de la région de la bouche

analyse de
co-inertie

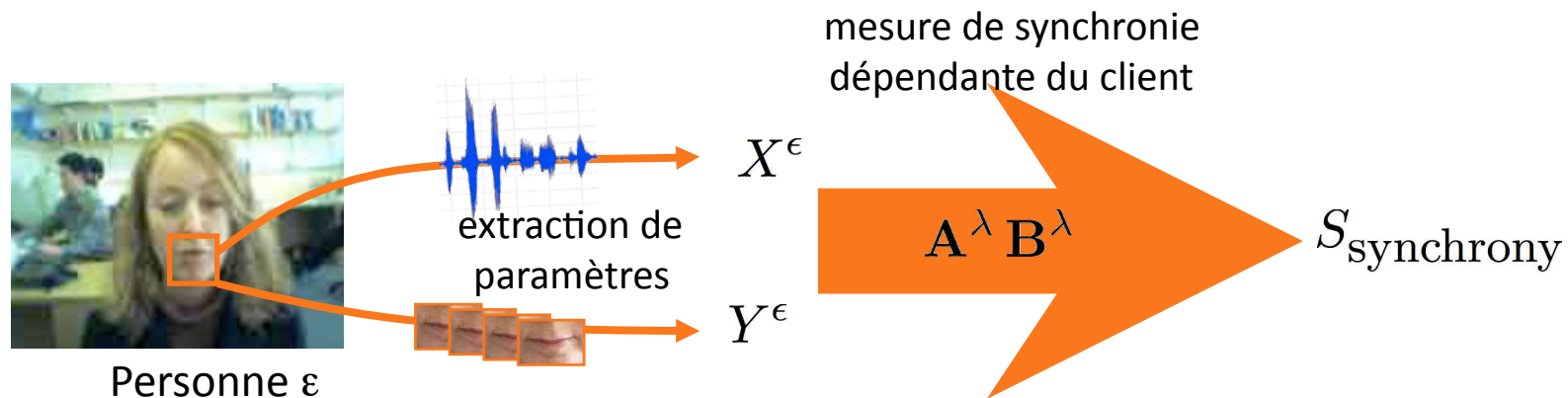
$$(\mathbf{a}_1^\lambda, \mathbf{b}_1^\lambda) = \operatorname{argmax}_{(a,b)} \operatorname{cov}(a^t X^\lambda, b^t Y^\lambda)$$

$$\mathbf{A}^\lambda = [\mathbf{a}_1^\lambda, \mathbf{a}_2^\lambda, \dots, \mathbf{a}_d^\lambda]$$

$$\mathbf{B}^\lambda = [\mathbf{b}_1^\lambda, \mathbf{b}_2^\lambda, \dots, \mathbf{b}_d^\lambda]$$

VÉRIFICATION AUDIOVISUELLE DE L'IDENTITÉ

Modalité voix x lèvres (test)



mesure de synchronie
dépendante du client

$$S_{\text{synchrony}} = \frac{1}{D} \sum_{k=1}^D \text{corr} \left(\mathbf{a}_k^{\lambda^t} X^\epsilon, \mathbf{b}_k^{\lambda^t} Y^\epsilon \right)$$

Accepté ($\epsilon=\lambda$) si $S_{\text{synchrony}} > \theta$, rejeté sinon.

identity
model λ

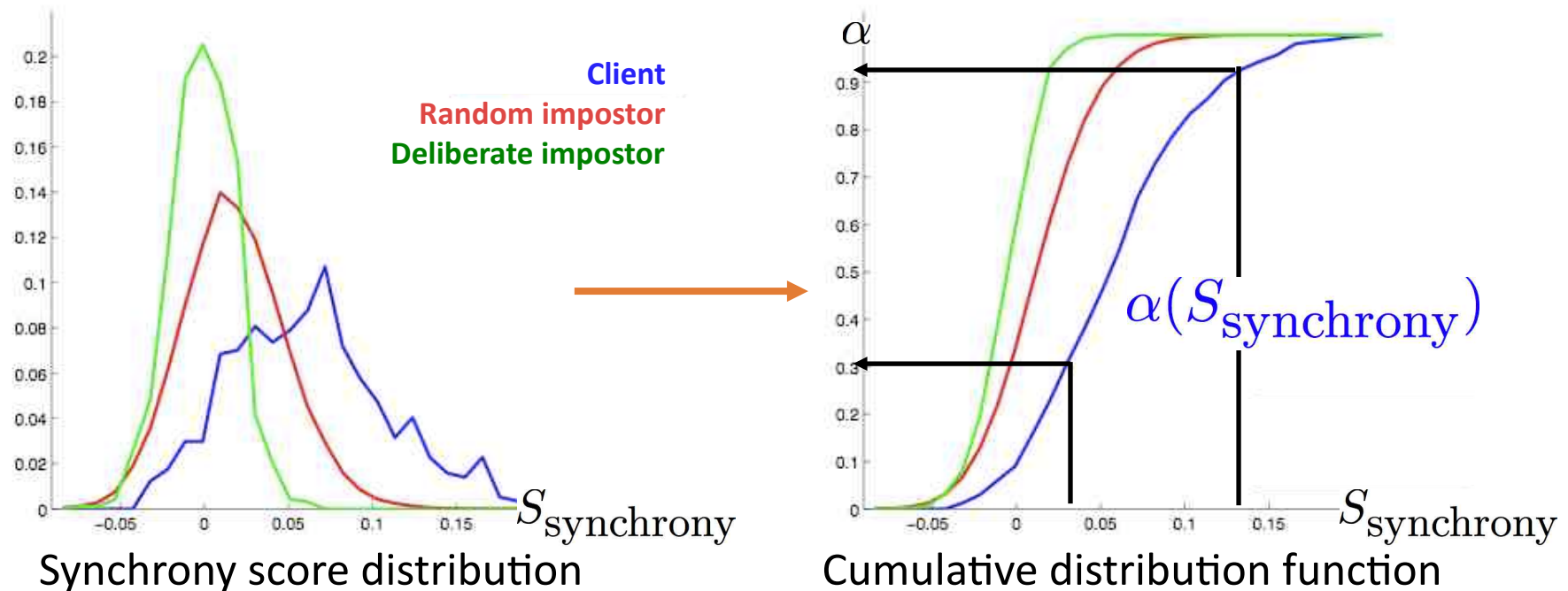
$$(\mathbf{a}_1^\lambda, \mathbf{b}_1^\lambda) = \text{argmax}_{(a,b)} \text{cov} (a^t X^\lambda, b^t Y^\lambda)$$

$$\mathbf{A}^\lambda = [\mathbf{a}_1^\lambda, \mathbf{a}_2^\lambda, \dots, \mathbf{a}_d^\lambda]$$

$$\mathbf{B}^\lambda = [\mathbf{b}_1^\lambda, \mathbf{b}_2^\lambda, \dots, \mathbf{b}_d^\lambda]$$

VÉRIFICATION AUDIOVISUELLE DE L'IDENTITÉ

Fusion adaptative



- Somme pondérée adaptative

$$S_{\text{final}} = \alpha(S_{\text{synchrony}}) \cdot S_{\text{baseline}} + \left[1 - \alpha(S_{\text{synchrony}})\right] \cdot S_{\text{synchrony}}$$

VÉRIFICATION AUDIOVISUELLE DE L'IDENTITÉ

Expériences (base BANCA)

| | Baseline (face + speaker) | New system (synchrony) | New fusion strategy |
|-------------------------|--|--|--|
| Random imposture | <div>DCF = 5.8%</div> <div>FAR = 2.1%</div> <div>FRR = 38%</div> | <div>DCF = 8.6%</div> <div>FAR = 1.0%</div> <div>FRR = 77%</div> | <div>DCF = 6.0%</div> <div>FAR = 0.6%</div> <div>FRR = 54%</div> |
| Deliberate imposture | <div>DCF = 97%</div> <div>FAR = 94%</div> <div>FRR = 38%</div> | <div>DCF = 7.6%</div> <div>FAR = 0%</div> <div>FRR = 77%</div> | <div>DCF = 7.3%</div> <div>FAR = 1.9%</div> <div>FRR = 54%</div> |

Performance brute du système de référence
&
Robustesses aux attaques

VÉRIFICATION AUDIOVISUELLE DE L'IDENTITÉ
Robustesse à l'imposture



Travaux de post-doctorat
à Dublin City University
(2008)

RÉSUMÉ AUTOMATIQUE DE SÉQUENCES AUDIOVISUELLES

RÉSUMÉ AUTOMATIQUE DE SÉQUENCES AUDIOVISUELLES

Contexte

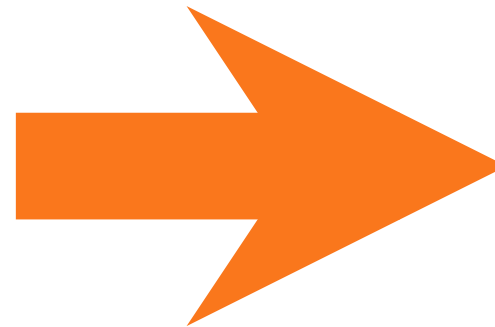
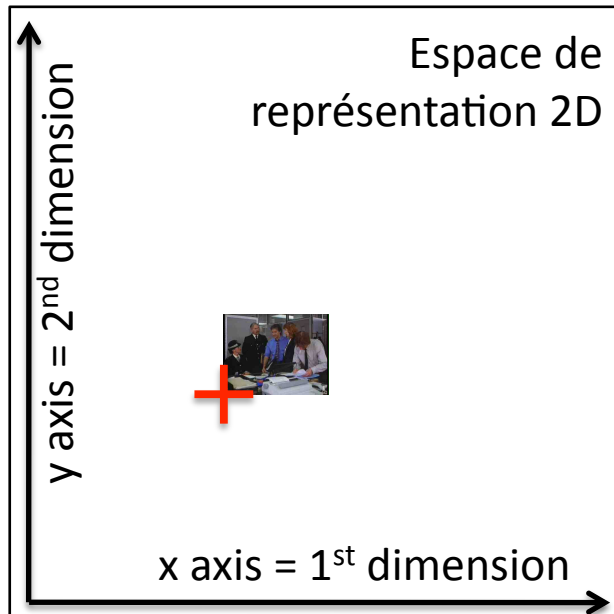
- Aide à la navigation dans les vidéos
Montage automatique
- Tâche
 - Entrée : document audiovisuel de durée D
 - Sortie : document audiovisuel de durée inférieure à 2% de D « *représentatif* » du document original
 - NIST : “*at least the summary should be usable by a professional*”
- @ DCU
 - D. Byrne, H. Lee, N. O’Connor, G. Jones

Principe général

- Segmentation en plan
- Description bas-niveau de chaque plan
- Sélection des plans représentatifs de la vidéo
- Sélection des segments représentatifs

RÉSUMÉ AUTOMATIQUE DE SÉQUENCES AUDIOVISUELLES

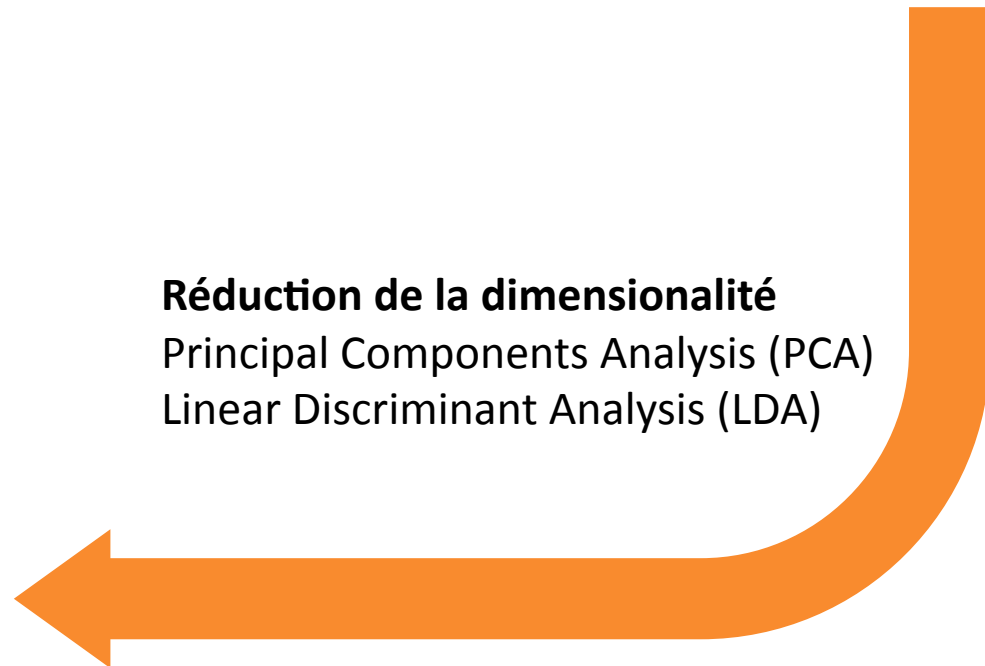
Description des trames



Histogramme RGB
(8x8x8 bins)

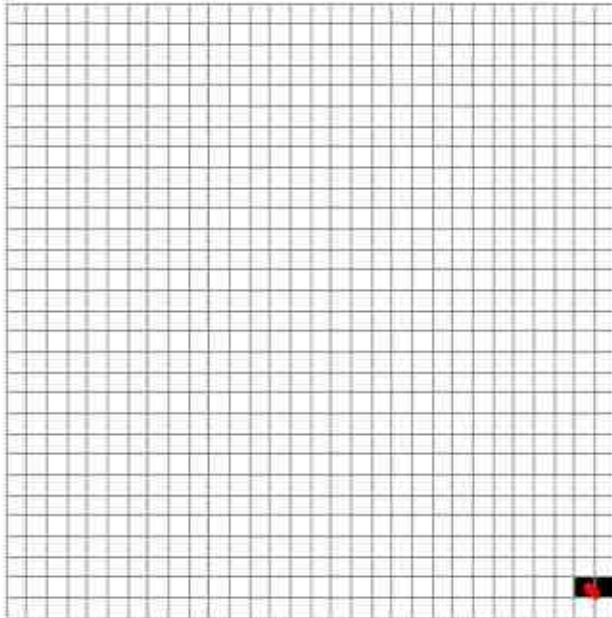
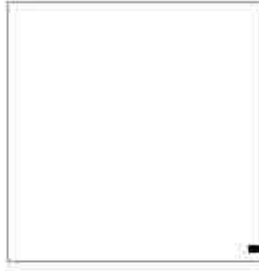
Espace de représentation
de dimension 512

Réduction de la dimensionalité
Principal Components Analysis (PCA)
Linear Discriminant Analysis (LDA)



RÉSUMÉ AUTOMATIQUE DE SÉQUENCES AUDIOVISUELLES

« Empreinte » visuelle



RÉSUMÉ AUTOMATIQUE DE SÉQUENCES AUDIOVISUELLES

Sélection des plans *représentatifs*

Step 1.

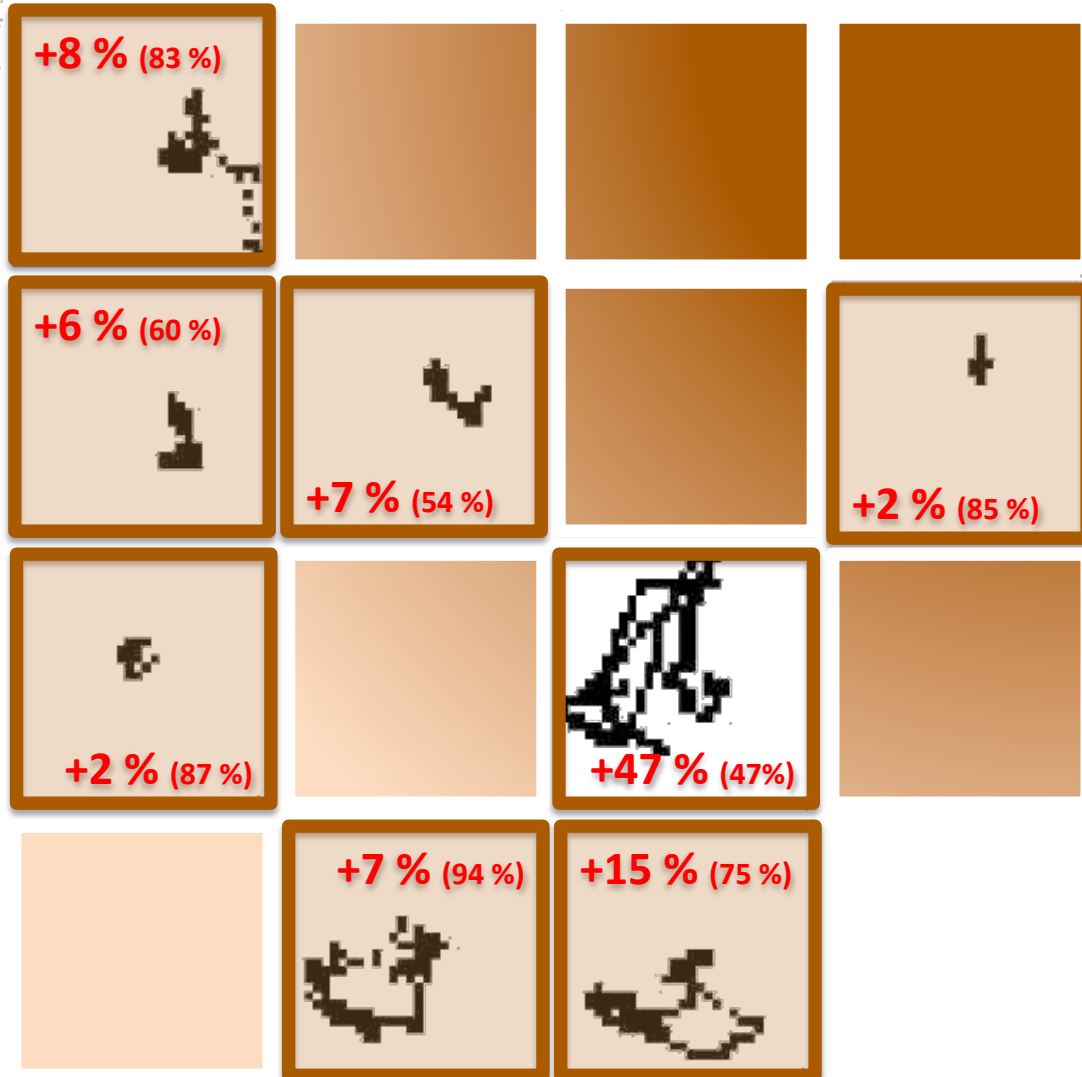
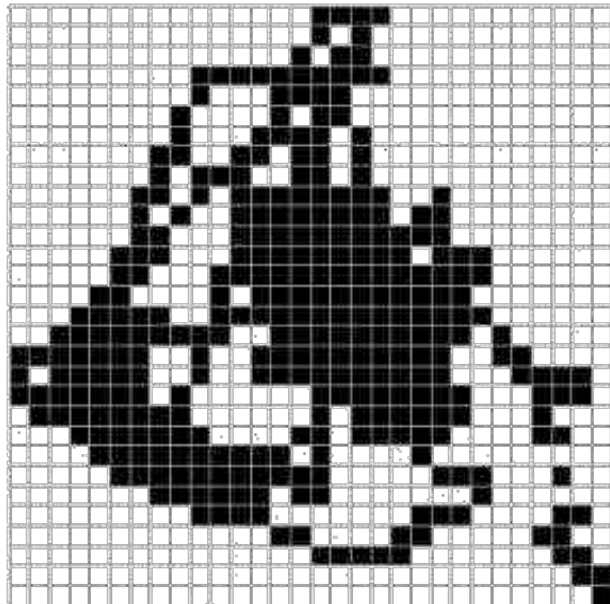
Select “largest” shot

Step 2.

Select shot with lowest
“relative” intersection

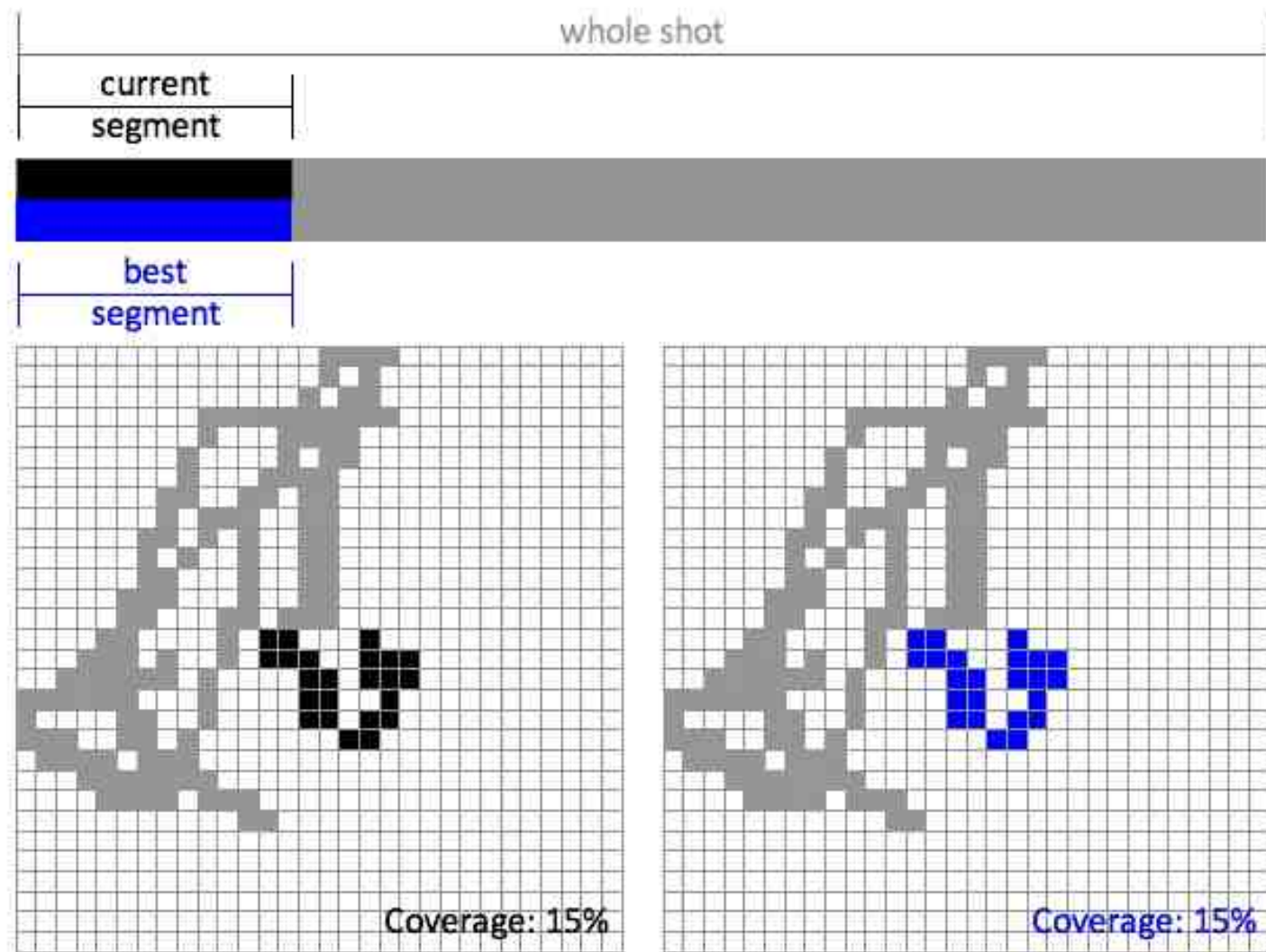
Step 3.

If sufficient coverage, stop.
Else go to step 2.



RÉSUMÉ AUTOMATIQUE DE SÉQUENCES AUDIOVISUELLES

Sélection du segment *représentatif*



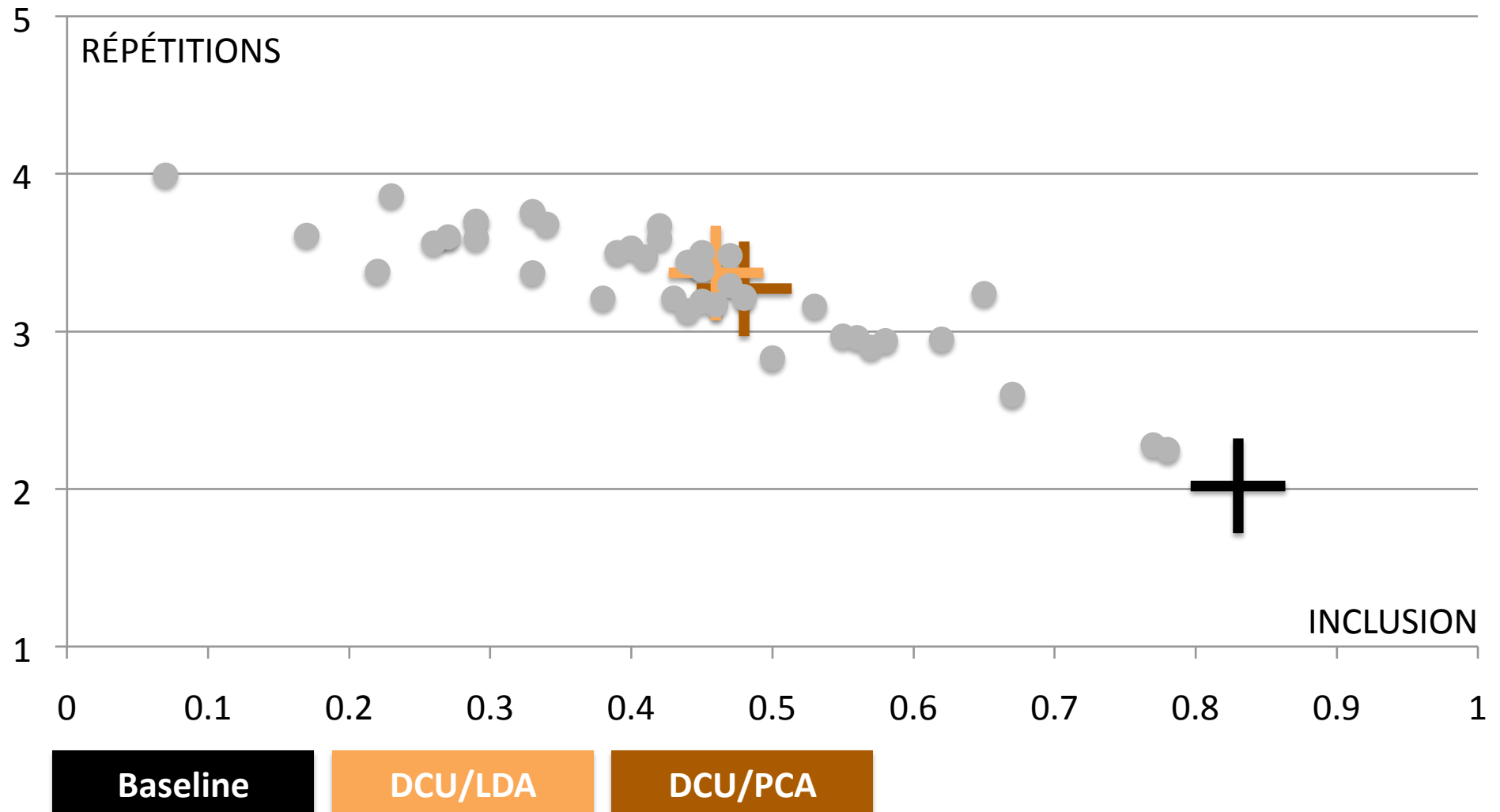
RÉSUMÉ AUTOMATIQUE DE SÉQUENCES AUDIOVISUELLES

Exemples de résumé



RÉSUMÉ AUTOMATIQUE DE SÉQUENCES AUDIOVISUELLES

TRECVID 2008



CNRS
IRIT (2008-2010)
LIMSI (2010-...)

INDEXATION PAR LE CONTENU DE DOCUMENTS AUDIOVISUELS

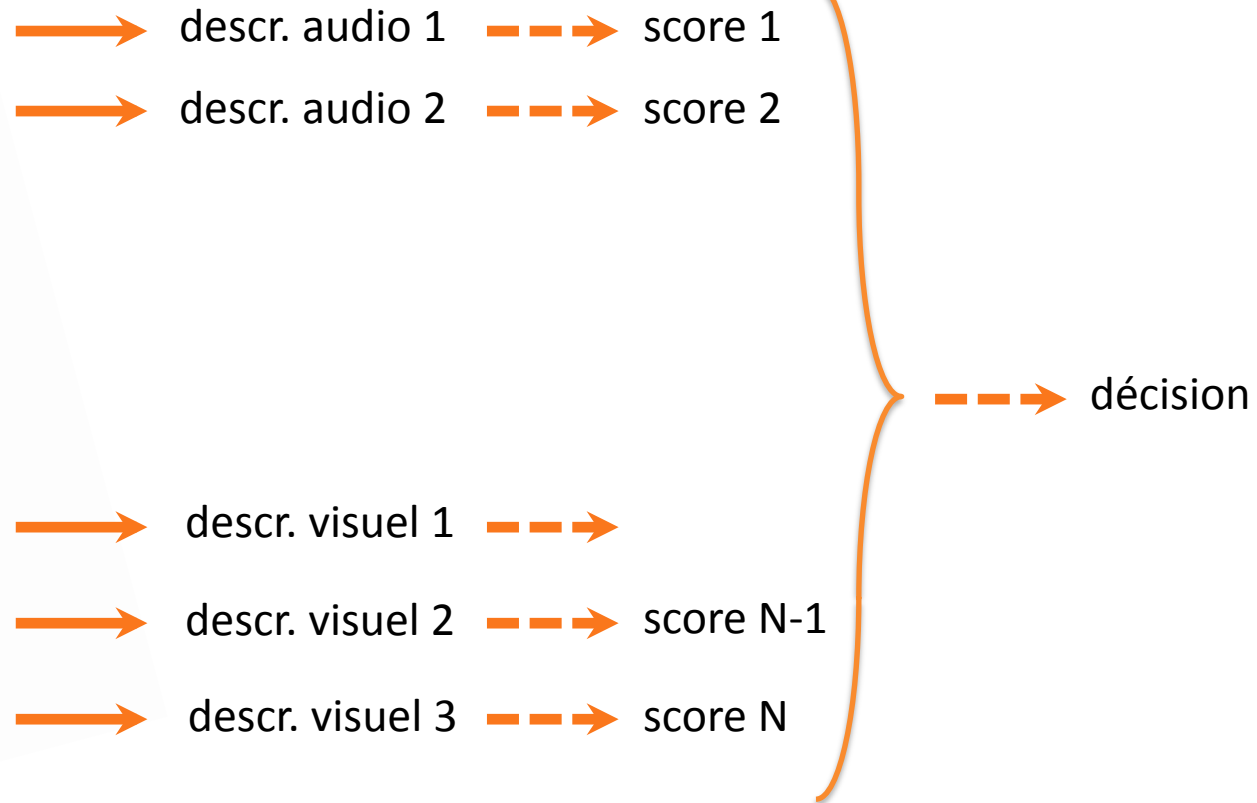
Détection de concepts sémantiques

- Problématique

– *le concept X apparaît-il dans la vidéo ?*

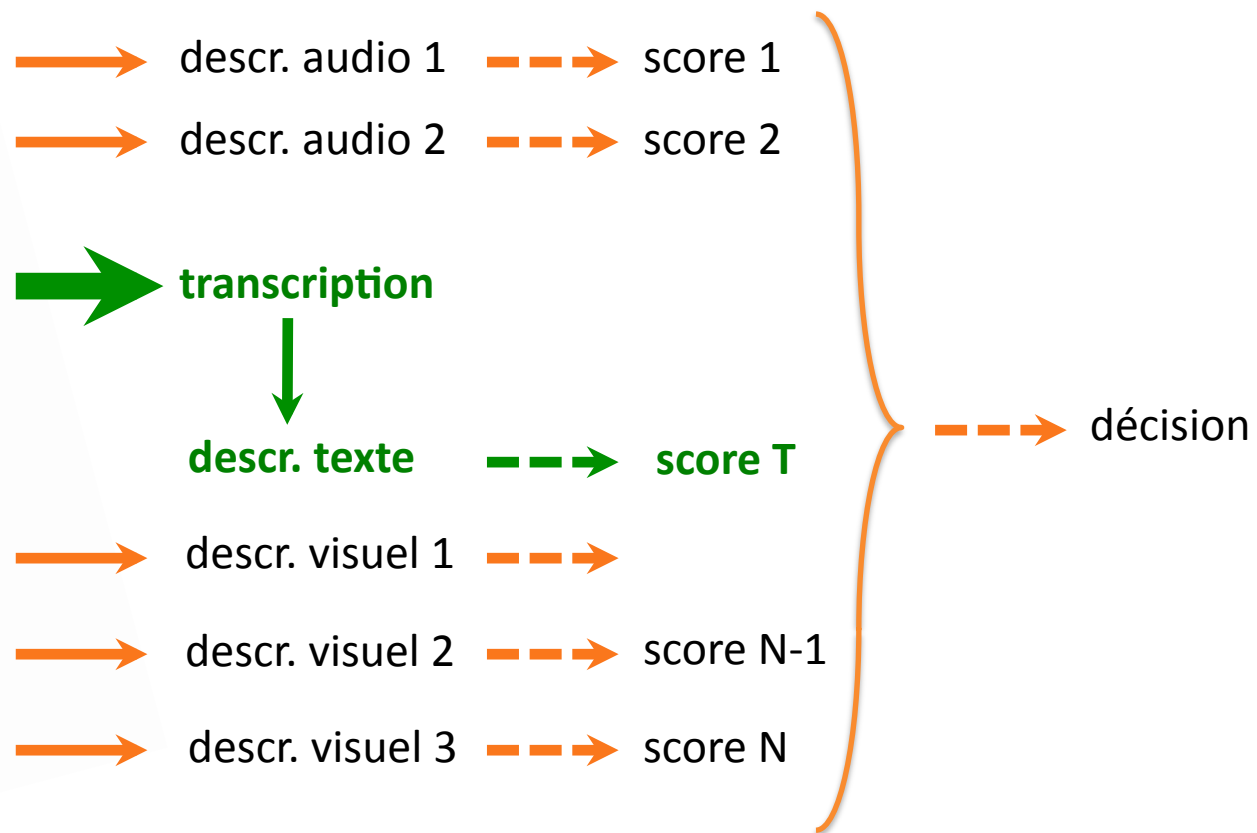
A word cloud of semantic concepts, including: stadium, dog, prisoner, chair, crowd, bridges, explosion, tent, vegetation, hospitals, athlete, indoor, laboratory, soccer player, walking, doorway, driver, eater, musical instrument, US flag, demonstration, roadway junction, tennis, scientist, singing, meeting, dancing, and waterscape.

Fusion tardive audio/vidéo



→ Extraction de descripteurs bas-niveau
--> Classification

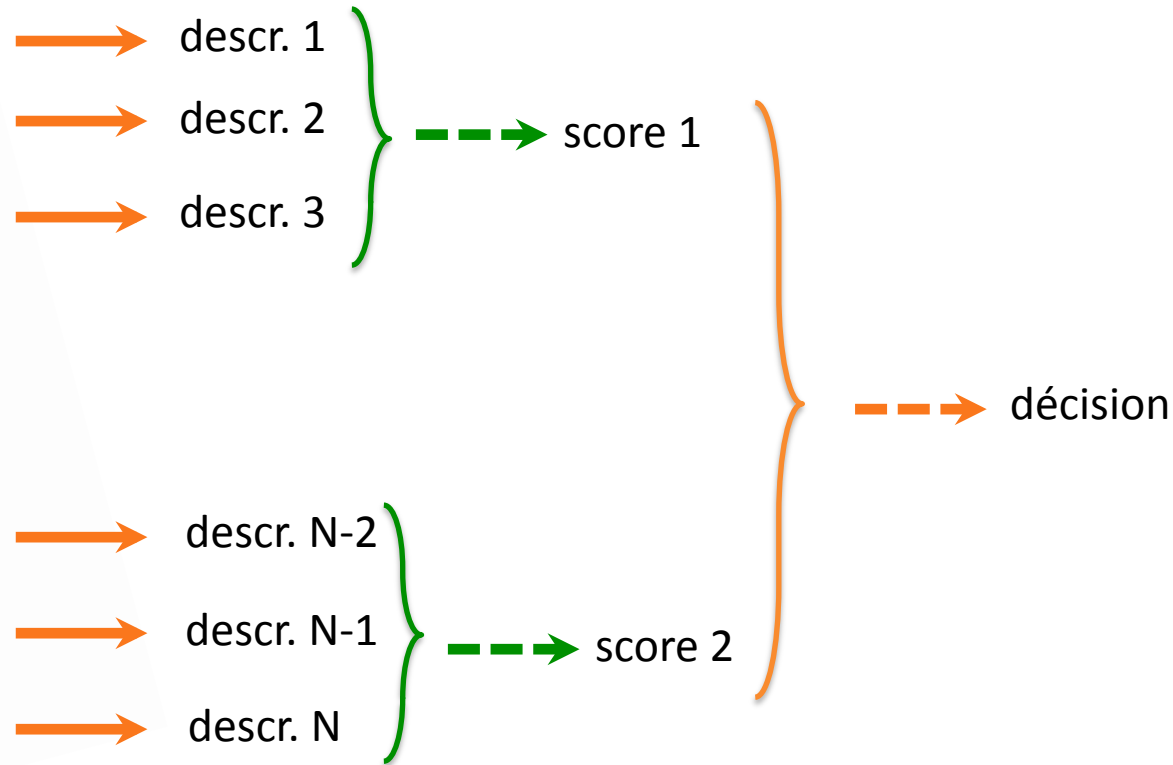
Fusion tardive audio/vidéo + transcription



—> Extraction de descripteurs bas-niveau
- - -> Classification

INDEXATION PAR LE CONTENU DE DOCUMENTS AUDIOVISUELS

Fusion **précoce** / hiérarchique



- Extraction de descripteurs bas-niveau
- - -→ Classification

LIMSI
Séminaire TLP
5 octobre 2010

THE END