

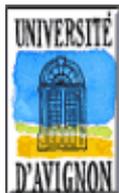
# Comment faire pour que l'opinion forgée à la sortie des urnes soit la bonne ? Application au défi DEFT 2007

**Juan Manuel Torres-Moreno**

Marc El-Bèze, Frédéric Béchet, Nathalie Camelin

3 juillet 2007

Laboratoire Informatique d'Avignon/UPRES 931  
Université d'Avignon et des Pays de Vaucluse



Laboratoire Informatique d'Avignon



# Corpus disponibles

- **Critiques** (films, livres, spectacles, bandes dessinées)
  - 3 460 critiques {0,1,2}
- **Jeux vidéo**
  - 4 231 tests {0,1,2}
- **Relectures**
  - 1 484 relectures d'articles scientifiques {0,1,2}
- **Débats**
  - 28 832 interventions parlementaires {0,1}
- 60% des corpus en **apprentissage** ; 40% **test**

# Evaluation

$$Fscore(\beta) = \frac{(\beta^2 + 1) \times \langle Précision \rangle \times \langle Rappel \rangle}{\beta^2 \times \langle Précision \rangle + \langle Rappel \rangle}$$

$$\langle Précision \rangle = \frac{\sum_{i=1}^n Précision_i}{n} ; \langle Rappel \rangle = \frac{\sum_{i=1}^n Rappel_i}{n}$$

$$Précision_i = \frac{\{\text{Nb de documents correctement attribués à la classe } i\}}{\{\text{Nb de documents attribués à la classe } i\}}$$

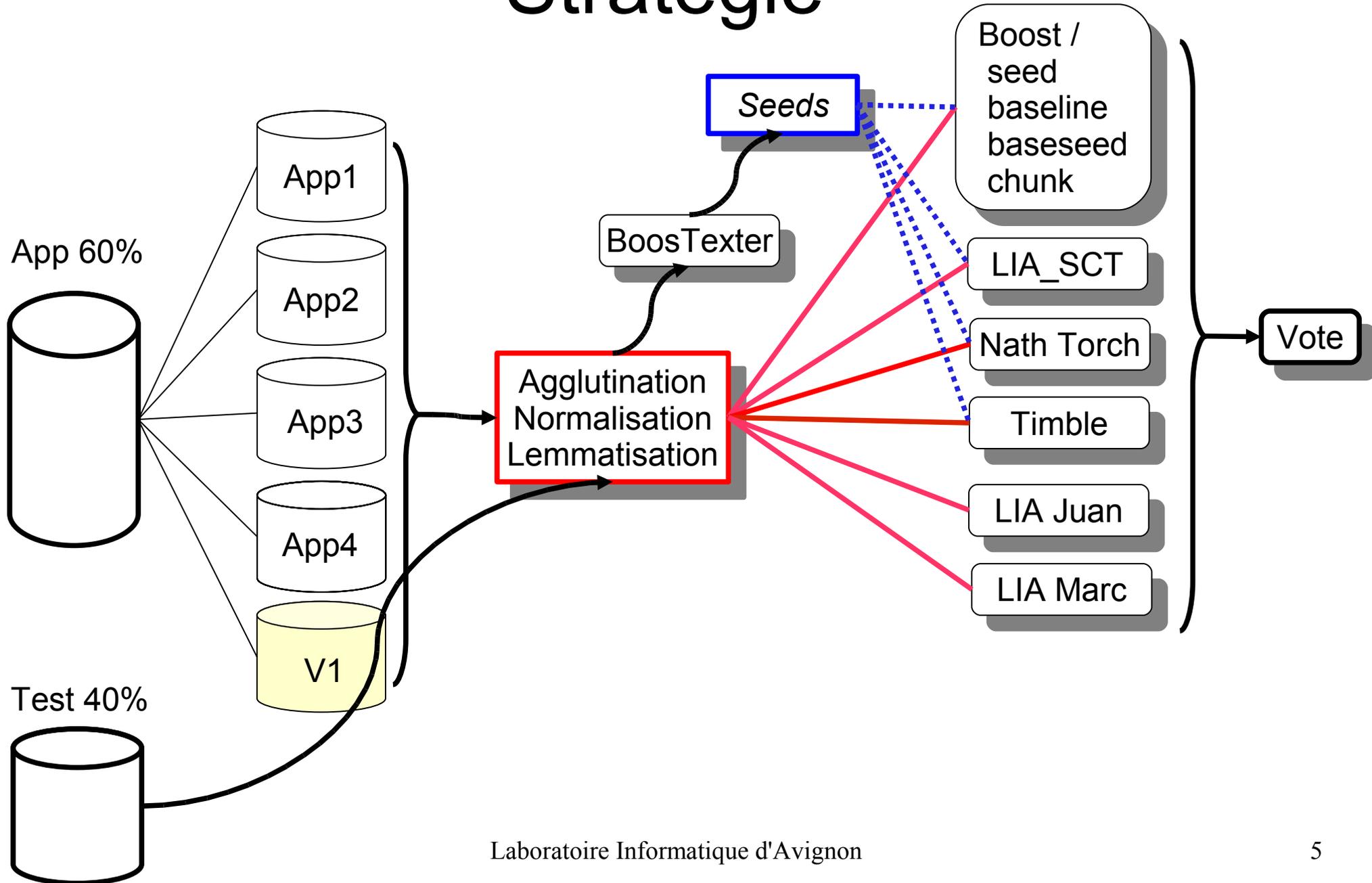
$$Rappel_i = \frac{\{\text{Nb de documents correctement attribués à la classe } i\}}{\{\text{Nb de documents appartenant à la classe } i\}}$$

Indice de confiance pondérée: probabilité d'un document d'appartenir à une classe d'opinion donnée

# Représentation de textes

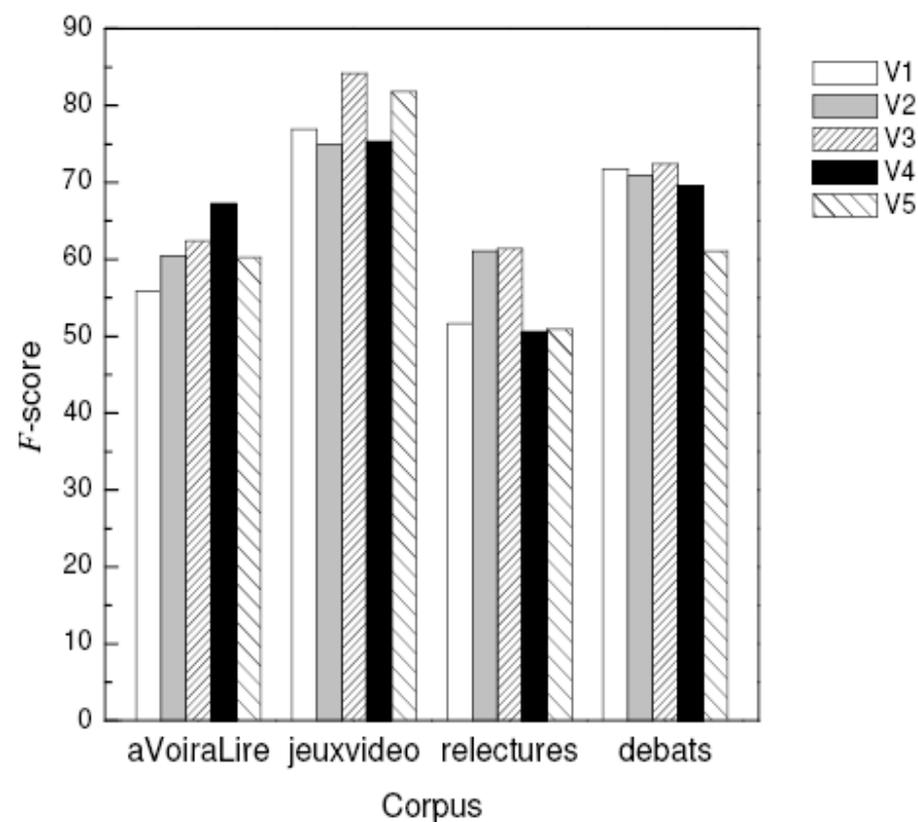
- Mots, POS,  $n$ -grammes,  $n$ -lemmes
- Etiquette supplémentaire: le *seed* (Wilson et al., 2005)
  - Polarité d'opinions
  - Liste manuelle
  - BoosTexter : *seed* remplacé par son lemme :  
lexique de 2000 *seeds*
- Normalisation/Agglutination/Lemmatisation
  - Avec ou sans

# Stratégie



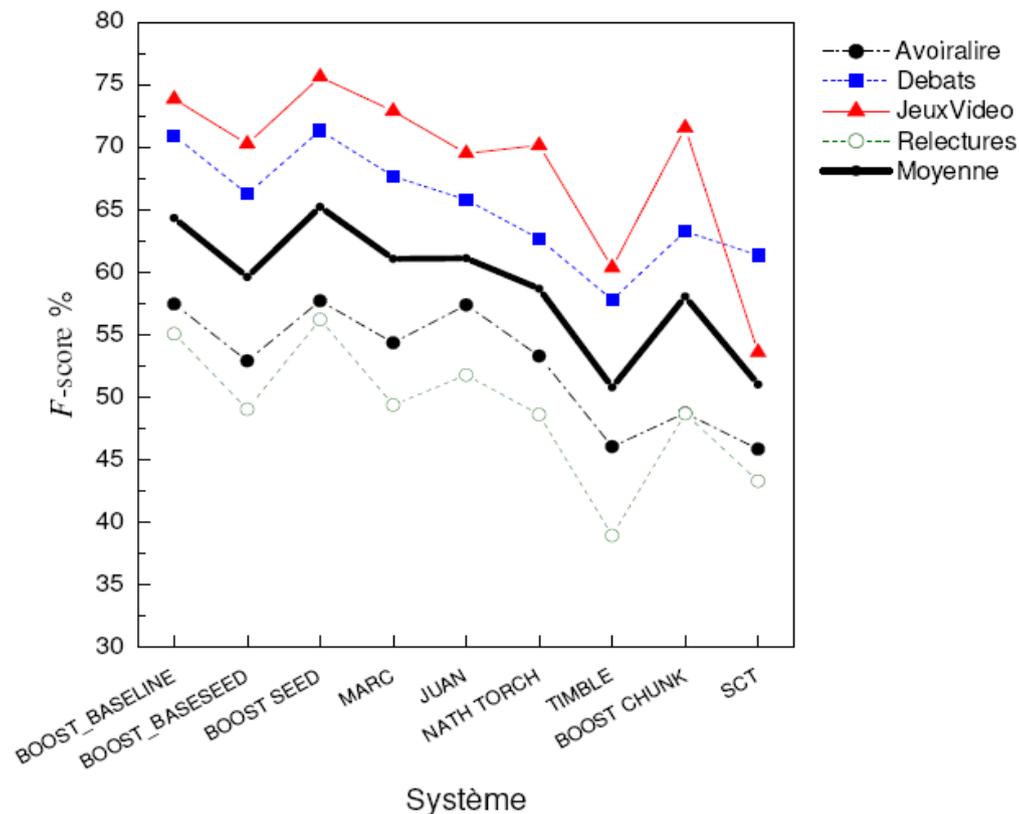
# Corpus de validation

Corpus	Précision	Rappel	$F$ -score	Correctes	Total
<b>aVoiraLire</b> (V)	0,6419	0,5678	0,6026	1 385	2 074
<b>jeuxvideo</b> (V)	0,8005	0,7730	0,7865	2 005	2 537
<b>relectures</b> (V)	0,5586	0,5452	0,5518	510	881
<b>debats</b> (V)	0,7265	0,7079	0,7171	12 761	17 299

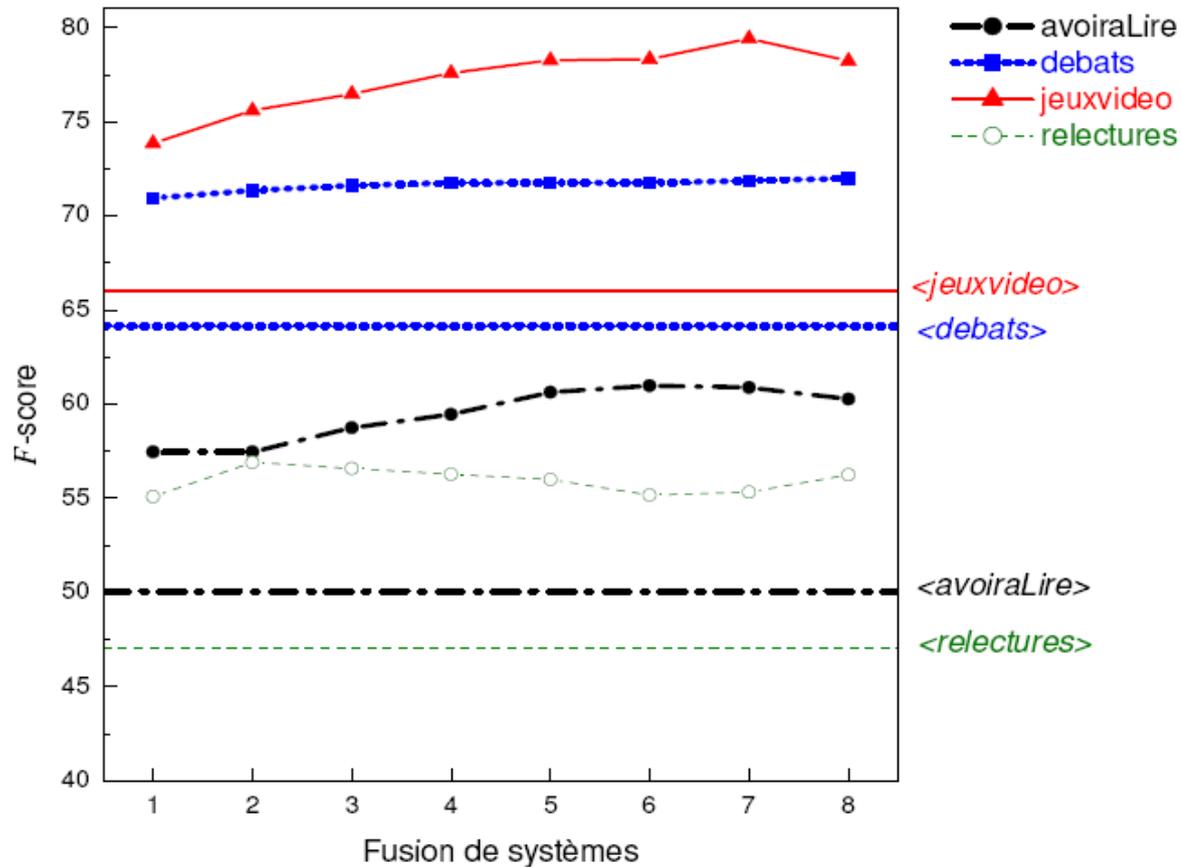


# Corpus de test (1)

Corpus	Précision	Rappel	$F$ -score	Correctes	Total
<b>aVoiraLire</b> (T)	0,6540	0,5590	0,6028	931	1 386
<b>jeuxvideo</b> (T)	0,8114	0,7555	0,7824	1 333	1 694
<b>relectures</b> (T)	0,5689	0,5565	0,5626	353	603
<b>debats</b> (T)	0,7307	0,7096	0,7200	8 403	11 533

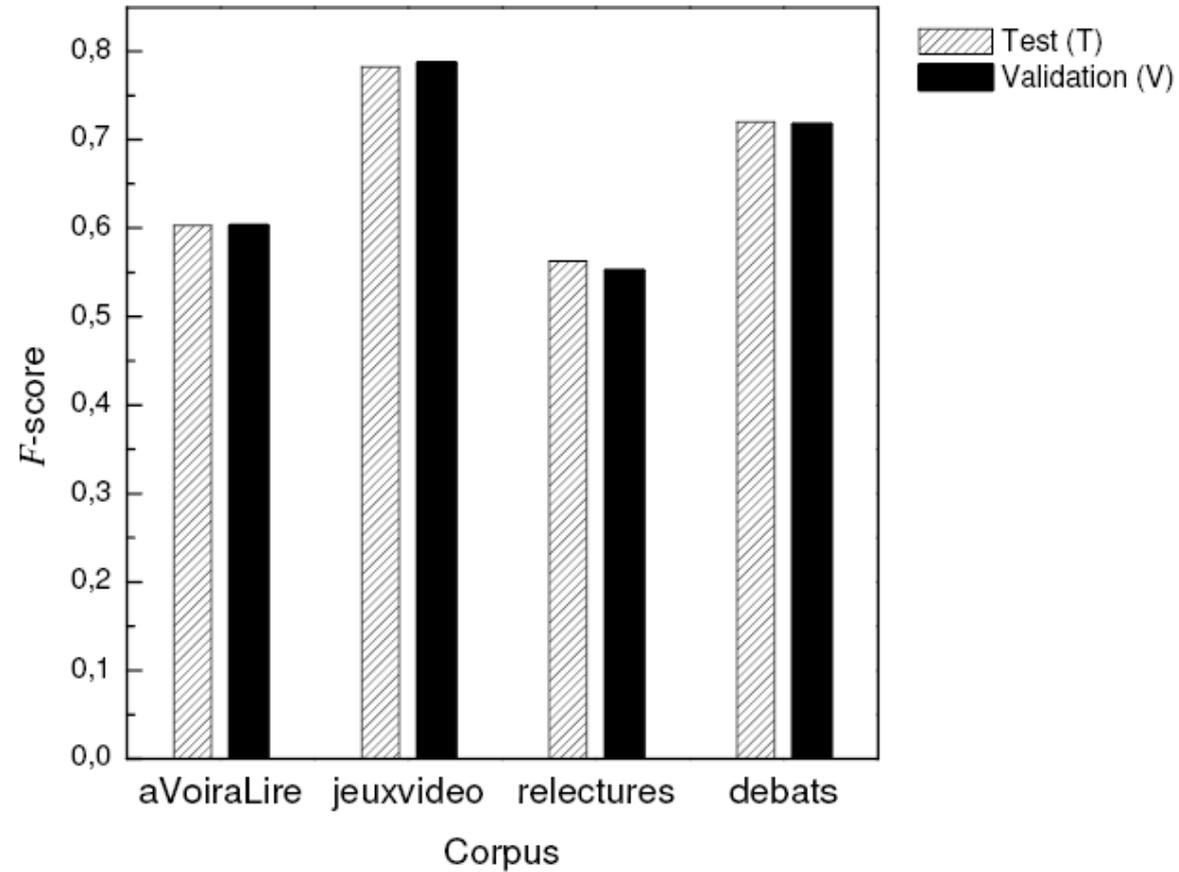


# Corpus de test (2)



- 1 BOOST\_BASELINE
- 2 BOOST\_BASELINE+ BOOST\_BASESEED + BOOST\_SEED
- 3 BOOST\_BASELINE+ BOOST\_BASESEED + BOOST\_SEED + MARC
- 4 BOOST\_BASELINE+ BOOST\_BASESEED + BOOST\_SEED + JUAN
- 5 BOOST\_BASELINE+ BOOST\_BASESEED + BOOST\_SEED + JUAN + NATH\_TORCH
- 6 BOOST\_BASELINE+ BOOST\_BASESEED + BOOST\_SEED + JUAN + NATH\_TORCH + TIMBLE
- 7 BOOST\_BASELINE+ BOOST\_BASESEED + BOOST\_SEED + JUAN + NATH\_TORCH + TIMBLE + BOOST\_CHUNK
- 8 BOOST\_BASELINE+ BOOST\_BASESEED + BOOST\_SEED + JUAN + NATH\_TORCH + TIMBLE + BOOST\_CHUNK + SCT

# Corpus de test (3)



# Discussion

## **3 :2 relectures**

*L'idée d'appliquer les méthodes de classification pour définir des classes homogènes de pages web est assez originale par contre, la méthodologie appliquée est classique. Je recommande donc un « weak accept » pour cet article.*

Classé **1** par le système (*accepté avec des modifications majeures*).  
De la lecture directe on pourrait en déduire que la classe est **1**...  
mais la référence est **2** (*accepté*)

## **3.6 relectures**

*Article trop court pour pouvoir être jugé. Je suggère de le mettre en POster si cela est prévu.*

Document trop court : les systèmes l'ont rejeté, mais il a été accepté

### **3 :9 relectures**

Question : comment est construit le réseau bayésien ? Un peu bref ici... Remarques de forme : page 2, 4ème ligne, « comprend » 5ème ligne : "annotées" ou "annoté" page 3 : revoir la phrase confuse précédant le tableau dernière ligne, répétition de "permet" page 5 : 7ème ligne accorder "diagnostiqué" et "visé" avec "états" ou avec "connaissances"

Rejeté mais il a été accepté... sur des remarques de forme uniquement !

### **aVoiraLire 1 :10**

Depuis trente-six ans, chaque nouvelle production de David Bowie est un événement. *Heathen*, ne fait pas exception à cette règle. On reconnaît instantanément la patte de son vieux compère Tony Visconti. La voix de Bowie est mise en avant. Agréable surprise, surtout qu'elle n'a rien perdu depuis ses débuts. Là, commence le voyage. Ambiance, mélange dosé des instruments. Dès l'ouverture de l'album avec *Sunday*, un sentiment étrange nous envahit. Comme si Bowie venait de rentrer d'un voyage expérimental au cœur même de la musique. Retour aux sources. L'ensemble du disque est rythmé par cette pulsation dont le rythme est le secret. Le tout saupoudré de quelques pincées d'électronique. Le groupe est réduit au minimum. Outre Bowie en chef d'orchestre et Visconti, David Torn ponctue les compositions de ses guitares aventureuses et Matt Chamberlain apporte de l'âme à la rythmique. Un quatuor à cordes fait une apparition, comme Pete Townshend (*The Who*) ou Dave Grohl (ex-batteur de Nirvana). Avec trois reprises réarrangées et neuf compositions originales, le 25e album de Bowie est à l'image d'une cohérence artistique retrouvée.

# Conclusions

- L'union fait l'*intelligence* : combinaison de plusieurs méthodes numériques/probabilistes
- Indépendance de la langue
- Indépendance du contexte
- Résultats honorables
- Adaptable à d'autres problématiques