

# Analyse de texte / du discours

Brigitte Grau  
Master 1 DAC - UPMC

Illustration de TextTiling: Rohit Kate, Univ. of Wisconsin

## Sommaire

- **Introduction**
- **Analyse thématique de texte**
  - Approche fondée sur les connaissances
  - Approche fondée sur une analyse de surface
    - Approches locales
      - TextTiling
      - Ajout de connaissances externes
      - Chaines lexicales
    - Approches globales
- **Structuration de texte**
- **Evaluation**
- **Conclusion**

2

## Discours

- Les productions en langue naturelle sont des successions de phrases
- Une phrase souvent non compréhensible de manière isolée

Marie est allée acheter un manteau. Elle en a trouvé un rouge qu'elle aimait vraiment. Quand elle l'a rapporté chez elle, elle a constaté qu'il allait parfaitement avec sa robe favorite.

### 2 questions :

- Q1 : Qu'est-ce que Marie a trouvé ?
- Q2 : Marie a-t-elle acheté quelque chose ?

3

## Discours

- Discours :
  - Ensemble cohérent et structuré de (groupes de) phrases
- Tâches pour analyser un discours
  - Analyse thématique (Segmentation en unités thématiques)
  - Détermination de relations de cohérence (ou relations discursives ou relations rhétoriques)
  - Résolution d'anaphores

4

## Analyse thématique de texte

### ■ Définition informelle

- Analyse thématique : rendre compte de ce dont "parlent" les textes
  - exemples : texte "parlant" de la crise du pétrole, ou d'une tentative d'assassinat de Martin Luther King

### ■ Par rapport à la langue

- analyse reposant sur une masse importante d'éléments extra-linguistiques
- statut non défini de la notion de thème / à la langue

5

## Décomposition de l'analyse thématique

### ■ 3 sous-problèmes en interrelation

- segmentation thématique
  - définir des segments de texte thématiquement homogènes
    - ➡ définition/délimitation des unités
- identification thématique
  - associer la représentation d'un thème à chaque segment
    - ➡ typage des unités
- suivi thématique
  - trouver les relations existant entre les segments (relations de changement de thème, de généralisation, ...)
    - ➡ rôle des unités ; définition d'une structure du texte

6

## Une illustration

Segment 1  
Séance  
de dédicace

Il y a quelques années, je me trouvais dans un grand magasin de Harlem, entouré de quelques centaines de personnes. J'étais en train de dédicacer des exemplaires de mon livre "Stride toward Freedom", qui relate le boycottage des autobus de Montgomery en 1955-56. Soudain, tandis que j'apposais ma signature sur une page,

-----  
changement de thème  
-----

Segment 2  
Attentat

je sentis quelque chose de pointu s'enfoncer brutalement dans ma poitrine. Je venais d'être poignardé à l'aide d'un coupe-papier, par une femme qui devait être reconnue folle par la suite. On me transporta d'urgence à l'Hôpital de Harlem

-----  
déviation de thème  
-----

Segment 3  
Hôpital

où je restai de longues heures sur un lit tandis qu'on faisait mille préparatifs pour extraire l'arme de mon corps

7

## Applications de l'analyse thématique

### ■ Segmentation

- structuration d'un flot textuel continu (ex.: transcriptions audios)
- RI et Extraction en contexte (ex. : détection d'événement, indexation, ...)
- routage de documents
- classification de documents

### ■ Suivi

- structuration des textes, repérage des éléments importants → résumé, filtrage, visualisation

8

## Approche fondée sur les connaissances

- **Années 80-90**
- **Modèle de discours (Grosz 1985)**
- **Travaux sur la compréhension de texte**
  - Travaux de Schank (1977: SAM/PAM), Dyer et al. (1983: Boris)
  - Définition de scripts, schémas représentant des situations
  - Exemple : aller au restaurant, tentative d'assassinat, réparer une voiture...
- **Structuration des connaissances :**
  - Un schéma fait référence à d'autres schémas pour des événements complexes
  - Hierarchisation des schémas

9

## Approche fondée sur les connaissances

### Mécanisme de compréhension

- **Un schéma** = un niveau de description d'un sujet
- **Mécanisme général de compréhension**
  - Identification du schéma, donc du thème
  - Remplissage par appariement avec les événements du texte
  - Recherche du lien entre différents schémas représentatifs de segments de discours différents
- **Illustration**
  - Programmes généraux de compréhension d'histoire
  - Les programmes SAM et PAM de Schank et Abelson (77)

10

## Approche fondée sur les connaissances

- **Bilan**
- **Intérêt**
  - Niveau de compréhension fin : segmentation, identification, suivi du thème.
  - Réponse à tout type de question: quoi, qui, comment, pourquoi
- **Inconvénient**
  - Difficile de modéliser ce type de connaissances
    - Application à des domaines restreints
- **Transformation de la tâche**
  - Extraction d'information à partir de « templates » ~ (~85)
  - Evaluations MUC
- **Connaît un renouveau**
  - Story telling

11

## Approche fondée sur des marques de surface

- **Segmentation d'un texte en différents segments thématiques**
- **Non supervisé :**
  - Pas de données d'entraînement
- **Critères**
  - Cohésion lexicale
  - Marques linguistiques
- **2 approches**
  - Globale : regrouper les phrases, passages en sous-thèmes qui forment un bloc cohésif
  - Locale : un changement de thème correspond à une rupture de cohésion
    - Travaux de Salton (1996), Hearst (1997)

12

## Définitions

### ■ Cohésion

- Lien entre unités textuelles dues à des procédés linguistiques (anaphores, cohésion lexicale, ellipses, conjonctions, etc.)

### ■ Cohésion lexicale :

- Utilisation des mêmes mots ou de mots similaires ou dans le même champ sémantique pour lier des unités de texte
  - Répétition des mots (distribution) : localiser un thème
  - Répartition des mots (position) : évaluer l'importance du mot/thème
- Exemple : extrait du [vin jaune](#)

13

## Segmentation thématique par analyse locale

### ■ Mécanisme d'analyse : TextTiling

- Hearst, 1997
- Principe : identification des changements de thème
- Association d'un vecteur de descripteurs à une zone de texte
  - Ex un descripteur : un mot plein
  - Ex. une valeur : nombre d'occurrences, pondéré par la répartition / différentes zones
  - Ex. zone = 20 mots

14

## Segmentation thématique par analyse locale

### ■ Mécanisme d'analyse : TextTiling

- Hearst, 1997
- Principe : identification des changements de thème
- Association d'un vecteur de descripteurs à une zone de texte (pseudo-phrase)
  - Ex un descripteur : un mot plein
  - Ex. une valeur : nombre d'occurrences, pondéré par la répartition / différentes zones
  - Ex. zone = 20 mots



15

## Segmentation thématique par analyse locale

### ■ TextTiling (suite)

### ■ Mesure d'un score entre vecteurs de zones consécutives

- ex. produit scalaire, cosinus : élevés quand mots cooccurrent

### ■ Score de cohésion lexicale

- Similarité des descripteurs avant et après le gap (ex. 10 pseudo-phrases avant et après)



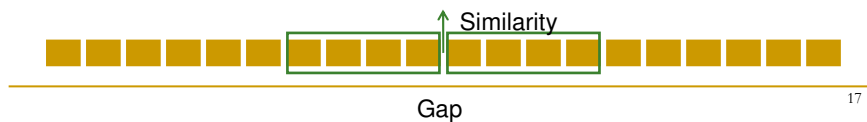
Gap

16



## Segmentation thématique par analyse locale

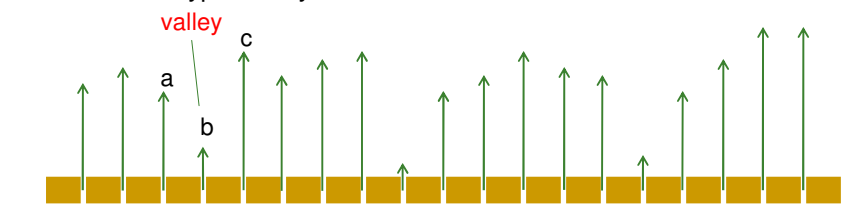
- TextTiling (suite)
- Mesure d'un score entre vecteurs de zones consécutives
  - ex. produit scalaire, cosinus : élevés quand mots cooccurrent
- Score de cohésion lexicale
  - Similarité des descripteurs avant et après le gap (ex. 10 pseudo-phrases avant et après)



17

## Segmentation thématique par analyse locale

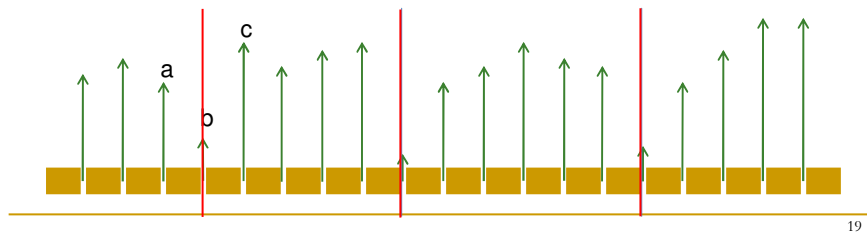
- TextTiling (suite)
- Construction du graphe des valeurs de similarités et calcul du score de rupture des « vallées »
  - $(a-b)+(c-b)$
- Un changement de thème
  - Si score de rupture > seuil
  - Ex. écart type > moyenne des scores



18

## Segmentation thématique par analyse locale

- TextTiling (suite)
- Construction du graphe des valeurs de similarités et calcul du score de rupture des « vallées »
  - $(a-b)+(c-b)$
- Un changement de thème
  - Si score de rupture > seuil
  - Ex. écart type > moyenne des scores



19

## Segmentation thématique par analyse locale

- Application
  - (Masson et al. 95)
- Un texte :
  - un ensemble de paragraphes
- Un paragraphe :
  - un vecteur de descripteurs
- Extraction des descripteurs :
  - balisage SGML, segmentation, étiquetage-lemmatisation, filtrage des catégories
  - Noms (simples et composés), verbes, adjectifs
- Pondération des vecteurs par le facteur **tf.idf**
  - renforce les descripteurs peu dispersés et inhibe ceux qui le sont
- Distance entre paragraphes adjacents :
  - coefficient de Dice (produit scalaire normalisé des vecteurs)

20

## Segmentation thématique par analyse locale

**Tf.idf :**  $w_{ij} = \text{tf}_{ij} \cdot \log \frac{N}{df_j}$

$w_{ij}$  : poids du terme  $T_j$  dans  $\S_i$   
 $\text{tf}_{ij}$  : nombre d'occurrences du terme  $T_j$  dans  $\S_i$   
 $N$  : nombre total de  $\S$  dans texte  
 $df_j$  : nombre de  $\S$  contenant  $T_j$

**Coefficient de Dice :**

$$\text{Dice} = \frac{2 \sum_{i=1}^t w(x_i) \cdot w(y_i)}{\sum_{i=1}^t w(x_i)^2 + \sum_{i=1}^t w(y_i)^2}$$

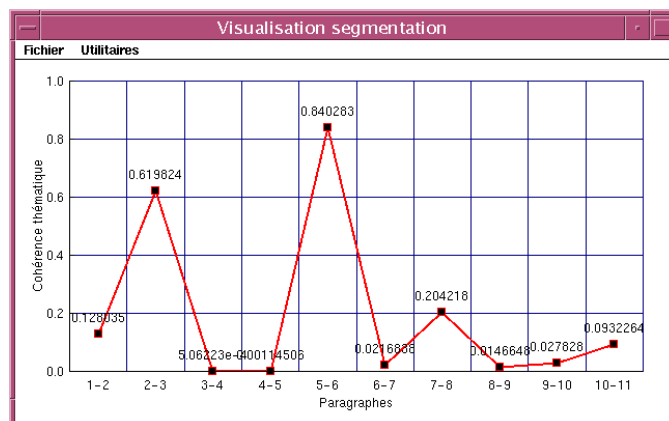
$X = (x_1, x_2, \dots, x_t)$

$Y = (y_1, y_2, \dots, y_t)$

21

## Segmentation thématique par analyse locale

### ■ Application : Résultats sur le texte “Vin jaune”



22

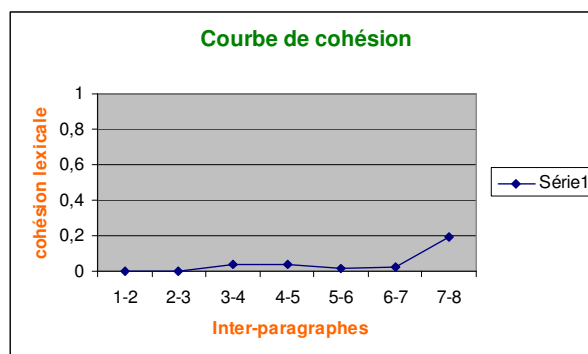
## Segmentation thématique par analyse locale

- **Application sur un texte de type journalistique**
  - Article paru dans Le Monde traitant d'un accident de train

23

## Segmentation thématique par analyse locale

- **Application : Résultats sur le texte journalistique**



Calcul de la cohésion inter-paragraphe du texte portant sur l'accident de train

24

## Segmentation thématique par analyse locale avec ajout de connaissances

### ■ Utilisation de marques de surface

- Litman et Passonneau 1997
  - Exemple : Ensuite, parce que, ...
- Problème : définition de ces marques

### ■ Utilisation de connaissances non dédiées

- Kozima (1993), Morris & Hirst (1991), Ferret (1998)
- Connaissances peu structurées :
  - ex. réseau de co-occurrences, thésaurus, dictionnaire
- Connaissances disponibles ou pouvant être construite automatiquement
- Même méthode que précédemment

25

## Segmentation thématique par analyse locale avec ajout de connaissances

### ■ But : améliorer le « rapprochement » des paragraphes thématiquement liés

### ■ 2 exemples de méthodes :

- **M1** : Utilisation des relations sémantiques et pragmatiques entre termes (Masson et al. 95)
  - Pour compléter la représentation vectorielle
    - Estimation des liaisons entre termes par comptage des cooccurrences dans un gros corpus
  - Pour calculer une valeur de cohésion à chaque pas du texte
- **M2** : Calcul de la cohésion lexicale sur réseau de mots avec des liaisons pondérées acquises automatiquement (Ferret et al. 98)

26

## Segmentation thématique par analyse locale avec ajout de connaissances

### Construction du réseau

- **Corpus** : 24 mois du journal « Le Monde » (entre 1990 et 1994)
- **Méthode**
  - Pré-traitement des textes
  - Comptage des cooccurrences dans une fenêtre glissante
    - Taille : 20 lemmes (environ 50 mots dans le texte original ⇔ 2 phrases)
    - Pas d'ordre :  $\text{coo}(m1, m2) + \text{coo}(m2, m1)$
    - Respect des frontières de texte
  - Filtrage des cooccurrences faibles ( $< 6$ )
  - Estimation de l'information mutuelle entre termes cooccurents



31 000 lemmes et 7 millions de relations

27

## Segmentation thématique par analyse locale avec ajout de connaissances

### Extrait du réseau

Lemme1	Lemme2	Nombre occurrences	Valeur de cohésion	Type de lien
imprimante	ordinateur	13	0,227	pragmatique
bateau	voilier	125	0,224	sémantique
prêtre	curé	44	0,209	sémantique
policier	cambriolage	41	0,190	pragmatique
chômage	emploi	1985	0,167	sémantique
prendre	racine	120	0,110	lexico-syntaxique
collision	franc	7	0,076	bruit

28

## Segmentation thématique par analyse locale avec ajout de connaissances

- **Méthode 1 : Utilisation du réseau lexical**
- **Principe : tenir compte des liaisons dans la représentation vectorielle**
  - $P1 = (kmot1, nmot2, 0, \dots)$  et  $P2 = (0, 0, tmot3, \dots)$
- **Renforcement du poids des descripteurs liés**
  - mot1 et mot2 liés par liaison w
  - $P1' = ((wn+k)mot1, (wk+n)mot2, 0, \dots)$
- **Ajout de descripteurs dans les vecteurs**
  - mot1 et mot3 liés par liaison z
  - $P1' = ((wn+k)mot1, (wk+n)mot2, zkmot3, \dots)$
  - $P2' = (ztmot1, 0, tmot3, \dots)$

29

## Segmentation thématique par analyse locale avec ajout de connaissances

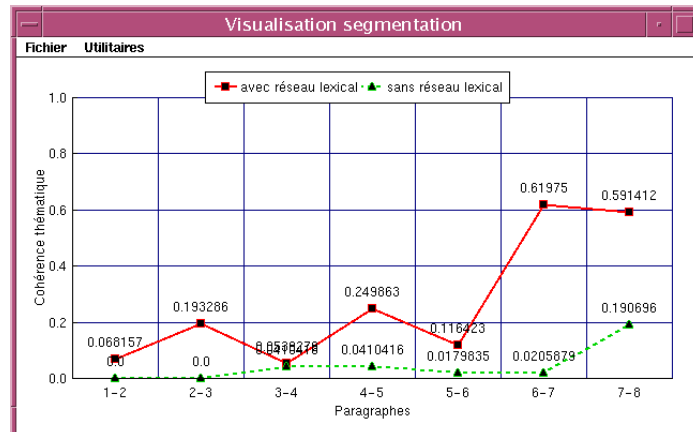
- **Méthode 1 (suite)**
- **Pré-traitement, vectorisation, renforcement et ajout de descripteurs**
- **Pondération de type tf.idf (élimination du bruit dû aux nombreux liens du réseau)**
- **Collocations avec « collision »**

Cooccurrences significatives		Cooccurrences peu significatives	
Accident	0.210329	Pays	0.0684458
Conducteur	0.19373	Président	0.0817475
Train	0.185325	Ministre	0.0890158
Gare	0.181422	Franc	0.0763806

30

## Segmentation thématique par analyse locale avec ajout de connaissances

### ■ Application sur le texte “Accident de train”



31

## Segmentation thématique par analyse locale avec ajout de connaissances

- Méthode 2 : variation sur le calcul de la cohésion lexicale
- Fenêtre glissante de 20 mots
  - Pas de découpage a priori en paragraphes
    - Evite le problème de tailles non homogènes
- Calcul de cohésion pour chaque mot au centre de la fenêtre

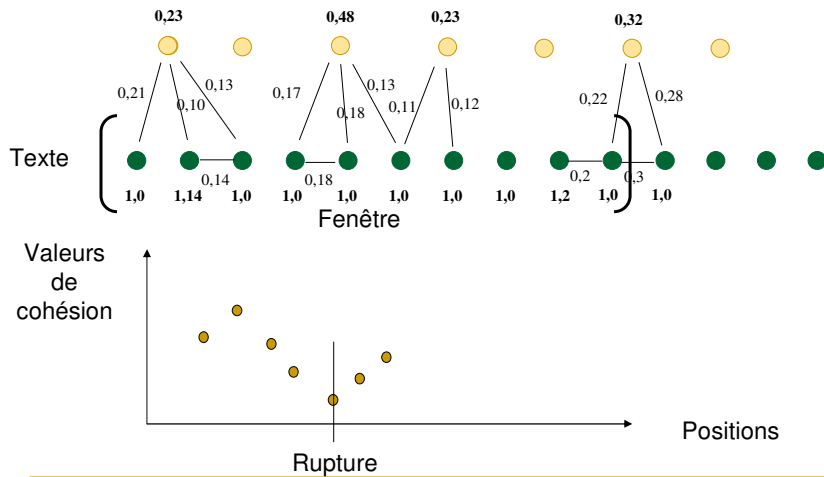
32



## Segmentation thématique par analyse locale avec ajout de connaissances

Méthode 2 : variation sur le calcul de la cohésion lexicale

Réseau de cooccurrences



33

## Segmentation thématique par analyse locale avec ajout de connaissances

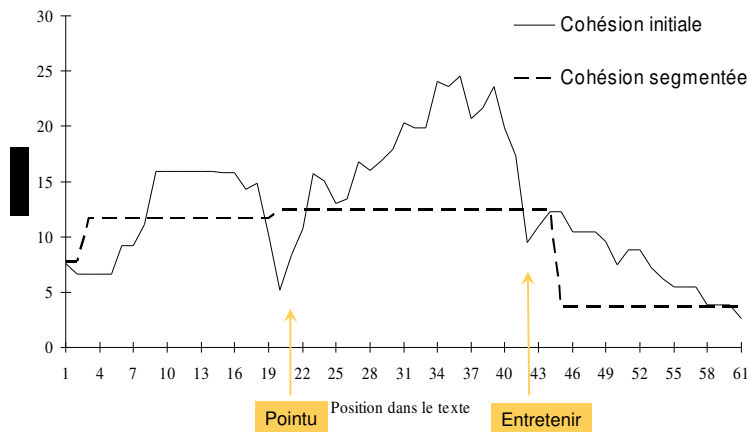
Exemple : Attentat de MLK

- Segment 1**  
Séance de dédicace
- Il y a quelques années, je me trouvais dans un grand magasin de Harlem, entouré de quelques centaines de personnes. J'étais en train de dédicacer des exemplaires de mon livre "Stride toward Freedom", qui relate le boycottage des autobus de Montgomery en 1955-56. Soudain, tandis que j'apposais ma signature sur une page,
- Segment 2**  
Attentat
- je sentis quelque chose de pointu s'enfoncer brutalement dans ma poitrine. Je venais d'être poignardé à l'aide d'un coupe-papier, par une femme qui devait être reconnue folle par la suite. On me transporta d'urgence à l'Hôpital de Harlem
- Segment 3**  
Hôpital
- où je restai de longues heures sur un lit tandis qu'on faisait mille préparatifs pour extraire l'arme de mon corps. Beaucoup plus tard, quand je fus en état de m'entretenir avec le Dr. Aubrey Maynard, le chirurgien en chef qui exécuta cette délicate et dangereuse intervention, j'appris les raisons de cette longue attente avant l'opération. La lame de l'instrument avait touché l'aorte et, pour l'extraire, il fallait ouvrir toute la cage thoracique.

34

## Segmentation thématique par analyse locale avec ajout de connaissances

Exemple : Attentat de MLK : courbes et segments obtenus



35

## Segmentation thématique par analyse locale avec ajout de connaissances

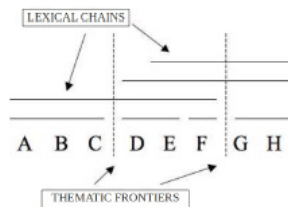
Exemple : Attentat de MLK : segments obtenus

- Segment 1**  
Séance de dédicace
- Il y a quelques années, je me trouvais dans un grand magasin de Harlem, entouré de quelques centaines de personnes. J'étais en train de dédicacer des exemplaires de mon livre "Stride toward Freedom", qui relate le boycottage des autobus de Montgomery en 1955-56. Soudain, tandis que j'apposais ma signature sur une page,
- Segment 2**  
Attentat
- je sentis quelque chose de **pointu** s'enfoncer brutalement dans ma poitrine. Je venais d'être poignardé à l'aide d'un coupe-papier, par une femme qui devait être reconnue folle par la suite. On me transporta d'urgence à l'Hôpital de Harlem
- Segment 3**  
Hôpital
- où je restai de longues heures sur un lit tandis qu'on faisait mille préparatifs pour extraire l'arme de mon corps. Beaucoup plus tard, quand je fus en état de **m'entretenir** avec le Dr. Aubrey Maynard, le chirurgien en chef qui exécuta cette délicate et dangereuse intervention, j'appris les raisons de cette longue attente avant l'opération. La lame de l'instrument avait touché l'aorte et, pour l'extraire, il fallait ouvrir toute la cage thoracique.

36

## Segmentation thématique par analyse locale : chaines lexicales

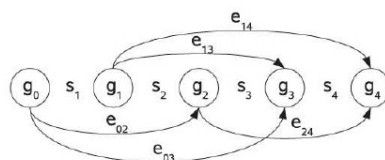
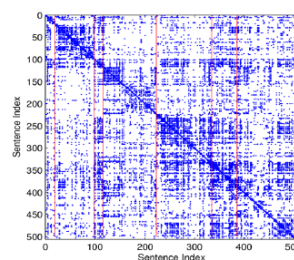
- **Chaine lexicale**
  - (Morris and Hirst, 1991)
  - Relie mêmes mots et mots sémantiquement liés
  - Liaison existe selon la position des mots
    - Distance entre 2 occurrences < seuil
- **Segmentation**
  - LCSeg (Galley et al., 2003) : outil disponible



37

## Segmentation thématique par analyse globale

- **Maximiser la valeur de cohésion lexicale de chaque segment obtenu**
  - Matrices de similarité des phrases
    - Dotplot (Reynar, 94)
  - Graphe
    - TextSeg (Utiyama & Isahara, 2001)
    - Arc = segment, nœud = frontière
    - Valeur de cohésion probabiliste = capacité d'un ML appris sur le segment à prédire les mots du segment



38

## Segmentation thématique

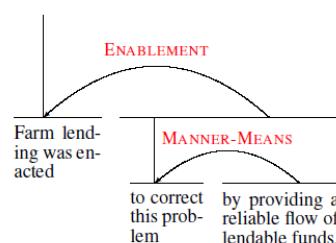
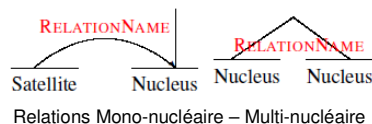
### Bilan

- **Niveau de compréhension**
  - ❑ Délimite les zones de texte relatives à différents thèmes
  - ❑ Pas de structuration (sauf par étude des valeurs de liaison entre toutes les zones)
  - ❑ Pas d'identification de thème
- **Intérêt**
  - ❑ Pas de connaissances spécifiques
  - ❑ Application à de nombreux textes indépendamment du domaine
  - ❑ Utilisation en résumé, recherche d'information
    - localiser les informations plus précisément dans le texte
- **Problèmes**
  - ❑ Résultats plus ou moins bons selon le type de texte : scientifiques ou dépêches
  - ❑ Résultats peu précis et peu fiables
  - ❑ Corpus de test

39

## Structuration de texte

- **Hierarchie de segments**
- **Structure par emboîtement de segments**
  - ❑ Ex. Projet REGAL : Structure « gros grain »
  - ❑ Exemple
- **Structure phrase par phrase : relations discursives**
  - ❑ Théorie RST (Rhetorical Structure Theory): Mann & Thomson



- **Ressources**
  - ❑ Penn Discourse Treebank : <http://www>
  - ❑ RST-DT (Carlson et al. 2001)

40

## Structuration de texte

### ■ Apprentissage des relations

- Système HILDA (Hernault et al. 2010)
  - Décomposition en EDU par apprentissage
    - $W \rightarrow [0, 1]$ ,  $W$  = mots, 1 = mot est frontière d'EDU
  - Etiquetage des relations (18)
    - 2 classifieurs : existence d'une relation + label
    - Traits : organisation du texte, marques de discours : 3-gram début et fin EDU, traits syntaxiques de dominance, ...
- Système DST
  - Structure « grain fin » (N. Hernandez, 2000)
  - Apprentissage des relations entre 2 phrases
    - Subordination, coordination, absence de relation
    - Critères :
      - marques linguistiques, cohésion lexicale, suivi thème-rhème, parallélisme syntaxique
  - Exemple

41

## Evaluation

### ■ TDT : Topic Detection and Tracking

#### ■ Trois tâches :

- Segmentation d'histoire : segmentation d'un flot continu
- Suivi de thème : association d'une histoire avec un thème connu du système
- Détection de thème : détection et suivi de thèmes inconnus du système

#### ■ Méthodes :

- Construction de modèle de langage pour représenter des thèmes à partir d'un corpus d'apprentissage

42

## Segmentation thématique

### Evaluation

- **Problème avec mesures de précision / rappel**
- **Introduction d'une métrique**
  - WindowDiff (Pevzner & Hearst, 2002)
  - Glissement d'une fenêtre de longueur k sur la référence (correcte) et sur la segmentation évaluée, et comptabilisation de nombre de coupures dans chaque
  - Métrique WindowDiff : moyenne des différences dans le nb de ruptures dans la fenêtre glissante

43

## Conclusion

- **Evaluation difficile**
- **Résultats assez moyens**
  - Place pour amélioration, mais ...
- **Questions ouvertes**
  - Définition de la notion de thème
  - Structurations différentes des textes
    - Structure thématique
    - Structure rhétorique
  - Tout passage appartient-il aux 2 structures ?
- **Lecture conseillée :**
  - Semantic structuring of video collections from speech : segmentation and hyperlinking, Anca Simon, thèse de Université de Rennes

44

<p3>

On soupçonna alors le 4,5 diméthyl 3 hydroxy 2(5H) furanone, ou **sotolon**, molécule construite autour d'un cycle de quatre **atomes** de carbone et d'un **atome** d'oxygène. Comme le **sotolon** et la **solérone** sont en concentrations minimales dans les vins de voile et, de surcroît, chimiquement instables, les chimistes dijonnais ont cherché à optimiser leur extraction afin de déterminer la **molécule** responsable du goût de jaune.

<p4>

L'analyse la plus directe d'extraits de vins est la chromatographie: on injecte un échantillon dans un solvant que l'on vaporise et on fait traverser au **mélange** une **colonne** revêtue intérieurement d'un polymère, qui retient les divers **composés** du **mélange** à des degrés divers; en bas de la **colonne**, on détecte la sortie des **composés** séparés. Le premier travail des chimistes fut la mise au point d'une variante de cette technique pour identifier les **composés** présents en quantités minimales dans des **mélanges** complexes.

<p5>

Les chromatogrammes d'échantillons de vin furent alors comparés à ceux de **solutions** pures de **sotolon** et de **solérone** de synthèse: le **sotolon** est ainsi présent entre 40 et 150 parties par milliard dans les sherrys; la **solérone** semble moins spécifique, et ses **concentrations** sont supérieures dans les sherrys, ce qui explique pourquoi on l'a d'abord trouvée dans ces vins.

<p6>

Enfin les dosages, complétés de tests sensoriels des fractions séparées, montrèrent que la **solérone**, aux **concentrations** trouvées dans du savagnin (le cépage à partir duquel on fabrique le vin jaune), n'était perçue par les consommateurs ni dans les vins, ni dans des **solutions** modèles: la **solérone** n'était pas la **molécule** caractéristique; le jugement était sans appel.

<p7>

En 1992, les chimistes se consacrèrent alors complètement au **sotolon**, qui avait été observé dans des molasses de canne à sucre, dans des graines de fenugrec, dans de la sauce de soja, dans du saké...

45

## Extrait de roman

- Extrait du livre *De la Terre à la Lune* de Jules Verne

46

## DST – type de relations

- (1) Les résumés par extraction sélectionnent des phrases d'un texte source selon leur importance.
- (2) Les critères d'importance incluent la présence de termes fréquents, des mots clefs tels que « en résumé », « meilleur », et la position de la phrase dans le texte.
- (3) Cette approche est illustrée par le système ADAM.
- (4) Un autre exemple est donné par [2].
- (5) Le problème de cette approche est que les phrases extraites ne constituent pas toujours un texte cohérent du fait d'anaphores ambiguës.

47

## DST – Types de relations

- (1) Les résumés par extraction sélectionnent des phrases d'un texte source selon leur importance.

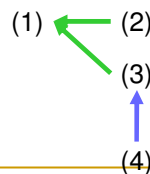
Subordination

- (2) Les critères d'importance incluent la présence de termes fréquents, des mots clefs tels que « en résumé », « meilleur », et la position de la phrase dans le texte.

Coordination

- (3) Cette approche est illustrée par le système ADAM.
- (4) Un autre exemple est donné par [2].

Structure construite



48



# REGAL - Principe de structuration

- Repérage de structures emboîtées [Masson, 1998]
  - Digressions, développements d'aspects particuliers
  - Fréquent dans textes expositifs

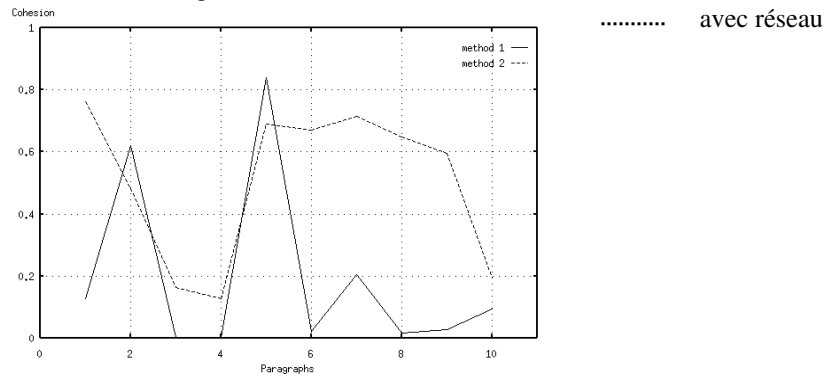
## Algorithme

- Recherche des 2 segments non-consécutifs les plus liés
- Application récursive pour les segments englobés ou non englobés restant

rg 101-102-103-104-105-106-107-108-109-110-111-112-113-114-115-116-117-118-119-120-121-122-123-124-125-126-127-128-129-130-131-132-133-134-135-136-137-138-139-140-141-142-143-144-145-146-147-148-149-150-151-152-153-154-155-156-157-158-159-160-161-162-163-164-165-166-167-168-169-170-171-172-173-174-175-176-177-178-179-180-181-182-183-184-185-186-187-188-189-190-191-192-193-194-195-196-197-198-199-200-201-202-203-204-205-206-207-208-209-210-211-212-213-214-215-216-217-218-219-220-221-222-223-224-225-226-227-228-229-230-231-232-233-234-235-236-237-238-239-240-241-242-243-244-245-246-247-248-249-250-251-252-253-254-255-256-257-258-259-260-261-262-263-264-265-266-267-268-269-270-271-272-273-274-275-276-277-278-279-280-281-282-283-284-285-286-287-288-289-290-291-292-293-294-295-296-297-298-299-300-301-302-303-304-305-306-307-308-309-310-311-312-313-314-315-316-317-318-319-320-321-322-323-324-325-326-327-328-329-330-331-332-333-334-335-336-337-338-339-340-341-342-343-344-345-346-347-348-349-350-351-352-353-354-355-356-357-358-359-360-361-362-363-364-365-366-367-368-369-370-371-372-373-374-375-376-377-378-379-380-381-382-383-384-385-386-387-388-389-390-391-392-393-394-395-396-397-398-399-400-401-402-403-404-405-406-407-408-409-410-411-412-413-414-415-416-417-418-419-420-421-422-423-424-425-426-427-428-429-430-431-432-433-434-435-436-437-438-439-440-441-442-443-444-445-446-447-448-449-450-451-452-453-454-455-456-457-458-459-460-461-462-463-464-465-466-467-468-469-470-471-472-473-474-475-476-477-478-479-480-481-482-483-484-485-486-487-488-489-490-491-492-493-494-495-496-497-498-499-500-501-502-503-504-505-506-507-508-509-510-511-512-513-514-515-516-517-518-519-520-521-522-523-524-525-526-527-528-529-530-531-532-533-534-535-536-537-538-539-540-541-542-543-544-545-546-547-548-549-550-551-552-553-554-555-556-557-558-559-560-561-562-563-564-565-566-567-568-569-570-571-572-573-574-575-576-577-578-579-580-581-582-583-584-585-586-587-588-589-590-591-592-593-594-595-596-597-598-599-600-601-602-603-604-605-606-607-608-609-610-611-612-613-614-615-616-617-618-619-620-621-622-623-624-625-626-627-628-629-630-631-632-633-634-635-636-637-638-639-640-641-642-643-644-645-646-647-648-649-650-651-652-653-654-655-656-657-658-659-660-661-662-663-664-665-666-667-668-669-670-671-672-673-674-675-676-677-678-679-680-681-682-683-684-685-686-687-688-689-690-691-692-693-694-695-696-697-698-699-700-701-702-703-704-705-706-707-708-709-710-711-712-713-714-715-716-717-718-719-720-721-722-723-724-725-726-727-728-729-730-731-732-733-734-735-736-737-738-739-740-741-742-743-744-745-746-747-748-749-750-751-752-753-754-755-756-757-758-759-760-761-762-763-764-765-766-767-768-769-770-771-772-773-774-775-776-777-778-779-780-781-782-783-784-785-786-787-788-789-790-791-792-793-794-795-796-797-798-799-800-801-802-803-804-805-806-807-808-809-810-811-812-813-814-815-816-817-818-819-820-821-822-823-824-825-826-827-828-829-830-831-832-833-834-835-836-837-838-839-840-841-842-843-844-845-846-847-848-849-850-851-852-853-854-855-856-857-858-859-860-861-862-863-864-865-866-867-868-869-870-871-872-873-874-875-876-877-878-879-880-881-882-883-884-885-886-887-888-889-890-891-892-893-894-895-896-897-898-899-900-901-902-903-904-905-906-907-908-909-910-911-912-913-914-915-916-917-918-919-920-921-922-923-924-925-926-927-928-929-930-931-932-933-934-935-936-937-938-939-940-941-942-943-944-945-946-947-948-949-950-951-952-953-954-955-956-957-958-959-960-961-962-963-964-965-966-967-968-969-970-971-972-973-974-975-976-977-978-979-980-981-982-983-984-985-986-987-988-989-990-991-992-993-994-995-996-997-998-999-1000-1001-1002-1003-1004-1005-1006-1007-1008-1009-1010-1011-1012-1013-1014-1015-1016-1017-1018-1019-1020-1021-1022-1023-1024-1025-1026-1027-1028-1029-1030-1031-1032-1033-1034-1035-1036-1037-1038-1039-1040-1041-1042-1043-1044-1045-1046-1047-1048-1049-1050-1051-1052-1053-1054-1055-1056-1057-1058-1059-1060-1061-1062-1063-1064-1065-1066-1067-1068-1069-1070-1071-1072-1073-1074-1075-1076-1077-1078-1079-1080-1081-1082-1083-1084-1085-1086-1087-1088-1089-1090-1091-1092-1093-1094-1095-1096-1097-1098-1099-1100-1101-1102-1103-1104-1105-1106-1107-1108-1109-1110-1111-1112-1113-1114-1115-1116-1117-1118-1119-1120-1121-1122-1123-1124-1125-1126-1127-1128-1129-1130-1131-1132-1133-1134-1135-1136-1137-1138-1139-1140-1141-1142-1143-1144-1145-1146-1147-1148-1149-1150-1151-1152-1153-1154-1155-1156-1157-1158-1159-1160-1161-1162-1163-1164-1165-1166-1167-1168-1169-1170-1171-1172-1173-1174-1175-1176-1177-1178-1179-1180-1181-1182-1183-1184-1185-1186-1187-1188-1189-1190-1191-1192-1193-1194-1195-1196-1197-1198-1199-1200-1201-1202-1203-1204-1205-1206-1207-1208-1209-1210-1211-1212-1213-1214-1215-1216-1217-1218-1219-1220-1221-1222-1223-1224-1225-1226-1227-1228-1229-1230-1231-1232-1233-1234-1235-1236-1237-1238-1239-1240-1241-1242-1243-1244-1245-1246-1247-1248-1249-1250-1251-1252-1253-1254-1255-1256-1257-1258-1259-1260-1261-1262-1263-1264-1265-1266-1267-1268-1269-1270-1271-1272-1273-1274-1275-1276-1277-1278-1279-1280-1281-1282-1283-1284-1285-1286-1287-1288-1289-1290-1291-1292-1293-1294-1295-1296-1297-1298-1299-1300-1301-1302-1303-1304-1305-1306-1307-1308-1309-1310-1311-1312-1313-1314-1315-1316-1317-1318-1319-1320-1321-1322-1323-1324-1325-1326-1327-1328-1329-1330-1331-1332-1333-1334-1335-1336-1337-1338-1339-1340-1341-1342-1343-1344-1345-1346-1347-1348-1349-1350-1351-1352-1353-1354-1355-1356-1357-1358-1359-1360-1361-1362-1363-1364-1365-1366-1367-1368-1369-1370-1371-1372-1373-1374-1375-1376-1377-1378-1379-1380-1381-1382-1383-1384-1385-1386-1387-1388-1389-1390-1391-1392-1393-1394-1395-1396-1397-1398-1399-1400-1401-1402-1403-1404-1405-1406-1407-1408-1409-1410-1411-1412-1413-1414-1415-1416-1417-1418-1419-1420-1421-1422-1423-1424-1425-1426-1427-1428-1429-1430-1431-1432-1433-1434-1435-1436-1437-1438-1439-1440-1441-1442-1443-1444-1445-1446-1447-1448-1449-1450-1451-1452-1453-1454-1455-1456-1457-1458-1459-1460-1461-1462-1463-1464-1465-1466-1467-1468-1469-1470-1471-1472-1473-1474-1475-1476-1477-1478-1479-1480-1481-1482-1483-1484-1485-1486-1487-1488-1489-1490-1491-1492-1493-1494-1495-1496-1497-1498-1499-1500-1501-1502-1503-1504-1505-1506-1507-1508-1509-1510-1511-1512-1513-1514-1515-1516-1517-1518-1519-1520-1521-1522-1523-1524-1525-1526-1527-1528-1529-1530-1531-1532-1533-1534-1535-1536-1537-1538-1539-1540-1541-1542-1543-1544-1545-1546-1547-1548-1549-1550-1551-1552-1553-1554-1555-1556-1557-1558-1559-1560-1561-1562-1563-1564-1565-1566-1567-1568-1569-1570-1571-1572-1573-1574-1575-1576-1577-1578-1579-1580-1581-1582-1583-1584-1585-1586-1587-1588-1589-1590-1591-1592-1593-1594-1595-1596-1597-1598-1599-1600-1601-1602-1603-1604-1605-1606-1607-1608-1609-1610-1611-1612-1613-1614-1615-1616-1617-1618-1619-1620-1621-1622-1623-1624-1625-1626-1627-1628-1629-1630-1631-1632-1633-1634-1635-1636-1637-1638-1639-1640-1641-1642-1643-1644-1645-1646-1647-1648-1649-1650-1651-1652-1653-1654-1655-1656-1657-1658-1659-1660-1661-1662-1663-1664-1665-1666-1667-1668-1669-1670-1671-1672-1673-1674-1675-1676-1677-1678-1679-1680-1681-1682-1683-1684-1685-1686-1687-1688-1689-1690-1691-1692-1693-1694-1695-1696-1697-1698-1699-1700-1701-1702-1703-1704-1705-1706-1707-1708-1709-1710-1711-1712-1713-1714-1715-1716-1717-1718-1719-1720-1721-1722-1723-1724-1725-1726-1727-1728-1729-1730-1731-1732-1733-1734-1735-1736-1737-1738-1739-1740-1741-1742-1743-1744-1745-1746-1747-1748-1749-1750-1751-1752-1753-1754-1755-1756-1757-1758-1759-1760-1761-1762-1763-1764-1765-1766-1767-1768-1769-1770-1771-1772-1773-1774-1775-1776-1777-1778-1779-1780-1781-1782-1783-1784-1785-1786-1787-1788-1789-1790-1791-1792-1793-1794-1795-1796-1797-1798-1799-1800-1801-1802-1803-1804-1805-1806-1807-1808-1809-1810-1811-1812-1813-1814-1815-1816-1817-1818-1819-1820-1821-1822-1823-1824-1825-1826-1827-1828-1829-1830-1831-1832-1833-1834-1835-1836-1837-1838-1839-1840-1841-1842-1843-1844-1845-1846-1847-1848-1849-1850-1851-1852-1853-1854-1855-1856-1857-1858-1859-1860-1861-1862-1863-1864-1865-1866-1867-1868-1869-1870-1871-1872-1873-1874-1875-1876-1877-1878-1879-1880-1881-1882-1883-1884-1885-1886-1887-1888-1889-1890-1891-1892-1893-1894-1895-1896-1897-1898-1899-1900-1901-1902-1903-1904-1905-1906-1907-1908-1909-1910-1911-1912-1913-1914-1915-1916-1917-1918-1919-1920-1921-1922-1923-1924-1925-1926-1927-1928-1929-1930-1931-1932-1933-1934-1935-1936-1937-1938-1939-1940-1941-1942-1943-1944-1945-1946-1947-1948-1949-1950-1951-1952-1953-1954-1955-1956-1957-1958-1959-1960-1961-1962-1963-1964-1965-1966-1967-1968-1969-1970-1971-1972-1973-1974-1975-1976-1977-1978-1979-1980-1981-1982-1983-1984-1985-1986-1987-1988-1989-1990-1991-1992-1993-1994-1995-1996-1997-1998-1999-2000-2001-2002-2003-2004-2005-2006-2007-2008-2009-2010-2011-2012-2013-2014-2015-2016-2017-2018-2019-2020-2021-2022-2023-2024-2025-2026-2027-2028-2029-2030-2031-2032-2033-2034-2035-2036-2037-2038-2039-2040-2041-2042-2043-2044-2045-2046-2047-2048-2049-2050-2051-2052-2053-2054-2055-2056-2057-2058-2059-2060-2061-2062-2063-2064-2065-2066-2067-2068-2069-2070-2071-2072-2073-2074-2075-2076-2077-2078-2079-2080-2081-2082-2083-2084-2085-2086-2087-2088-2089-2090-2091-2092-2093-2094-2095-2096-2097-2098-2099-2100-2101-2102-2103-2104-2105-2106-2107-2108-2109-2110-2111-2112-2113-2114-2115-2116-2117-2118-2119-2120-2121-2122-2123-2124-2125-2126-2127-2128-2129-2130-2131-2132-2133-2134-2135-2136-2137-2138-2139-2140-2141-2142-2143-2144-2145-2146-2147-2148-2149-2150-2151-2152-2153-2154-2155-2156-2157-2158-2159-2160-2161-2162-2163-2164-2165-2166-2167-2168-2169-2170-2171-2172-2173-2174-2175-2176-2177-2178-2179-2180-2181-2182-2183-2184-2185-2186-2187-2188-2189-2190-2191-2192-2193-2194-2195-2196-2197-2198-2199-2200-2201-2202-2203-2204-2205-2206-2207-2208-2209-2210-2211-2212-2213-2214-2215-2216-2217-2218-2219-2220-2221-2222-2223-2224-2225-2226-2227-2228-2229-2230-2231-2232-2233-2234-2235-2236-2237-2238-2239-2240-2241-2242-2243-2244-2245-2246-2247-2248-2249-2250-2251-2252-2253-2254-2255-2256-2257-2258-2259-2260-2261-2262-2263-2264-2265-2266-2267-2268-2269-2270-2271-2272-2273-2274-2275-2276-2277-2278-2279-2280-2281-2282-2283-2284-2285-2286-2287-2288-2289-2290-2291-2292-2293-2294-2295-2296-2297-2298-2299-2300-2301-2302-2303-2304-2305-2306-2307-2308-2309-2310-2311-2312-2313-2314-2315-2316-2317-2318-2319-2320-2321-2322-2323-2324-2325-2326-2327-2328-2329-2330-2331-2332-2333-2334-2335-2336-2337-2338-2339-2340-2341-2342-2343-2344-2345-2346-2347-2348-2349-2350-2351-2352-2353-2354-2355-2356-2357-2358-2359-2360-2361-2362-2363-2364-2365-2366-2367-2368-2369-2370-2371-2372-2373-2374-2375-2376-2377-2378-2379-2380-2381-2382-2383-2384-2385-2386-2387-2388-2389-2390-2391-2392-2393-2394-2395-2396-2397-2398-2399-2400-2401-2402-2403-2404-2405-2406-2407-2408-2409-2410-2411-2412-2413-2414-2415-2416-2417-2418-2419-2420-2421-2422-2423-2424-2425-2426-2427-2428-2429-2430-2431-2432-2433-2434-2435-2436-2437-2438-2439-2440-2441-2442-2443-2444-2445-2446-2447-2448-2449-2450-2451-2452-2453-2454-2455-2456-2457-2458-2459-2460-2461-2462-2463-2464-2465-2466-2467-2468-2469-2470-2471-2472-2473-2474-2475-2476-2477-2478-2479-2480-2481-2482-2483-2484-2485-2486-2487-2488-2489-2490-2491-2492-2493-2494-2495-2496-2497-2498-2499-2500-2501-2502-2503-2504-2505-2506-2507-2508-2509-2510-2511-2512-2513-2514-2515-2516-2517-2518-2519-2520-2521-2522-2523-2524-2525-2526-2527-2528-2529-2530-2531-2532-2533-2534-2535-2536-2537-2538-2539-2540-2541-2542-2543-2544-2545-2546-2547-2548-2549-2550-2551-2552-2553-2554-2555-2556-2557-2558-2559-2560-2561-2562-2563-2564-2565-2566-2567-2568-2569-2570-2571-2572-2573-2574-2575-2576-2577-2578-2579-2580-2581-2582-2583-2584-2585-2586-2587-2588-2589-2590-2591-2592-2593-2594-2595-2596-2597-2598-2599-2600-2601-2602-2603-2604-2605-2606-2607-2608-2609-2610-2611-2612-2613-2614-2615-2616-2617-2618-2619-2620-2621-2622-2623-2624-2625-2626-2627-2628-2629-2630-2631-2632-2633-2634-2635-2636-2637-2638-2639-2640-2641-2642-2643-2644-2645-2646-2647-2648-2649-2650-2651-2652-2653-2654-2655-2656-2657-2658-2659-2660-2661-2662-2663-2664-2665-2666-2667-2668-2669-2670-2671-2672-2673-2674-2675-2676-2677-2678-2679-2680-2681-2682-2683-2684-2685-2686-2687-2688-2689-2690-2691-2692-2693-2694-26

## Problèmes posés par l'utilisation du réseau

Liens entre des mots non significatifs masquant la réitération de mots significatifs absents du réseau



**Vin jaune:** regroupement non souhaité des paragraphes 1 et 2, 6 et 7, 8 et 9, 9 et 10, 10 et 11

51

## Vers une combinaison des méthodes

Appliquer la méthode fondée sur la récurrence des mots jusqu'à la phase de pondération des termes

**Si**  $x$  % des termes jugés significatifs (termes n'étant pas présents dans tous les paragraphes ; pondération *tf.idf*) ne sont pas dans le réseau de cooccurrences

**alors**

continuer avec la méthode fondée sur la récurrence des mots

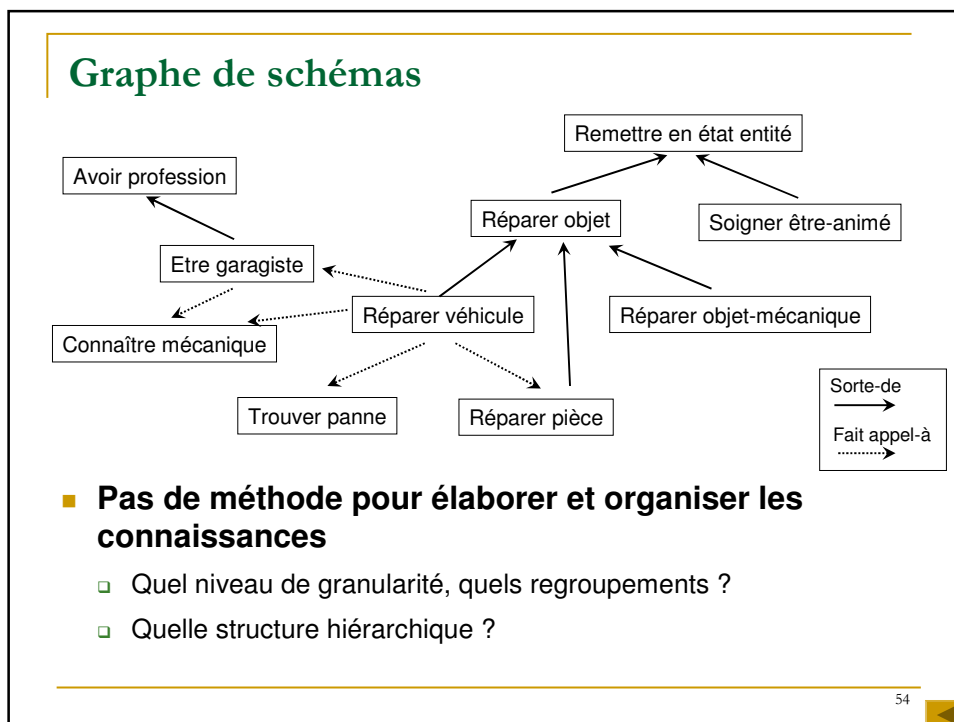
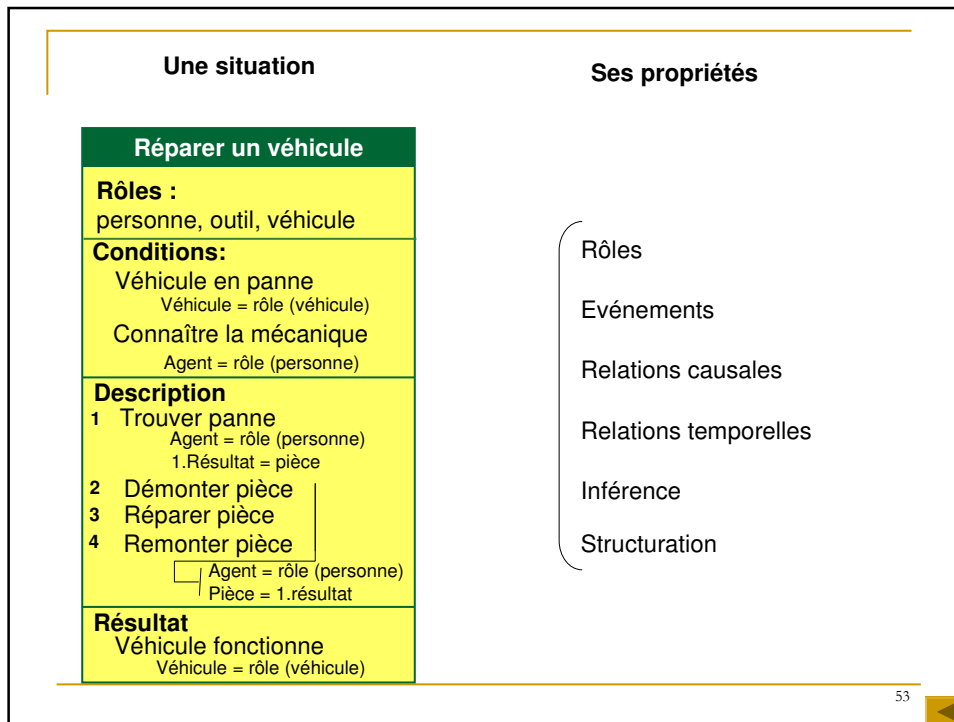
**Sinon**

appliquer la méthode utilisant le réseau de cooccurrences

**Fin\_si**

estimation de  $x$  : 25 %

52



**Thème général : vin jaune**

[illegible]

**Thème général : vin jaune**

[illegible]

28

## Exemple d'exploration d'un texte

### Thème général : vin jaune

<p>... les vins de Jura, dont le vin jaune, sont des vins de France. Ils sont produits dans la région de Jura, en France, et sont connus pour leur qualité et leur caractère unique. Le vin jaune est un vin blanc, produit à partir de raisins de la variété Chardonnay, qui sont cultivés dans la région de Jura. Le vin jaune est un vin blanc, produit à partir de raisins de la variété Chardonnay, qui sont cultivés dans la région de Jura. Le vin jaune est un vin blanc, produit à partir de raisins de la variété Chardonnay, qui sont cultivés dans la région de Jura.</p> <p>... les vins de Jura, dont le vin jaune, sont des vins de France. Ils sont produits dans la région de Jura, en France, et sont connus pour leur qualité et leur caractère unique. Le vin jaune est un vin blanc, produit à partir de raisins de la variété Chardonnay, qui sont cultivés dans la région de Jura. Le vin jaune est un vin blanc, produit à partir de raisins de la variété Chardonnay, qui sont cultivés dans la région de Jura. Le vin jaune est un vin blanc, produit à partir de raisins de la variété Chardonnay, qui sont cultivés dans la région de Jura.</p> <p>... les vins de Jura, dont le vin jaune, sont des vins de France. Ils sont produits dans la région de Jura, en France, et sont connus pour leur qualité et leur caractère unique. Le vin jaune est un vin blanc, produit à partir de raisins de la variété Chardonnay, qui sont cultivés dans la région de Jura. Le vin jaune est un vin blanc, produit à partir de raisins de la variété Chardonnay, qui sont cultivés dans la région de Jura. Le vin jaune est un vin blanc, produit à partir de raisins de la variété Chardonnay, qui sont cultivés dans la région de Jura.</p> <p>... les vins de Jura, dont le vin jaune, sont des vins de France. Ils sont produits dans la région de Jura, en France, et sont connus pour leur qualité et leur caractère unique. Le vin jaune est un vin blanc, produit à partir de raisins de la variété Chardonnay, qui sont cultivés dans la région de Jura. Le vin jaune est un vin blanc, produit à partir de raisins de la variété Chardonnay, qui sont cultivés dans la région de Jura. Le vin jaune est un vin blanc, produit à partir de raisins de la variété Chardonnay, qui sont cultivés dans la région de Jura.</p> <p>... les vins de Jura, dont le vin jaune, sont des vins de France. Ils sont produits dans la région de Jura, en France, et sont connus pour leur qualité et leur caractère unique. Le vin jaune est un vin blanc, produit à partir de raisins de la variété Chardonnay, qui sont cultivés dans la région de Jura. Le vin jaune est un vin blanc, produit à partir de raisins de la variété Chardonnay, qui sont cultivés dans la région de Jura. Le vin jaune est un vin blanc, produit à partir de raisins de la variété Chardonnay, qui sont cultivés dans la région de Jura.</p>	<p><b>Thème global :</b></p> <p>vin</p> <p><b>Thème local :</b></p> <p>mélange, composé</p> <p><b>Meta-descripteurs :</b></p> <p>analyse, technique</p>	<p>&lt;seg&gt;</p> <p>L'analyse la plus directe d'extraits de vins est la chromatographie : on injecte un échantillon dans un solvant que l'on vaporise et on fait traverser au mélange une colonne revêtue intérieurement d'un polymère, qui retient les divers composés du mélange à des degrés divers ; en bas de la colonne, on détecte la sortie des composés séparés. Le premier travail des chimistes fut la mise au point d'une variante de cette technique pour identifier les composés présents en quantités minimales dans des mélanges complexes.</p> <p>&lt;/seg&gt;</p>
	<p><b>Thème global :</b></p> <p>vin</p> <p><b>Thèmes local :</b></p> <p>goût, noisette</p> <p><b>Méta-descripteurs :</b></p> <p>caractéristique</p>	