Séminaire INALCO "Corpus parallèles", 17 mai 2010

ALIGNEMENT AUTOMATIQUE DES PROPOSITIONS FRANÇAIS-JAPONAIS

YAYOI NAKAMURA-DELLOYE

Équipe ALPAGE, INRIA-Rocquencour <u>yayoi@yayoi.fr</u>

Introduction

INTRODUCTION

TEXTES PARALLÈLES ET ALIGNEMENT 2/2

Corpus alignés (bi- ou multitexte)

Différentes applications (e.g. mémoires de traduction, dictionnaires)

- Corpus alignés anglais-japonais disponibles sur le site du NICT (National Institute of Information and Communications Technology)
- Corpus aligné français-japonais dans le package d'évaluation ARCADE II
- Systèmes d'alignement
 - ✓ disponibles sur Internet (ex. GIZA++)
 - ✓ Système d'alignement phrastique Fr-Jp : AlALeR

(Nakamura-Delloye, 2005)

INTRODUCTION

TEXTES PARALLÈLES ET ALIGNEMENT 1/2

- ** Textes parallèles = ensemble de textes de langues différentes, constitué d'un texte original et de ses traductions
 - Corpus compilés anglais-japonais de taille importante
 - Textes parallèles français-japonais
 - Le Monde Diplomatique
 - Manuels de logiciels libres
 - Textes G7, G8, etc.
 - Textes littéraires (Aozora Bunko : http://www.aozora.gr.jp/)

2

Yayoi NAKAMURA - DELLOYE

Introduct

INTRODUCTION

POURQUOI LA PROPOSITION ? 1/2

- Nécessité d'une unité de traitement plus petite que la phrase
 - → Proposition = bon candidat
 - la documentation technique;
 - l'analyse discursive ;
 - ...

Introduction

INTRODUCTION

POURQUOI LA PROPOSITION ? 2/2

- Dans le cas de l'alignement : ex. mémoire de traduction
 - Par rapport à la **phrase** : réutilisabilité plus importante
 - Par rapport aux **unités inférieures** : portabilité plus élevée de correspondance
 - mot français « **compte** » = 10 définitions constituées de noms japonais différents non interchangeables
 - « tenir compte » = « 考慮する »

5

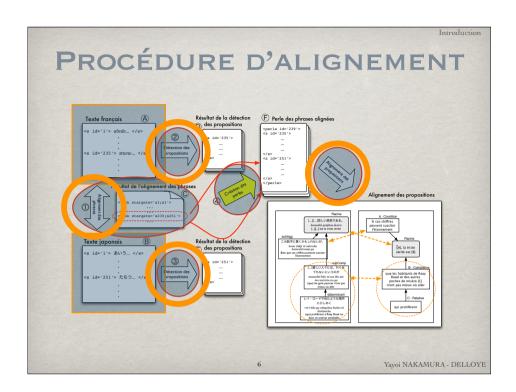
Yayoi NAKAMURA - DELLOYE

あられ 霰【arare】n.

1. Perle de glace. 2. Petit biscuit de riz. 3. inform. AlALER (système d'Alignement Autonome, Léger et Robuste) Aligneur adapté au traitement du japonais caractérisé par l'absence d'utilisation d'analyseur morphologique et de dictionnaire.

ALALER

SYSTÈME D'ALIGNEMENT DES PHRASES

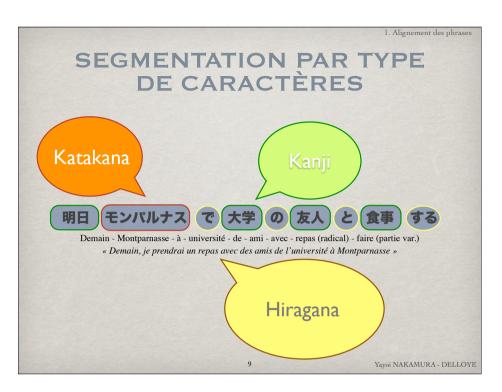


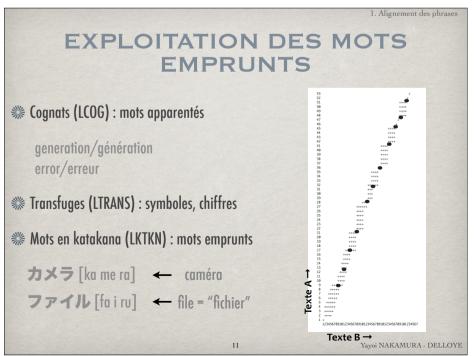
1. Alignement des phras

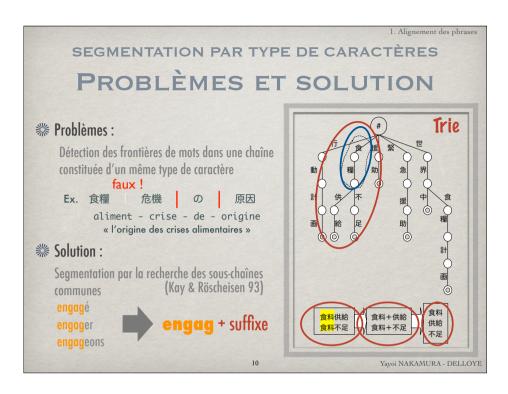
MÉTHODES D'ALIGNEMENT DES PHRASES

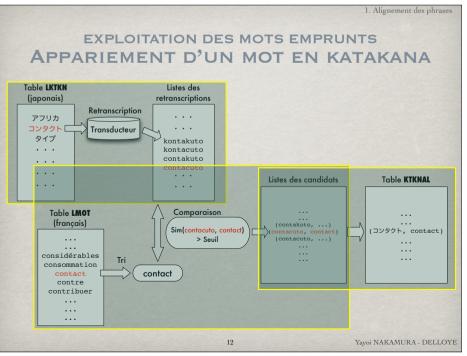
- * Premières méthodes dans le cadre de travaux sur la TA
 - Méthodes basées sur la distribution lexicale (Kay & Röscheisen 93)
 - Méthodes basées sur la corrélation des longueurs (Brown et al. 91, Gale & Church 93)
 - Amélioration par l'introduction de la notion de cognat (Simard et al. 92, Langlais 97, Kraif 01)
- Problèmes pour la conception d'un système autonome capable de traiter le japonais :
 - Segmentation sans analyseur morphologique
 - → Amélioration de la méthode de segmentation par type de caractères
 - Détermination d'ancrages sûrs sans dictionnaire bilingue
 - **→ Exploitation des mots emprunts**

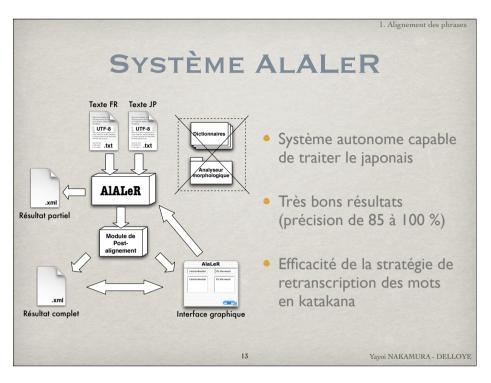
8

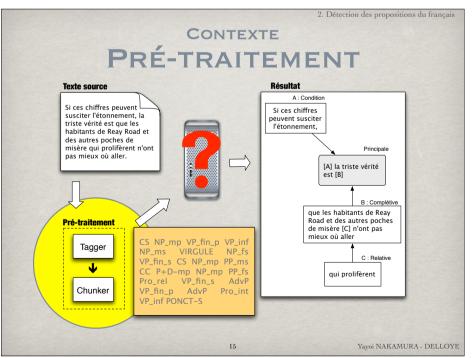












しぐれ 時雨【ʃigɯre】n.

1. Brève averse. **2.** INFORM. SIGLÉ (Système d'Identification de propositions avec Grammaire Légère) système réalisant la détection des propositions françaises caractérisé par l'utilisation d'une grammaire hors contexte écrite dans un formalisme DCG et par une implémentation en langage PROLOG.

SIGLÉ FR

Système d'identification des propositions du français

PLAN

Études linguistiques sur la proposition du français

Réalisation

Évaluation

Perspectives

Études linguistiques

- Typologie des
- Classement des subordonnées

Réalisation

Évaluation

Perspectives

ÉTUDES LINGUISTIQUES

₩ But:

définir une grammaire permettant de détecter des propositions à partir d'un résultat de chunker

- Typologie des propositions, basée sur des critères formels
- Classement des subordonnées selon la position d'apparition

17

Yayoi NAKAMURA - DELLOYE

2. Détection des propositions du français

Études linguistiques

- Typologie des propositions
- Classement des subordonnées

Réalisation

Évaluation

Perspectives

NOTION PRINCIPALE: PROPOSITION

Proposition = sujet + prédicat

Subordonnée

- · Il était déjà rentré quand je suis arrivé.
- · Je pense qu'il viendra.
- Je me demande s'il est parti.
- · La peinture qui m'a fascinée.
- La déception du père quand il a entendu cette nouvelle.
- · Je voterai pour qui me promettra moins d'impôts.
- Où il y a de la gène, il n'y a pas de plaisir.
- Que le gouvernement propose une nouvelle loi, l'opposition crie au scandale.
- Il n'a pas pu lire cette lettre comme sa mère l'a deviné.
- Tu peux poser ton manteau où tu veux.
- Je pars, que cela vous plaise ou non.
- Le crocodile n'eut pas le temps de se demander ce que lui voulait ce lourdaud, que Gropopotin s'était déjà assis sur son dos.
- La maison est restée aussi conviviale qu'elle l'était avant.

Yayoi NAKAMURA - DELLOYE

Détection des propositions du français

- Typologie de propositions
- Classement des subordonnées

Réalisation

Évaluation

Perspectives

NOTION PRINCIPALE:

PROPOSITION

Proposition = suiet + prédicat

- Types de proposition :
- 1. Racine (Principale)
- 2. Coordonnée

Mon père est professeur et ma mère travaille à la banque.

3. Subordonnée

Il était déjà rentré quand je suis arrivé.

4. Incidente

Il s'en est, me semble-t-il, bien sorti.

+ Éléments extra-prédicatifs (Charolles 1997) (Combettes 1998)

L'autre jour, ... Cette affaire étant réglée, ... En ce qui concerne X, ...

18

Yayoi NAKAMURA - DELLOYE

Études

- Typologie des propositions
- Classement des subordonnées

Réalisation

Evaluation

Perspectives

PROBLÈMES DES TYPOLOGIES USUELLES

- Définitions des subordonnées souvent selon la nature du connecteur qui les introduit
 - ... MAIS
 - → Étiquetage automatique très difficile;

qui	pronoms relatif/interrogatif
que / qu'	pronoms relatif/interrogatif, adverbe, adverbe exclamatif, conjonction de subordination
quand	adverbe interrogatif, conjonction de subordination
comme	adverbe exclamatif, conjonction de subordination, conjonction de coordination, préposition
si/s'	adverbe, conjonction de subordination, affirmation, clitique

20

2. Détection des propositions du françai

Études linguistiques

- Typologie des propositions
- Classement des

Réalisation

Évaluation

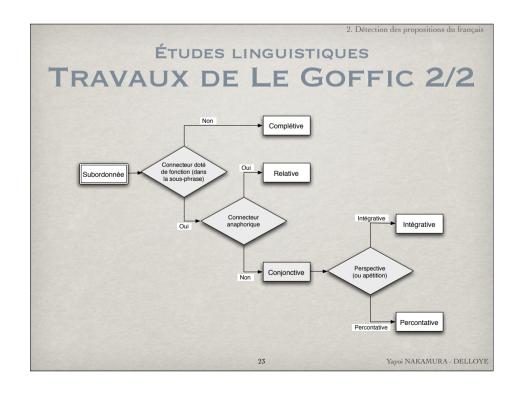
Perspectives

PROBLÈMES DES TYPOLOGIES USUELLES

- Comportements syntaxiques différents de termes appartenant à la même catégorie
- 1. Le gouvernement a retiré sa proposition dont <u>la conformité à</u> la Constitution avait été remise en cause.
- 2. À cette occasion, se sont réunis huit représentants dont <u>notre</u> <u>Président</u>.
 - → Qu'est-ce qu'une locution conjonctive (LC) ? Les LCs constituent-elles une liste fermée ?

2

Yayoi NAKAMURA - DELLOYE



2. Détection des propositions du français

ÉTUDES LINGUISTIQUES

TRAVAUX DE LE GOFFIC 1/2

(Le Goffic 1992, 1993, 2002)

** Termes en «qu-» = seuls connecteurs du français

Une vieille famille indo-européenne en «Kw-»:

- pronoms : qui, que, quoi, lequel ;
- adjectif: quel;
- adverbes: où, quand, comme, comment, combien, que, dont, pourquoi.
- Subordonnées introduites par une LC = des GAdv ou GPrép comprenant une subordonnée constituée par un connecteur en «qu-»
- → un traitement unifié et homogène des types de subordonnées

22

Yayoi NAKAMURA - DELLOYE

ÉTUDES LINGUISTIQUES

TRAVAUX DE LE GOFFIC 2/2

Complétive : complétive

Je crois qu'il va pleuvoir ...

Relative : relative avec antécédent

Le médecin qui est venu / la maison où je suis né...

- Intégrative :
 - Pronominale : relative sans antécédent

Qui dort dîne / embrassez qui vous voulez...

• Adverbiale : circonstancielle en qu- ou si

Quand on veut, on peut / si vous avez fini, vous pouvez sortir / il est à peine sorti qu'il a commencé à pleuvoir...

Percontative : interrogative/exclamative indirecte

Je sais qui a gagné / où il est allé / comment il l'a fait...

Le classement de Le Goffic (comme beaucoup d'autres classements proposés) présuppose une analyse correcte de la nature du connecteur, très difficile à réaliser de manière automatique

ÉTUDES LINGUISTIQUES

AUTRES TYPOLOGIES

* Typologie selon la catégorie du mot simple

équivalent : (Le Bon Usage, 11ème éd.) (Biskri et Desclés 2005)

- Substantive : Je pense qu'il viendra / Que tu m'aimes me réjouit
- · Adjective : La femme que tu vois / la ville où i'babite
- Adverbiale : Il était déjà rentré quand je suis arrivé
- *Typologie selon la fonction de la subordonnée dans la principale : (Chevalier et al. 1964)(Grevisse 1969)(Wilmet 1997)
 - Sujet : Que je sois malade ne l'a jamais effleuré
 - · Attribut : La triste vérité est qu'il est fou
 - Objet : Marie sait que Paul viendra
 - Circonstancielle: Il était déjà rentré quand je suis arrivé
 - Complément de nom : la certitude que son but était atteint
 - ... etc.

Yavoi NAKAMURA - DELLOYE

Postverbale:

Initiale/ Finale:

Autres

SN:

POSITION POST-VERBALE (SUBORDONNÉE COMPLÉMENT : SUDO)

Subordonnées substantives :

- Complétives
 - Je pense qu'il viendra.
- * Intégratives pronominales (relatives sans ant.)
 - Il a le droit d'embrasser qui il veut.
- Percontatives (ou interrogatives)
 - Je me demande s'il est parti.
 - Il ne m'a pas dit quand il rentrerait.
 - Voyez comme c'est facile.

ÉTUDES LINGUISTIQUES

TYPOLOGIE DES SUBORDONNÉES SELON LA POSITION

- ***** Typologie selon la catégorie
 - · Substantive, adverbiale, adjective
- ***** Typologie selon la position
 - Initiale/Finale:
 - Il était déjà rentré quand je suis arrivé
 - Post-verbale:
 - Je pense qu'il viendra
 - Autres positions SN:
 - Que tu m'aimes me réjouit
 - Post-nominale:
 - La femme que tu vois / la ville où j'habite
 - Post-adj. et -adv. :
 - De même que ... / bien que ... / aujourd'hui que ...
 - → Description systématique de chaque type (Le Goffic, 2000)

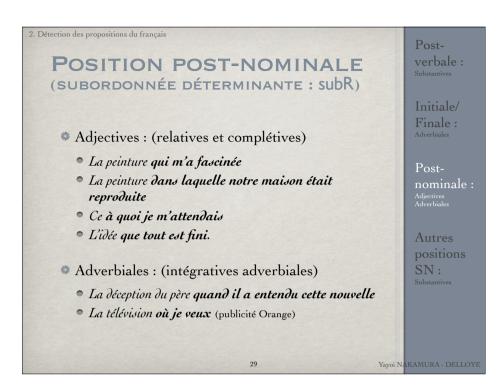
Yavoi NAKAMURA - DELLOYE

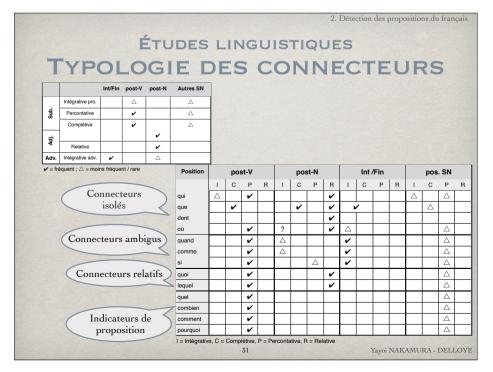
POSITION INITIALE/FINALE (SUBORDONNÉE CIRCONSTANCIELLE: SUBP)

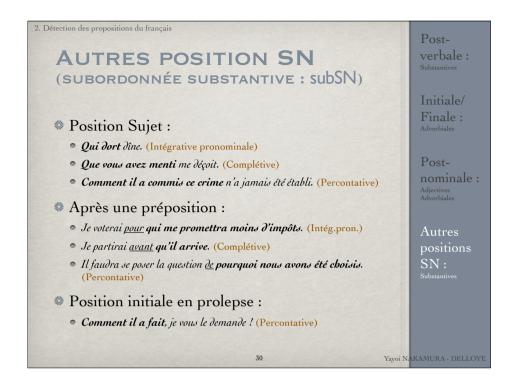
- Quand je suis arrivé, il était déjà rentré.
- Si tu ne manges pas, tu ne guériras pas.
- Comme elle est écrite en chinois, il n'a pas pu lire cette lettre.
- Où il y a de la gène, il n'y a pas de plaisir.
- Que le gouvernement propose une nouvelle loi, l'opposition crie au scandale.
- Il n'a pas pu lire cette lettre comme sa mère l'a deviné.
- Tu peux poser ton manteau où tu veux.
- Je pars, que cela vous plaise ou non.
- Le crocodile n'eut pas le temps de se demander ce que lui voulait ce lourdaud, que Gropopotin s'était déjà assis sur son dos.
- La maison est restée aussi conviviale qu'elle l'était avant.

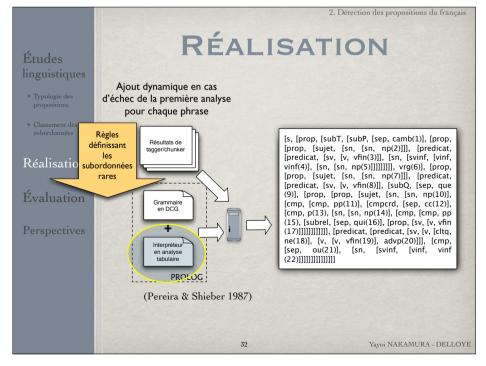
Finale:

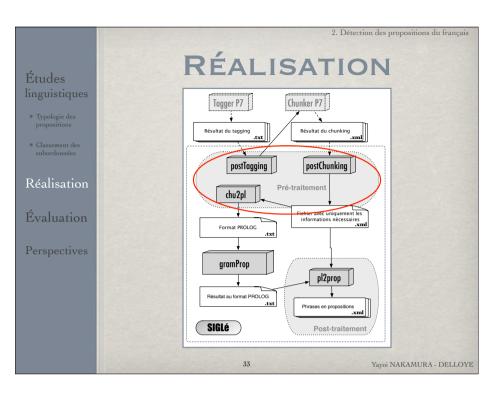
Autres positions SN:

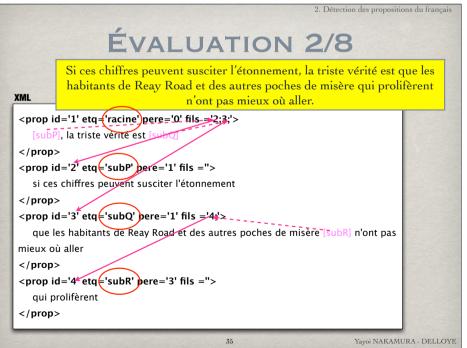




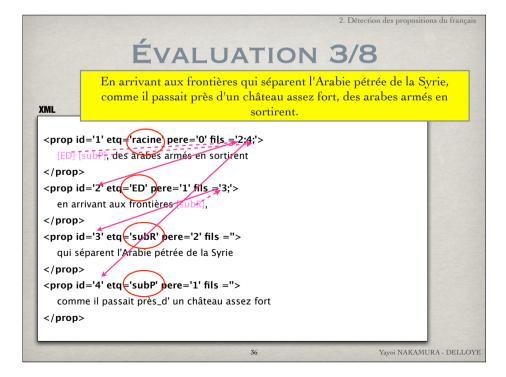








ÉVALUATION 1/8 Études Préc. 2 Nb de phr Rappel Préc.1 Préc. T 96,2 G8 53 98.0 100,0 98.0 274 Unicode 81.4 96.2 97.8 94,1 Zadig 1206 88.6 92.8 95,3 88.4 1713 LMD* 84.9 89.2 98,0 87,4 *Le Monde Diplomatique Analyse linéaire Analyse structurale A: Nombre total de phrases B: Nombre de phrases dont l'analyse a abouti Rappel (%) = $\frac{B}{\Delta}$ × 100 C : Nombre de phrases dont les frontières de Perspectives propositions sont correctement détectées Précision 1 (%) = $\frac{C}{R}$ × 100 D : Nombre de phrases dont les relations des propositions sont correctement analysées Précision 2 (%) = $\frac{D}{C}$ × 100 Précision T (%) = $\frac{\text{Préc. 1} \times \text{Préc.2}}{100} = \frac{D}{B} \times 100$ Yavoi NAKAMURA - DELLOYE



ÉVALUATION 4/8

Tout ce qui passe sur mes terres est à moi, dit -il, aussi bien que ce que je trouve sur les terres des autres

XML

```
prop id='1' etq='racine')pere='0' fils=='2:3;4:'>
  tout ce [subR] est à moi [lact aussi bien que ce [subR]
```

</prop>

qui passe sur mes terres

</prop>

prop id='3*eta='Inc' pere='1' fils ="> . dit -il.

</prop>

que je trouve sur les terres des autres

</prop>

Yavoi NAKAMURA - DELLOYE

ÉVALUATION 5/8

- * Erreurs de l'analyse linéaire :
 - ambiguïté de la position post-prépositionnelle intégrative pronominale (subSN) ⇔ subordonnée déterminante (subR)
- 1. [...] le pouvoir est seulement entre les mains (1) 🗢 🖰 subR |de|qui détient des armes à feu (2) subSN |, de | qui possède les richesses (3)

À comparer :

il admet aussi un Être supérieur (1) I, de qui la forme et la matière dépendent (2) © subR

ÉVALUATION 5/8

- * Erreurs de l'analyse linéaire :
 - étiquettes erronées attribuées par le pré-traitement ;
 - mauvaise interprétation de connecteurs : introducteur de proposition ⇔ intr. de syntagme (structure réduite sans verbe)

Et dire qu'au moment de son apogée, dans les années 1950, Cockerill employait encore plus de 25 000 personnes, que la ville de Seraing était toujours noire de fumée, de bruit, de monde, de travail.

Yavoi NAKAMURA - DELLOYE

ÉVALUATION 6/8

- * Trois types d'erreurs de l'analyse structurale
 - Coordination de subordonnées sans connecteur De son côté, Taikong Corp. explique (1) que la firme n'a pas encore le droit de les vendre en France (2) , mais peut les exposer (3)
 - Ambiguïté de la virgule précédant une subordonnée : Subordonnée coordonnée ? ⇔ subR simple ?

Personne ne m'a expliqué (1) | qu'il s'agissait de la première étape de l'expansion prétendument bienveillante d'une nation nouvelle (2) , mais que cette expansion signifiait en réalité l'expulsion violente des Indiens de la totalité du continent (3) 1, qu'elle serait jalonnée d'atrocités indicibles (4) à l'issue desquelles on parquerait les survivants dans des réserves (5)

- Ambiguïté des positions
 - 1. C'est facile à dire (1) | quand on n'est pas concerné dans sa chair (2)
 - 2. Le pontife trouva dans son coeur (1) que cela valait beaucoup (2)

Relations ambiguës

Paris avait estimé, à l'époque, (1)

|| qu'une référence aux valeurs religieuses n'était pas acceptable (2)

|| car elle soulevait des problèmes politiques et constitutionnels en France. (5)

(Paris avait estimé X) + (car Y)?

Paris avait estimé (une référence aux valeurs religieuses n'était pas acceptable + car Y)?

41

Yayoi NAKAMURA - DELLOYE

2. Détection des propositions du français

Études linguistiques

- Typologie des propositions
- Classement des subordonnées

Réalisation

Évaluation

Perspectives

PISTES D'AMÉLIORATION

- Amélioration du pré-traitement
 - → amélioration des modules de pré-traitement
 - → utilisation d'autres systèmes : système de segmentation en super-chunks (Blanc et al., 2007)
- Introduction de plus d'informations
 - risque de multiplication des calculs
- Affinement des étiquettes attribution d'étiquettes syntactico-sémantiques

ÉVALUATION 8/8

* Fréquence des subordonnées

occurrence(%)

		Int/Fin	post-V	post-N	Autres SN	
	Intégrative pro.		△ 0	0	△ 0,4	0,
Sub.	Percontative		✓ 2	4	△ 0	0
	Complétive		✓ 20	27	△ 0,2	0
Adj.				✓ 2	0,3	
Ā	Relative		V	✓ 60	57	
Adv.	Intégrative adv.	✓ 15	11	△ 0	0	

✓ = fréquent ; △ = moins fréquent / rare

42

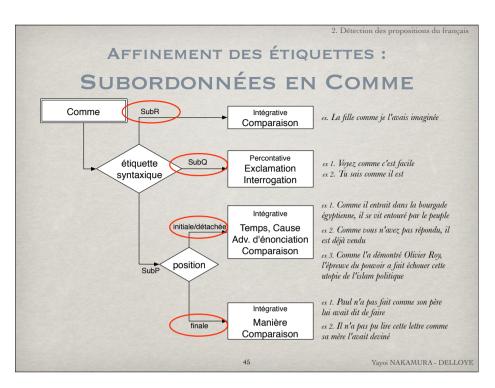
Yayoi NAKAMURA - DELLOYE

2. Détection des propositions du français

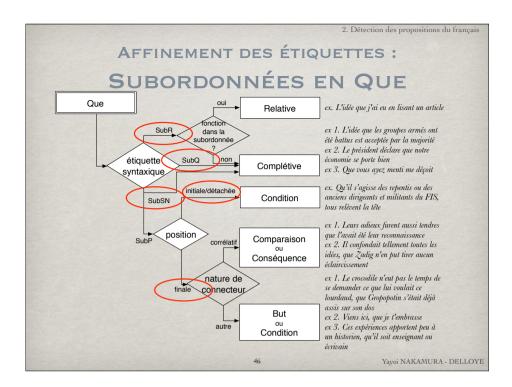
PISTES D'AMÉLIORATION AFFINEMENT DES ÉTIQUETTES

Position	sition post-V				post-N			Int /Fin				pos. SN				
	1	С	Р	R	1	С	Р	R	1	С	Р	R	1	С	Р	R
qui	Δ		~					~					Δ		Δ	
que		~				~		~	١	/				Δ		
dont								~								
où			~		?			~	Δ						Δ	
quand			~		Δ				~						Δ	
comme			~		Δ				~						Δ	
si			~				Δ		~						Δ	
quoi			~					~							Δ	
lequel			~		İ			~							Δ	
quel			~												Δ	
combien			V		İ				İ				İ		Δ	
comment			~		Ì				Ì				Ì		Δ	
pourquoi			V		İ				İ				İ		Δ	

I = Intégrative, C = Complétive, P = Percontative, R = Relativ







PLAN

- Études linguistiques sur la phrase et la proposition du japonais
- * Réalisation
- # Évaluation

Études linguistiques

Structure de la

 Classement des propositions

Réalisation

Évaluation

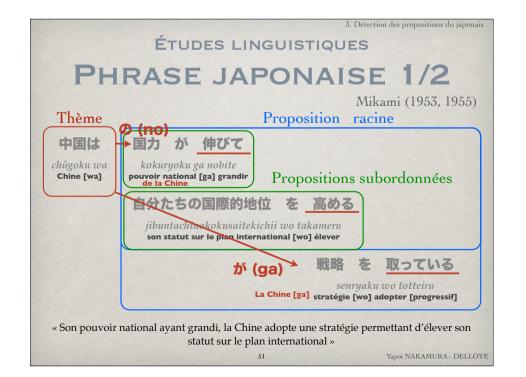
ÉTUDES LINGUISTIQUES

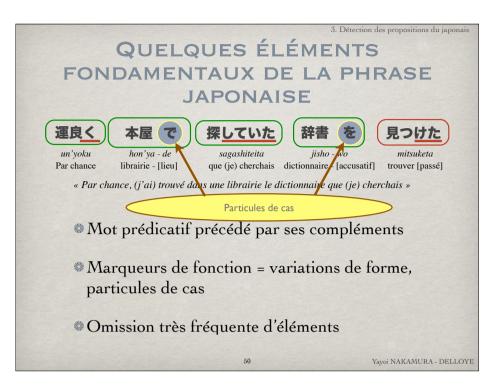
₩ But:

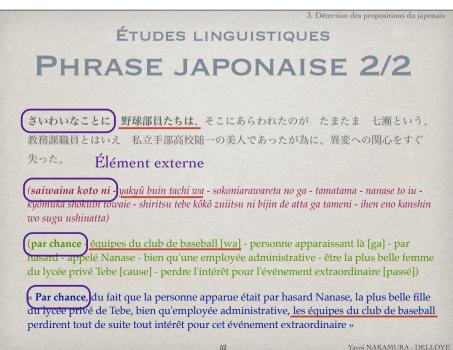
définir les propositions à l'aide uniquement des critères formels pour pouvoir les identifier automatiquement

- 1. Structure de la phrase japonaise, basée sur l'opposition thème-rhème
- 2. Classement des propositions japonaises

49







Études linguistiques

phrase

Réalisation

Évaluation

TOU OCIT DEC

TYPOLOGIE DES SUBORDONNÉES

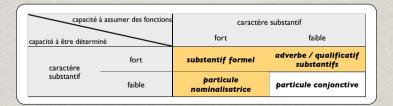
Classement selon uniquement des critères morphosyntaxiques (Teramura 82)

- Subordonnées sans connecteur
 - 1. Subordonnée neutre
 - 2. Subordonnée de condition
- 3. Subordonnée déterminante sans connecteur
 - Subordonnées avec connecteur
- 1. Subordonnée avec particule conjonctive
- 2. Subordonnée avec connecteur agglutinant
 - 3. Subordonnée de citation
 - 4. Subordonnée déterminante avec connecteur

53

Yayoi NAKAMURA - DELLOYE

ÉTUDES LINGUISTIQUES PROBLÈMES LIÉS AUX CONNECTEURS



Mot agglutinant (Sakuma 1940)

ÉTUDES LINGUISTIQUES PROBLÈMES LIÉS AUX CONNECTEURS * Subordonnée déterminante sans connecteur 友達 Substantif nihon e - itta - tomodachi Japon [e] - aller/partir [passé] - ami « ami qui est parti au Japon » * Subordonnée avec particule conjonctive Particule conjonctive nihon e - itta - aa Japon [e] - aller/partir [passé] - [opposition] « bien que (je sois/tu sois/il soit...) parti au Japon » * Subordonnée avec connecteur agglutinant 日本へ 行った 時 Substantif formel nihon e - itta - toki Japon [e] - aller/partir [passé] - temps « quand (je suis/tu es/il est...) parti au Japon »

Études linguistiques

 Structure de l phrase

 Classement des propositions

Réalisation

Évaluation

RÉALISATION DU SYSTÈME

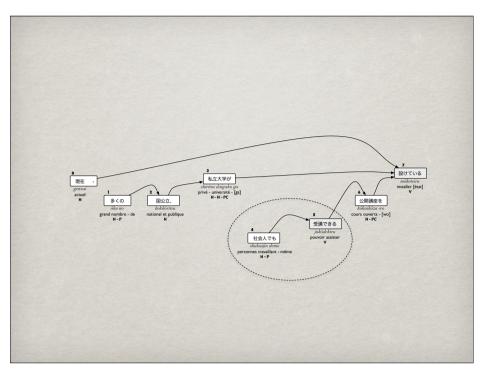
Méthode existante : CBAP (Maruyama et al. 2004) :

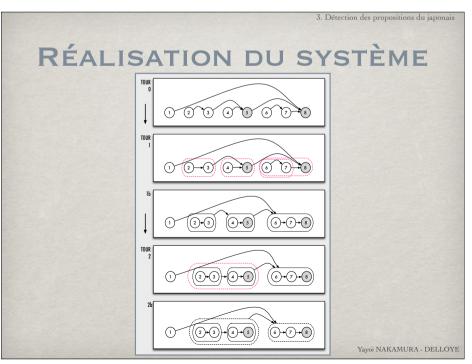
incapable de traiter les structures imbriquées

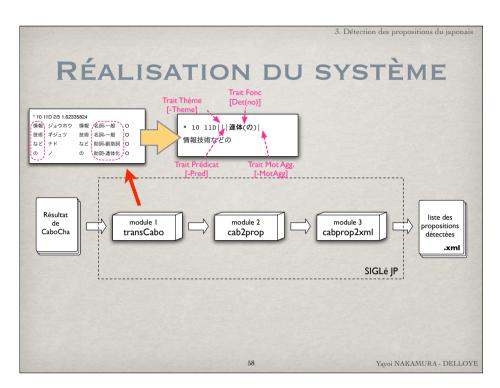
→ Utilisation d'un analyseur des relations de dépendance (Kudo & Matsumoto 2002)

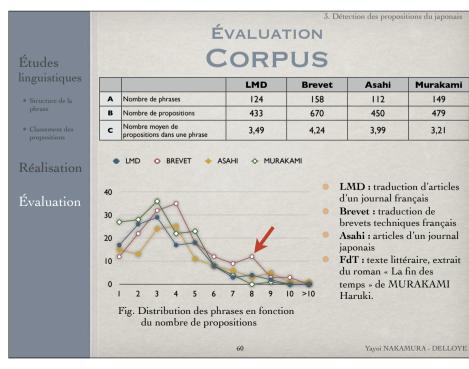
informations sur le chunk
informations sur chaque
unité constituant le chunk

Yavoi NAKAMURA - DELLOYE









3. Détection des propositions du japonai

ÉVALUATION RÉSULTATS 1/2 ÉVALUATION RÉSULTATS 2/2

		LMD	Brevet	Asahi	Murakami	<s id="62"> <prop id="1" pere="</th"></prop></s>
Α	Nombre de phrases	124	158	112	149	<txt1>「在目朝</txt1>
В	Nombre de propositions	433	670	450	479	<pre><pre></pre> <pre><pre>id='2' pere</pre></pre></pre>
С	Nombre moyen de propositions dans une phrase	3,49	4,24	3,99	3,21	<pre><pre><pre><txt1>03年に指</txt1></pre></pre></pre>
D	Nombre de propositions détectées	444	672	453	490	<pre><pre><pre><pre><pre><pre><txt1>科協系企業</txt1></pre></pre></pre></pre></pre></pre>
E	Nombre de propositions détectées correctement	389	589	391	426	<pre><pre></pre></pre>
F	Rappel (= E/B)	0,898	0,879	0,869	0,891	<txt1>[中立]北朝</txt1>
G	Précision (= E/D)	0,876	0,876	0,863	0,869	<pre><pre><pre>op id='5' pere=</pre></pre></pre>
н	Analyse linéaire = nombre de phrases correctement analysées (H/A %)	99 (80%)	119 (75%)	85 (76%)	120 (81%)	<pre><txt1>[連体]都内 <pre>cprop id='6' peres</pre></txt1></pre>
ı	Analyse structurale = nombre de phrases correctement analysées (I/H %)	94 (95%)	107 (90%)	79 (93%)	111 (93%)	<txt1>[吸節(一方</txt1>

科協系企業が関与して kakyôkei kigyô ga kan'yo shite le(s) filiale(s) de l'Association Kakyô y participe(nt)

Yayoi NAKAMURA - DELLOYE

Yavoi NAKAMURA - DELLOYE ÉVALUATION RÉSULTATS 2/2 [吸節(一方)+で][吸節(など)]技術・物資流出への 「在日朝鮮人科学技術者の親睦(しんぼく) 団体」とされる一方で 関与が指摘されてきた。 [sub.agg(ippô)+de] "zainichi chôsenjin kagaku gijutsusha no shinboku (shinboku) [sub.agg(ippô)+de] [sub.agg(nado)] gijutsu · busshi ryûshutsu eno dantai" to sareru ippô de kanyo ga shitekisarete kita Alors qu'elle est considérée comme un organisme amical on faisait remarquer son implication dans des fuites techniques et matérielles [...] [...] des scientifiques nord-coréens demeurant au Japon [sub.agg(nado)] [連体]都内メーカーのミサイル関連機器不正輸出事件では [吸節(こと)+が]判明するなど、 [déterminant] tonai mêkâ no misairu kanren kiki husei vushutsu iiken dew [sub.agg(koto)+ga] hanmeisuru nado comme par exemple, dans l'affaire d'exportation clandestine d'équipements pour missiles par un fabriquant de Tokyo [...], [...] a été découvert [déterminant] [sub.agg(koto)+ga] [中立]北朝鮮にも送っていたことが 0.3年に摘発された [neutre] kitachôsen nimo okutteita koto ga 03 nen ni tekihatsusareta le fait qu'(on les) envoyait aussi en Corée du Nord [nominatif] (qui) a été dénoncé en 2003

63

<s id='62'>

foid='1' pere='6' fils='' kakari='吸節(一方)+で'>

<p

Détection des propositions du japonai

Yavoi NAKAMURA - DELLOYE

ÉVALUATION RÉSULTATS 2/2

- * Présence d'erreurs dues aux mauvaises analyses fournies par le système de pré-traitement
- Résultat des extractions expérimentales des thèmes en wa et des éléments externes démontrant la nécessité d'une étude linguistique plus poussée

6

みぞれ 霙【midzore】n.

1. grésil, neige fondue. **2.** dessert en glaçon râpé au sirop. **3.** radis blanc râpé. **4.** inform. *MIZOLé* système réalisant l'alignement des propositions sur la base de l'approche spectrale de l'alignement des graphes ou de la méthode inspirée de la classification ascendante hiérarchique.

MIZOLÉ

SYSTÈME D'ALIGNEMENT DES PROPOSITIONS

Problèmes et éléments de solution

Trois méthodes

Évaluation

4. Alignement des propositions

PROBLÈMES

- Peu de travaux (Boutsis & Piperidis 1998) (Wong & Ren 2005)
- Impossibilité d'une simple application des méthodes classiques d'alignement des phrases pour l'alignement fr-jp due au nonparallélisme
 - → Alignement à l'aide des graphes

PLAN

Problèmes et éléments de solution

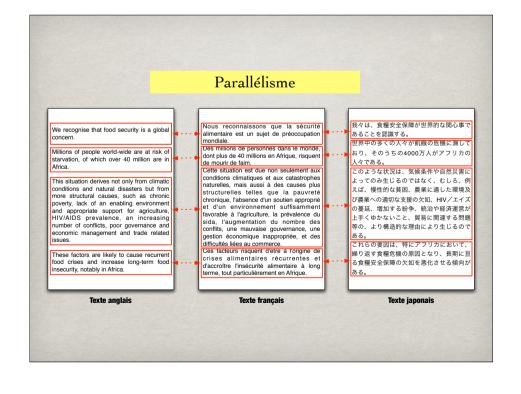
Trois méthodes réalisées

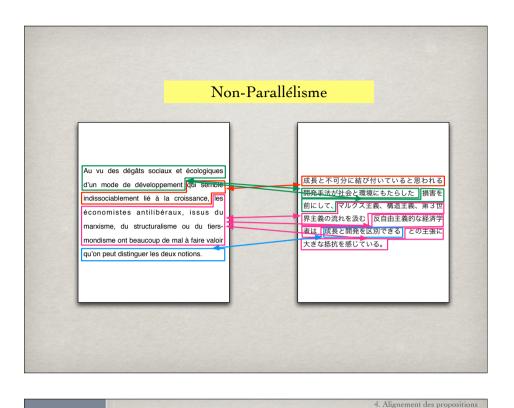
Évaluation

Constats sur les corpus utilisés

Résultats

Analyse des erreurs





ALGORITHME 1

MÉTHODES SPECTRALES

Appariements des graphes inexacts (Kosinov & Caelli 2002, 2004)

- Combinaison des avantages des techniques de décomposition spectrale, de projection et de classification

- Projeter les nœuds sur un sous-espace propre et regrouper les points projetés avec un algorithme de classification

ELÉMENTS DE SOLUTION → Alignement à l'aide des graphes Qui serricle Qui serricle Qui serricle Pastive Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave Restave

4. Alignement des proposition

Yavoi NAKAMURA - DELLOYE

ALGORITHME 1

MÉTHODES SPECTRALES 2/2

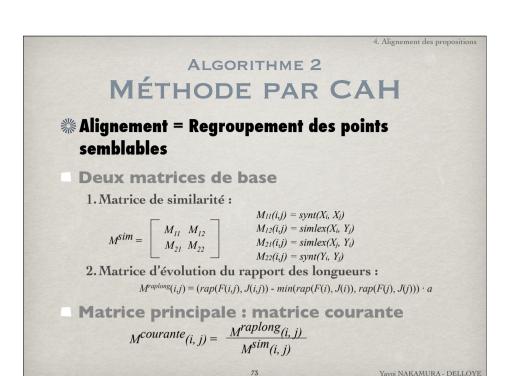
* Amélioration pour des graphes valués (Lerallut 2006)

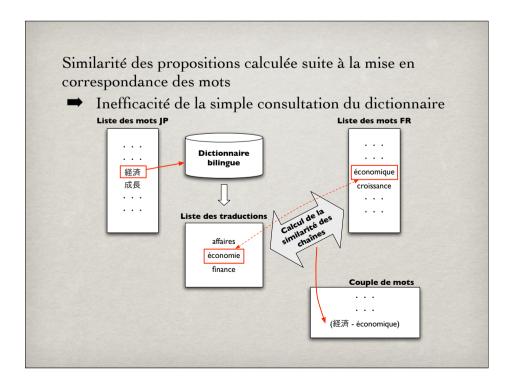
$$M^{final} = a \frac{M^{couleur}}{\max(M^{couleur})} + (1 - a) \frac{M^{topo}}{\max(M^{topo})}$$

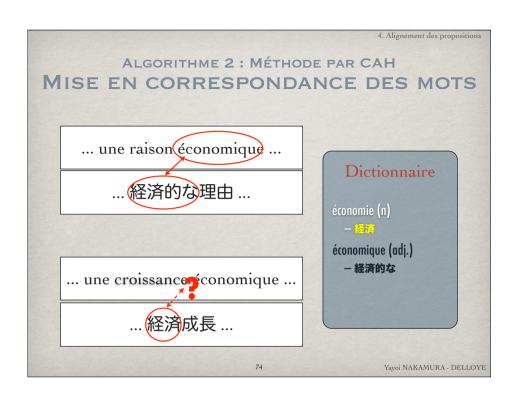
→ Adaptation à l'alignement des propositions

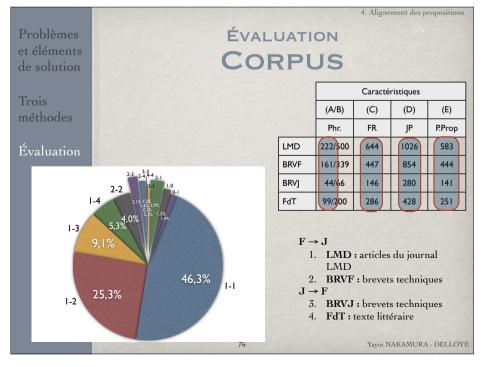
- 1. Alignement avec la topologie (Kosinov)
- 2. Alignement des graphes valués par les types de propositions (Lerallut)

72









ÉVALUATION RÉSULTAT

Trois méthodes:

MI: Topologique

M2: Graphe valué

M3: CAH

		Exact (F)		Partiel (G)				
	MI	M2	M3	MI	M2	M3		
LMD	0,127	0,200	0,591	0,643	0,784	0,951		
BRVF	0,081	0,158	0,706	0,619	0,705	0,977		
BRVJ	0,048	0,078	0,537	0,663	0,689	0,990		
FdT	0,138	0,151	0,464	0,670	0,659	0,932		

Résultats non satisfaisants de MI et de M2 :

informations insuffisantes pour l'alignement

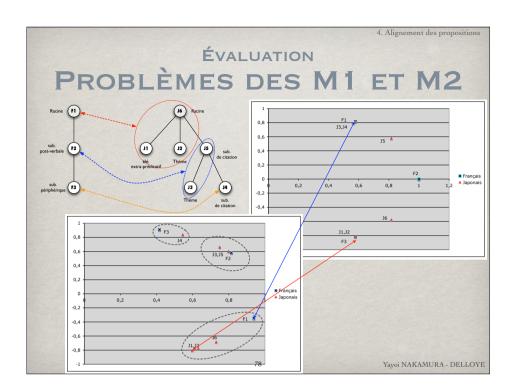
77

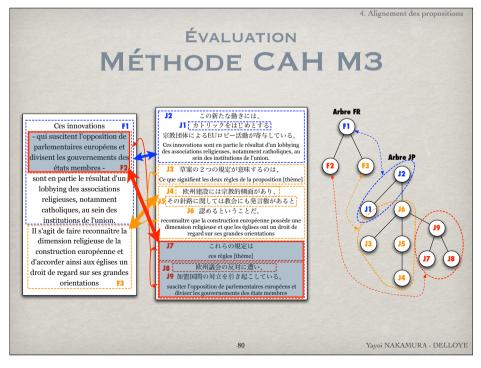
Yayoi NAKAMURA - DELLOYE

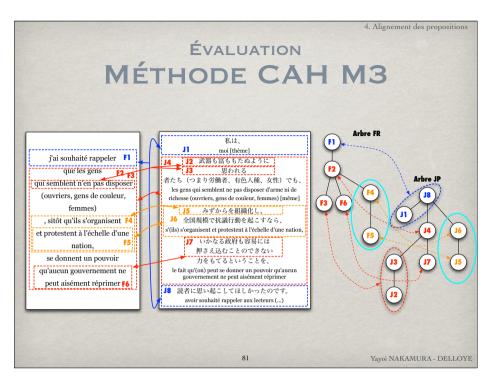
Yayoi NAKAMURA - DELLOYE

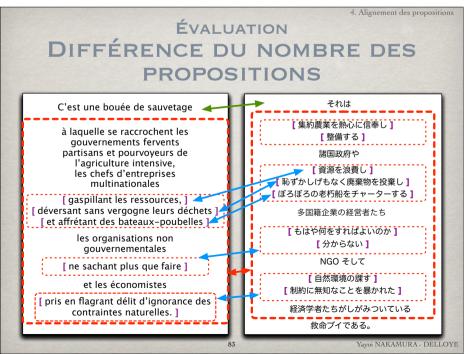
ÉVALUATION MÉTHODE CAH M3

L'introduction des informations lexicales a permis d'aligner correctement des phrases pour lesquelles la topologie et les informations sur les types des propositions ne suffisaient pas.

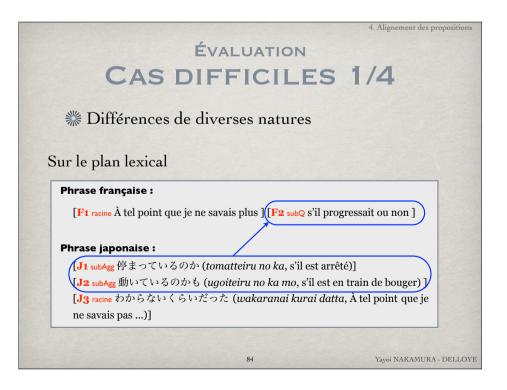








ÉVALUATION MÉTHODE CAH M3 * Méthode prometteuse mais résultat non encore tout à fait satisfaisant * Mauvais résultat du calcul de similarité lexicale: — Étude sur la réorganisation du dictionnaire — Causes plus complexes: absence d'unités dans le texte japonais * Deux types de sources d'erreurs — divergence de la définition des propositions — alignement fondamentalement difficile





Sur le plan syntaxique

Phrase française:

[F1 racine c'est notamment lors des débats sur les programmes d'aide aux pays du sud][F2 subR que les questions de la contraception et du statut de la famille sont abordées]

Phrase japonaise:

[J1 thème 避妊と家族の地位という問題は (hinin to kazoku no chii toiu mondai wa, les questions de la contraception et du statut de la famille [thème]] [J2 racine 特に開発途上国援助プログラムをめぐる議論の中で大きく取り上げられた (tokuni kaihatsu tojô koku enjo puroguramu wo meguru giron no nakade ôkiku toriagerareta, être abordé, notamment lors des débats sur les programmes d'aide aux pays en voie de développement)]

85

Yayoi NAKAMURA - DELLOYE

ÉVALUATION CAS DIFFICILES 4/4

Sur le plan rhétorique

Phrase française:

[F1 ED En y réfléchissant,]([F2 ED les trucages,] [F3 racine je n'étais pas près de les découvrir :)] [F4 ED déjà,] [F5 propord je ne savais pas] F6 subQ si l'ascenseur marchait ou non]

Phrase japonaise:

[**J1** subCond 考えてみれば (kangaete mireba, si (je) réfléchis)] [**J2** subCit たねどころか私には (tane dokoro ka watashi ni wa, sans aller jusqu'aux trucages, à moi) [**J3** subAgg エレベーターが動いているのか (erebêtâ ga ugoiteiru no ka, si l'ascenseur est en train de bouger)] [**J4** subAgg 停まっているのかさえ (tomatteiru no ka sae, s'il est arrêté)] わからないのだ (wakaranai no da, je ne sais pas...)]

Si je réfléchis bien, je ne sais, sans aller jusqu'aux trucages, même pas si l'ascenseur est en train de bouger ou s'il est arrêté EVALUATION
CAS DIFFICILES 3/4

Sur le plan rhétorique

Phrase française:

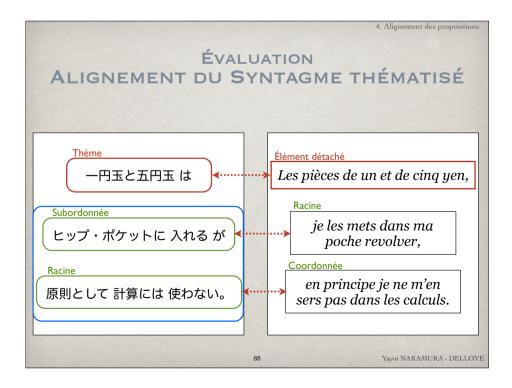
[F1 racine Je n'arrivais pas à croire] [F2 subQ que c'était moi] [F3 subR qui avais émis ce bruit]

Phrase japonaise:

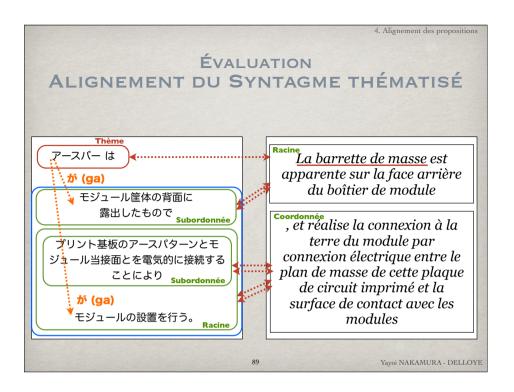
[J1 racine 私には (watashi ni wa, à moi) (J2 subCit それが (sore ga, ceci [ga]) [J3) subR 自分の体から発せられた (jibun no karada kara hasserareta, émis de mon corps)] 音だとは (oto da to wa, être un son/bruit [citation+wa])]とうしても思えなかった (dôshitemo omoenakatta, je n'arrivais pas à croire...)]

Je n'arrivais pas à croire que c'était un bruit émis de mon corps.

Yavoi NAKAMURA - DELLOYE



8



Conclusion & perspectives

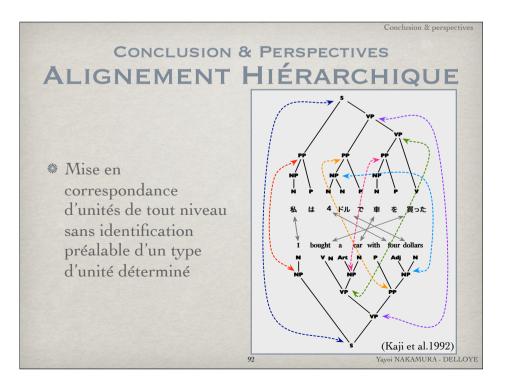
CONCLUSION & PERSPECTIVES 1/2

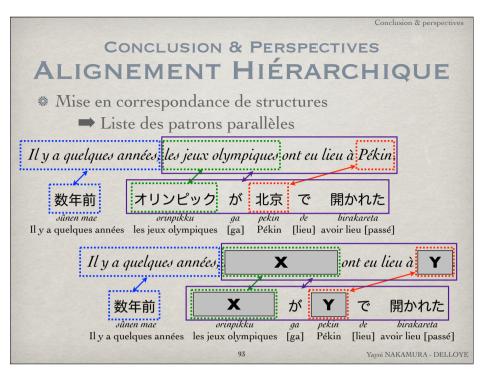
- * Indépendance de chaque système : différentes pistes d'amélioration pour chacun
- * Alignement de textes comparables

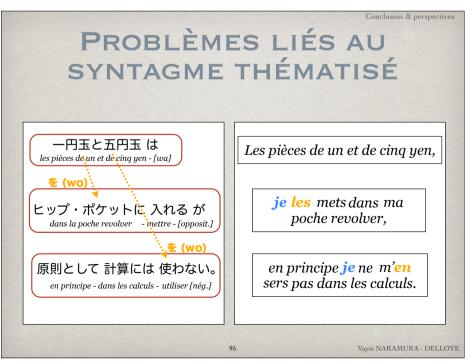
Textes comparables ≠ textes parallèles un ensemble de textes sur un même sujet dont aucun n'est traduction de l'un d'entre eux

· Ré-examen des unités à aligner (alignement hiérarchique)

CONCLUSION







CONCLUSION & PERSPECTIVES 2/2 Fort investissement dans les études linguistiques, mais un grand nombre de questions en suspens dans les travaux linguistiques du japonais Problèmes liés au traitement des syntagmes thématisés

