

ALIGNEMENT AUTOMATIQUE DES PROPOSITIONS FRANÇAIS-JAPONAIS

Y A Y O I N A K A M U R A - D E L L O Y E

Université Paris VII, LATTICE - CNRS UMR 8094

Université Paris X, MoDyCo - CNRS UMR 7114

yayoi@free.fr

INTRODUCTION TEXTES PARALLÈLES ET ALIGNEMENT 1/2

- ☼ **Alignement** = mise en correspondance entre des unités de **textes parallèles**
- ☼ **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
 - ✓ Label France
 - ✓ Manuels de logiciels libres
 - ✓ Textes G7, G8, etc.
 - ✓ Textes littéraires (Aozora Bunko : <http://www.aozora.gr.jp/>)

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 : **AIALeR**
(Nakamura-Delloye, 2005)

- ☼ Possibilité d'alignement à différents niveaux

INTRODUCTION POURQUOI LA PROPOSITION ? 1/2

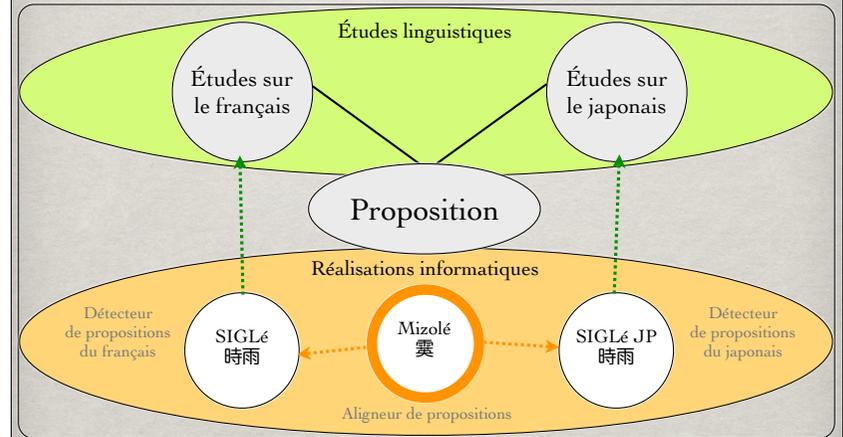
- ☼ Nécessité d'une unité de traitement plus petite que la phrase
 - ➔ Proposition = bon candidat
 - la documentation technique ;
 - l'analyse discursive ;
 - ...

INTRODUCTION 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 la correspondance
 - mot français « **compte** » = 10 définitions constituées de noms japonais différents non interchangeables
 - « **tenir compte** » = « 考慮する »

SCHÉMA DES TRAVAUX



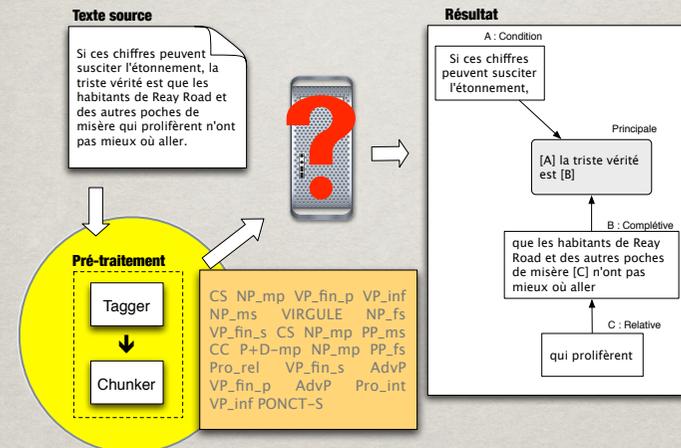
しぐれ 時雨 【figure】 n.

1. Brève averse. **2. INFORM. SIGLÉ** (*Systeme 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

CONTEXTE PRÉ-TRAITEMENT



PLAN

- ✿ Études linguistiques sur la proposition du français
- ✿ Réalisation
- ✿ Évaluation
- ✿ Perspectives

ÉTUDES LINGUISTIQUES

Études linguistiques

- Typologie des propositions
- Classement des subordonnées

Réalisation

Évaluation

Perspectives

- ✿ **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**

NOTION PRINCIPALE : PROPOSITION

Proposition = sujet + 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, ...

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 gêne, 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.*

Études linguistiques

- Typologie des propositions
- Classement des subordonnées

Réalisation

Évaluation

Perspectives

PROBLÈMES DES TYPOLOGIES USUELLES

☉ Définitions des subordonnées souvent selon la nature du connecteur qui les introduit

... MAIS

- ➔ Étiquetage automatique très difficile ;
- ➔ Comportements syntaxiques différents de termes appartenant à la même catégorie
- ➔ Qu'est-ce qu'une locution conjonctive (LC) ? Les LCs constituent-elles une liste fermée ?

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

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...

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 j'habite*
 - 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.

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)

POSITION POST-VERBALE (SUBORDONNÉE COMPLÉMENT : subQ)

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.*

Post-verbale :
Substantives

Initiale/
Finale :
Adverbiales

Post-nominale :
Adjectives
Adverbiales

Autres
positions
SN :
Substantives

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 gêne, 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.*

Post-verbale :
Substantives

Initiale/
Finale :
Adverbiales

Post-nominale :
Adjectives
Adverbiales

Autres
positions
SN :
Substantives

POSITION POST-NOMINALE (SUBORDONNÉE DÉTERMINANTE : subR)

- ☀ Adjectives : (relatives et complétives)
 - *La peinture qui m'a fascinée*
 - *La peinture dans laquelle notre maison était reproduite*
 - *Ce à quoi je m'attendais*
 - *L'idée que tout est fini.*
- ☀ Adverbiales : (intégratives adverbiales)
 - *La déception du père quand il a entendu cette nouvelle*
 - *La télévision où je veux* (publicité Orange)

Post-verbale :
Substantives

Initiale/
Finale :
Adverbiales

Post-nominale :
Adjectives
Adverbiales

Autres
positions
SN :
Substantives

AUTRES POSITION SN (SUBORDONNÉE SUBSTANTIVE : subSN)

Position Sujet :

- *Qui dort dine.* (Intégrative pronominale)
- *Que vous avez menti me déçoit.* (Complétive)
- *Comment il a commis ce crime n'a jamais été établi.* (Percontative)

Après une préposition :

- *Je voterai pour qui me promettra moins d'impôts.* (Intég.pron.)
- *Je partirai avant qu'il arrive.* (Complétive)
- *Il faudra se poser la question de pourquoi nous avons été choisis.* (Percontative)

Position initiale en prolepse :

- *Comment il a fait, je vous le demande !* (Percontative)

Post-verbale :
Substantives

Initiale/
Finale :
Adverbiales

Post-nominale :
Adjectives
Adverbiales

Autres positions
SN :
Substantives

ÉTUDES LINGUISTIQUES NOTRE TYPOLOGIE DES CONNECTEURS

| | Int/Fin | post-V | post-N | Autres SN |
|------|------------------|--------|--------|-----------|
| Sub. | Intégrative pro. | | △ | △ |
| | Percontative | | ✓ | △ |
| Adj. | Complétive | ✓ | | △ |
| | Relative | | | ✓ |
| Adv. | Intégrative adv. | ✓ | △ | |

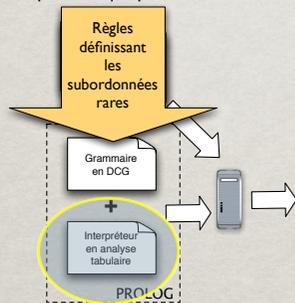
✓ = fréquent ; △ = moins fréquent / rare

| Position | post-V | | | post-N | | | Int / Fin | | | pos. SN | | |
|----------|--------|---|---|--------|---|---|-----------|---|---|---------|---|---|
| | I | C | P | I | C | P | R | I | C | P | R | |
| qui | △ | | ✓ | | | | ✓ | | | | △ | △ |
| que | ✓ | | | | ✓ | | ✓ | ✓ | | | △ | |
| dort | | | | | | | ✓ | | | | | |
| où | | ✓ | | ? | | | ✓ | △ | | | | △ |
| quand | | | | | △ | | | ✓ | | | | △ |
| comme | | ✓ | | △ | | | | ✓ | | | | △ |
| si | | ✓ | | | | △ | | | | | | △ |
| quoi | | ✓ | | | | | ✓ | | | | | △ |
| lequel | | ✓ | | | | | ✓ | | | | | △ |
| quel | | ✓ | | | | | | | | | | △ |
| combien | | ✓ | | | | | | | | | | △ |
| comment | | ✓ | | | | | | | | | | △ |
| pourquoi | | ✓ | | | | | | | | | | △ |

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

RÉALISATION

Ajout dynamique en cas d'échec de la première analyse pour chaque phrase



(Pereira & Shieber 1987)

```

<prop id='1' etq='principale' pere='0' fils ='2;3;'>
  [subP], la triste vérité est [subQ]
</prop>
<prop id='2' etq='subP' pere='1' fils ='>
  si ces chiffres peuvent susciter l' étonnement
</prop>
<prop id='3' etq='subQ' pere='1' fils ='4;'>
  que les habitants de Reay Road et des autres
  poches de misère [subR]n' ont pas mieux où aller
</prop>
<prop id='4' etq='subR' pere='3' fils ='>
  qui prolifèrent
</prop>
    
```

Études linguistiques

- Typologie des propositions
- Classement des subordinées

Réalisation

Évaluation

Perspectives

ÉVALUATION 1/8

| | Nb de phr | Rappel | Préc.1 | Préc. 2 | Préc. T |
|---------|-----------|--------|--------|---------|---------|
| G8 | 53 | 96,2 | 98,0 | 100,0 | 98,0 |
| Unicode | 274 | 81,4 | 96,2 | 97,8 | 94,1 |
| Zadig | 1206 | 88,6 | 92,8 | 95,3 | 88,4 |
| LMD* | 1713 | 84,9 | 89,2 | 98,0 | 87,4 |

* Le Monde Diplomatique

- A : Nombre total de phrases
- B : Nombre de phrases dont l'analyse a abouti
- C : Nombre de phrases dont les frontières de propositions sont correctement détectées
- D : Nombre de phrases dont les relations des propositions sont correctement analysées

$$\text{Rappel (\%)} = \frac{B}{A} \times 100$$

$$\text{Précision 1 (\%)} = \frac{C}{B} \times 100$$

$$\text{Précision 2 (\%)} = \frac{D}{C} \times 100$$

$$\text{Précision T (\%)} = \frac{\text{Préc. 1} \times \text{Préc. 2}}{100} = \frac{D}{B} \times 100$$

Études linguistiques

- Typologie des propositions
- Classement des subordinées

Réalisation

Évaluation

Perspectives

ÉVALUATION 2/8

Si ces chiffres peuvent susciter l'étonnement, la triste vérité est que les habitants de Reay Road et des autres poches de misère qui prolifèrent n'ont pas mieux où aller.

XML

```
<prop id='1' etq='racine' pere='0' fils ='2;3;4;'>
  [subP], la triste vérité est [subQ]
</prop>
<prop id='2' etq='subP' pere='1' fils =">
  si ces chiffres peuvent susciter l'étonnement
</prop>
<prop id='3' etq='subQ' pere='1' fils ='4;'>
  que les habitants de Reay Road et des autres poches de misère [subR] n'ont pas
  mieux où aller
</prop>
<prop id='4' etq='subR' pere='3' fils =">
  qui prolifèrent
</prop>
```

ÉVALUATION 3/8

En arrivant aux frontières qui séparent l'Arabie pétérée de la Syrie, comme il passait près d'un château assez fort, des arabes armés en sortirent.

XML

```
<prop id='1' etq='racine' pere='0' fils ='2;4;'>
  [ED] [subP], des arabes armés en sortirent
</prop>
<prop id='2' etq='ED' pere='1' fils ='3;'>
  en arrivant aux frontières [subR],
</prop>
<prop id='3' etq='subR' pere='2' fils =">
  qui séparent l'Arabie pétérée de la Syrie
</prop>
<prop id='4' etq='subP' pere='1' fils =">
  comme il passait près_d' un château assez fort
</prop>
```

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 [inc] aussi bien que ce [subR]
</prop>
<prop id='2' etq='subR' pere='1' fils =">
  qui passe sur mes terres
</prop>
<prop id='3' etq='inc' pere='1' fils =">
  , dit -il,
</prop>
<prop id='4' etq='subR' pere='1' fils =">
  que je trouve sur les terres des autres
</prop>
```

ÉVALUATION 5/8

✪ Erreurs de l'analyse linéaire :

- étiquettes erronées attribuées par le pré-traitement ;
- mauvaise interprétation de connecteurs :
introduceur de proposition
⇔ intr. de syntagme (structure réduite sans verbe)
- ambiguïté de la position post-prépositionnelle
intégrative pronominale (subSN)
⇔ subordonnée déterminante (subR)

Et dire [complétive 1] qu'au moment de son apogée, dans les années 1950,
Cockerill employait encore plus de 25 000 personnes [complétive 2] que la ville de
Seraing était toujours noire de fumée, de bruit, de monde, de travail.]

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) **!**, **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 cœur (1) **!** **que** cela valait beaucoup (2)

ÉVALUATION 7/8

(Paris avait estimé X) + (car Y) ?

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

☼ 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. (3)

ÉVALUATION 8/8

☼ Fréquence des subordonnées

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

✓ = fréquent ; △ = moins fréquent / rare

PISTES D'AMÉLIORATION

Études linguistiques

- Typologie des propositions
- Classement des subordonnées

Réalisation

Évaluation

Perspectives

■ 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

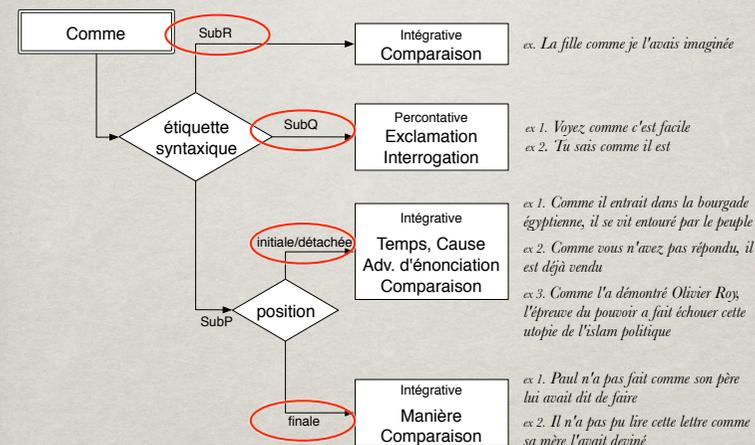
attribution d'étiquettes syntactico-sémantiques

PISTES D'AMÉLIORATION AFFINEMENT DES ÉTIQUETTES

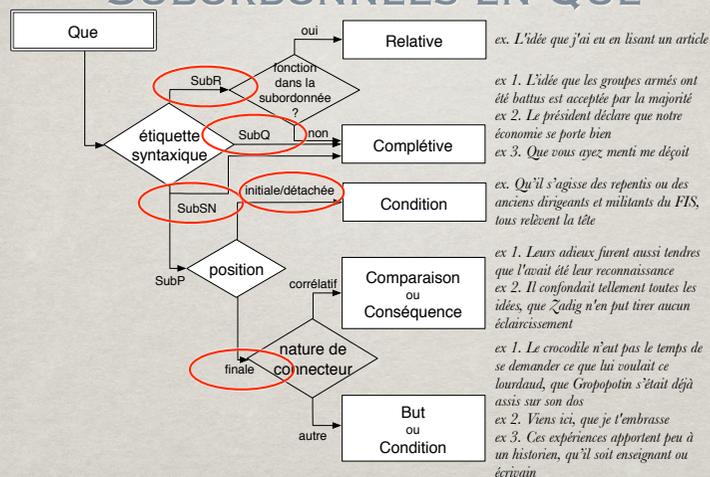
| Position | post-V | | | | post-N | | | | Int /Fin | | | | pos. SN | | | |
|----------|--------|---|---|---|--------|---|---|---|----------|---|---|---|---------|---|---|---|
| | I | C | P | R | I | C | P | R | I | C | P | R | I | C | P | R |
| qui | △ | | ✓ | | | | | ✓ | | | | | △ | | △ | |
| que | | ✓ | | | | ✓ | | ✓ | ✓ | | | | | △ | | |
| dont | | | | | | | | | | | | | | | | |
| dont | | | | | | | | | | | | | | | | |
| où | | | ✓ | | | ? | | ✓ | △ | | | | | | △ | |
| quand | | | ✓ | | | △ | | | | | | | | | △ | |
| comme | | | ✓ | | | △ | | | | | | | | | △ | |
| si | | | ✓ | | | | △ | | ✓ | | | | | | △ | |
| quoi | | | | | | | | | | | | | | | △ | |
| lequel | | | | | | | | ✓ | | | | | | | △ | |
| quel | | | | | | | | | | | | | | | △ | |
| combien | | | ✓ | | | | | | | | | | | | △ | |
| comment | | | ✓ | | | | | | | | | | | | △ | |
| pourquoi | | | ✓ | | | | | | | | | | | | △ | |

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

AFFINEMENT DES ÉTIQUETTES : SUBORDONNÉES EN COMME



AFFINEMENT DES ÉTIQUETTES : SUBORDONNÉES EN QUE



SIGLÉ JP

SYSTÈME D'IDENTIFICATION DES PROPOSITIONS DU JAPONAIS

PLAN

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

Études linguistiques

- Structure de la phrase
- 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**

QUELQUES ÉLÉMENTS FONDAMENTAUX DE LA PHRASE JAPONAISE

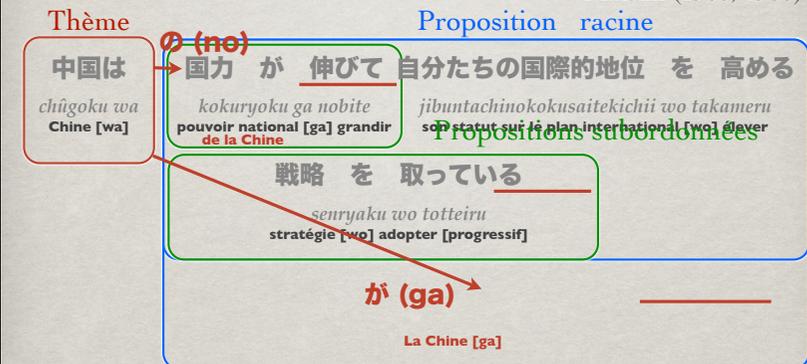
| | | | | |
|---------------------|---------------------|-----------------------|----------------------------|----------------------|
| 運 <small>よく</small> | 本屋 <small>で</small> | 探 <small>していた</small> | 辞書 <small>を</small> | 見 <small>つけた</small> |
| un'yoku | hon'ya - de | sagashiteita | jisho - wo | mitsuketa |
| Par chance | librairie - [lieu] | que (je) cherchais | dictionnaire - [accusatif] | trouver [passé] |

« Par chance, (j'ai) trouvé dans une librairie le dictionnaire que (je) cherchais »

- ✿ Mot prédicatif précédé par ses compléments
- ✿ Marqueurs de fonction = variations de forme, particules de cas
- ✿ Omission très fréquente d'éléments

ÉTUDES LINGUISTIQUES PHRASE JAPONAISE 1/2

Mikami (1953, 1955)



« Son pouvoir national ayant grandi, la Chine adopte une stratégie permettant d'élever son statut sur le plan international »

ÉTUDES LINGUISTIQUES PHRASE JAPONAISE 2/2

さいわいなことに 野球部員たちは、そこにあらわれたのが たまたま 七瀬という、
教務課職員とはいえ 私立手部高校随一の美人であったが為に、異変への関心をすぐ
失った。 **Élément externe**

(*saiwaina koto ni - yakyū buin tachi wa - sokoniarawareta no ga - tamatama - nanase to iu -
kyōmuka shokutū towaie - shiritsu tebe kōkō zuiitsu ni bijin de atta ga tameni - ihen eno kanshin
wo sugu ushinatta*)

(*par chance - équipes du club de baseball [wa] - personne apparaissant là [ga] - par
hasard - appelé Nanase - bien qu'une employée administrative - être la plus belle femme
du lycée privé Tebe [cause] - perdre l'intérêt pour l'événement extraordinaire [passé]*)

« **Par chance** du fait que la personne apparue était par hasard Nanase, la plus belle fille
du lycée privé de Tebe, bien qu'employée administrative, les équipes du club de baseball
perdirent tout de suite tout intérêt pour cet événement extraordinaire »

TYPOLOGIE DES SUBORDONNÉES

Classement selon uniquement des critères morpho-
syntaxiques (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

ÉTUDES LINGUISTIQUES PROBLÈMES LIÉS AUX CONNECTEURS

- Subordonnée déterminante sans connecteur

日本へ **行った** 友達 **Substantif**

nihon e - itta - tomodachi
Japon [e] - aller/partir [passé] - ami
« **ami qui est parti au Japon** »

- Subordonnée avec particule conjonctive

日本へ **行った** **が** **Particule conjonctive**

nihon e - itta - ga
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** **Mot agglutinant**

nihon e - itta - toki
Japon [e] - aller/partir [passé] - temps
« **quand (je suis/tu es/il est...) parti au Japon** »

(Sakuma 1940)

ÉTUDES LINGUISTIQUES

RÉALISATION DU SYSTÈME

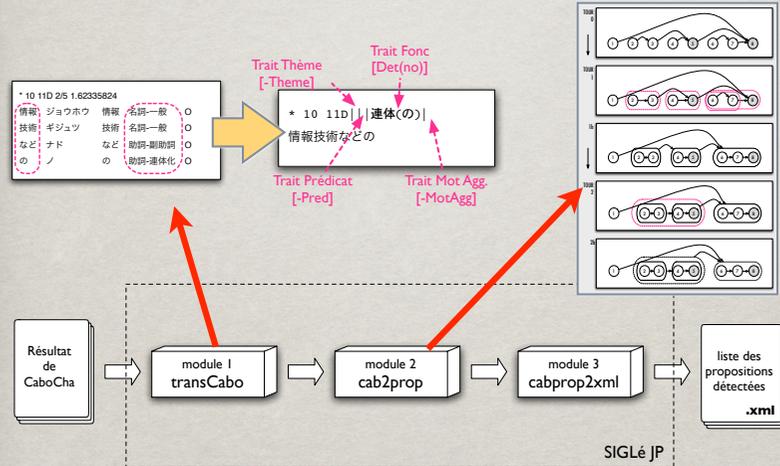
- Méthode existante : CBAP (Maruyama et al. 2004) :
incapable de traiter les structures imbriquées
- ➔ Utilisation d'un analyseur des relations de dépendance (Kudo & Matsumoto 2002)

informations sur le chunk

informations sur chaque unité constituant le chunk

| | |
|-------------------------|---|
| 現在クワンサイ現在名詞-動詞可能 | 0 |
| 記号-読点 | 0 |
| 1 2D 0/1 0.14490557 | 0 |
| 多くオオク 多く名詞-動詞可能 | 0 |
| の ノ の 助詞-連体化 | 0 |
| * 2 3D 0/0 0.11489553 | 0 |
| 国公立 コッコウリツ 国公立 名詞-一般 | 0 |
| 記号-読点 | 0 |
| * 3 3D 1/2 0.95327759 | 0 |
| 私立シリンツ 私立名詞-一般 | 0 |
| 大学ダイガク 大学名詞-一般 | 0 |
| が ガ が 助詞-格助詞-一般 | 0 |
| * 4 3D 1/2 1.64770347 | 0 |
| 社会シャカイ 社会名詞-一般 | 0 |
| 人ジンズ人 名詞-接尾-一般 | 0 |
| も モ も 助詞-格助詞 | 0 |
| * 5 6D 1/1 1.33629884 | 0 |
| 家庭ケイゴウ 家庭名詞-他家接尾 | 0 |
| できる デキル できる 動詞-自立 一段基本形 | 0 |
| * 6 7D 1/2 0.00000000 | 0 |
| 公開コウカイ 公開名詞-他家接尾 | 0 |
| 講座コウザイ 講座名詞-一般 | 0 |
| をヲを 助詞-格助詞-一般 | 0 |
| * 7 10 0/2 0.00000000 | 0 |
| 受けモツケ 受ける 動詞-自立 一段連用形 | 0 |
| て テ て 助詞-接続助詞 | 0 |
| いるイルいる 動詞-非自立 一段基本形 | 0 |
| 記号-句点 | 0 |
| EOS | 0 |

RÉALISATION DU SYSTÈME



ÉVALUATION CORPUS

Études linguistiques

- Structure de la phrase
- Classement des propositions

Réalisation

Évaluation

| | | LMD | Brevet | Asahi | Murakami |
|----------|--|------|--------|-------|----------|
| A | Nombre de phrases | 124 | 158 | 112 | 149 |
| B | Nombre de propositions | 433 | 670 | 450 | 479 |
| C | Nombre moyen de propositions dans une phrase | 3,49 | 4,24 | 3,99 | 3,21 |

● LMD ○ BREVET ◆ ASAHI ◇ MURAKAMI

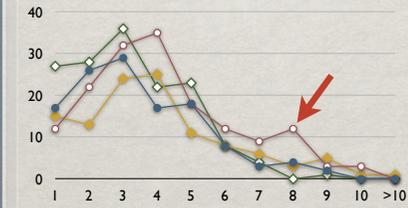


Fig. Distribution des phrases en fonction du nombre de propositions

- LMD : traduction d'articles d'un journal français
- Brevet : traduction de brevets techniques français
- Asahi : articles d'un journal japonais
- PdT : texte littéraire, extrait du roman « La fin des temps » de MURAKAMI Haruki.

ÉVALUATION RÉSULTATS 1/2

| | | LMD | Brevet | Asahi | Murakami |
|----------|--|----------|-----------|----------|-----------|
| A | Nombre de phrases | 124 | 158 | 112 | 149 |
| B | Nombre de propositions | 433 | 670 | 450 | 479 |
| C | Nombre moyen de propositions dans une phrase | 3,49 | 4,24 | 3,99 | 3,21 |
| D | Nombre de propositions détectées | 444 | 672 | 453 | 490 |
| E | Nombre de propositions détectées correctement | 389 | 589 | 391 | 426 |
| F | Rappel (= E/B) | 0,898 | 0,879 | 0,869 | 0,891 |
| G | Précision (= E/D) | 0,876 | 0,876 | 0,863 | 0,869 |
| H | Analyse linéaire = nombre de phrases correctement analysées (H/A %) | 99 (80%) | 119 (75%) | 85 (76%) | 120 (81%) |
| I | Analyse structurale = nombre de phrases correctement analysées (I/H %) | 94 (95%) | 107 (90%) | 79 (93%) | 111 (93%) |

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

みぞれ 霰【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

PLAN

- ☼ Problèmes et éléments de solution
- ☼ Trois méthodes réalisées
- ☼ Procédure d'alignement
- ☼ Évaluation
 - Constats sur les corpus utilisés
 - Résultats
 - Analyse des erreurs

Problèmes et éléments de solution

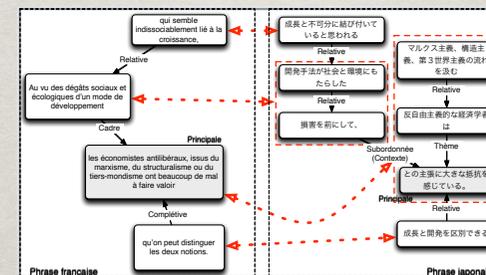
PROBLÈMES

- ☼ Peu de travaux (Boutsis & Piperidis 1998) (Wang & Ren 2005)
- ☼ Impossibilité d'une simple application des méthodes classiques d'alignement des phrases pour l'alignement fr-jp due au non-parallélisme

➔ **Alignement à l'aide des graphes**

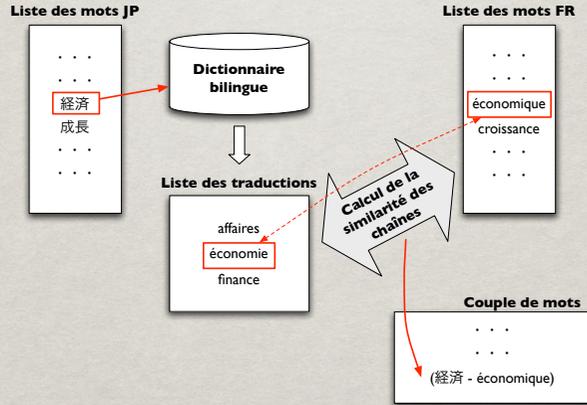
ÉLÉMENTS DE SOLUTION

▶ Alignement à l'aide des graphes



- ➔ **Méthodes spectrales d'appariement des graphes**
- ➔ **Classification Ascendante Hiérarchique (CAH)**

ALGORITHME 2 : MÉTHODE PAR CAH MISE EN CORRESPONDANCE DES MOTS



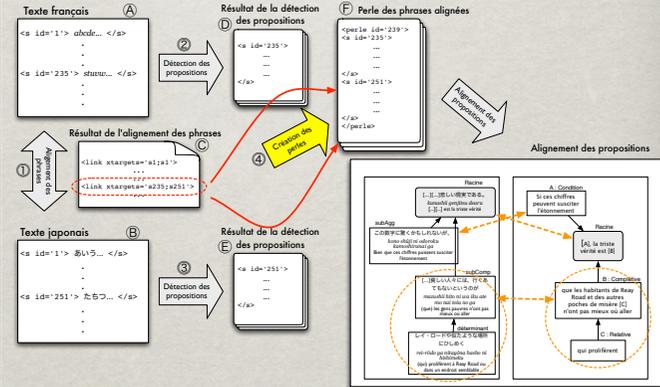
PROCÉDURE D'ALIGNEMENT

Problèmes et éléments de solution

Trois méthodes

Procédure

Évaluation



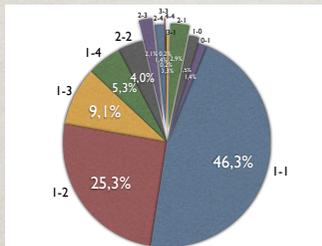
ÉVALUATION CORPUS

Problème et éléments de solution

Trois méthodes

Procédure

Évaluation



| | Caractéristiques | | | |
|------|------------------|-----|------|-------|
| | (A/B) | (C) | (D) | (E) |
| | Phr. | FR | JP | PProp |
| LMD | 222/500 | 644 | 1026 | 583 |
| BRVF | 161/339 | 447 | 854 | 444 |
| BRVJ | 44/66 | 146 | 280 | 141 |
| FdT | 99/200 | 286 | 428 | 251 |

F → J

1. LMD : articles du journal LMD
2. BRVF : brevets techniques

J → F

3. BRVJ : brevets techniques
4. FdT : texte littéraire

ÉVALUATION RÉSULTAT

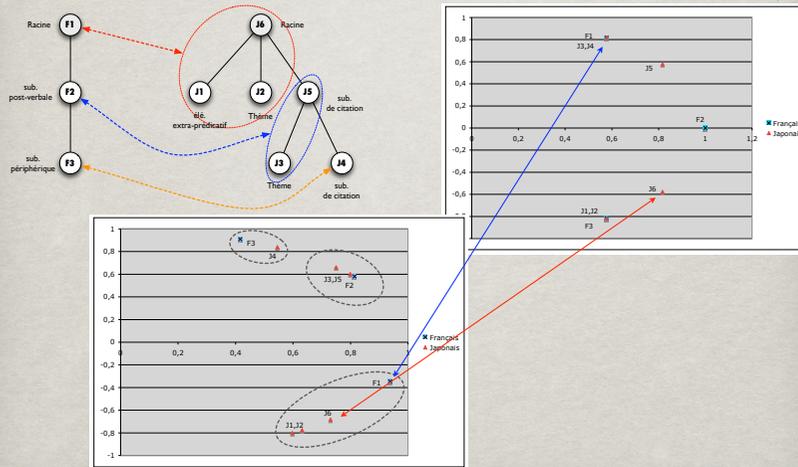
Trois méthodes :

- M1 : Topologique
- M2 : Graphe valué
- M3 : CAH

| | Exact (F) | | | Partiel (G) | | |
|------|-----------|-------|-------|-------------|-------|-------|
| | M1 | M2 | M3 | M1 | 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 M1 et de M2 : informations insuffisantes pour l'alignement

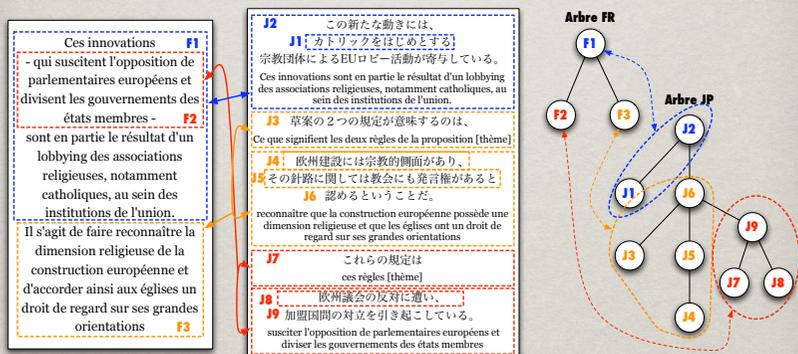
ÉVALUATION PROBLÈMES DES M1 ET M2



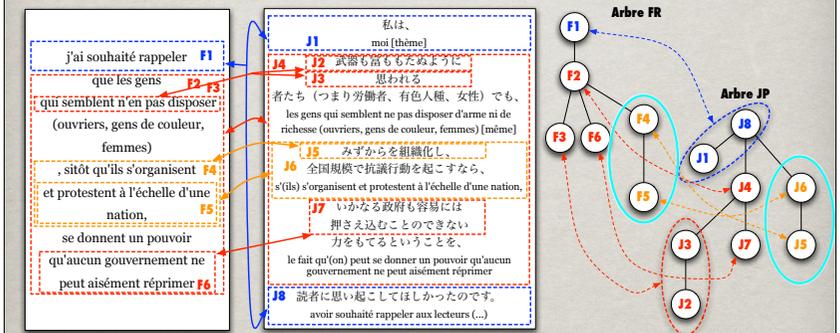
ÉVALUATION MÉTHODE CAH M3

✂ 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



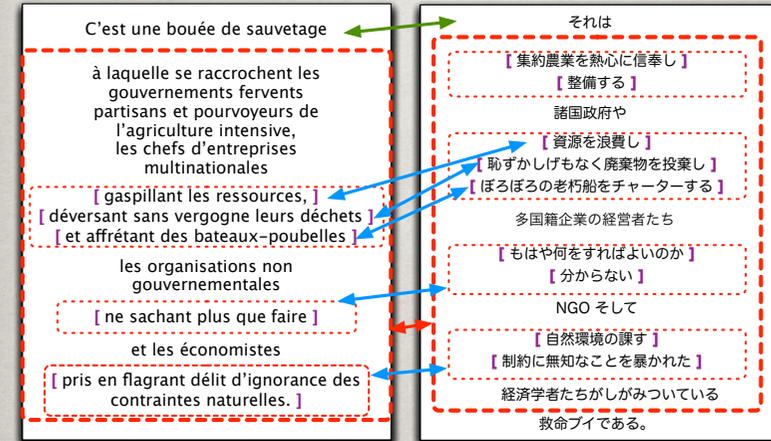
ÉVALUATION MÉTHODE CAH M3



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
 - ➡ Recherche d'une meilleure méthode de mise en correspondance
- Causes peut-être plus complexes de cette inefficacité de mise en correspondance des mots par dictionnaire
 - ➡ À cause de l'absence d'unités dans le texte japonais ?

ÉVALUATION DIFFÉRENCE DU NOMBRE DES PROPOSITIONS



ÉVALUATION CAS DIFFICILES 1/4

☀ Différences de diverses natures

Sur le plan lexical

Phrase française :

[F1 racine À tel point que je ne savais plus] [F2 subQ s'il progressait ou non]

Phrase japonaise :

- [J1 subAgg 停まっているのか (tomatteiru no ka, s'il est arrêté)]
- [J2 subAgg 動いているのかも (ugoiteiru no ka mo, s'il est en train de bouger)]
- [J3 racine わからないくらいだった (wakaranai kurai datta, À tel point que je ne savais pas ...)]

ÉVALUATION CAS DIFFICILES 2/4

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 kaihatu 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)]

ÉVALUATION 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])] どうしても思えなかった (dôshitomo omoenakatta, je n'arrivais pas à croire...)]

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

ÉVALUATION CAS DIFFICILES 4/4

Sur le plan rhétorique

Phrase française :

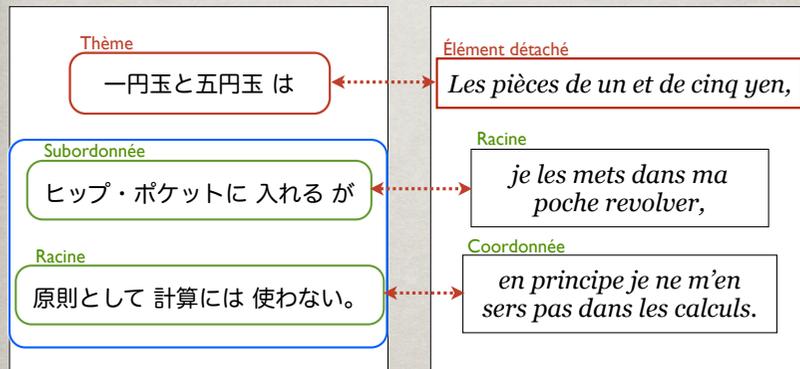
[F1 ED En y réfléchissant,] [F2 ED les trucages,] [F3 racine je n'étais pas près de les découvrir :] [F4 ED déjà,] [F5 proprd je ne savais pas] [F6 subQ si l'ascenseur marchait ou non]

Phrase japonaise :

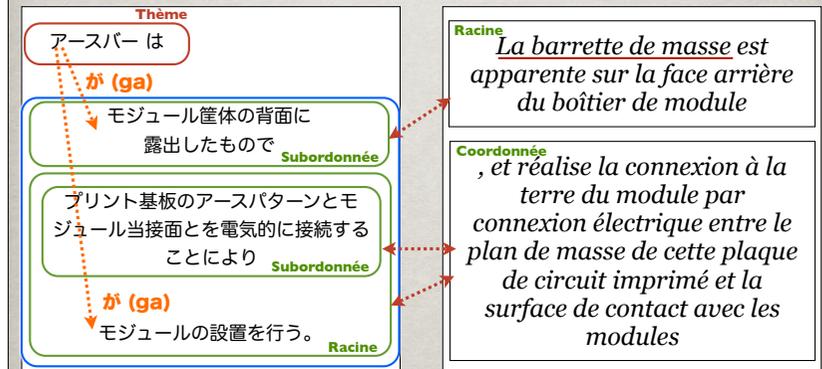
[J1 subCond 考えてみれば (kangaete mireba, si (je) réfléchis)] [J2 subCit たねどころか私には (tane dokoro ka watashi ni wa, sans aller jusqu'aux trucages, à moi)] [J3 subAgg エレベーターが動いているのか (erebêta ga ugoiteiru no ka, si l'ascenseur est en train de bouger)] [J4 subAgg 停まっているのかさえ (tomatteiru no ka sae, s'il est arrêté)]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

ÉVALUATION ALIGNEMENT DU SYNTAGME THÉMATISÉ



ÉVALUATION ALIGNEMENT DU SYNTAGME THÉMATISÉ



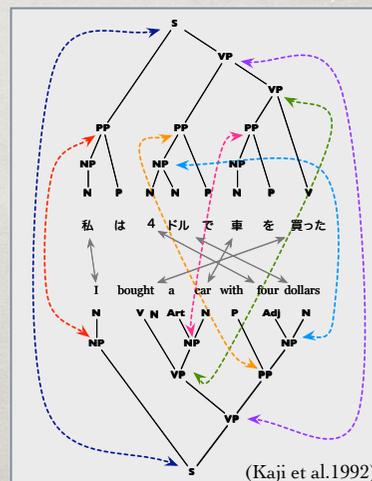
CONCLUSION

CONCLUSION & PERSPECTIVES 1/2

- Indépendance de chaque système : différentes pistes d'amélioration pour chacun
- Amélioration du concept même du système d'alignement des propositions
- Ré-examen des unités à aligner (alignement hiérarchique)

CONCLUSION & PERSPECTIVES ALIGNEMENT HIÉRARCHIQUE

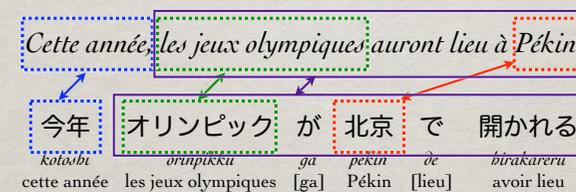
- Mise en correspondance d'unités de tout niveau sans identification préalable d'un type d'unité déterminé



(Kaji et al.1992)

CONCLUSION & PERSPECTIVES ALIGNEMENT HIÉRARCHIQUE

- Mise en correspondance de structures
➡ Liste des patrons parallèles



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

PROBLÈMES LIÉS AU SYNTAGME THÉMATISÉ

