

# Traitement “good-enough” du pronom “on” : Vers une modélisation de la coréférence floue.

M. DELABORDE & F. LANDRAGIN

LATTICE (CNRS, ENS, Université de la Sorbonne Nouvelle, USPC, PSL)

Linguistic and Psycholinguistic Approaches to Text Structuring  
19/01/2018

## Définitions : référence

- **Expression référentielle** = maillon = mention : marqueur linguistique qui réfère à une entité qui existe dans le monde [Frege, 1892] ou que l'on peut se représenter [Charolles, 2002]
- **Référence** : acte de langage consistant à désigner un objet extralinguistique (par une expression référentielle)
- **Référence floue** : acte de langage consistant à désigner un objet extralinguistique **imprécis** (par une expression référentielle)
- **Coréférence** : relation entre plusieurs expressions qui réfèrent à la même entité

## Référence floue et chaînes de coréférences

- **Chaîne de coréférence** [Kleiber, 1994], [Schnecker and Landragin, 2014]: chaîne composée d'expressions coréférentes (au moins 3 maillons pour [Corblin, 1995], 2 suffisent en TAL)
- Référent parfois difficile à identifier par un lecteur (référence floue)  
Exemple : "On a besoin de travailler" → personne spécifique ou générique?
- Problématique pour les chaînes de coréférences :
  - Expressions non prises en compte
  - Perte d'une partie du sens

## Traitement “good-enough” des expressions

- Expériences psycholinguistiques : mouvance “good-enough”  
[Ferreira et al., 2002]

- Mauvaise compréhension des phrases passives  
[Ferreira and Stacey, 2000]
- Mauvaise compréhension des phrases à garden-path  
[Christianson et al., 2001]

→ Sujets capables de traiter certaines expressions de manière incomplète

⇒ Repenser la manière de modéliser les chaînes de coréférences.

# Approche

- **Approche descriptive**

- recueil d'exemples et description de leur fonctionnement (détails et paramètres de la résolution de la référence)

→ Collection d'occurrences de "on"

- **Volonté de formalisation**

- "traduction" de ce fonctionnement en forme logique ou dans un formalisme tel que la DRT (Discourse Representation Theory)

- **Proposition d'instructions / suggestions en TAL**

⇒ Partie descriptive avancée, pistes de formalisation

# Problématique

- **Approches classiques** : construction d'une chaîne quand il y a coréférence stricte.
  - Chaîne habituelle :
    - Description
    - Corpus
    - Manuel d'annotation
    - Schéma d'annotation
    - Outils d'annotation et d'exploration des données annotées
    - Outils de TAL
- Mais pourquoi ne pas construire une chaîne même quand le référent est **vague** ou **indéterminé**?
- [Charolles, 2014] : Profondeur de traitement dans les annotations des expressions référentielles.

# Le cas de “on”

- **Terminologie :**

- pronom personnel [Charaudeau, 1992]
- pronom indéfini [Sandfeld, 1970], [Grevisse and Goosse, 2002]
- pronom personnel indéfini [TLFI, 2002]
- pronom impersonnel [Cabredo Hofherr, 2008]

- **Selon les grammaires du français** [Riegel et al., 1994] :

- Toujours nominal
- Sujet
- Désigne des référents humains animés

- **Identification souvent difficile du référent :**

→ inclusif/exclusif? spécifique/générique?

- **Illustration du phénomène de coréférence floue par le pronom “on”**

→ Interprétation particulièrement propice au phénomène “good-enough”

## Identification du référent

- **“On”** [Viollet, 1988], [Rabatel, 2001], [Fløttum et al., 2008]  
3 façons de l'identifier dans le contexte d'une chaîne de référence :
  - 1 Identification claire
  - 2 Désambiguïstation par d'autres expressions
  - 3 Identification incomplète, floue

## Exemple 1 : Identification claire

(1)

“Comme ils ont aussi une oeuvre de Kiefer, j'ai écrit que je souhaitais qu'ils constituent une salle avec **Kiefer, Kantor et moi**. Parce qu'**on** parle de la même histoire, chacun depuis sa place.”

BOLTANSKI & GRENIER *La vie possible de Christian Boltanski*, 2007

→ Groupe bien identifié.

## Exemple 2 : Désambiguïsation par d'autres expressions

(2)

“**On** travaillait, **leur père** travaillait, **je** travaillais.”

LAGARCE, *Juste la fin du monde*, 1990

→ “leur père” et “je” permettent une désambiguïsation de “on”

→ Aucun autre référent candidat possible

## Exemple 3 : Identification incomplète, floue

(3)

“La nuit il rêvait d'une longue galerie de mine, noire, froide, et par terre, une ligne blanche, lumineuse, **on** suivait, **on** longeait la ligne, et l'**on** marchait, toujours plus nombreux, droit devant, mais soudain, la ligne se dérobaît, disparaissait, elle dérapait et s'esquivaît, de gauche, de droite, se recoupant sans cesse, ce n'était plus qu'un gribouillis au sol, éblouissant, et **les camarades** à côté appelaient au secours, hurlaient, se perdaient dans les ténèbres et l'**on** entendait les cris de ceux qui pour toujours perdaient la trace, happés par le noir comme par le glougloutement du doute.”

FILIPPETTI, *Les derniers jours de la classe ouvrière*, 2003

→ Désignation du groupe de manière floue

## Critères de construction de chaîne

- **Hypothèse** [Landragin and Tanguy, 2014] :
  - “On” = marqueur de **cohésion** et de **cohérence** au niveau du discours
  - Catégorisation selon **différents degrés de coréférence**
    - Proposé pour le concept de near-identity par [Recasens et al., 2011]
- Construction rationnelle des chaînes → définition de **critères** :
  - ① Au moins un référent doit être inclus dans l'interprétation de deux occurrences  

Attention : pas une unique chaîne avec toutes les mentions floues!  
→ Référent commun inclus de manière stricte dans le groupe flou
  - ② Tous les référents strictement inclus dans l'interprétation de deux occurrences doivent être identiques

## Différents degrés de coréférence

- Degré 1 : **“Coréférence stricte”**  
→ Référents identiques
- Degré 2 : **“Coréférence inclusive”**  
→ Un référent en inclut un autre
- Degré 3 : **“Coréférence floue”**  
→ Deux référents sont des groupes flous dont l'intersection est possible tout en restant potentiellement floue elle-même

# Théories de modélisation

- **Théorie des ensembles flous** [Zadeh, 1965]
  - Idée d'appartenance partielle à une classe
  - **Fonction d'appartenance** = Degré d'appartenance à des variables (entre 0 et 1)
- Généralisation de la théorie des ensembles avec des situations intermédiaires
- Champ "degré d'appartenance"
  - Choix libre de la valeur
  - Discrétisé dans le manuel

# Théories de modélisation

- **Théorie des possibilités** [Zadeh, 1978], [Dubois and Prade, 2012]

- Utilise les outils de la théorie des ensembles flous
- Mesure de **possibilité** d'appartenance à une classe
- Deux approches possibles :
  - Ordinale
  - Numérique

→ Champ à ajouter dans un schéma et un manuel d'annotation

⇒ Outils mathématiques permettant la gestion de l'**imprécision** et de l'**incertitude**

## Avantages par rapport aux méthodes actuelles

- Prise en compte du phénomène good-enough pour les chaînes de (co)référence :
  - Prendre en compte différents degrés de coréférence
    - Récupération des informations sémantiques
    - Meilleure précision des systèmes de TAL

# Conséquences en linguistique de corpus

- Revoir les **procédures d'annotation** de corpus
- [Charolles, 2014] : **schémas d'annotations** plus permissifs  
→ relâchements de contraintes
- Relâcher les métriques d'**accord inter-annotateurs**  
(erreur sur du flou : moins de poids)
- Métriques à définir, il existe déjà :
  - **Kappa** : plus ancienne
  - **Gamma** : plus récente mais peu d'applications,  
distinction entre deux types d'erreurs (unitizing et categorizing)

# Conséquences en TAL

- Procédure : évaluation des systèmes avec des métriques [Désoyer et al., 2015], pour le moment :
  - MUC : unité de base = liens de coréférence, erreurs d'insertion et de suppression prises en compte individuellement
  - B3 : unité de base = mention, singletons pris en compte
  - CEAF : temps de calcul plus long, comparaison de chaînes de coréférence
  - BLANC : plus récente, liens de coréférence et de non coréférence
- Premier cas : **systèmes classiques à base de règles**  
→ Plus grande complexité des algorithmes  
(règles spécifiques à prévoir)
- Vers une multiplication des identifications de flou  
(même quand ce n'est pas pertinent)

# Conséquences en TAL

- Deuxième cas : **système d'apprentissage** (quelle que soit la technique)
  - Reste envisageable avec des méthodes comme le **deep learning**
  - Mais → identifier des **features** qui aident le système à trouver des paramètres pour détecter des situations de flou
    - La simple forme de surface "on" = un de ces paramètres
    - Mais il en faut d'autres pour l'ensemble des expressions référentielles potentiellement floues

## Conclusion

- **Spéculations issues d'un travail de réflexions**
- Flou référentiel : peu étudié
  - La méthodologie reste à bâtir
  - Linguistique théorique et descriptive → applications TAL
- **Perspectives :**
  - Typologie plus complète de la référence floue  
Pour **toutes** les expressions référentielles
  - Schéma + manuel d'annotation
    - 3 degrés de coréférence
    - fonction d'appartenance
  - Test de la procédure
    - Corpus
    - Exploitation linguistique / statistique / TAL

- Cabredo Hofherr, P. (2008). Les pronoms impersonnels humains - syntaxe et interprétation. *Modèles linguistiques XXIX (1)*, 57:35–56.
- Charaudeau, P. (1992). *Grammaire du sens et de l'expression*. Hachette.
- Charolles, M. (2002). *La référence et les expressions référentielles en français*. Editions Ophrys.
- Charolles, M. (2014). Annotation des expressions référentielles et profondeur de traitement. In Fossard, M. and Béguelin, M.-J., editors, *Nouvelles perspectives sur l'anaphore : points de vue linguistique, psycholinguistique et acquisitionnel*, number vol. 111 in Sciences pour la communication, pages 55–98. Peter Lang.
- Christianson, K., Hollingworth, A., Halliwell, J., and Ferreira, F. (2001). Thematic roles assigned along the garden path linger. *Cognitive Psychology*, 42(4):368–407.
- Corblin, F. (1995). *Les formes de reprise dans le discours. Anaphores et chaînes de référence*. Presses Universitaires de Rennes.
- Dubois, D. and Prade, H. (2012). *Possibility theory: an approach to computerized processing of uncertainty*. Springer Science & Business Media.
- Désoyer, A., Landragin, F., Tellier, I., Lefeuvre, A., and Antoine, J.-Y. (2015). Les coréférences à l'oral : une expérience d'apprentissage automatique sur le corpus ANCOR. 55(2):97–121.
- Ferreira, F., Bailey, K., and Ferraro, V. (2002). Good-enough representations in language comprehension. 11(1):11–15.
- Ferreira, F. and Stacey, J. (2000). The misinterpretation of passive sentences.
- Fløttum, K., Jonasson, K., and Norén, C. (2008). *On : Pronom à facettes*. Duculot.

- Frege, G. (1892). Über sinn und bedeutung. *Zeitschrift für Philosophie und philosophische Kritik*.
- Grevisse, M. and Goosse, A. (2002). *Le bon usage*. De Boeck/Duculot.
- Kleiber, G. (1994). *Anaphores et pronoms*. Duculot.
- Landragin, F. and Tanguy, N. (2014). Référence et coréférence du pronom indéfini on. 195(3):99.
- Rabatel, A. (2001). La valeur de "on" pronom indéfini/pronom personnel dans les perceptions représentées. 88(1):28–32.
- Recasens, M., Hovy, E., and Martí, A. (2011). Identity, non-identity, and near-identity : Addressing the complexity of coreference. 121:1138–1152.
- Riegel, M., Pellat, J.-C., and Rioul, R. (1994). *Grammaire méthodique du français*. Presses Universitaires de France - PUF.
- Sandfeld, K. (1970). *Syntaxe du français contemporain. I. Les pronoms*. Honoré Champion.
- Schneedecker, C. and Landragin, F. (2014). Les chaînes de référence : présentation. (195):3–22.
- TLFI (2002). Trésor de la langue française informatisé.
- Viollet, C. (1988). Mais qui est on ? 18(1):67–75.
- Zadeh, L. (1965). Fuzzy sets. *Information and Control*, 8(3):338–353.
- Zadeh, L. (1978). Fuzzy sets as a basis for a theory of possibility. *Fuzzy Sets and Systems*, 1:3–28.

Coréférence inclusive :

” Je vous ennuie, j’ennuie tout le monde avec ça, les enfants, on croit être intéressante. (Lagarce, Juste la fin du monde, 1990 : 14) “