Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

# Méthodologie de la recherche

## Cours de Traitement Automatique des Langues
## M1 DAC
## UPMC

Anne-Laure Ligozat

2015/2016

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

1. Ce qui vous sera demandé dans ce cours

2. Qu'est-ce qu'un article scientifique ?

3. Comment faire la synthèse bibliographique ?

4. Application

5. Thèmes et sujets d'application

**Ce qui vous sera demandé dans ce cours**
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

### À rendre

- Synthèse bibliographique
- Application

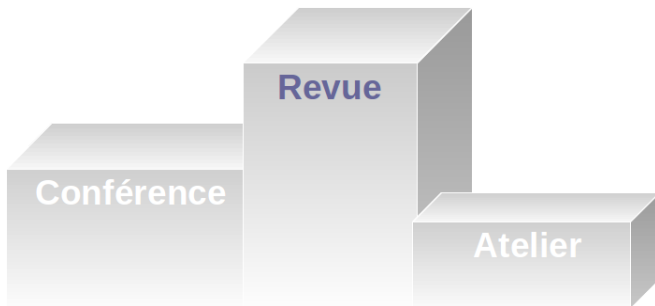- un thème à choisir (cf dernière section du cours)
- en binôme

### Objectifs du cours

- Scientifiques

  Découverte d'un **thème du TAL**

- Techniques

  Manipulation d'**outils de TAL**

- Méthodologiques

  Rédaction d'une **synthèse bibliographique**

Ce qui vous sera demandé dans ce cours
**Qu'est-ce qu'un article scientifique ?**
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

Où trouver les articles ?
Comment un article est-il organisé ?

Ce qui vous sera demandé dans ce cours
**Qu'est-ce qu'un article scientifique ?**
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

Où trouver les articles ?
Comment un article est-il organisé ?

# Nature des articles

- **Revue** (*journal*) internationale ou nationale
  - processus de sélection (plus ou moins) strict et long
  - environ 20/30 pages

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

Où trouver les articles ?
Comment un article est-il organisé ?

# Nature des articles

- **Conférence**
    - comité de lecture
    - annuelles ou bisannuelles
    - environ 10 pages

Ce qui vous sera demandé dans ce cours
**Qu'est-ce qu'un article scientifique ?**
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

Où trouver les articles ?
Comment un article est-il organisé ?

## Nature des articles

- **Atelier** (*workshop*), working notes de campagnes d'évaluation
  - intérêt : rencontre des spécialistes du domaines
  - description précise de systèmes

Ce qui vous sera demandé dans ce cours
**Qu'est-ce qu'un article scientifique ?**
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

Où trouver les articles ?
Comment un article est-il organisé ?

# Principales revues et conférences en TAL

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

Où trouver les articles ?
Comment un article est-il organisé ?

# Accès aux articles : Google scholar

Travaux de recherche de tous domaines

Ce qui vous sera demandé dans ce cours
**Qu'est-ce qu'un article scientifique ?**
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

**Où trouver les articles ?**
Comment un article est-il organisé ?

# Accès aux articles : anthologie ACL

Conférences et revues internationales en TAL

Ce qui vous sera demandé dans ce cours
**Qu'est-ce qu'un article scientifique ?**
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

Où trouver les articles ?
Comment un article est-il organisé ?

# Organisation d'un article

- Introduction
  - problème dans son **contexte** général
  - à quelle **problématique** les auteurs veulent-ils répondre ?
  - très rapide **état de l'art**
  - **apport** de l'article
  - **organisation** de l'article
- Corps de l'article
  - **état de l'art**
  - sinon, pas de format "type" en informatique, mais en général
    - **hypothèse**
    - algorithme/**méthode**
    - **résultats** et analyse d'erreurs
    - **comparaison** à l'existant
- Conclusion
  - **rappel** des idées fortes et résultats principaux
  - **discussion**
  - proposition de futures **directions** de recherche

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

Où trouver les articles ?
Comment un article est-il organisé ?

# Exemple d'article

## À vous de jouer : annotez un article

- Event Extraction as Dependency parsing
- David McClosky, Mihai Surdeanu, and Christopher D. Manning
- ACL
- 2011

Ce qui vous sera demandé dans ce cours
**Qu'est-ce qu'un article scientifique ?**
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

Où trouver les articles ?
Comment un article est-il organisé ?

Résumé de l'état de l'art

Problématique

Contexte

Apports

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

Où trouver les articles ?
Comment un article est-il organisé ?

Ce qui vous sera demandé dans ce cours
**Qu'est-ce qu'un article scientifique ?**
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

Où trouver les articles ?
**Comment un article est-il organisé ?**

Méthode

Ce qui vous sera demandé dans ce cours
**Qu'est-ce qu'un article scientifique ?**
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

Où trouver les articles ?
Comment un article est-il organisé ?

Résultats

Analyse d'erreurs

Ce qui vous sera demandé dans ce cours
**Qu'est-ce qu'un article scientifique ?**
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

Où trouver les articles ?
Comment un article est-il organisé ?

Discussion

Pistes

Rappel

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
**Comment faire la synthèse bibliographique ?**
Application
Thèmes et sujets d'application

Qu'est-ce qu'une synthèse ?
Outils pour la gestion de la bibliographie

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
**Comment faire la synthèse bibliographique ?**
Application
Thèmes et sujets d'application

Qu'est-ce qu'une synthèse ?
Outils pour la gestion de la bibliographie

# Objectif d'une synthèse

### Idée

Présenter le domaine en donnant les **principaux axes** des travaux effectués

- $=$ **état de l'art** ($\neq$ analyse de l'existant)
- présenter également les **limites** des méthodes actuelles

### Attention: théorique vs technique

Ex : formalisme d'analyse syntaxique

- structure en dépendances ou en constituants $=$ théorique
- arbre ou format XML ou parenthésé $=$ technique

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
**Comment faire la synthèse bibliographique ?**
Application
Thèmes et sujets d'application

Qu'est-ce qu'une synthèse ?
Outils pour la gestion de la bibliographie

# Aspects techniques vs théoriques : exemple

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
**Comment faire la synthèse bibliographique ?**
Application
Thèmes et sujets d'application

Qu'est-ce qu'une synthèse ?
Outils pour la gestion de la bibliographie

# Plan d'une synthèse (1/2)

### Introduction

- présenter le **sujet**
  - Ex : Cette synthèse aborde la problématique de l'extraction de relations.
- donner les **définitions** des termes principaux
  - Ex : extraction d'information, extraction de relations
- expliquer les **difficultés**
  - ambiguïtés de rattachement, variations...
- présenter l'**historique** du domaine et les articles fondateurs

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
**Comment faire la synthèse bibliographique ?**
Application
Thèmes et sujets d'application

Qu'est-ce qu'une synthèse ?
Outils pour la gestion de la bibliographie

# Plan d'une synthèse (2/2)

## Corps

- organiser en fonction des **axes de recherche** actuels (et non du type d'application, des équipes...)
- citer des **travaux représentatifs** de chaque axe
- **évaluation** dans le domaine

## Conclusion

- **pistes** de recherche futures

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
**Comment faire la synthèse bibliographique ?**
Application
Thèmes et sujets d'application

Qu'est-ce qu'une synthèse ?
Outils pour la gestion de la bibliographie

## En détails

- Attention au **vocabulaire** utilisé: reprendre les termes du domaine (quitte à se répéter)
  - traduire les termes : tokenization = segmentation en mots
  - être précis : terme $\neq$ mot $\neq$ entité
- Donner des **exemples**
- Donner des **résultats** de systèmes (beaucoup de campagnes d'évaluation en TAL)
- Attention au format des **références**

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
**Comment faire la synthèse bibliographique ?**
Application
Thèmes et sujets d'application

Qu'est-ce qu'une synthèse ?
Outils pour la gestion de la bibliographie

# Exemple de plan de synthèse

Extraction de relations

- Introduction
- Approche à base de patrons
    - surfaciques
    - syntaxiques
    - acquisition automatique de patrons
- Approche par apprentissage
    - avec attributs vectoriels
    - sur structure arborescente
- Évaluation
- Conclusion

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
**Comment faire la synthèse bibliographique ?**
Application
Thèmes et sujets d'application

Qu'est-ce qu'une synthèse ?
Outils pour la gestion de la bibliographie

# Bibliographie : références

- Citer ses sources !
- Les citer correctement
  - Éléments indispensables dans la référence :
    - **titre** de l'article
    - noms des **auteurs**
    - titre de la **ressource** : nom de la conférence ou de la revue
    - **année** de publication
  - Exemple :
    - Anja Belz, Michael White, Dominic Espinosa, Eric Kow, Deirdre Hogan, and Amanda Stent. 2011. The first surface realisation shared task: Overview and evaluation results. In Proceedings of the 13th European Workshop on Natural Language Generation (ENLG), Nancy, France

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
**Comment faire la synthèse bibliographique ?**
Application
Thèmes et sujets d'application

Qu'est-ce qu'une synthèse ?
Outils pour la gestion de la bibliographie

# Bibliographie : format ACL

- Citations dans le texte
  - Citations within the text appear in parentheses as (Gusfield, 1997) or, if the author's name appears in the text itself, as Gusfield (1997). Append lowercase letters to the year in cases of ambiguity. Treat double authors by using both authors' last names (e.g., (Aho and Ullman, 1972), but use et al. when more than two authors are involved. (e.g. (Chandra et al., 1981)) Collapse multiple citations (e.g., (Gusfield, 1997; Aho and Ullman, 1972).)

- Dans la partie Références
  - Alfred V. Aho and Jeffrey D. Ullman. 1972. The Theory of Parsing, Translation and Compiling, volume 1. Prentice-Hall, Englewood Cliffs, NJ.
  - Ashok K. Chandra, Dexter C. Kozen, and Larry J. Stockmeyer. 1981. Alternation. Journal of the Association for Computing Machinery, 28(1):114–133.
  - Dan Gusfield. 1997. Algorithms on Strings, Trees and Sequences. Cambridge University Press, Cambridge, UK.

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
**Comment faire la synthèse bibliographique ?**
Application
Thèmes et sujets d'application

Qu'est-ce qu'une synthèse ?
Outils pour la gestion de la bibliographie

# BibTex : format des références

## BibTex : gestion et traitement des données bibliographiques

forme : clé = valeur

```
@inproceedings{LanglaisEACL2009,
  author = {Langlais, Philippe and Yvon, François and
            Zweigenbaum, Pierre},
  title = {Improvements in Analogical Learning:
           Application to Translating Multi-Terms
           of the Medical Domain},
  booktitle = {Proceedings of the 12th Conference of
               the European Chapter of the Association
               for Computational Linguistics (EACL 2009)},
  publisher = {Association for Computational Linguistics},
  year = {2009},
  pages = {487-495},
  url = {http://www.aclweb.org/anthology/E09-1056}
}
```

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
**Comment faire la synthèse bibliographique ?**
Application
Thèmes et sujets d'application

Qu'est-ce qu'une synthèse ?
Outils pour la gestion de la bibliographie

# BibTex : utilisation dans LaTeX (pdflatex)

- dans le .tex :
  \cite{LanglaisEACL2009} montrent qu'il est également possible de traduire...
- dans le .pdf :
  (Langlais et al., 2009) montrent qu'il est également possible de traduire..

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
**Comment faire la synthèse bibliographique ?**
Application
Thèmes et sujets d'application

Qu'est-ce qu'une synthèse ?
**Outils pour la gestion de la bibliographie**

# JabRef

- gestion graphique des références bibliographiques
- exports BibTeX, texte, OpenOffice (plugin)...

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
**Comment faire la synthèse bibliographique ?**
Application
Thèmes et sujets d'application

Qu'est-ce qu'une synthèse ?
Outils pour la gestion de la bibliographie

# Caractéristiques des références

- références anciennes si besoin (articles fondateurs) et récentes (cinq dernières années)
- articles d'auteurs, laboratoires et pays variés

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
**Application**
Thèmes et sujets d'application

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
**Application**
Thèmes et sujets d'application

# Objectif de l'application

- Objectif : aborder le domaine du TAL de la synthèse d'un point de vue pratique
- Application = mise en œuvre d'un algorithme, évaluation/comparaison d'outils/méthodes...

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
**Application**
Thèmes et sujets d'application

## Choix à faire

- **Corpus de test**
  - quelle source ? quel format ? prétraitement nécessaires ? → cette année, corpus unique
- **Entrée/sorties** du système
  - quel format en entrée, quel format en sortie ?
- **Tests** de l'application
  - cas nominal, limites, erreurs
- Mode d'**évaluation**
  - choisir dès le départ une évaluation standard (ex : rappel, précision, f-mesure)

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
**Application**
Thèmes et sujets d'application

# Rapport sur l'application

- **Expliciter** et **justifier** les choix
- Donner des **exemples** précis des résultats du système

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
**Thèmes et sujets d'application**

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

# Analyse morphologique

- notions : lemmatisation, stemming, familles morphologiques
- articles : étiquetage du français et construction de familles morphologiques
- application : Comparer stemming, lemmatisation et familles morphologiques (en utilisant les derivationally related forms de Wordnet ou avec un outil de segmentation morphologique comme Morfessor) pour de la sélection de passages répondant à des questions

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
Thèmes et sujets d'application

# Terminologie : extraction de termes et de collocations

- notions : termes (Multi Word Unit - MWU ou MultiWord Expression - MWE), mesures de cooccurences/collocations
- articles : reconnaissance d'acronymes, reconnaissance de termes
- application : Constitution automatique d'une base d'acronymes avec leur signification et annotation des acronymes dans les textes; évaluation de l'apport à la recherche de passages

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
**Thèmes et sujets d'application**

## Terminologie, variations de termes

- notions : termes, expressions, variations linguistiques
- articles : reconnaissance de variantes, validation de relations entre termes
- application : Validation en contexte de variations morpho-sémantiques de mots (en utilisant les informations de synonymie et morphologiques de Wordnet) à partir des cooccurrents (issus par exemple de la base de cooccurrences Wortschatz)

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
**Thèmes et sujets d'application**

## Variations, paraphrase

- notions : paraphrase, implication textuelle
- articles : reconnaissance de paraphrases, implication textuelle
- application : Utilisation d'un outil d'implication textuelle (exemple : EDITS, ou BIUTEE) ou une banque de paraphrases (exemple : PPDB) et évaluation sur QA4MRE

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
**Thèmes et sujets d'application**

# Entités nommées

- notions : reconnaissance d'entités nommées, désambiguïsation et résolution d'entités nommées
- articles : typage non supervisé d'entités; résolution d'entités nommées
- application : Suivi d'entités nommées

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
**Thèmes et sujets d'application**

# Analyse syntaxique

- notions : analyse syntaxique en constituants et en dépendances, arbres syntaxiques
- articles : génération de questions, analyse syntaxique du français, correction d'analyse morpho-syntaxique
- application 1 : Génération d'hypothèses et validation
- application 2 : Correction d'analyse syntaxique de questions

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
**Thèmes et sujets d'application**

# Sémantique : synonymie, structuration des connaissances

- notions : sens, ambiguité, relation sémantiques, évaluation, construction de ressources
- articles : synonymie, construction de ressource

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
**Thèmes et sujets d'application**

# Anaphore et coréférence

- notions : coréférence
- article : apprentissage et coreference; anaphore
- application : Évaluation d'un système d'annotation de coréférence et analyse du corpus

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
**Thèmes et sujets d'application**

## Analyse thématique

- notions : cohésion lexicale, distribution des mots dans le texte et dans blocs, segmentation thématique
- article : segmentation par ressources segmentation thématique
- application : Étude de la segmentation thématique pour la sélection de passages (2 segmenteurs)

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
**Thèmes et sujets d'application**

# Résumé automatique

- notions : résumé par extraction, critères de sélection de phrases importantes
- articles : résumé par ordonnancement résumé multidocuments
- application : Identification de thèmes

Ce qui vous sera demandé dans ce cours
Qu'est-ce qu'un article scientifique ?
Comment faire la synthèse bibliographique ?
Application
**Thèmes et sujets d'application**

# Analyse du discours

- notions : relations du discours, structure logique
- articles : analyse par règles apprentissage de relations implicites
- application 1 : Segmentation automatique en phrases et reconnaissance des titres
- application 2 : Reconnaissance de relations du discours