

# Vers une simplification automatique de la parole en français

---

*Les enjeux de l'extraction des données d'apprentissage  
pour la simplification linguistique*



**UNIVERSITÉ  
DE GENÈVE**

**FACULTÉ DE TRADUCTION  
ET D'INTERPRÉTATION**



**Lucía Ormaechea**  
*Doctorante, assistante*  
Université de Genève et  
Université Grenoble-Alpes

---

*Colloque International l'Association for French Language Studies (AFLS)*

*8 septembre de 2023, Université de Lille, Villeneuve-d'Ascq*

# Plan



## 1. Introduction

- 1.1. Qu'est-ce que la simplification automatique de textes ?
- 1.2. SAT en français : une tâche à faibles ressources
- 1.3. *Bridging the gap* : une nouvelle méthode pour extraire des paires *complexes-simples*

## 2. Acquisition de données : Wikipédia et Vikidia

## 3. Méthode de filtrage en deux étapes pour la SAT

- 3.1. Étape I : Préservation du sens original
- 3.2. Étape II : Gain en simplicité linguistique

## 4. Résultats : *Wikipedia-Vikidia Corpus* (WiViCo)

## 5. Conclusions

# 1. Introduction

- 1.1. Qu'est-ce que la simplification automatique de textes ?
- 1.2. SAT en français : une tâche à faibles ressources
- 1.3. *Bridging the gap* : une nouvelle méthode pour extraire des paires *complexes-simples*

La Simplification Automatique de Textes (SAT) est un domaine du TAL qui vise à **réduire** automatiquement la **complexité linguistique** des textes, **sans pour autant perdre les informations** et la **signification originales** [Horn et al., 2014; Stajner, 2021].

ORIGINAL	La Ligue nationale de hockey ou LNH est une <b>association sportive</b> professionnelle <del>nord-américaine regroupant des franchises</del> de hockey sur glace <b>du Canada et des États-Unis</b> .
SIMPLIFIÉ	La Ligue nationale de hockey ou <b>tout simplement</b> LNH est une <b>ligue sportive</b> professionnelle de hockey sur glace <b>d'Amérique du Nord</b> .

# 1. Introduction

- 1.1. Qu'est-ce que la simplification automatique de textes ?
- 1.2. SAT en français : une tâche à faibles ressources
- 1.3. *Bridging the gap* : une nouvelle méthode pour extraire des paires complexes-simples

La Simplification Automatique de Textes (SAT) est un domaine du TAL qui vise à **réduire** automatiquement la **complexité linguistique** des textes, **sans pour autant perdre les informations** et la **signification originales** [Horn et al., 2014; Stajner, 2021].

ORIGINAL	La Ligue nationale de hockey ou LNH est une <b>association sportive</b> professionnelle <del>nord-américaine regroupant des franchises</del> de hockey sur glace <b>du Canada et des États-Unis</b> .
SIMPLIFIÉ	La Ligue nationale de hockey ou <b>tout simplement</b> LNH est une <b>ligue sportive</b> professionnelle de hockey sur glace <b>d'Amérique du Nord</b> .

SAT : une analogie avec la relation *forme-substance* :



FORME

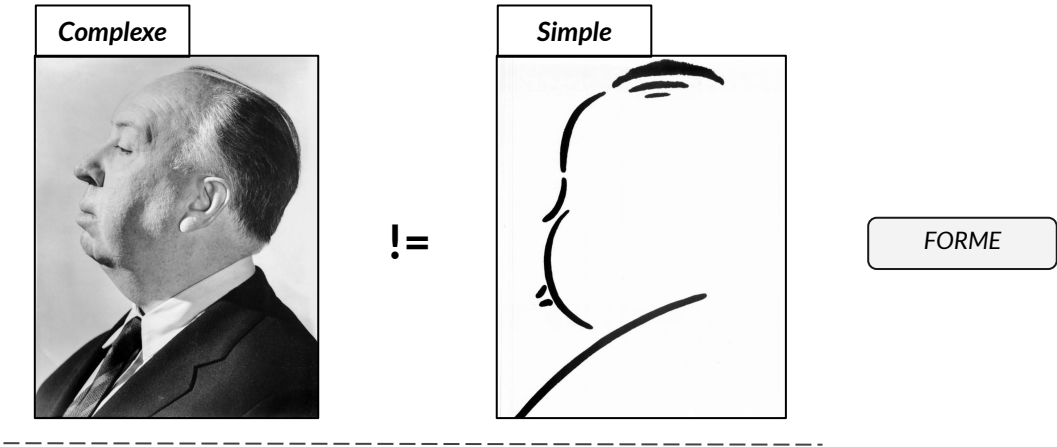
# 1. Introduction

- 1.1. Qu'est-ce que la simplification automatique de textes ?
- 1.2. SAT en français : une tâche à faibles ressources
- 1.3. *Bridging the gap* : une nouvelle méthode pour extraire des paires complexes-simples

La Simplification Automatique de Textes (SAT) est un domaine du TAL qui vise à **réduire** automatiquement la **complexité linguistique** des textes, **sans pour autant perdre les informations** et la **signification originales** [Horn et al., 2014; Stajner, 2021].

ORIGINAL	La Ligue nationale de hockey ou LNH est une <b>association sportive</b> professionnelle <del>nord-américaine regroupant des franchises</del> de hockey sur glace <b>du Canada et des États-Unis</b> .
SIMPLIFIÉ	La Ligue nationale de hockey ou <b>tout simplement</b> LNH est une <b>ligue sportive</b> professionnelle de hockey sur glace <b>d'Amérique du Nord</b> .

SAT : une analogie avec la relation *forme-substance* :



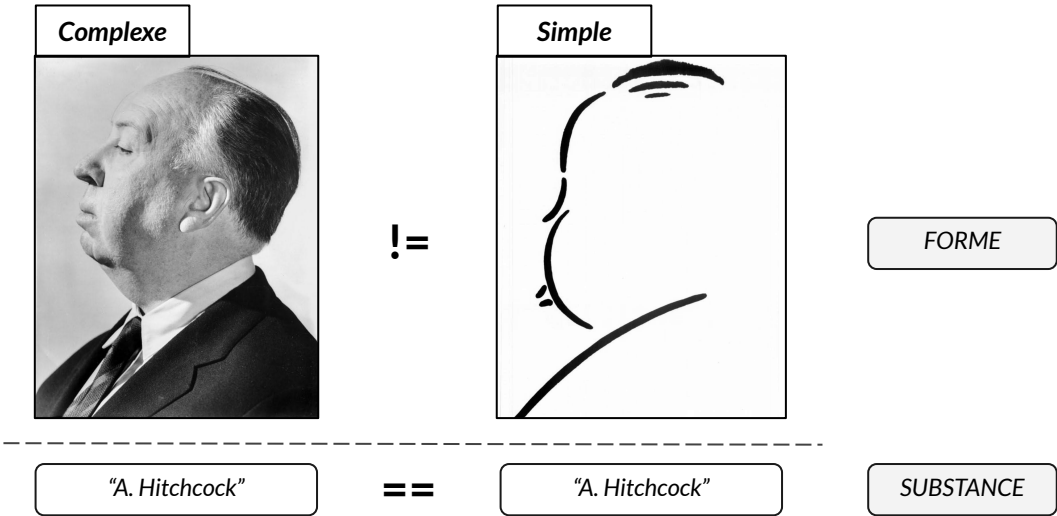
# 1. Introduction

- 1.1. Qu'est-ce que la simplification automatique de textes ?
- 1.2. SAT en français : une tâche à faibles ressources
- 1.3. *Bridging the gap* : une nouvelle méthode pour extraire des paires complexes-simples

La Simplification Automatique de Textes (SAT) est un domaine du TAL qui vise à **réduire** automatiquement la **complexité linguistique** des textes, **sans pour autant perdre les informations** et la **signification originales** [Horn et al., 2014; Stajner, 2021].

ORIGINAL	La Ligue nationale de hockey ou LNH est une <b>association sportive</b> professionnelle <del>nord-américaine regroupant des franchises</del> de hockey sur glace <del>du Canada et des États-Unis</del> .
SIMPLIFIÉ	La Ligue nationale de hockey ou <b>tout simplement</b> LNH est une <b>ligue sportive</b> professionnelle de hockey sur glace <del>d'Amérique du Nord</del> .

SAT : une analogie avec la relation *forme-substance* :



# 1. Introduction

- 1.1. Qu'est-ce que la simplification automatique de textes ?
- 1.2. SAT en français : une tâche à faibles ressources
- 1.3. *Bridging the gap* : une nouvelle méthode pour extraire des paires *complexes-simples*



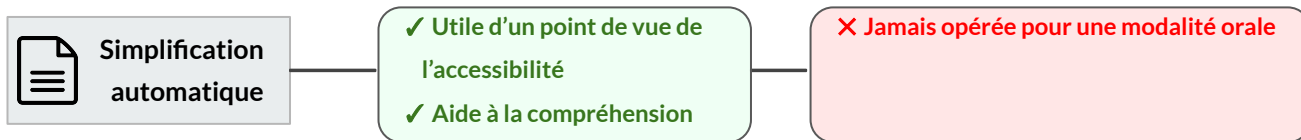
Simplification  
automatique

- ✓ Utile d'un point de vue de l'accessibilité
- ✓ Aide à la compréhension

ORIGINAL	La Ligue nationale de hockey ou LNH est une <b>association sportive</b> professionnelle <del>nord-américaine regroupant des franchises</del> de hockey sur glace <b>du Canada et des États-Unis</b> .
SIMPLIFIÉ	La Ligue nationale de hockey ou <b>tout simplement</b> LNH est une <b>ligue sportive</b> professionnelle de hockey sur glace <b>d'Amérique du Nord</b> .

# 1. Introduction

- 1.1. Qu'est-ce que la simplification automatique de textes ?
- 1.2. SAT en français : une tâche à faibles ressources
- 1.3. *Bridging the gap* : une nouvelle méthode pour extraire des paires *complexes-simples*

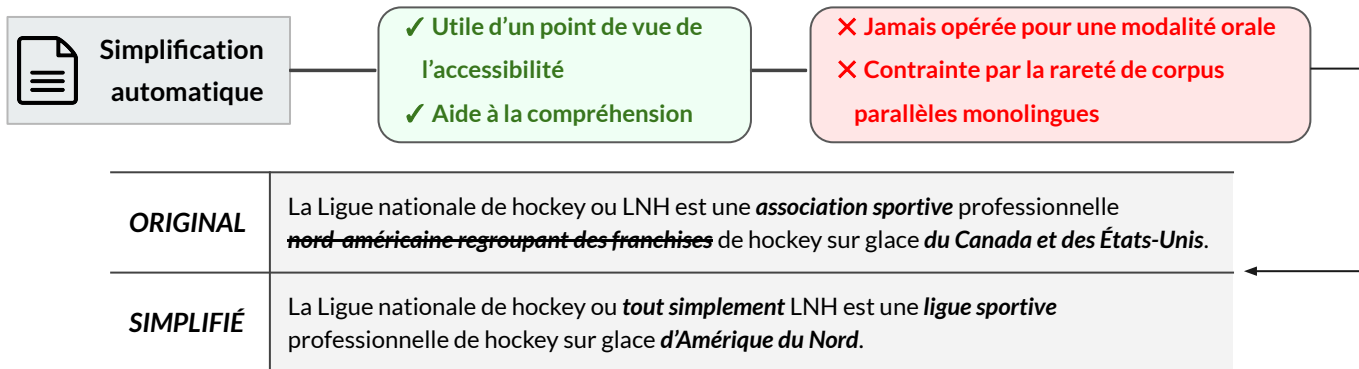


ORIGINAL	La Ligue nationale de hockey ou LNH est une <b>association sportive</b> professionnelle <del>nord-américaine regroupant des franchises</del> de hockey sur glace <b>du Canada et des États-Unis</b> .
SIMPLIFIÉ	La Ligue nationale de hockey ou <b>tout simplement</b> LNH est une <b>ligue sportive</b> professionnelle de hockey sur glace <b>d'Amérique du Nord</b> .



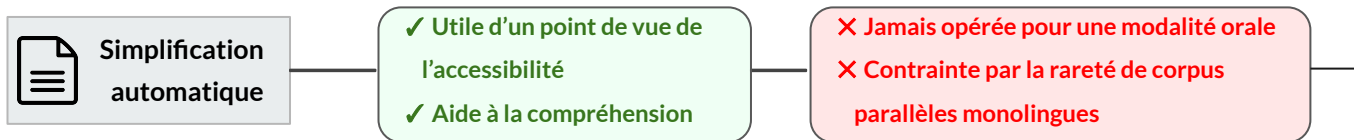
# 1. Introduction

- 1.1. Qu'est-ce que la simplification automatique de textes ?
- 1.2. SAT en français : une tâche à faibles ressources
- 1.3. *Bridging the gap* : une nouvelle méthode pour extraire des paires *complexes-simples*



# 1. Introduction

- 1.1. Qu'est-ce que la simplification automatique de textes ?
- 1.2. SAT en français : une tâche à faibles ressources
- 1.3. *Bridging the gap* : une nouvelle méthode pour extraire des paires *complexes-simples*



ORIGINAL	La Ligue nationale de hockey ou LNH est une <b>association sportive</b> professionnelle <del>nord-américaine regroupant des franchises</del> de hockey sur glace <b>du Canada et des États-Unis</b> .
SIMPLIFIÉ	La Ligue nationale de hockey ou <b>tout simplement</b> LNH est une <b>ligue sportive</b> professionnelle de hockey sur glace <b>d'Amérique du Nord</b> .

Dans des langues moins dotées que l'anglais :

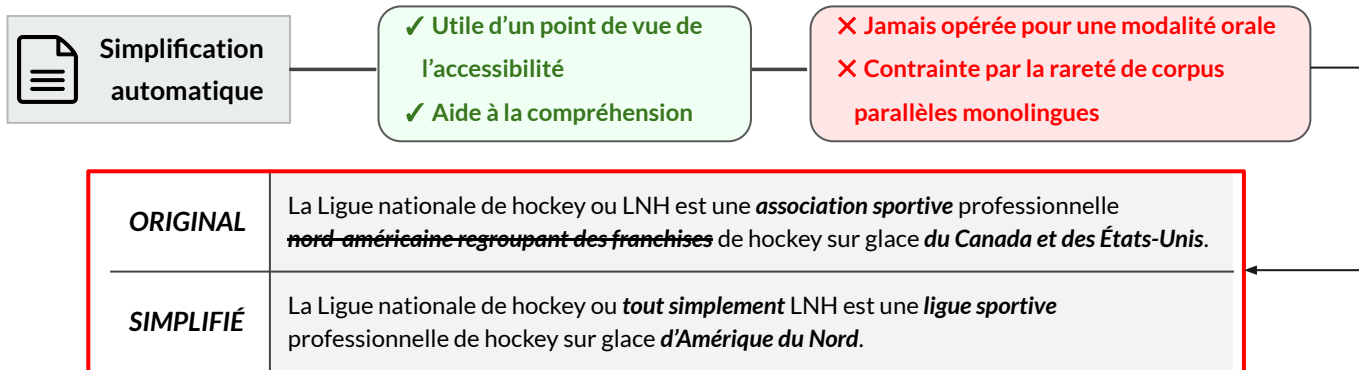


Le seul jeu de données monolingue disponible → **ALECTOR** [Gala et al., 2020]:

- Créé manuellement par des éditeurs professionnels.
- Comprenant des textes scolaires, **simplifiés** pour un **public enfant**.

# 1. Introduction

- 1.1. Qu'est-ce que la simplification automatique de textes ?
- 1.2. SAT en français : une tâche à faibles ressources
- 1.3. *Bridging the gap* : une nouvelle méthode pour extraire des paires *complexes-simples*



Dans des langues moins dotées que l'anglais :



Le seul jeu de données monolingue disponible → **ALECTOR** [Gala et al., 2020]:

- **Créé manuellement** par des éditeurs professionnels.
- Comprenant des textes scolaires, **simplifiés** pour un **public enfant**.

**Jeux de données créés manuellement**

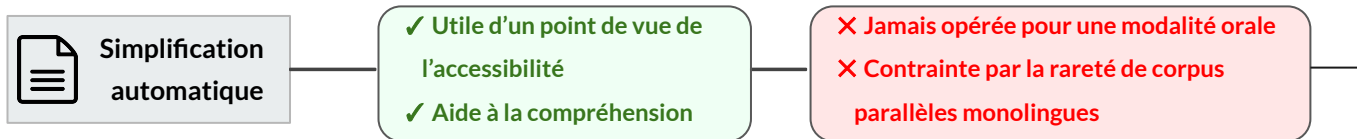
✓ **Simplifications fiables**

✗ **Création coûteuse**

✗ **Très petite taille**

# 1. Introduction

- 1.1. Qu'est-ce que la simplification automatique de textes ?
- 1.2. SAT en français : une tâche à faibles ressources
- 1.3. *Bridging the gap* : une nouvelle méthode pour extraire des paires *complexes-simples*



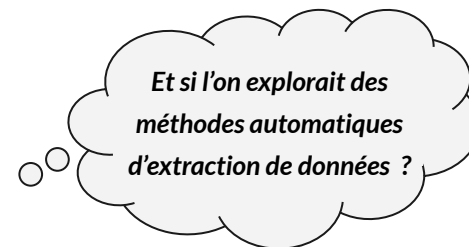
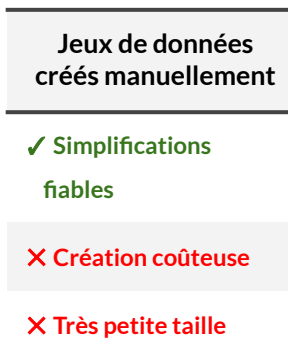
ORIGINAL	La Ligue nationale de hockey ou LNH est une <b>association sportive</b> professionnelle <del>nord-américaine regroupant des franchises</del> de hockey sur glace <del>du Canada et des États-Unis</del> .
SIMPLIFIÉ	La Ligue nationale de hockey ou <b>tout simplement</b> LNH est une <b>ligue sportive</b> professionnelle de hockey sur glace <b>d'Amérique du Nord</b> .

Dans des langues moins dotées que l'anglais :



Le seul jeu de données monolingue disponible → ALECTOR [Gala et al., 2020]:

- Créé manuellement par des éditeurs professionnels.
- Comprenant des textes scolaires, **simplifiés** pour un **public enfant**.



# 1. Introduction

- 1.1. Qu'est-ce que la simplification automatique de textes ?
- 1.2. SAT en français : une tâche à faibles ressources
- 1.3. *Bridging the gap* : une nouvelle méthode pour extraire des paires *complexes-simples*

## PRINCIPALE CONTRIBUTION

- Présentation d'une **nouvelle méthode** pour **extraire** des paires de **phrases complexes-simples** à partir des **corpus comparables**.
- Conçue spécifiquement pour la **tâche en aval**.



*“Création d’un corpus parallèle en français  
de paires de phrases complexes-simples”*

# 1. Introduction

- 1.1. Qu'est-ce que la simplification automatique de textes ?
- 1.2. SAT en français : une tâche à faibles ressources
- 1.3. *Bridging the gap* : une nouvelle méthode pour extraire des paires *complexes-simples*

## PRINCIPALE CONTRIBUTION

- Présentation d'une **nouvelle méthode** pour **extraire** des paires de **phrases complexes-simples** à partir des **corpus comparables**.
- Conçue spécifiquement pour la **tâche en aval**.



*“Création d'un corpus parallèle en français  
de paires de phrases complexes-simples”*

## Approches automatiques pour la compilation des données en SAT :

- Basées sur des **corpus comparables** ayant des **registres distincts** :
  - **Wikipédia** → Édition française.
  - **Vikidia** → Version adaptée et simplifiée de la première.



# 1. Introduction

- 1.1. Qu'est-ce que la simplification automatique de textes ?
- 1.2. SAT en français : une tâche à faibles ressources
- 1.3. *Bridging the gap* : une nouvelle méthode pour extraire des paires *complexes-simples*

## PRINCIPALE CONTRIBUTION

- Présentation d'une **nouvelle méthode** pour **extraire** des paires de **phrases complexes-simples** à partir des **corpus comparables**.
- Conçue spécifiquement pour la **tâche en aval**.



*“Création d'un corpus parallèle en français de paires de phrases complexes-simples”*

## Approches automatiques pour la compilation des données en SAT :

- Basées sur des **corpus comparables** ayant des **registres distincts** :
  - **Wikipédia** → Édition française.
  - **Vikidia** → Version adaptée et simplifiée de la première.
- Normalement fondés sur des **mesures de similarité sémantique**, mais **ignorent** si la **phrase cible** constitue effectivement une **simplification**.



# 1. Introduction

- 1.1. Qu'est-ce que la simplification automatique de textes ?
- 1.2. SAT en français : une tâche à faibles ressources
- 1.3. *Bridging the gap* : une nouvelle méthode pour extraire des paires *complexes-simples*

## PRINCIPALE CONTRIBUTION

- Présentation d'une **nouvelle méthode** pour **extraire** des paires de **phrases complexes-simples** à partir des **corpus comparables**.
- Conçue spécifiquement pour la **tâche en aval**.



*“Création d'un corpus parallèle en français de paires de phrases complexes-simples”*

## Approches automatiques pour la compilation des données en SAT :

- Basées sur des **corpus comparables** ayant des **registres distincts** :
  - **Wikipédia** → Édition française.
  - **Vikidia** → Version adaptée et simplifiée de la première.
- Normalement fondés sur des **mesures de similarité sémantique**, mais **ignorent** si la **phrase cible** constitue effectivement une **simplification**.



## C'est pourquoi, nous proposons :

Étapes	Deux conditions à respecter	Examinées de manière séquentielle
1	Préservation du sens original	Filtrage par similarité sémantique
2	Gain en simplicité linguistique	Filtrage par gain de simplicité



## 2. Acquisition de données : Wikipédia et Vikidia



*“Obtenir des données différenciées par registre de langue”*

Extraction du contenu textuel des deux **encyclopédies** :

- Récupération des textes **standards** (Wikipédia) et leurs versions **simplifiées** (Vikidia).
- Accent mis sur les **lead sections** (résumés introductifs).

## 2. Acquisition de données : Wikipédia et Vikidia



*“Obtenir des données différenciées par registre de langue”*

Extraction du contenu textuel des deux **encyclopédies** :

- Récupération des textes **standards** (Wikipédia) et leurs versions **simplifiées** (Vikidia).
- Accent mis sur les **lead sections** (résumés introductifs).
- Sous l'hypothèse que → **Plus de chances de trouver des phrases alignées** :

1. Tous deux partagent un **style définitionnel**.
2. Les **extraits communs** sont **plus susceptibles d'apparaître au début**.

## 2. Acquisition de données : Wikipédia et Vikidia

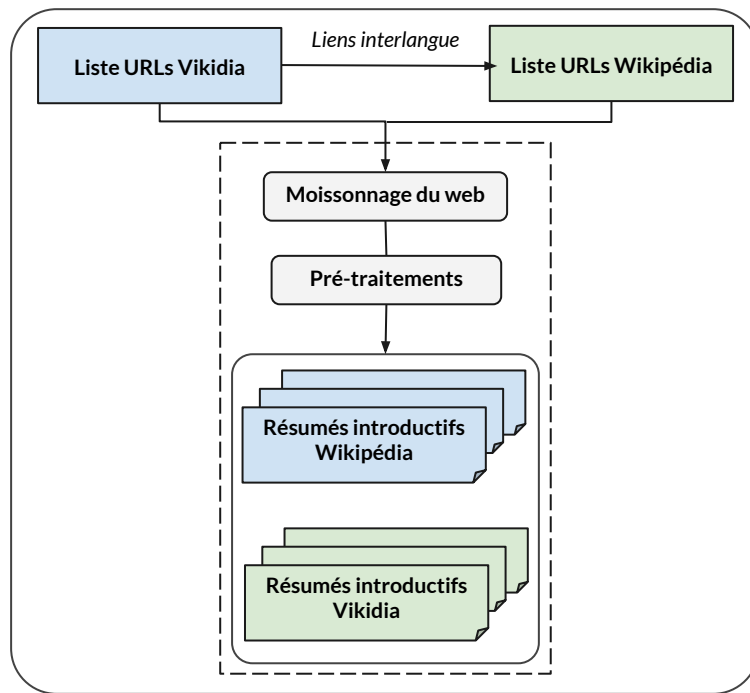


*“Obtenir des données différenciées par registre de langue”*

Extraction du contenu textuel des deux **encyclopédies** :

- Récupération des textes **standards** (Wikipédia) et leurs versions **simplifiées** (Vikidia).
- Accent mis sur les **lead sections** (résumés introductifs).
- Sous l'hypothèse que → **Plus de chances de trouver des phrases alignées** :

1. Tous deux partagent un **style définitionnel**.
2. Les **extraits communs** sont **plus susceptibles d'apparaître au début**.

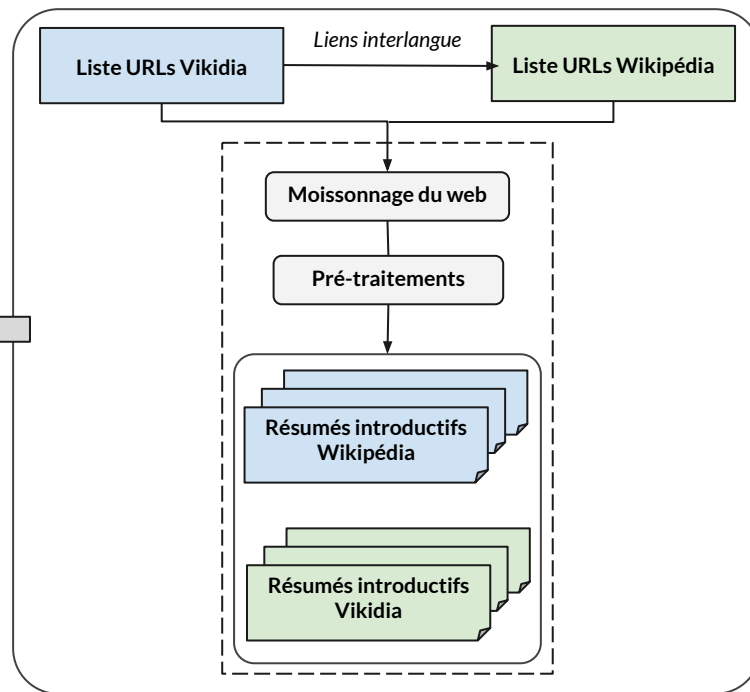


## 2. Acquisition de données : Wikipédia et Vikidia



*“Obtenir des données différenciées par registre de langue”*

Acquisition de données	Wiki-textes	Viki-textes
# documents	34,806	
# phrases	165,806	134,348
# tokens	4,030,148	2,373,045
# types	294,979	195,791
Type/token ratio	7.32	8.25
Longueur moyenne mot	5.32	5.08
Longueur moyenne phrase	25.22	18.16



### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Identifier les paires complexes-simples qui sont appropriées pour la SAT”*

Acquisition de données	Wiki-textes	Viki-textes
# documents	34,806	
# phrases	165,806	134,348
# tokens	4,030,148	2,373,045
...		

#### Alignement automatique des phrases :

Utilisation de *multilingual sentence transformers* (SBERT) <sup>\*</sup> :

- Générant des **vecteurs de sens** à taille fixe, au niveau de la **phrase**.
- **Similarité cosinus** entre les phrases **Wiki:Viki** [-1,1].

<sup>\*</sup> <https://huggingface.co/sentence-transformers/paraphrase-xlm-r-multilingual-v1>



### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Identifier les paires complexes-simples qui sont appropriées pour la SAT”*

Mais quelles paires sont suffisamment cohérentes d'un point de vue sémantique ?

Acquisition de données	Wiki-textes	Viki-textes
# documents	34,806	
# phrases	165,806	134,348
# tokens	4,030,148	2,373,045
...		

#### Alignement automatique des phrases :

Utilisation de *multilingual sentence transformers* (SBERT) :

- Générant des **vecteurs de sens** à taille fixe, au niveau de la **phrase**.
- **Similarité cosinus** entre les phrases **Wiki:Viki** [-1,1].

\* <https://huggingface.co/sentence-transformers/paraphrase-xlm-r-multilingual-v1>

### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Extraire les paires Wiki et Viki qui présentent un chevauchement sémantique élevé”*

#### 3 exemples de paires de phrases Wiki:Viki :

Label	Wiki-texte	Viki-texte
Valable	<i>L'expression « Maison-Blanche » est souvent employée pour désigner, par métonymie, l'administration du président.</i>	<i>Par métonymie, la Maison-Blanche désigne aussi le gouvernement américain et son entourage.</i>
Partiellement valable	<i>Neal McDonough est un acteur et producteur américain né le 13 février 1966 à Dorchester (Massachusetts).</i>	<i>Neal McDonough est un acteur américain.</i>
Non valable	<i>L'information désigne à la fois le message à communiquer et les symboles utilisés pour l'écrire.</i>	<i>Les écrits, les sons, les images, les odeurs ou les goûts contiennent de l'information.</i>

### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Extraire les paires Wiki et Viki qui présentent un chevauchement sémantique élevé”*

#### 3 exemples de paires de phrases Wiki:Viki :

Label	Wiki-texte	Viki-texte
Valable	<i>L'expression « Maison-Blanche » est souvent employée pour désigner, par métonymie, l'administration du président.</i>	<i>Par métonymie, la Maison-Blanche désigne aussi le gouvernement américain et son entourage.</i>
Partiellement valable	<i>Neal McDonough est un acteur <b>et producteur</b> américain <b>né le 13 février 1966 à Dorchester (Massachusetts).</b></i>	<i>Neal McDonough est un acteur américain.</i>
Non valable	<i>L'information désigne à la fois le message à communiquer et les symboles utilisés pour l'écrire.</i>	<i>Les écrits, les sons, les images, les odeurs ou les goûts contiennent de l'information.</i>

### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Extraire les paires Wiki et Viki qui présentent un chevauchement sémantique élevé”*

#### 3 exemples de paires de phrases Wiki:Viki :

Label	Wiki-texte	Viki-texte
Valable	<i>L'expression « Maison-Blanche » est souvent employée pour désigner, par métonymie, l'administration du président.</i>	<i>Par métonymie, la Maison-Blanche désigne aussi le gouvernement américain et son entourage.</i>
Partiellement valable	<i>Neal McDonough est un acteur et producteur américain né le 13 février 1966 à Dorchester (Massachusetts).</i>	<i>Neal McDonough est un acteur américain.</i>
Non valable	<i>L'information <b>désigne à la fois le message à communiquer et les symboles utilisés pour l'écrire.</b></i>	<i>Les écrits, les sons, les images, les odeurs ou les goûts contiennent de l'information.</i>



### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Extraire les paires Wiki et Viki qui présentent un chevauchement sémantique élevé”*

#### Recours à une annotation manuelle (1/2) :

- Sélection aléatoire de **500 paires de phrases** de l'ensemble de données initial.
- **Deux annotateurs** ont évalué dans quelle mesure chaque paire de phrases **Wiki:Viki** véhiculait le même sens :
  1. **Valable** → L'information et le sens véhiculés par **Wiki** et **Viki** sont équivalents.
  2. **Partiellement valable** → Des informations sont partiellement perdues de **Wiki** à **Viki** ou *vice versa*.
  3. **Non-valable** → Les informations entre **Wiki** et **Viki** sont divergentes.
- **Accord inter-annotateur élevé** (Cohen's kappa = **0.87**).

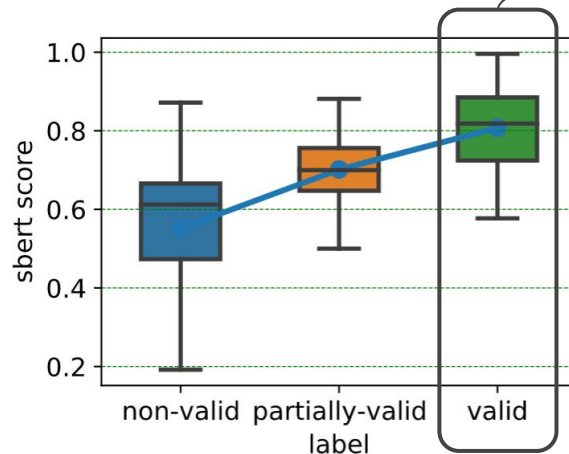
### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Extraire les paires Wiki et Viki qui présentent un chevauchement sémantique élevé”*

Recours à une annotation manuelle (2/2) :



Score moyen pour le label valide = 0.81

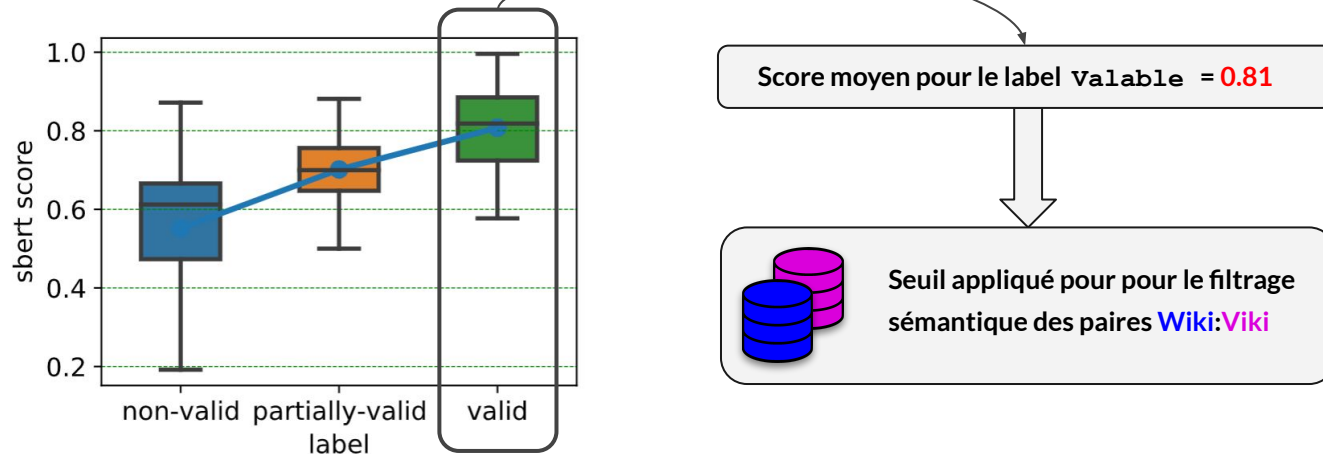
### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Extraire les paires Wiki et Viki qui présentent un chevauchement sémantique élevé”*

Recours à une annotation manuelle (2/2) :



### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique

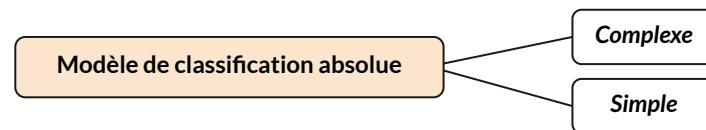




### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique

Très souvent, la simplicité est évaluée automatiquement avec...



### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique

Très souvent, la simplicité est évaluée automatiquement avec...

Modèle de classification absolue

Complexe

Simple

✗ Approche insuffisante pour la SAT

La **simplification** est une **opération intrinsèquement relative** → On transforme une phrase en une version **relativement plus simple** (ce qui n'équivaut pas nécessairement à **simple**).

### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique

Très souvent, la simplicité est évaluée automatiquement avec...

Modèle de classification absolue

Complexe

Simple

✗ Approche insuffisante pour la SAT

La **simplification** est une **opération intrinsèquement relative** → On transforme une phrase en une version **relativement plus simple** (ce qui n'équivaut pas nécessairement à **simple**).

C'est pourquoi :

- **Complémenter** deux modèles au classification absolu ; basés sur l'**affinage** des modèles **FLAUBERT** [Le et al., 2020].
- **Objectif** : proposer une **méthode plus fine** pour trouver des **paires de phrases** pertinentes pour la **SAT**.

### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique

Très souvent, la simplicité est évaluée automatiquement avec...

Modèle de classification absolue

Complexe

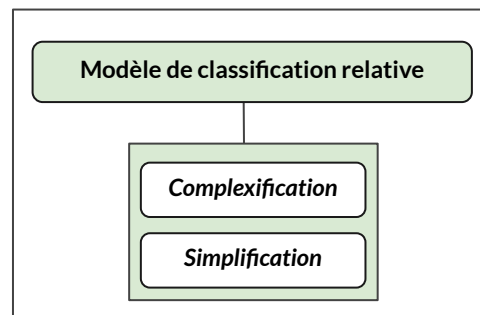
Simple

✗ Approche insuffisante pour la SAT

La **simplification** est une **opération intrinsèquement relative** → On transforme une phrase en une version **relativement plus simple** (ce qui n'équivaut pas nécessairement à **simple**).

C'est pourquoi :

- Complémenter deux modèles au classement absolu ; basés sur l'affinage des modèles FLAUBERT [Le et al., 2020].
- Objectif : proposer une méthode plus fine pour trouver des paires de phrases pertinentes pour la SAT.



### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique

Très souvent, la simplicité est évaluée automatiquement avec...

Modèle de classification absolue

Complexe

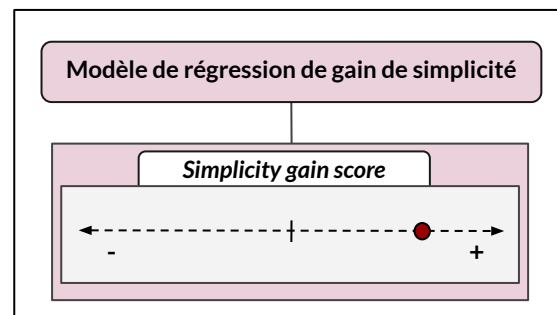
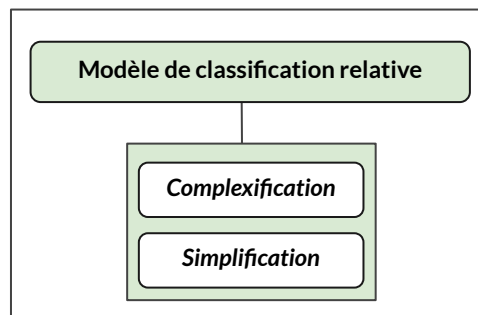
Simple

✗ Approche insuffisante pour la SAT

La **simplification** est une **opération intrinsèquement relative** → On transforme une phrase en une version **relativement plus simple** (ce qui n'équivaut pas nécessairement à **simple**).

C'est pourquoi :

- Complémenter deux modèles au classification absolu ; basés sur l'affinage des modèles FLAUBERT [Le et al., 2020].
- Objectif : proposer une **méthode plus fine** pour trouver des **paires de phrases** pertinentes pour la SAT.



### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique

Très souvent, la simplicité est évaluée automatiquement avec...

Modèle de classification absolue

Complexe

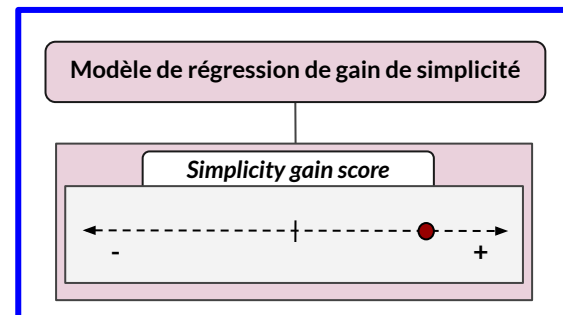
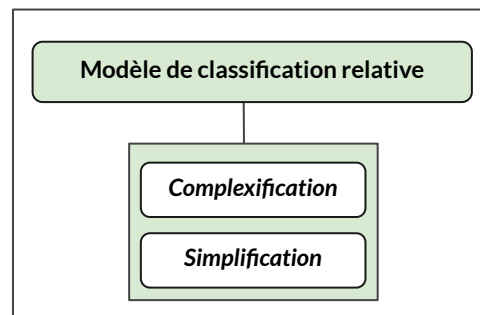
Simple

✗ Approche insuffisante pour la SAT

La **simplification** est une **opération intrinsèquement relative** → On transforme une phrase en une version **relativement plus simple** (ce qui n'équivaut pas nécessairement à **simple**).

C'est pourquoi :

- Complémenter deux modèles au classement absolu ; basés sur l'affinage des modèles FLAUBERT [Le et al., 2020].
- Objectif : proposer une **méthode plus fine** pour trouver des **paires de phrases** pertinentes pour la SAT.



### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Déterminer quelles paires constituent des simplifications valables”*

**GAIN** → **Différence absolue** entre la valeur d'un trait dans la phrase source et la phrase cible, y compris sa **polarité**.

- Ensemble de **traits**, divisés en trois groupes : **structurels**, **lexicaux** et **syntactiques**, suivant la littérature précédente [Tanguy & Tulechki, 2009 ; Brunato et al., 2022].
- Définissant un **score** de gain de simplicité.

### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Déterminer quelles paires constituent des simplifications valables”*

**GAIN** → Différence absolue entre la valeur d'un trait dans la phrase source et la phrase cible, y compris sa **polarité**.

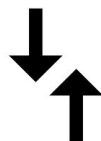
- Ensemble de **traits**, divisés en trois groupes : **structurels**, **lexicaux** et **syntactiques**, suivant la littérature précédente [Tanguy & Tulechki, 2009 ; Brunato et al., 2022].
- Définissant un **score** de gain de simplicité.

Trait
Nombre de mots

8 mots



Wiki



6 mots



Viki

=



gain = +2

modèle de régression

pairs	feature gains			
$wv_1$	$g_{11}$	$g_{12}$	$g_{13}$	...
$wv_2$	$g_{21}$	$g_{22}$	$g_{23}$	...
...				



### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Déterminer quelles paires constituent des simplifications valables”*

Tâche	Modèle de classification absolue		Modèle de classification relative		Modèle de régression	
Eval. dataset	Corpus test (précision %)		Corpus test (précision %)		Corpus test (score MSE)	
Modèle transformer	untuned	tuned	untuned	tuned	untuned	tuned
<b>flaubert-small</b>	49,54	<b>70,11</b>	49,78	92,99	1,89	0,39
<b>flaubert-base</b>	50,97	69,82	49,88	93,82	1,18	0,35
<b>flaubert-large</b>	52,29	69,19	52,18	<b>94,16</b>	4,59	<b>0,23</b>

### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Déterminer quelles paires constituent des simplifications valables”*

Wiki

Viki

*En France, ce lézard est strictement protégé par la loi.*

*En France, il est protégé par la loi.*

Modèle de classification absolue

Modèle de classification relative

Modèle de régression de gain de simplicité

### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Déterminer quelles paires constituent des simplifications valables”*

Wiki

Viki

*En France, ce lézard est strictement protégé par la loi.*

*En France, il est protégé par la loi.*

Modèle de classification absolue

Complexe

Simple

Modèle de classification relative

Modèle de régression de gain de simplicité

### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Déterminer quelles paires constituent des simplifications valables”*

Wiki

Viki

*En France, ce lézard est strictement protégé par la loi.*

*En France, il est protégé par la loi.*

Modèle de classification absolue

Complexe

Simple

Modèle de classification relative

Simplification

Modèle de régression de gain de simplicité

### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Déterminer quelles paires constituent des simplifications valables”*

Wiki

Viki

*En France, ce lézard est strictement protégé par la loi.*

*En France, il est protégé par la loi.*

Modèle de classification absolue

Complexe

Simple

Modèle de classification relative

Simplification

Modèle de régression de gain de simplicité

+0,84

### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Déterminer quelles paires constituent des simplifications valables”*

Wiki

Viki

*Praticien précoce et représentant éminent du concept français de la haute gastronomie, il est considéré comme le fondateur de ce style grandiose, recherché à la fois par les cours royales et les nouveaux riches de Paris*

*Il est considéré comme l'un des pionniers, sinon le fondateur, de la gastronomie française.*

Modèle de classification absolue

Modèle de classification relative

Modèle de régression de gain de simplicité

### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Déterminer quelles paires constituent des simplifications valables”*

Wiki

Viki

*Praticien précoce et représentant éminent du concept français de la haute gastronomie, il est considéré comme le fondateur de ce style grandiose, recherché à la fois par les cours royales et les nouveaux riches de Paris*

*Il est considéré comme l'un des pionniers, sinon le fondateur, de la gastronomie française.*

Modèle de classification absolue

Complexe

Complexe

Modèle de classification relative

Modèle de régression de gain de simplicité

### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Déterminer quelles paires constituent des simplifications valables”*

Wiki

Viki

*Praticien précoce et représentant éminent du concept français de la haute gastronomie, il est considéré comme le fondateur de ce style grandiose, recherché à la fois par les cours royales et les nouveaux riches de Paris*

*Il est considéré comme l'un des pionniers, sinon le fondateur, de la gastronomie française.*

Modèle de classification absolue

Complexe

Complexe

Modèle de classification relative

Simplification

Modèle de régression de gain de simplicité

+1,95



### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Déterminer quelles paires constituent des simplifications valables”*

Wiki

Viki

*Makassar ou Macassar est une ville d'Indonésie et la capitale de la province de Sulawesi du Sud.*

*Macassar ou Makassar est une ville d'Indonésie, située sur l'île de Sulawesi (ou Célèbes), en bordure du détroit du même nom.*

Modèle de classification absolue

Modèle de classification relative

Modèle de régression de gain de simplicité

### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Déterminer quelles paires constituent des simplifications valables”*

Wiki

Viki

*Makassar ou Macassar est une ville d'Indonésie et la capitale de la province de Sulawesi du Sud.*

*Macassar ou Makassar est une ville d'Indonésie, située sur l'île de Sulawesi (ou Célèbes), en bordure du détroit du même nom.*

Modèle de classification absolue

Simple

Complexe

Modèle de classification relative

Modèle de régression de gain de simplicité

### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Déterminer quelles paires constituent des simplifications valables”*

Wiki

Viki

*Makassar ou Macassar est une ville d'Indonésie et la capitale de la province de Sulawesi du Sud.*

*Macassar ou Makassar est une ville d'Indonésie, située sur l'île de Sulawesi (ou Célèbes), en bordure du détroit du même nom.*

Modèle de classification absolue

Simple

Complexe

Modèle de classification relative

Complexification

Modèle de régression de gain de simplicité

### 3. Méthode de filtrage en deux étapes pour la SAT

- 1.1. Étape I : Préservation du sens original
- 1.2. Étape II : Gain en simplicité linguistique



*“Déterminer quelles paires constituent des simplifications valables”*

Wiki

Viki

*Makassar ou Macassar est une ville d'Indonésie et la capitale de la province de Sulawesi du Sud.*

*Macassar ou Makassar est une ville d'Indonésie, située sur l'île de Sulawesi (ou Célèbes), en bordure du détroit du même nom.*

Modèle de classification absolue

Simple

Complexe

Modèle de classification relative

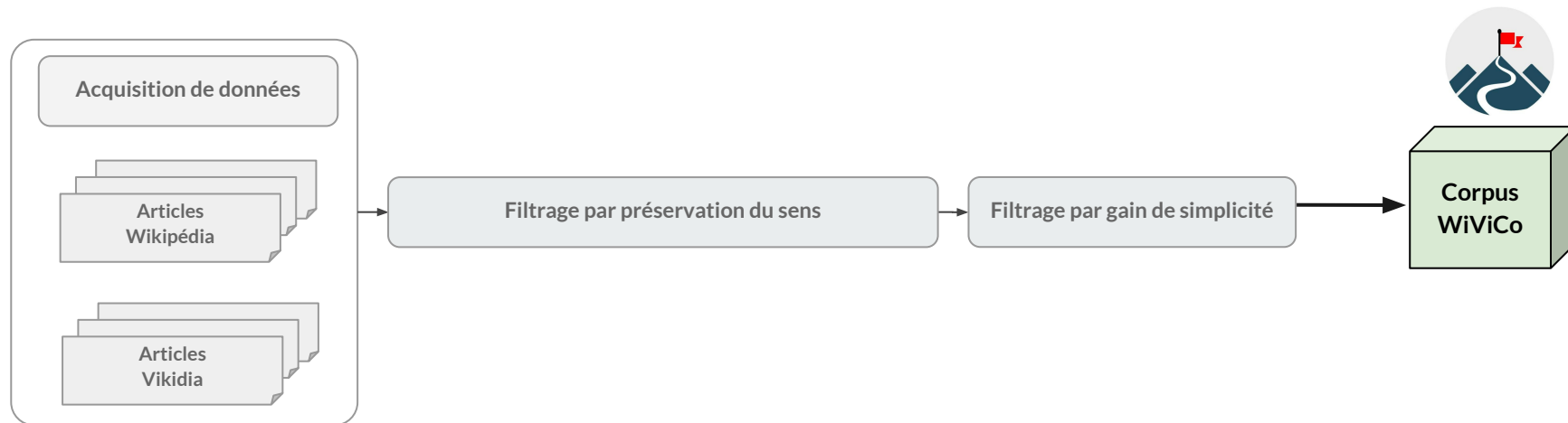
Complexification

Modèle de régression de gain de simplicité

-2,65

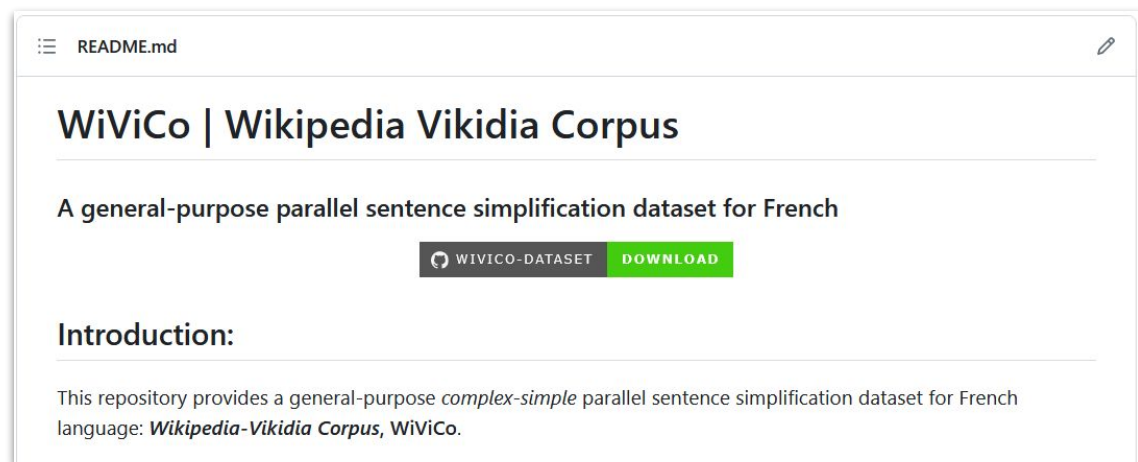
#### 4. Résultats : *Wikipedia-Vikidia Corpus (WiViCo)*

Suite à l'application de ce filtrage en deux étapes...



## 4. Résultats : Wikipedia-Vikidia Corpus (WiViCo)

À l'issue de l'application de cette méthode de filtrage :



<https://github.com/lormaechea/wivico>

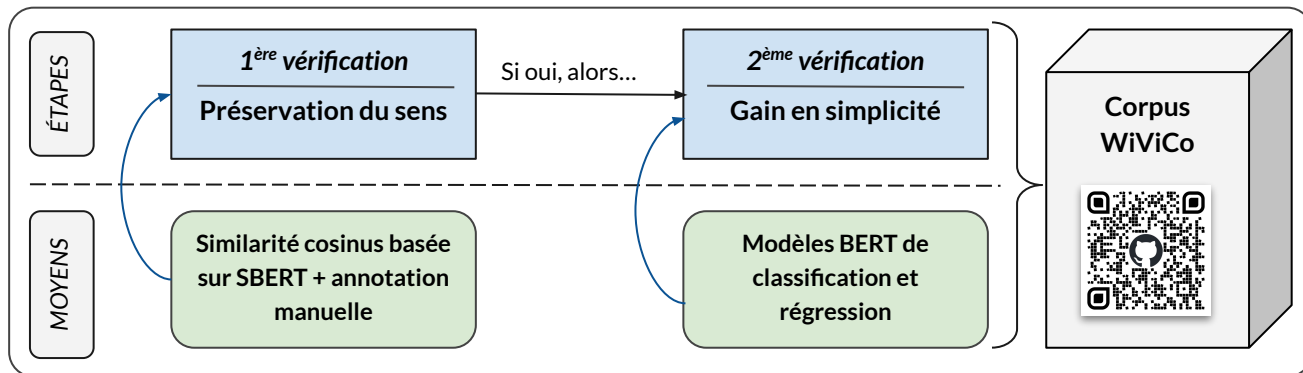
- Jeu de données en français associant des paires de **phrases *complexes-simples***.
- Version actuelle : presque **65k phrases parallèles**.

## 5. Conclusions

### Principale contribution :

Introduction d'une nouvelle **méthode** pour **exploiter des corpus comparables** :

- Spécifiquement **ciblée** pour la **tâche** de **SAT**.
- **Approche séquentielle** pour satisfaire les **2 conditions principales** qu'un **texte simplifié** doit remplir pour être considéré comme étant **valable** :



## 5. Conclusions

### Poursuites :

#### RÉGLAGE

- Amélioration de l'interprétabilité du score de gain de simplicité.
- Besoin d'un contraste avec des jugements humains.



## 5. Conclusions

### Poursuites :

#### RÉGLAGE

- Amélioration de l'interprétabilité du score de gain de simplicité.
- Besoin d'un contraste avec des jugements humains.

#### GÉNÉRATION

- Entraînement des modèles de SAT en français.
- Affinage d'un grand modèle de langage pré-entraîné → Pour notre tâche en aval.

## 5. Conclusions

### Poursuites :

#### RÉGLAGE

- Amélioration de l'interprétabilité du score de gain de simplicité.
- Besoin d'un contraste avec des jugements humains.

#### GÉNÉRATION

- Entraînement des modèles de SAT en français.
- Affinage d'un grand modèle de langage pré-entraîné → Pour notre tâche en aval.

#### ADAPTATION

- Étendre l'application de modèles de SAT à une modalité orale.
- Examiner comment les adapter pour les rendre robustes au bruit induit par les caractéristiques de la langue parlée.

---

# Merci de votre attention !

---



[Lucia.OrmaecheaGrijalba@unige.ch](mailto:Lucia.OrmaecheaGrijalba@unige.ch)



<https://luciaormaechea.com/fr/>

# Références

---

[Par ordre d'apparition]

- **C. Horn, C. Manduca, and D. Kauchak.** Learning a lexical simplifier using wikipedia. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 458–463. Association for Computational Linguistics, **2014**. URL <http://aclweb.org/anthology/P14-2075>
- **S. Stajner.** Automatic text simplification for social good: Progress and challenges. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2637–2652. Association for Computational Linguistics, **2021**. URL <https://aclanthology.org/2021.findings-acl.233>
- **N. Gala, A. Tack, L. Javourey-Drevet, T. François, and J. C. Ziegler.** Alector: A Parallel Corpus of Simplified French Texts with Alignments of Misreadings by Poor and Dyslexic Readers. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, Marseille, France. European Language Resources Association, **2020**. URL <https://aclanthology.org/2020.lrec-1.169>
- **X. Zhang and M. Lapata.** Sentence Simplification with Deep Reinforcement Learning. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 584–594, Copenhagen, Denmark. Association for Computational Linguistics, **2017**. URL <https://aclanthology.org/D17-1062/>