Extracting sentence simplification pairs from French comparable corpora using a two-step filtering method

Swiss Text Analytics Conference 2023

June 13th, University of Applied Sciences and Arts of Western Switzerland, Neuchâtel

Lucía Ormaechea and Nikos Tsourakis

Department of Translation Technology





- 1.1. What is Automatic Text Simplification?
- 1.2. ATS in French: a particularly low-resourced task
- 1.3. Bridging the gap: a new method to extract simplification pairs

2. Data acquisition: scraping Wikipedia and Vikidia

3. Two-step ATS-targeted filtering method

- 3.1. Stage I: Sentence semantic similarity filtering
- 3.2. Stage II: Simplicity gain filtering
- 4. Results: Wikipedia-Vikidia Corpus (WiViCo)
- 5. Conclusions

1.1. What is Automatic Text Simplification?

- 1.2. ATS in French: a particularly low-resourced task
- 1.3. Bridging the gap: a new method to extract simplification pairs

Automatic Text Simplification (ATS) is an area of NLP that aims at automatically converting texts into simpler variants, by reducing their linguistic complexity, albeit preserving their original meaning [Horn *et al.*, 2014; Stajner, 2021].

ORIGINAL	The second largest city of China and one of the world's major cities , Beijing has played a vital role in Chinese history.
SIMPLIFIED	Beijing is the second biggest city of China. Beijing has played an important role in Chinese history.

1.1. What is Automatic Text Simplification?

- 1.2. ATS in French: a particularly low-resourced task
- 1.3. Bridging the gap: a new method to extract simplification pairs

Automatic Text Simplification (ATS) is an area of NLP that aims at automatically converting texts into simpler variants, by reducing their linguistic complexity, albeit preserving their original meaning [Horn *et al.*, 2014; Stajner, 2021].

ORIGINAL	The second <i>largest</i> city of China and one of the world's major cities , Beijing has played a vital role in Chinese history.
SIMPLIFIED	Beijing is the second biggest city of China. Beijing has played an important role in Chinese history.

ATS: an analogy with the form-substance relation:





1.1. What is Automatic Text Simplification?

- 1.2. ATS in French: a particularly low-resourced task
- 1.3. Bridging the gap: a new method to extract simplification pairs

Automatic Text Simplification (ATS) is an area of NLP that aims at automatically converting texts into simpler variants, by reducing their linguistic complexity, albeit preserving their original meaning [Horn *et al.*, 2014; Stajner, 2021].

ORIGINAL	The second <i>largest</i> city of China and one of the world's major cities , Beijing has played a vital role in Chinese history.
SIMPLIFIED	Beijing is the second biggest city of China. Beijing has played an important role in Chinese history.

ATS: an analogy with the form-substance relation:



1.1. What is Automatic Text Simplification?

- 1.2. ATS in French: a particularly low-resourced task
- 1.3. Bridging the gap: a new method to extract simplification pairs

Automatic Text Simplification (ATS) is an area of NLP that aims at automatically converting texts into simpler variants, by reducing their linguistic complexity, albeit preserving their original meaning [Horn *et al.*, 2014; Stajner, 2021].

	ORIGINAL	The second largest city of China and one of the world's major cities , Beijing has played a vital role in Chinese history.
_	SIMPLIFIED	Beijing is the second biggest city of China. Beijing has played an important role in Chinese history.

ATS: an analogy with the form-substance relation:



Text Simplification

ORIGINAL	The second largest city of China and one of the world's major cities , Beijing has played a vital role in Chinese history.
SIMPLIFIED	Beijing is the second biggest city of China. Beijing has played an important role in Chinese history.

- 1.1. What is Automatic Text Simplification?
- 1.2. ATS in French: a particularly low-resourced task
- 1.3. Bridging the gap: a new method to extract simplification pairs

Text Simplification

✓ Text accessibility

✓ Comprehensibility aid

X Scarcity of parallel monolingual data

X Bottleneck for the advancement of ATS

ORIGINAL	The second <i>largest</i> city of China and one of the world's major cities , Beijing has played a vital role in Chinese history.	•
SIMPLIFIED	Beijing is the second biggest city of China. Beijing has played an important role in Chinese history.	-

- 1.1. What is Automatic Text Simplification?
- 1.2. ATS in French: a particularly low-resourced task
- 1.3. Bridging the gap: a new method to extract simplification pairs

✓ Text accessibility

✓ Comprehensibility aid

× Scarcity of parallel monolingual data

X Bottleneck for the advancement of ATS

ORIGINAL	The second largest city of China and one of the world's major cities , Beijing has played a vital role in Chinese history.	∢ —
SIMPLIFIED	<i>Beijing</i> is the second <i>biggest</i> city of China. Beijing has played <i>an important</i> role in Chinese history.	

say.... French

In less resource-rich languages than English: \bigcirc^{\bigcirc}

Only available monolingual parallel dataset \rightarrow **ALECTOR** [Gala *et al.*, 2020]:

Text Simplification

- Manually created by professional editors.
- Comprising schoolbook texts, simplified for child audiences.

- 1.1. What is Automatic Text Simplification?
- 1.2. ATS in French: a particularly low-resourced task
- 1.3. Bridging the gap: a new method to extract simplification pairs

✓ Text accessibility

✓ Comprehensibility aid

× Scarcity of parallel monolingual data

X Bottleneck for the advancement of ATS

ORIGINAL	The second <i>largest</i> city of China and one of the world's major cities , Beijing has played a vital role in Chinese history.	•
SIMPLIFIED	<i>Beijing</i> is the second <i>biggest</i> city of China. Beijing has played <i>an important</i> role in Chinese history.	

say.... French

In less resource-rich languages than English: \bigcirc^{\bigcirc}

Only available monolingual parallel dataset \rightarrow **ALECTOR** [Gala *et al.*, 2020]:

Text Simplification

- Manually created by professional editors.
- Comprising schoolbook texts, simplified for child audiences.

Manually crafted datasets

 High-quality and reliable simplifications

X Costly compilation, both

economically and time-wise

× Compact size, preventing the implementation of ML algorithms

- 1.1. What is Automatic Text Simplification?
- 1.2. ATS in French: a particularly low-resourced task
- 1.3. Bridging the gap: a new method to extract simplification pairs

✓ Text accessibility

✓ Comprehensibility aid

× Scarcity of parallel monolingual data

X Bottleneck for the advancement of ATS

ORIC	GINAL	The second largest city of China and one of the world's major cities , Beijing has played a vital role in Chinese history.
SIM	PLIFIED	Beijing is the second biggest city of China. Beijing has played an important role in Chinese history.

In less resource-rich languages than English: \bigcirc

Only available monolingual parallel dataset \rightarrow **ALECTOR** [Gala *et al.*, 2020]:

Text Simplification

- Manually created by professional editors.
- Comprising schoolbook texts, simplified for child audiences.

Manually crafted datasets

say.... French

- ✓ High-quality and reliable simplifications
- × Costly compilation, both economically and time-wise
- X Compact size, preventing the implementation of ML algorithms



What about exploring automatic methods?

- 1.1. What is Automatic Text Simplification?
- 1.2. ATS in French: a particularly low-resourced task
- 1.3. Bridging the gap: a new method to extract simplification pairs

MAIN CONTRIBUTION • Introduce a new method to mine comparable corpora to extract complex-simple pairs.

• Properly **adapted** to the **ATS** task.



"Create a parallel French corpus of complex-simple pairs for ATS"

- 1.1. What is Automatic Text Simplification?
- 1.2. ATS in French: a particularly low-resourced task
- 1.3. Bridging the gap: a new method to extract simplification pairs

• Introduce a new method to mine comparable corpora to extract complex-simple pairs. MAIN CONTRIBUTION

"Create a parallel French corpus of complex-simple pairs for ATS"

• Properly adapted to the ATS task.

Introduction

1.1. What is Automatic Text Simplification?

1.

- 1.2. ATS in French: a particularly low-resourced task
- 1.3. Bridging the gap: a new method to extract simplification pairs

Automatic approaches for ATS data compilation:

- Typically based on register-diversified comparable corpora:
 - **Wikipedia** \rightarrow French edition. 0
 - **Vikidia** \rightarrow Simplified version of the former. 0



MAIN CONTRIBUTION
Introduce a new method to mine comparable corpora to extract complex-simple pairs.
Properly adapted to the ATS task.

"Create a parallel French corpus of complex-simple pairs for ATS"

1. Introduction

- 1.1. What is Automatic Text Simplification?
- 1.2. ATS in French: a particularly low-resourced task
- 1.3. Bridging the gap: a new method to extract simplification pairs

Automatic approaches for ATS data compilation:

- Typically based on register-diversified comparable corpora:
 - $\circ \qquad \textbf{Wikipedia} \rightarrow French \ edition.$
 - \circ Vikidia \rightarrow Simplified version of the former.
- Usually rely on semantic similarity measures, but are oblivious of whether the target sentence constitutes indeed a simplification.



• Introduce a new method to mine comparable corpora to extract complex-simple pairs. MAIN CONTRIBUTION

"Create a parallel French corpus of complex-simple pairs for ATS"

• Properly adapted to the ATS task.

Introduction 1.

- 1.1. What is Automatic Text Simplification?
- ATS in French: a particularly 1.2. low-resourced task
- 1.3. Bridging the gap: a new method to extract simplification pairs

Automatic approaches for ATS data compilation:

- Typically based on register-diversified comparable corpora:
 - **Wikipedia** \rightarrow French edition. 0
 - **Vikidia** \rightarrow Simplified version of the former. 0
- Usually rely on semantic similarity measures, but are oblivious of whether the target sentence constitutes indeed a simplification.

That's why we propose...

	Stages	Two primary con	ditions within ATS	Tackled in a sequential manner
	1Retention of the original meaning and information		Sentence semantic similarity filtering	
2 Linguistic simplicity gain with respect to the reference		Simplicity gain filtering		





Extracting textual content from both encyclopedias:

- Scrape standard texts (Wikipedia) and their simplified versions (Vikidia).
- Special focus on *lead sections* (initial summaries).



Extracting textual content from both encyclopedias:

- Scrape standard texts (Wikipedia) and their simplified versions (Vikidia).
- Special focus on *lead sections* (initial summaries).
- Assumption that → Greater chances of finding aligned sentences:
 - 1. Both share a typically **definitional style**.
 - 2. **Common excerpts** are potentially **more likely** to appear in the **beginning**.



Extracting textual content from both encyclopedias:

- Scrape standard texts (Wikipedia) and their simplified versions (Vikidia).
- Special focus on *lead sections* (initial summaries).
- Assumption that → Greater chances of finding aligned sentences:
 - 1. Both share a typically **definitional style**.
 - 2. Common excerpts are potentially more likely to appear in the beginning.





Data acquisition	Wiki-texts	Viki-texts
# documents	34,806	
# sentences	165,806	134,348
# tokens	4,030,148	2,373,045
# types	294,979	195,791
Type/token ratio	7.32	8.25
Avg. word length	5.32	5.08
Avg. sent. length	25.22	18.16



3. Two-step ATS-targeted filtering method

- 1.1. Stage I: Sentence semantic similarity filtering
- 1.2. Stage II: Simplicity gain filtering



- 3. Two-step ATS-targeted filtering method
 - 1.1. Stage I: Sentence semantic similarity filtering
 - 1.2. Stage II: Simplicity gain filtering



- 3. Two-step ATS-targeted filtering method
 - 1.1. Stage I: Sentence semantic similarity filtering
 - 1.2. Stage II: Simplicity gain filtering

"Determine which complex-simple pairs are suitable for ATS"



Data acquisition	Wiki-texts	Viki-texts
# documents	34	I,806
# sentences	165,806	134,348
# tokens	4,030,148	2,373,045
	•••	-

Automatic sentence alignment:

- **SBERT** (Sentence-BERT) using **multilingual** sentence transformers^{*}:
 - Generate fixed-length sentence **embeddings**.
 - The output is a **768-dimensional** dense vector representation.
 - **Cosine similarity** between complex and simple sentences.

https://huggingface.co/sentence-transformers/paraphrase-xlm-r-multilingual-v1

- 3. Two-step ATS-targeted filtering method
 - 1.1. Stage I: Sentence semantic similarity filtering
 - 1.2. Stage II: Simplicity gain filtering



Three Wiki: Viki example sentences

Label	Wikipedia sentence	Vikidia sentence
Valid	The term "White House" is often used as a metonym for the president's administration.	By metonymy, the White House also refers to the US government and its entourage.
Partially valid	Neal McDonough is an American actor and producer born on February 13, 1966 in Dorchester, Massachusetts.	Neal McDonough is an American actor.
Non-valid	The information refers both to the message to be communicated and the symbols used to write it.	The written words, sounds, images, smells or tastes contain information.

- 3. Two-step ATS-targeted filtering method
 - 1.1. Stage I: Sentence semantic similarity filtering
 - 1.2. Stage II: Simplicity gain filtering



Three Wiki: Viki example sentences

	Label	Wikipedia sentence	Vikidia sentence
	Valid	The term "White House" is often used as a metonym for the president's administration.	By metonymy, the White House also refers to the US government and its entourage.
>	Partially valid	Neal McDonough is an American actor and producer born on February 13, 1966 in Dorchester, Massachusetts.	Neal McDonough is an American actor.
	Non-valid	The information refers both to the message to be communicated and the symbols used to write it.	The written words, sounds, images, smells or tastes contain information.

- 3. Two-step ATS-targeted filtering method
 - 1.1. Stage I: Sentence semantic similarity filtering
 - 1.2. Stage II: Simplicity gain filtering



Three Wiki: Viki example sentences

	Label	Wikipedia sentence	Vikidia sentence
	Valid	The term "White House" is often used as a metonym for the president's administration.	By metonymy, the White House also refers to the US government and its entourage.
	Partially valid	Neal McDonough is an American actor and producer born on February 13, 1966 in Dorchester, Massachusetts.	Neal McDonough is an American actor.
>	Non-valid	The information refers both to the message to be communicated and the symbols used to write it.	The written words, sounds, images, smells or tastes contain information.

- 3. Two-step ATS-targeted filtering method
 - 1.1. Stage I: Sentence semantic similarity filtering
 - 1.2. Stage II: Simplicity gain filtering



Manual annotation:

- Randomly selecting **500 samples** from the initial dataset.
- **Two annotators** determined to which extent each pair of *Wiki:Viki* sentences conveyed the same meaning using three judgement labels:
 - **Valid**, where the meaning and information from **Wiki** to **Viki** is fully preserved.
 - **Partially valid**, where information is partially lost from **Wiki** to **Viki** or vice versa.
 - **Non-valid**, where information between **Wiki** to **Viki** is dissimilar.
- High agreement (Cohen's kappa = 0.87).

- 3. Two-step ATS-targeted filtering method
 - 1.1. Stage I: Sentence semantic similarity filtering
 - 1.2. Stage II: Simplicity gain filtering

"Extract the Wiki- and Viki-pairs that exhibit a high semantic overlap"

Manual annotation:

The mean score for the valid case was equal to 0.81, which we consider as the cutoff threshold for the semantic filtering of Wiki to Viki pairs.





3. Two-step ATS-targeted filtering method

- 1.1. Stage I: Sentence semantic similarity filtering
- 1.2. Stage II: Simplicity gain filtering



- 3. Two-step ATS-targeted filtering method
 - 1.1. Stage I: Sentence semantic similarity filtering

8 words

Wiki

1.2. Stage II: Simplicity gain filtering

Feature

Number of words



gain = -2

Gain: Signifies the **absolute difference** between a feature value in the source and target sentence, including polarity.

6 words

Viki

pairs	feature gains			
wv ₁	9 ₁₁	g ₁₂	g ₁₃	
wv ₂	9 ₂₁	g ₂₂	g ₂₃	
wv ₃	9 ₃₁	9 ₃₂	9 ₃₃	
wv ₄	9 ₄₁	9 ₄₂	g ₄₃	

. . .

3. Two-step ATS-targeted filtering method

- 1.1. Stage I: Sentence semantic similarity filtering
- 1.2. Stage II: Simplicity gain filtering

Feature group	Feature	Description
Structural	Sentence Length (SL)	Difference in the number of characters between the target and source sentences.
	Number of Words (NW)	Difference in the number of words between the target and source sen- tences.
	Word Error Rate (WER) BLEU score	Word-based similarity between the source and target sentences. <i>n</i> -gram overlap via precision of the target sentence with its correspond- ing source sentence.
Lexical	Number of Named Entities (NNE)	Difference in the number of named entities (organizations, people, places, <i>etc.</i>) between the target and source sentences.
		words multiplied by their complexity weight value.
	Maximum Depth Tree (MDT)	Difference in the maximum depth of the dependency tree between the target and source sentences.
Syntactic	Incomplete Dependency Theory (IDT)	Within a phrase, the average number of incomplete dependencies be- tween the current and next token.
	Dependency Locality Theory (DLT)	For every head token in a sentence, the number of discourse referents starting from the current token and ending to its longest leftmost dependent [32]. Values are then combined using an average function.
	Combined IDT+DLT	Sum of IDT+DLT metrics for all tokens in a sentence. Resulting values are then combined using an average function.
	Left Embeddedness (LE)	Within a sentence, the number of tokens on the left-hand-side of the root verb that are not verbs.
_	Noun Nested Distance (NND)	The average nested distance of all nouns within a phrase that have as ancestor another noun in the dependency tree.

- 3. Two-step ATS-targeted filtering method
 - 1.1. Stage I: Sentence semantic similarity filtering
 - 1.2. Stage II: Simplicity gain filtering



"Determine which pairs constitute valid simplifications"

- Trained a simplicity gain classifier on these features.
 - Feed-Forward Neural Network (FFNN) Four hidden layers of 256 nodes each.
- Exploited ~180k example pairs from WikiLarge (EN).
 - Resorted to Google Translate (FR).
- Extracted values of features independently for each article and combined them for a single pair.
- Split the dataset into **two halves** interchanging the source/target.
 - Acquired simplification/complexification pairs.
 - Calculated the gain values for each pair.

- 3. Two-step ATS-targeted filtering method
 - 1.1. Stage I: Sentence semantic similarity filtering
 - 1.2. Stage II: Simplicity gain filtering

Feature Selection

Feature	Description	
SL	Sentence Length	
CEFR	Common European Framework of Reference score	
NW	Number of Words	
IDT	Incomplete Dependency Theory	
IDTDLT	Combined IDT+DLT	
MDT	Maximum Depth Tree	
LE	Left Embeddedness	
NNE	Number of Named Entities	
BLEU	BLEU score	
WER	Word Error Rate	
NND	Noun Nested Distance	

Classification Performance



L. Ormaechea and N. Tsourakis



4. Results: Wikipedia-Vikidia Corpus (WiViCo)

- With the implementation of the proposed two-step filtering method we created the **WiViCo monolingual** parallel corpus.
- Utilizing the sigmoid output layer of the classifier we provide simplification pairs based on **lenient** or **stricter thresholds**.
- Based on different **cutoff probability thresholds**, we enumerate all **Wiki:Viki** samples in each class:

	Label		
Probability	0 (non-simplified)	1 (simplified)	
>0.9	44,049	20,692	
>0.8	22,556	42,185	
>0.7	18,642	46,099	
>0.6	11,087	53,654	
>0.5	7,302	57,439	

README.md WiViCo | Wikipedia Vikidia Corpus A general-purpose parallel sentence simplification dataset for French WIVICO-DATASET DOWNLOAD Introduction: This repository provides a general-purpose complex-simple parallel sentence simplification dataset for French language: Wikipedia-Vikidia Corpus, WiViCo.

https://github.com/lormaechea/wivico

Main contribution:

5. Conclusions

Introduction of a **new method** to **mine comparable corpora**:

- Specifically targeted to ATS.
- Sequential approach to satisfy the 2 main conditions for a simplified text to be valid:



Goals to be explored:

5. Conclusions

Once the data compilation is done, move on to the **generation stage**, by:

- Training general-domain seq2seq ATS models for French.
- Fine-tuning LLMs for our downstream task.

Goals to be explored:

Once the data compilation is done, move on to the **generation stage**, by:

- Training general-domain seq2seq ATS models for French.
- Fine-tuning LLMs for our downstream task.

Goals being currently explored:

• Introducing a **comprehensive approach** for **assessing sentence simplicity**.



Simplification is inherently a relative process \rightarrow A given text is transformed into a relatively simpler version, which does not necessarily equate to simple.

• Examining a **BERT-based fine-tuning approach** to **qualify** and **quantify sentence simplicity**.

5. Conclusions

Thanks for your attention!

Any questions?





References

[In order of appearance]

- **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 <u>http://aclweb.org/anthology/P14-2075</u>
- 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 <u>https://aclanthology.org/2020.lrec-1.169</u>