# MUSS: Multilingual Unsupervised Sentence Simplification by Mining Paraphrases

**Louis Martin**, Angela Fan, Éric de la Clergerie, Antoine Bordes, Benoît Sagot

Meta AI Paris & Inria Paris

FACEBOOK AI

# Access to Information is Hard

# … and many People struggle with Reading Difficulties

- Intellectual Disabilities

- Low literacy

- Non-native speakers

How can we make information easier to read and comprehend for each and everyone?

# Automatic Sentence Simplification

**Goal**: Simplify a sentence while preserving its meaning

# A Typical Human Simplification

Source

Simplification

The second largest city of Russia and one of the world's major cities, St . Petersburg has played a vital role in Russian history.
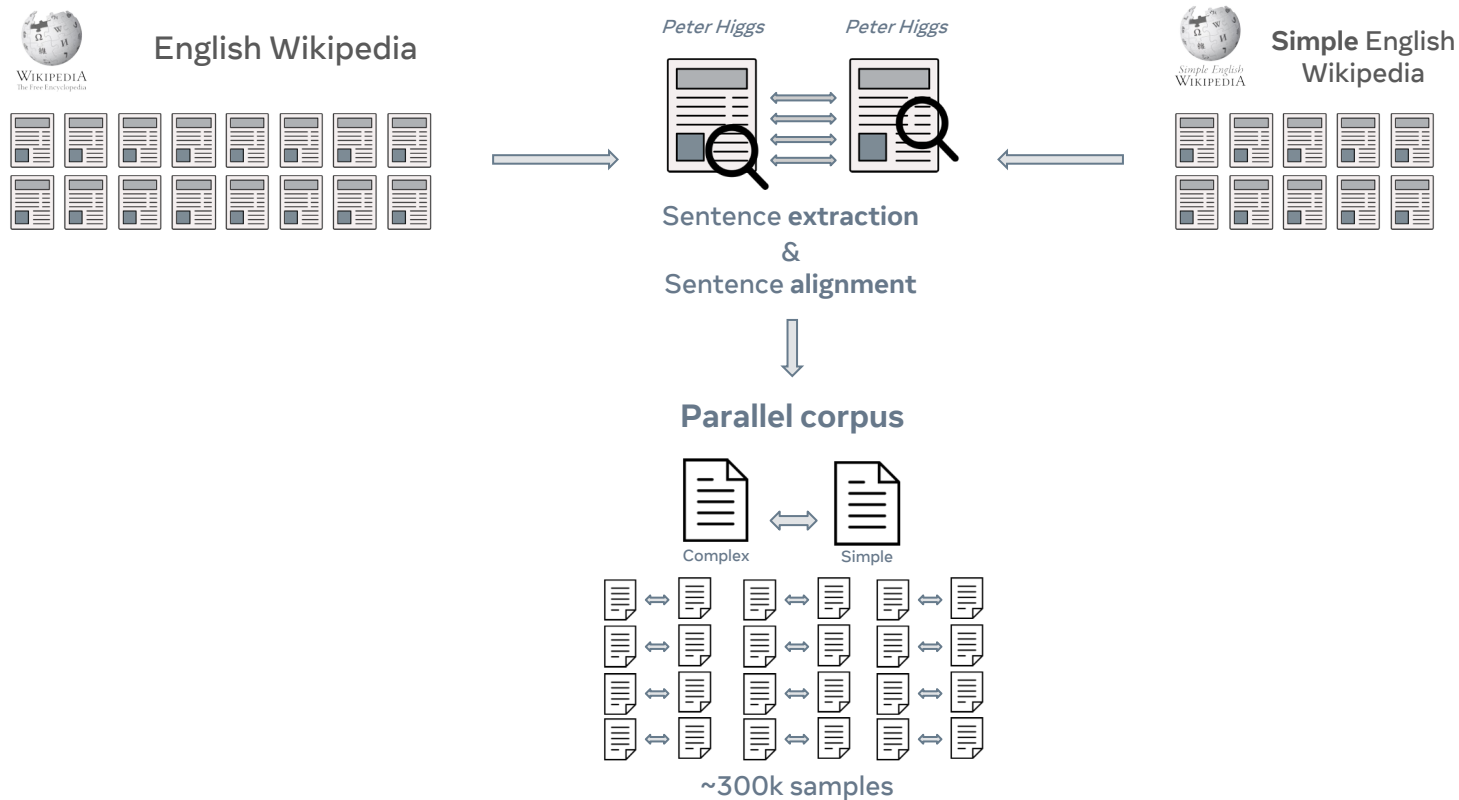
Sentence Splitting

St. Petersburg is the second biggest city of Russia.

St. Petersburg has played an important role in Russian history.

- **Lexical Simplification** - Replace uncommon words

- **Syntactic Simplification** - Simplify complex syntactic structures

- **Compression** - Retain key information only

# How to train sentence simplification models?

# Traditional Simplification Datasets



English Wikipedia

*Peter Higgs*  *Peter Higgs*

**Simple** English Wikipedia

Sentence **extraction**
&
Sentence **alignment**

**Parallel corpus**

Complex    Simple

~300k samples

# Problems with Simple English Wikipedia Alignment

- Contains alignment errors.

- Encyclopedic Domain only.

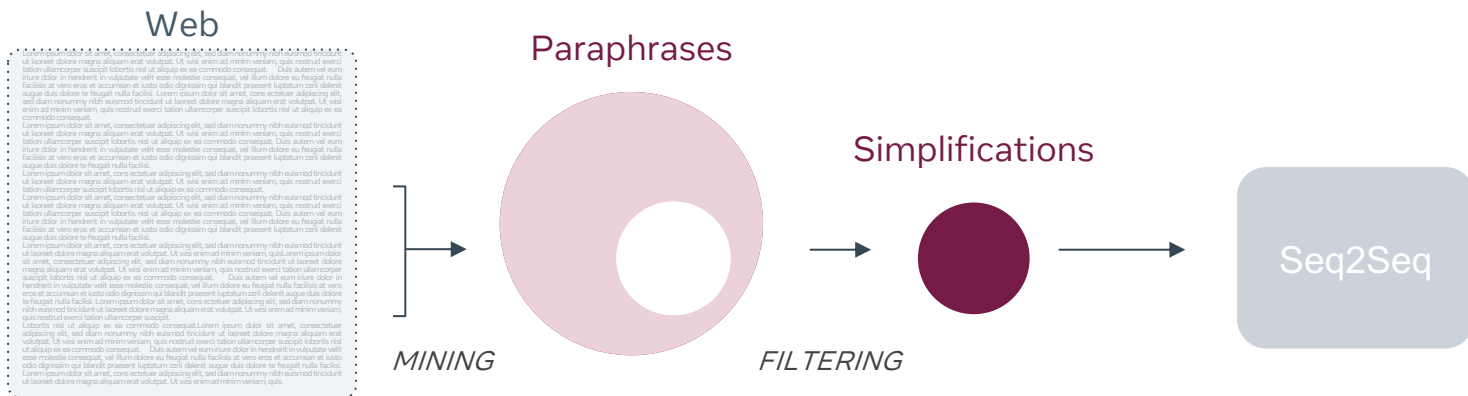- Simple English Wikipedia **only in English**.

# Can we find Parallel Sentences in any Language on the **Web**?

# MUSS: Multilingual Unsupervised Sentence Simplification by Mining Paraphrases

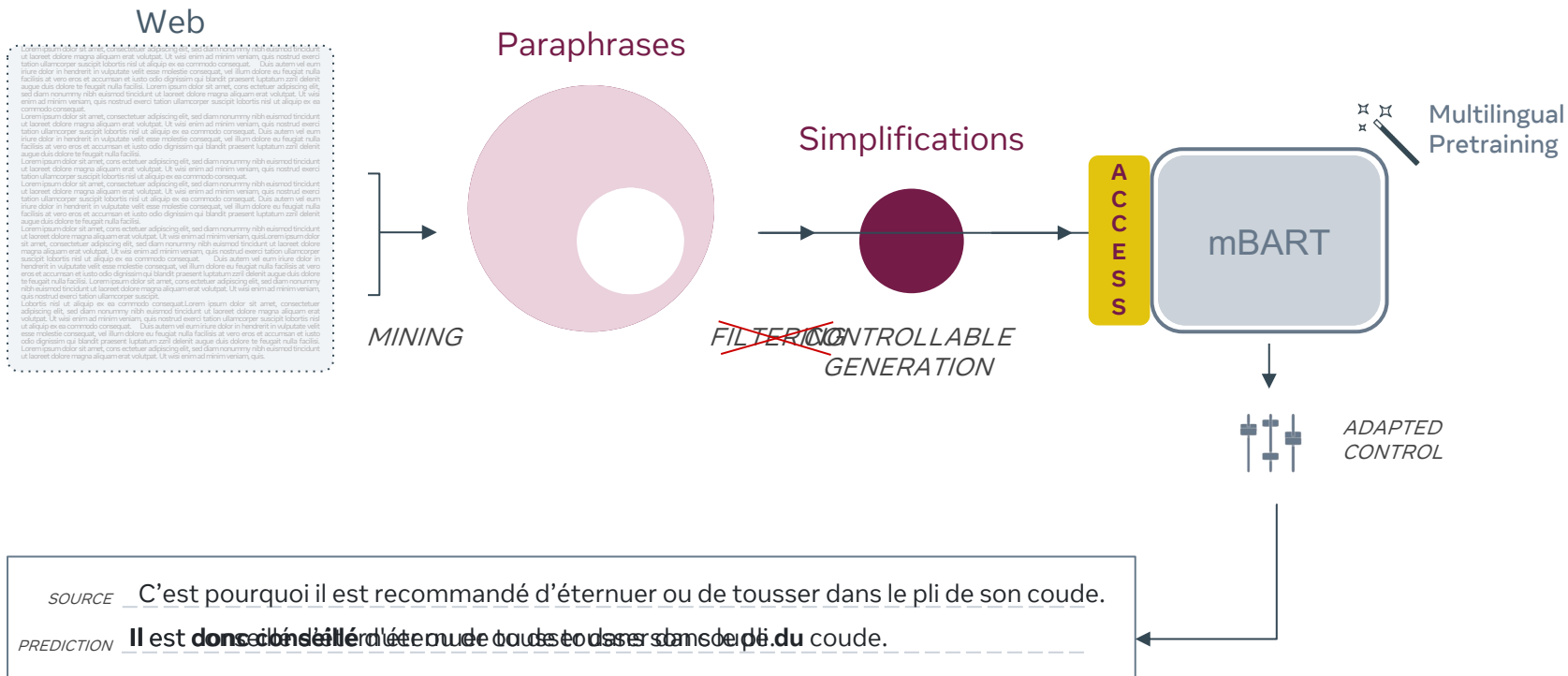**Louis Martin**, Angela Fan, Éric de la Clergerie, Antoine Bordes, Benoît Sagot
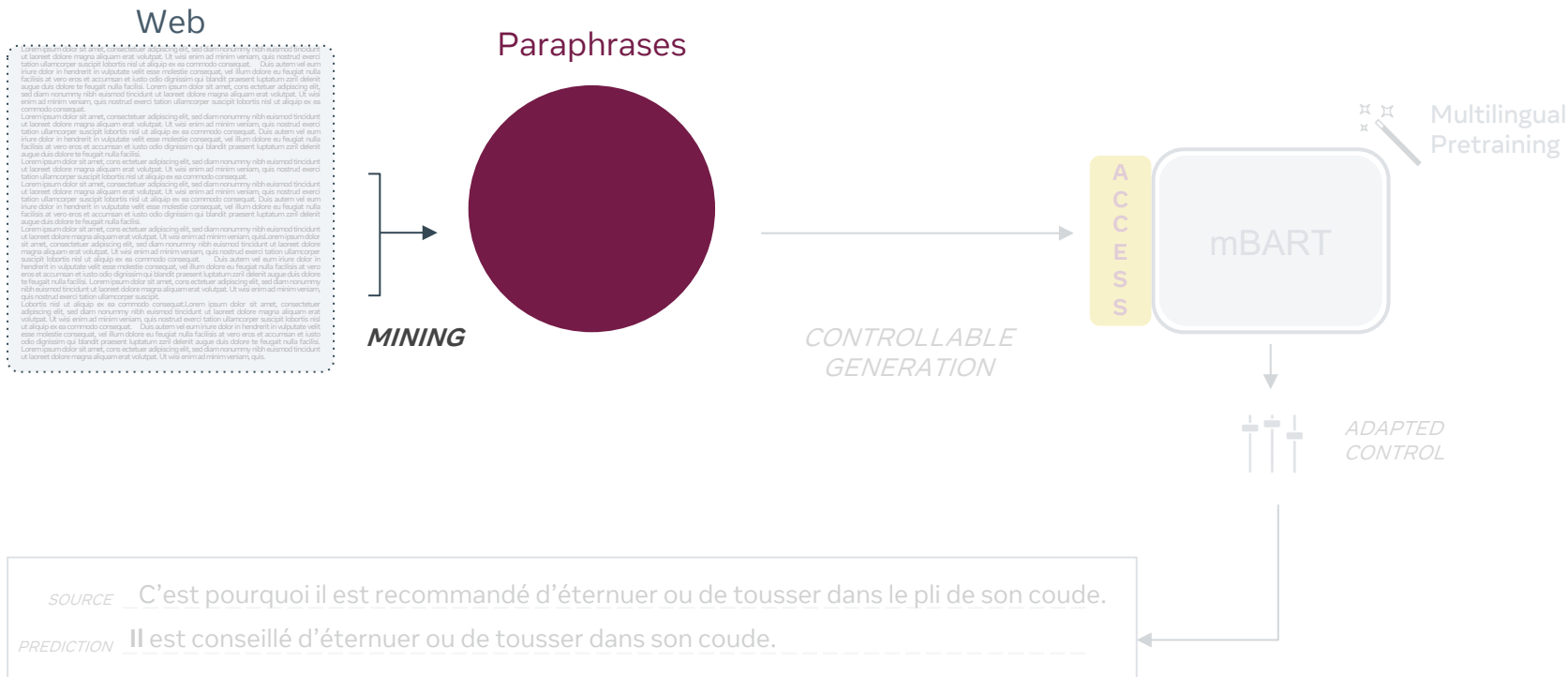
FACEBOOK AI

# Mining Simplifications on the Web

Web

Paraphrases

Simplifications

Seq2Seq

*MINING*

*FILTERING*

What if we didn't need this extra filtering step?

# The MUSS Approach



Web

Paraphrases

Simplifications

Multilingual Pretraining

A C C E S S

mBART

*MINING*

*FILTERING* *CONTROLLABLE GENERATION*

*ADAPTED CONTROL*

| | |
|---|---|
| *SOURCE* | C'est pourquoi il est recommandé d'éternuer ou de tousser dans le pli de son coude. |
| *PREDICTION* | Il est **donc conseillé** d'éternuer ou de tousser dans le pli **du** coude. |

# Mine Paraphrases

Web

Paraphrases

MINING

CONTROLLABLE
GENERATION

ACCESS

mBART

Multilingual
Pretraining

ADAPTED
CONTROL

SOURCE   C'est pourquoi il est recommandé d'éternuer ou de tousser dans le pli de son coude.

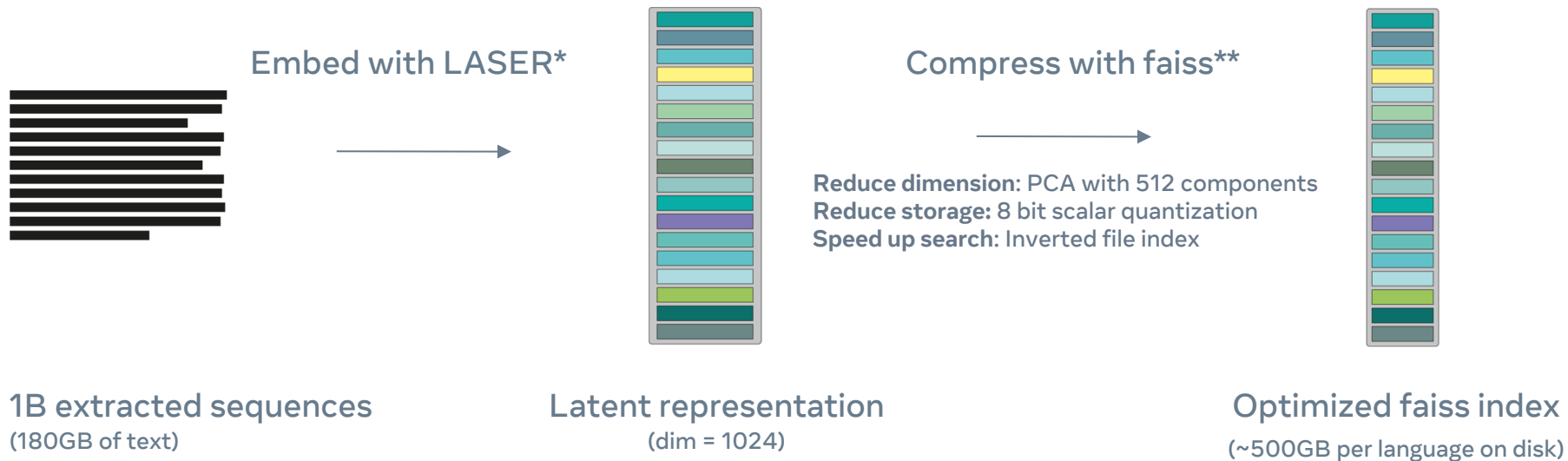PREDICTION   Il est conseillé d'éternuer ou de tousser dans son coude.

# Mine Paraphrases

**Idea**: Mine Paraphrases using Sentence Embeddings

**Paraphrases** = **Nearest Neighbours** in Embedding Space

# Index Creation

Embed with LASER*

Compress with faiss**

**Reduce dimension**: PCA with 512 components
**Reduce storage**: 8 bit scalar quantization
**Speed up search**: Inverted file index

1B extracted sequences
(180GB of text)

Latent representation
(dim = 1024)

Optimized faiss index
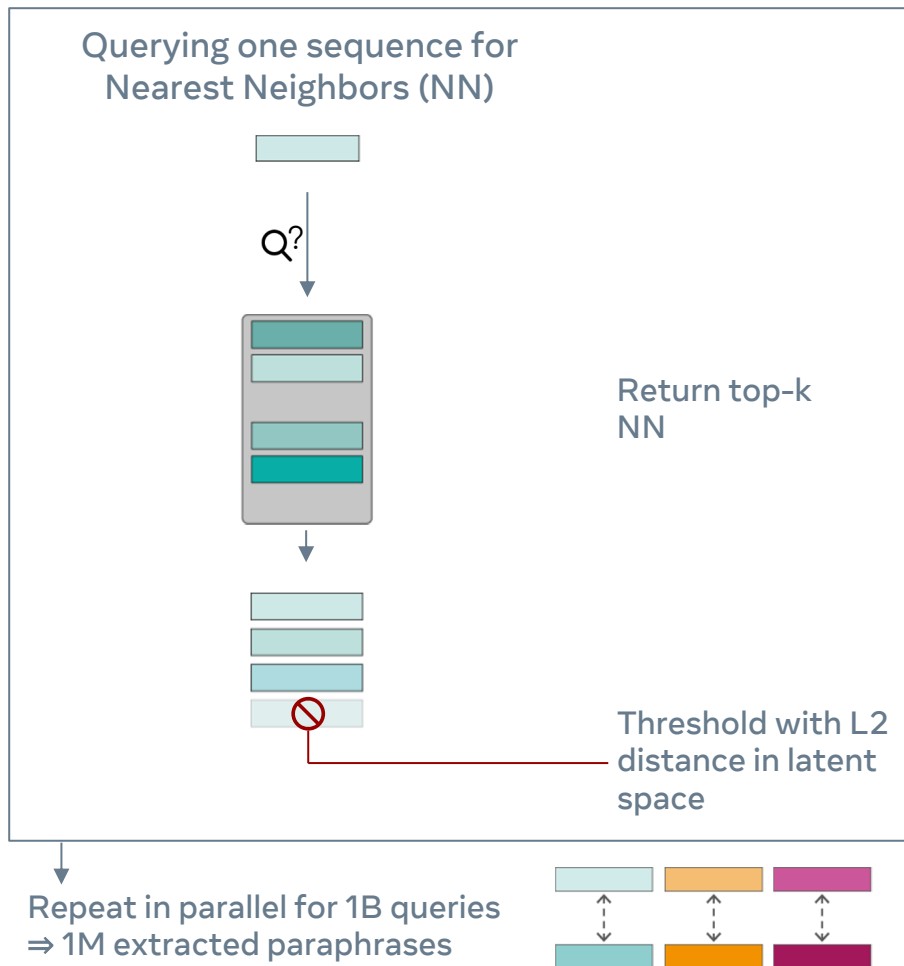(~500GB per language on disk)

*LASER: Multilingual sentence embeddings model
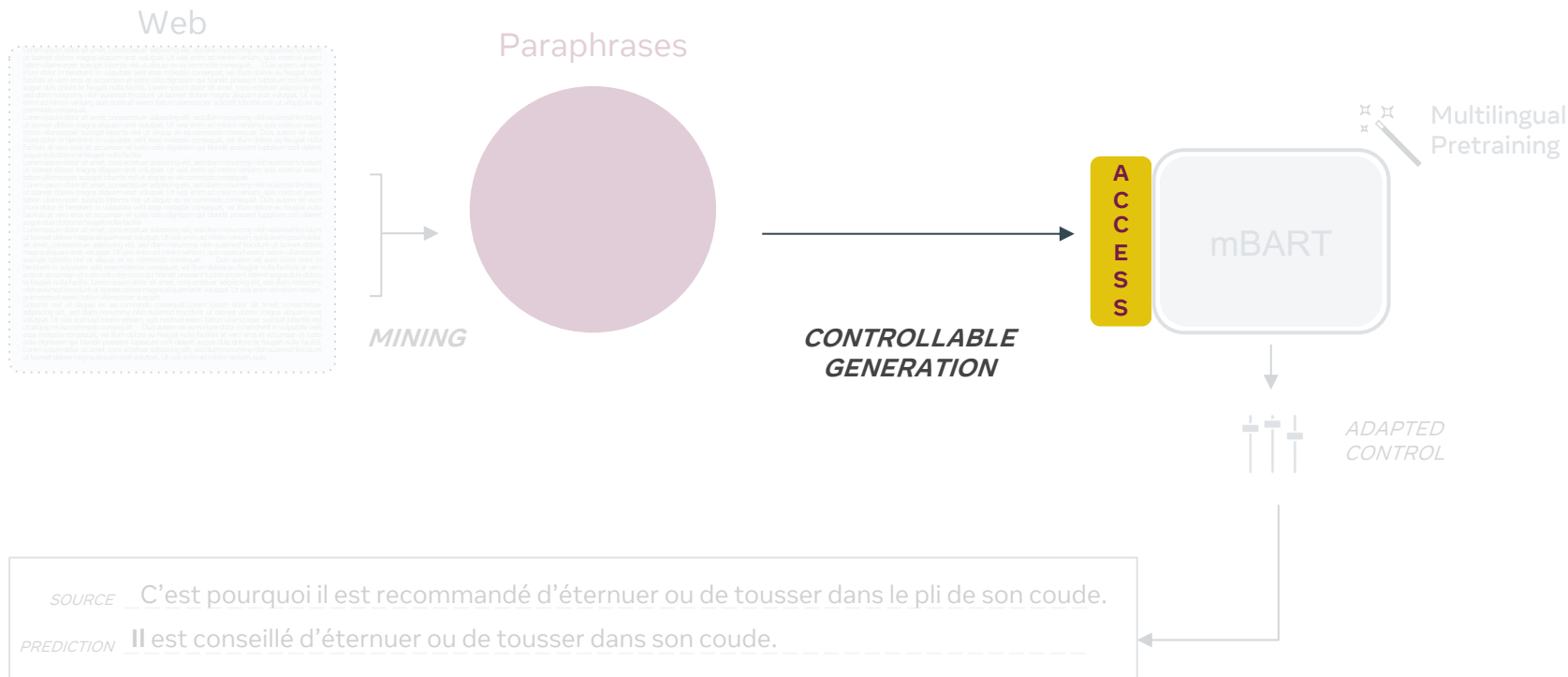**faiss: Fast nearest neighbour search library

# Paraphrase Mining

Optimized faiss index

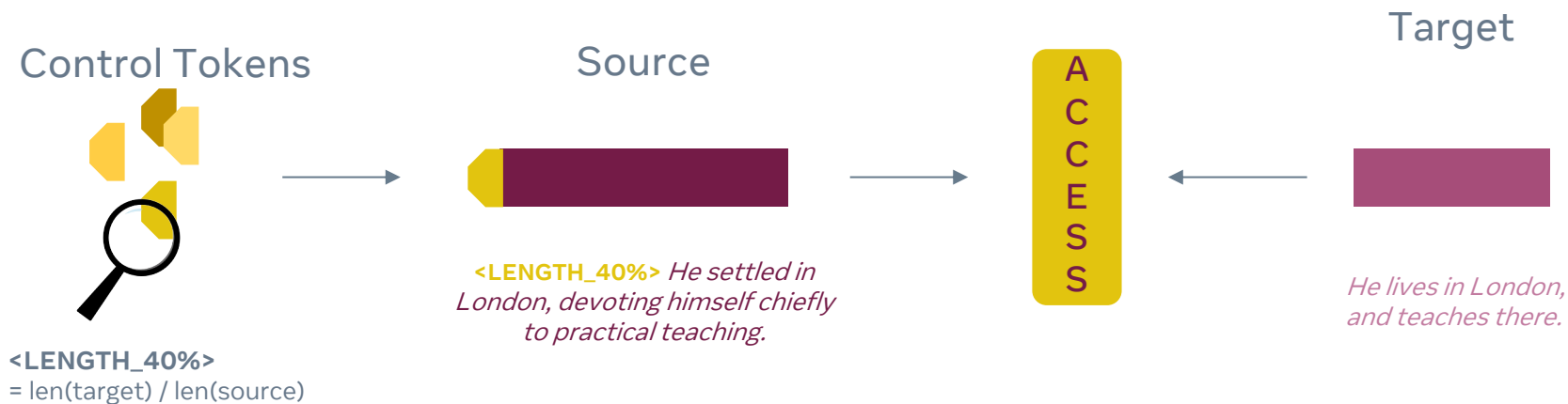(~500GB per language on disk)

Querying one sequence for Nearest Neighbors (NN)

Q?

Return top-k NN

Threshold with L2 distance in latent space

Repeat in parallel for 1B queries ⇒ 1M extracted paraphrases

# Use ACCESS for Controllable Generation

Web

Paraphrases

MINING

CONTROLLABLE
GENERATION

Multilingual
Pretraining

ACCESS

mBART

ADAPTED
CONTROL

SOURCE    C'est pourquoi il est recommandé d'éternuer ou de tousser dans le pli de son coude.

PREDICTION    Il est conseillé d'éternuer ou de tousser dans son coude.

# Conditioning on Control Tokens during Training

Control Tokens

Source

Target

A
C
C
E
S
S

**<LENGTH_40%>** *He settled in London, devoting himself chiefly to practical teaching.*

*He lives in London, and teaches there.*

**<LENGTH_40%>**
= len(target) / len(source)

# Choose Desired Length at Test Time

Tokens

Sentence

Prediction

A
C
C
E
S
S

**<LENGTH_25%>** *He settled in London, devoting himself chiefly to practical teaching.*

*He lives and teaches in London.*
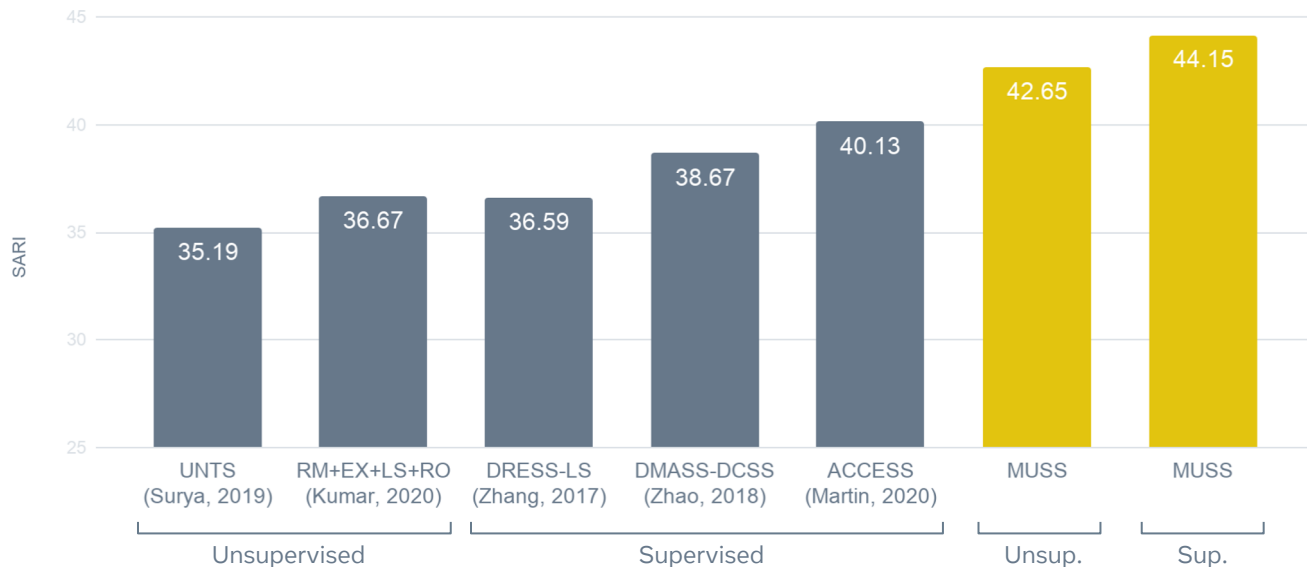
In practice: **find ratio** with best results on valid

# Condition on Many Attributes

- Length

- Lexical Complexity

- Syntactic Complexity

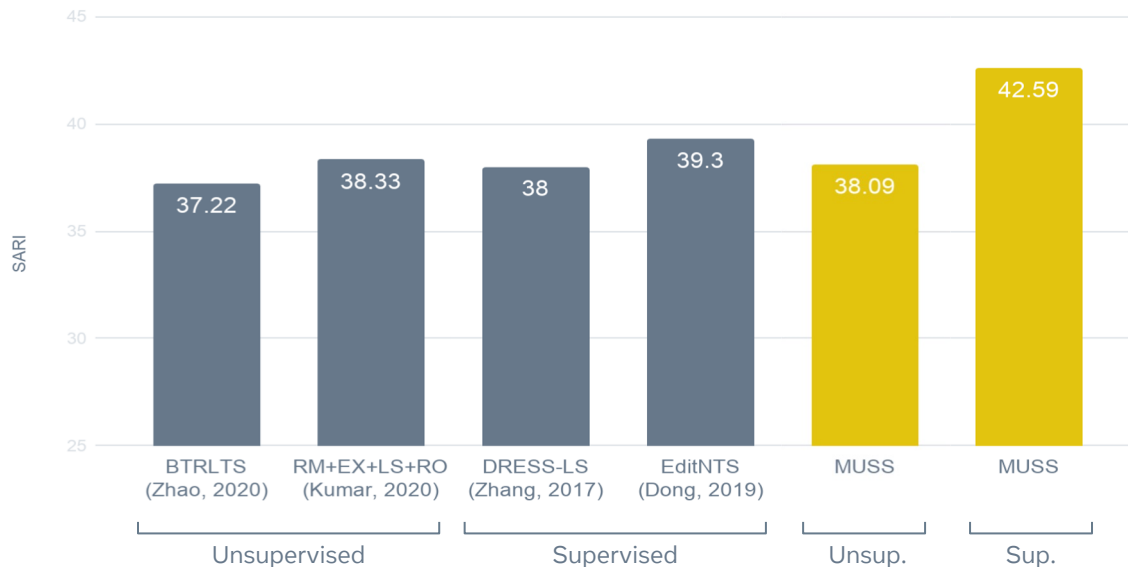- Amount of Paraphrasing

# The MUSS Approach



Web

Paraphrases

Multilingual Pretraining

*MINING*

*CONTROLLABLE GENERATION*

ACCESS

mBART

*ADAPTED CONTROL*

*SOURCE*    C'est pourquoi il est recommandé d'éternuer ou de tousser dans le pli de son coude.

*PREDICTION*    **Il** est conseillé d'éternuer ou de tousser dans son coude.

# English Results - ASSET



- MUSS improves over previous methods

- Incorporating labelled data improves further.

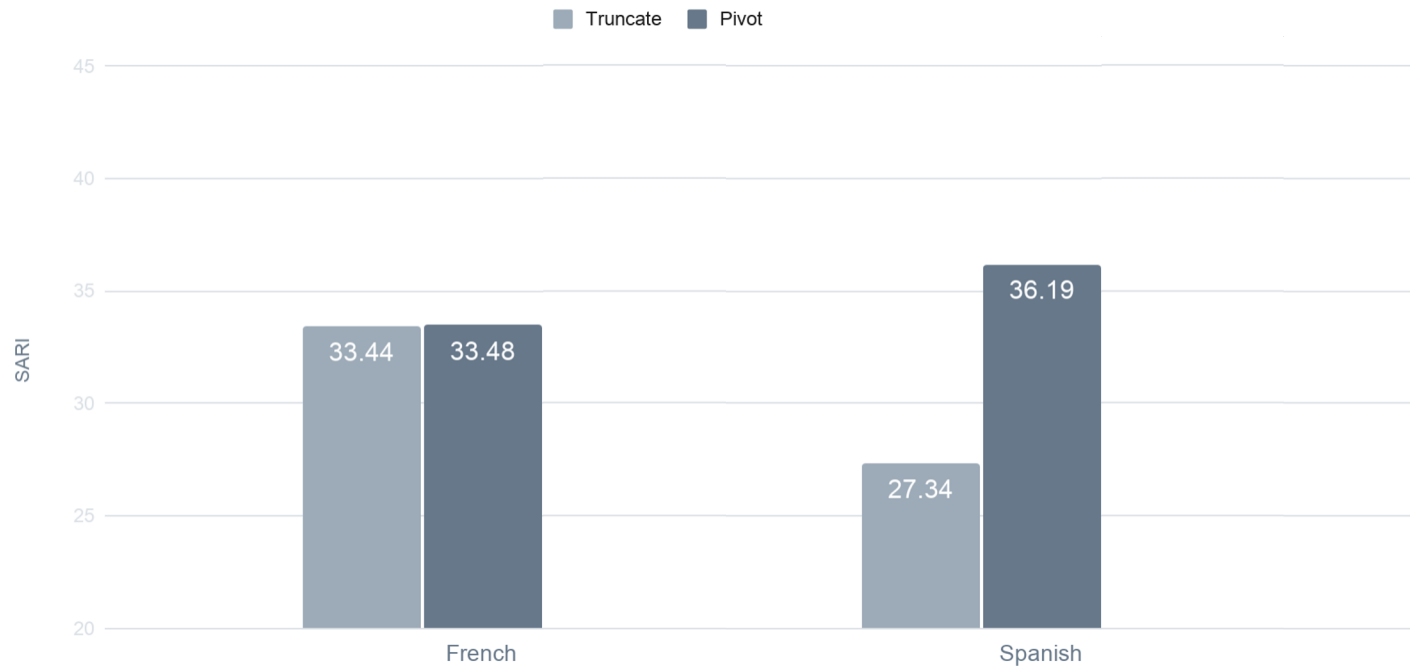# English Results - Newsela



- In-domain data important for professional News Corpus

# Multilingual Results

Baselines

- **Truncate**

    ○ Drop last 20% tokens.

- **Pivot**

    ○ **Fr**⇒**En** Translation - **En**⇒**En** Simplification - **En**⇒**Fr** Translation

# Multilingual Results



- Good results compared to strong baselines…

- But benchmarks are still imperfect
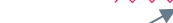
# MUSS Simplifications Example

## Correct Simplification

It is situated at the coast of the Baltic Sea, where it encloses the city of Stralsund.
It is located on the Baltic Sea. The city of Stralsund is located in it.

## Simplification Error

In 1998, Culver ran for Iowa Secretary of State and was victorious.
In 1998, Culver ran for Governor of Iowa and won.

**Different meaning**

# Conclusion

**MUSS**: Paraphrase Mining + Controllable Generation

- Fully unsupervised sentence simplification

- Can be applied in any language

# Perspectives

Towards document simplification

- How to mine full documents?

- Will controls be as effective?

Apply MUSS to other text rewriting tasks

- Paraphrasing

- Style transfer

- Summarization

# Thank You