

# An Interactive Machine Translation Framework for Modernizing the Language of Historical Documents

Miguel Domingo, Francisco Casacuberta

`midobal@prhlt.upv.es`, `fcn@prhlt.upv.es`

Pattern Recognition and Human Language Technology Research Center  
Universitat Politècnica de València

IbPRIA 2022

Aveiro, May 4, 2022

# Outline

1. Introduction
2. Language Modernization
3. Interactive Machine Translation
4. Experimental Framework
5. Results
6. Conclusions and Future Work

# Outline

1. Introduction
2. Language Modernization
3. Interactive Machine Translation
4. Experimental Framework
5. Results
6. Conclusions and Future Work

# Introduction

- Historical documents are part of our cultural heritage.
- However, due to their linguistic characteristics they are mostly limited to scholars.
- Language modernization helps non-experts to understand their content, but it is not error-free.
- The interactive machine translation framework can help scholars to generate error-free modernization.

# Outline

1. Introduction
2. Language Modernization
3. Interactive Machine Translation
4. Experimental Framework
5. Results
6. Conclusions and Future Work

# Language Modernization

**Goal:** make historical documents more accessible to a general audience.

# Language Modernization

**Goal:** make historical documents more accessible to a general audience.

## Original

To be, or not to be? That is the question  
Whether tis nobler in the mind to suffer  
The slings and arrows of outrageous fortune,  
Or to take arms against a sea of troubles,  
And, by opposing, end them?

## Modernized

The question is: is it better to be alive or dead?  
Is it nobler to put up  
with all the nasty things that luck throws your way,  
or to fight against all those troubles  
by simply putting an end to them once and for all?

# Language Modernization

## Approaches:

$$\hat{y}_1^I = \arg \max_{y_1^I} Pr(y_1^I | x_1^J)$$

- Statistical machine translation (SMT).
- Neural machine translation (NMT).
  - ▶ Recurrent neural networks with long short-term memory units (LSTM).
  - ▶ Transformer.

# Outline

1. Introduction
2. Language Modernization
3. Interactive Machine Translation
4. Experimental Framework
5. Results
6. Conclusions and Future Work

# Interactive Machine Translation

**Goal:** collaborative framework in which human and machine work together to produce the final high-quality translations.

# Interactive Machine Translation

Prefix-based interactive machine translation (IMT)

# Interactive Machine Translation

## Prefix-based interactive machine translation (IMT)

**Source:** Ealle ing he foresceawa and wt, and ealra eoda gereord he cann.

**Target translation:** All things he foresees and knows, and he understands the tongues of all nations.

All things he foresees and knows, and  
of all nations language he understands.

# Interactive Machine Translation

## Prefix-based interactive machine translation (IMT)

**Source:** Ealle ing he foresceawa and wt, and ealra eoda gereord he cann.

**Target translation:** All things he foresees and knows, and he understands the tongues of all nations.



All things he foresees and knows, and  
of all nations language he understands.

# Interactive Machine Translation

## Prefix-based interactive machine translation (IMT)

**Source:** Ealle ing he foresceawa and wt, and ealra eoda gereord he cann.

**Target translation:** All things he foresees and knows, and he understands the tongues of all nations.



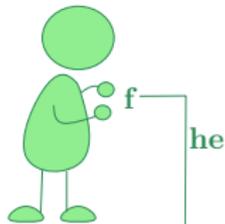
All things he foresees and knows, and  
[of]all nations language he understands.

# Interactive Machine Translation

## Prefix-based interactive machine translation (IMT)

**Source:** Ealle ing he foresceawa and wt, and ealra eoda gereord he cann.

**Target translation:** All things he foresees and knows, and he understands the tongues of all nations.



All things he foresees and knows, and

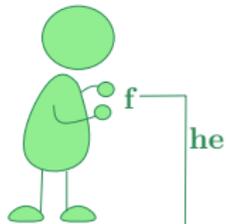
→ of all nations language he understands.

# Interactive Machine Translation

## Prefix-based interactive machine translation (IMT)

**Source:** Ealle ing he foresceawa and wt, and ealra eoda gereord he cann.

**Target translation:** All things he foresees and knows, and he understands the tongues of all nations.



All things he foresees and knows, and  
 of all nations language he understands.



All things he foresees and knows, and  
 he understands of all nations language.

# Interactive Machine Translation

## Prefix-based interactive machine translation (IMT)

Suffix generation:

$$\hat{y}_{i+1}^I = \arg \max_{I, y_{i+1}^I} Pr(y_{i+1}^I \mid x_1^J, f = \tilde{y}_1^i) = \arg \max_{I, y_{i+1}^I} Pr(\tilde{y}_1^i y_{i+1}^I \mid x_1^J)$$

# Interactive Machine Translation

## Prefix-based interactive machine translation (IMT)

Suffix generation:

$$\hat{y}_{i+1}^I = \arg \max_{l, y_{i+1}^I} Pr(y_{i+1}^I | x_1^J, f = \tilde{y}_1^I) = \arg \max_{l, y_{i+1}^I} Pr(\tilde{y}_1^I y_{i+1}^I | x_1^J)$$

MT fundamental equation:

$$\hat{y}_1^I = \arg \max_{l, y_1^I} Pr(y_1^I | x_1^J)$$

# Interactive Machine Translation

## Segment-based IMT

# Interactive Machine Translation

## Segment-based IMT

**Source:** Ealle ing he foresceawa and wt, and ealra eoda gereord he cann.

**Target translation:** All things he foresees and knows, and he understands the tongues of all nations.

All things he foresees and knows, and  
of all nations language he understands.

# Interactive Machine Translation

## Segment-based IMT

**Source:** Ealle ing he foresceawa and wt, and ealra eoda gereord he cann.

**Target translation:** All things he foresees and knows, and he understands the tongues of all nations.



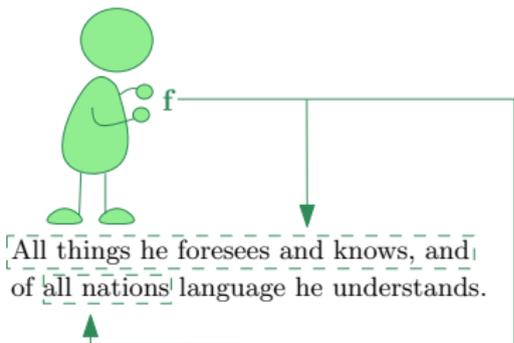
All things he foresees and knows, and  
of all nations language he understands.

# Interactive Machine Translation

## Segment-based IMT

**Source:** Ealle ing he foresceawa and wt, and ealra eoda gereord he cann.

**Target translation:** All things he foresees and knows, and he understands the tongues of all nations.

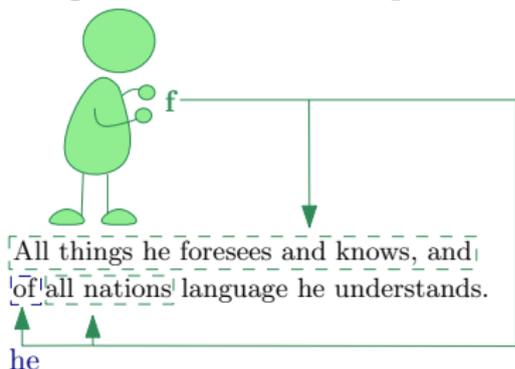


# Interactive Machine Translation

## Segment-based IMT

**Source:** Ealle ing he foresceawa and wt, and ealra eoda gereord he cann.

**Target translation:** All things he foresees and knows, and he understands the tongues of all nations.

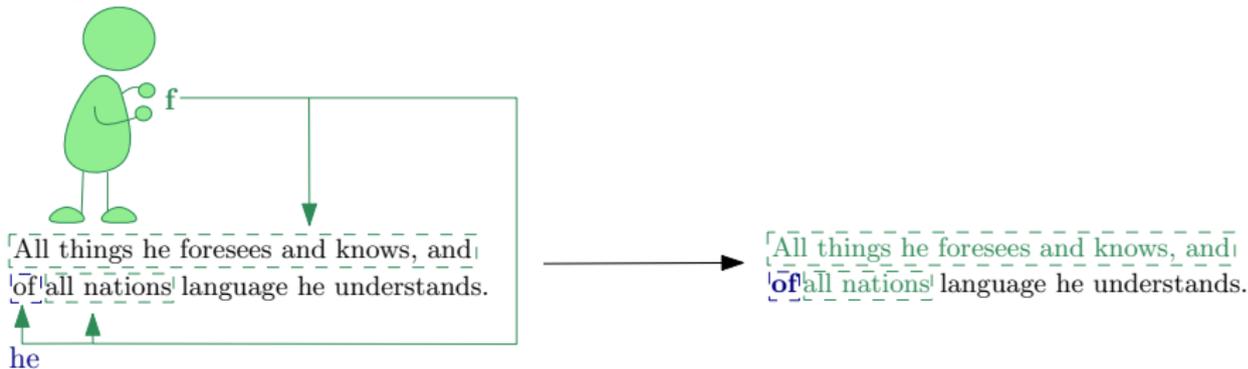


# Interactive Machine Translation

## Segment-based IMT

**Source:** Ealle ing he foresceawa and wt, and ealra eoda gereord he cann.

**Target translation:** All things he foresees and knows, and he understands the tongues of all nations.



# Outline

1. Introduction
2. Language Modernization
3. Interactive Machine Translation
4. Experimental Framework
5. Results
6. Conclusions and Future Work

## Corpora

- Dutch Bible (17<sup>th</sup> century Dutch; 30K segments).
- El Quijote (17<sup>th</sup> century Spanish; 10K segments).
- OE-ME (11<sup>th</sup> century English; 3K segments).

## Metrics

- Word Stroke Ratio (WSR).
- Mouse Action Ratio (MAR).
- Translation Error Rate (TER).
- BiLingual Evaluation Understudy (BLEU).
  
- We applied approximate randomization tests<sup>1</sup>, with 10,000 repetitions and using a  $p$ -value of 0.05.

---

<sup>1</sup>Stefan Riezler and John T Maxwell (2005). "On some pitfalls in automatic evaluation and significance testing for MT". In: *Proceedings of the workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pp. 57–64.

# Outline

1. Introduction
2. Language Modernization
3. Interactive Machine Translation
4. Experimental Framework
5. Results
6. Conclusions and Future Work

## Results

Corpus	Approach	Modernization quality		Prefix-based		Segment-based	
		TER [↓]	BLEU [↑]	WSR [↓]	MAR [↓]	WSR [↓]	MAR [↓]
Dutch Bible	SMT	11.5	77.5	14.3	4.4	<b>9.0</b>	<b>10.8</b>
	NMT <sub>LSTM</sub>	50.7 <sup>†</sup>	43.4	42.6 <sup>‡</sup>	9.2	42.6 <sup>‡</sup>	50.9
	NMT <sub>Transformer</sub>	50.3 <sup>†</sup>	35.8	49.2 <sup>‡</sup>	10.4	49.2 <sup>‡</sup>	48.3
El Quijote	SMT	30.7	58.3	38.8	10.9	<b>22.0</b>	<b>19.7</b>
	NMT <sub>LSTM</sub>	42.9	50.4	68.9 <sup>‡</sup>	11.8	68.9 <sup>‡</sup>	47.8
	NMT <sub>Transformer</sub>	47.3	46.1	73.2 <sup>‡</sup>	13.4	73.2 <sup>‡</sup>	50.5
OE-ME	SMT	39.6	39.6	58.2	15.5	<b>28.2</b>	<b>26.1</b>
	NMT <sub>LSTM</sub>	56.4	30.3	72.1 <sup>‡</sup>	12.8 <sup>†</sup>	72.1 <sup>‡</sup>	59.5
	NMT <sub>Transformer</sub>	58.9	28.2	73.5 <sup>‡</sup>	13.3 <sup>†</sup>	73.5 <sup>‡</sup>	49.5

All results are significantly different between all approaches except those denoted with <sup>†</sup>. Given the same approach, all results are significantly different between the different IMT protocols except those denoted with <sup>‡</sup>.

## Examples

**source (x):** Ealle ing he foresceawa and wt, and ealra eoda gereord he cann.

**target translation (y):** All things he foresees and knows, and he understands the tongues of all nations.

IT-0	<b>System</b>	All things he foresceawa and knows, and of all nations language he understands.
IT-1	<b>User</b>	All things he foresees and knows, and of all nations language he understands.
	<b>System</b>	All things he foresees and knows, and of all nations language he understands.
IT-2	<b>User</b>	All things he foresees and knows, and he all nations language he understands.
	<b>System</b>	All things he foresees and knows, and he understands of all nations language.
IT-3	<b>User</b>	All things he foresees and knows, and he understands the all nations language.
	<b>System</b>	All things he foresees and knows, and he understands the beginning of all nations language.
IT-4	<b>User</b>	All things he foresees and knows, and he understands the tongues of all nations language.
	<b>System</b>	All things he foresees and knows, and he understands the tongues all.
IT-5	<b>User</b>	All things he foresees and knows, and he understands the tongues of
	<b>System</b>	All things he foresees and knows, and he understands the tongues of all.
IT-6	<b>User</b>	All things he foresees and knows, and he understands the tongues of all nations.
	<b>System</b>	All things he foresees and knows, and he understands the tongues of all nations.
END	<b>User</b>	All things he foresees and knows, and he understands the tongues of all nations.

Example of a prefix-based IMT session.

## Examples

**source (x):** Ealle ing he foresceawa and wt, and ealra eoda gereord he cann.

**target translation (y):** All things he foresees and knows, and he understands the tongues of all nations.

IT-0	System	All things he foresceawa and knows, and of all nations language he understands.
IT-1	User	All things he foresees and knows, and of all nations language he understands.
	System	All things he foresceawa foresees and knows, and language he understand of all nations .
IT-2	User	All things he foresees and knows, and he understand the of all nations .
	System	All things he foresees and knows, and he understand the language of all nations .
IT-3	User	All things he foresees and knows, and he understand the tongues of all nations .
	System	All things he foresees and knows, and he understand the tongues of all nations .
END	User	All things he foresees and knows, and he understands the tongues of all nations.

Example of a segment-based IMT session.

# Outline

1. Introduction
2. Language Modernization
3. Interactive Machine Translation
4. Experimental Framework
5. Results
6. Conclusions and Future Work

## Conclusions and Future Work

- We deployed the interactive framework into language modernization to help scholar to generate error-free modernizations.
- We deployed two different protocols: SMT and NMT.
- While IMT always succeeds in reducing the human effort, the SMT approach yielded the best results.
- We evaluated our proposal under simulated conditions. In a future work, we would like to perform a human evaluation with the help of scholars.