

SEGMENT-BASED INTERACTIVE MACHINE TRANSLATION

Authors: Miguel Domingo, Álvaro Peris Supervisor: Francisco Casacuberta

Ph.D. Program: Ph.D. in Computer Science

PRHLT Research Center {midobal, lvapeab, fcn}@prhlt.upv.es

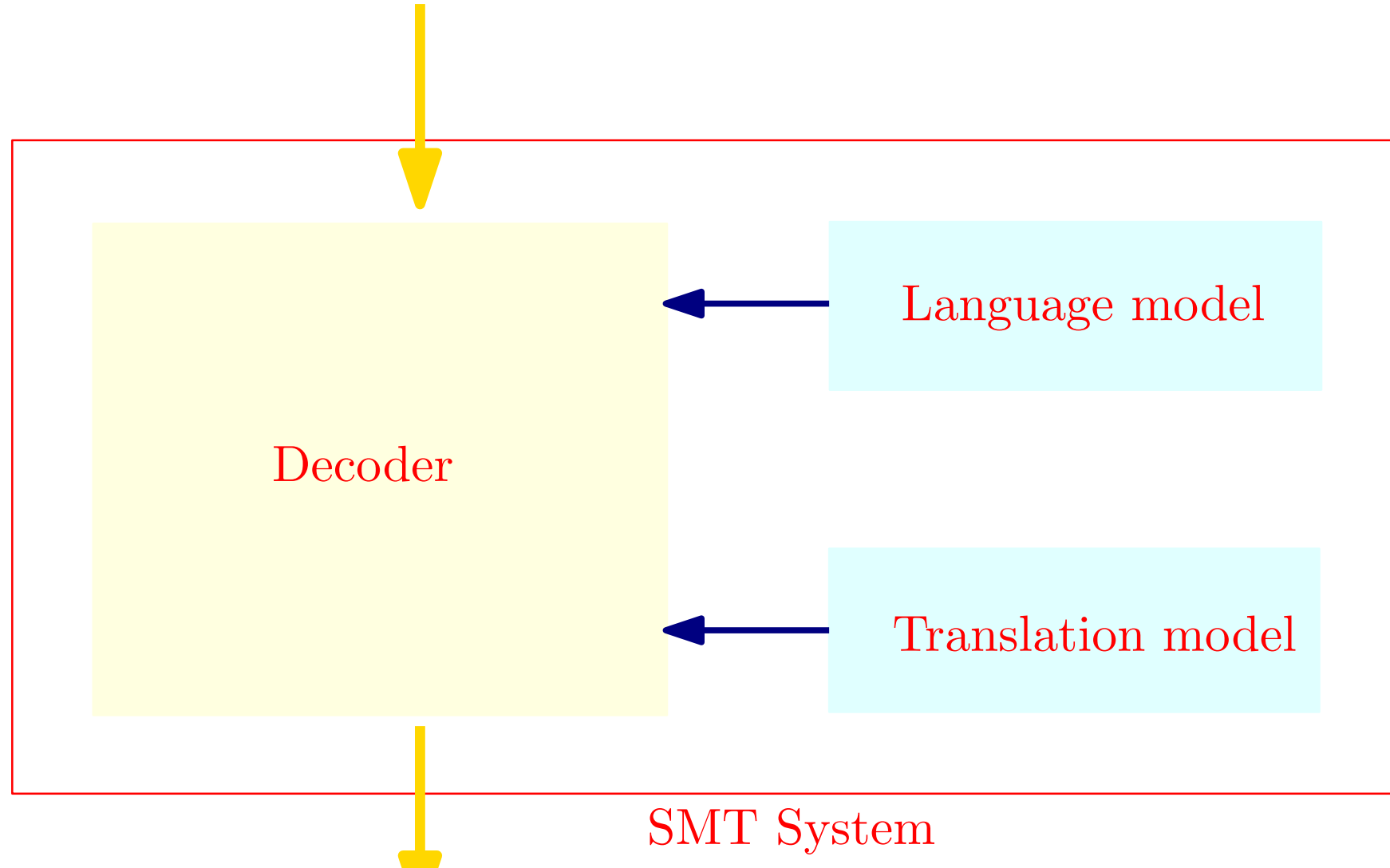


INTRODUCTION

- Statistical Machine Translation (SMT) is still not good enough.
- Human correction is needed to improve translation quality.
- Interactive Machine Translation (IMT) aims at reducing the effort of this process.
- Prefix-based IMT was an interesting contribution to the field.

STATISTICAL MACHINE TRANSLATION

Source: Et la question n' a pas encore été évaluée chez les patients atteints de cancer gastrique



Translation: And the issue has not been investigated among patients with gastric cancer

SMT main equation:

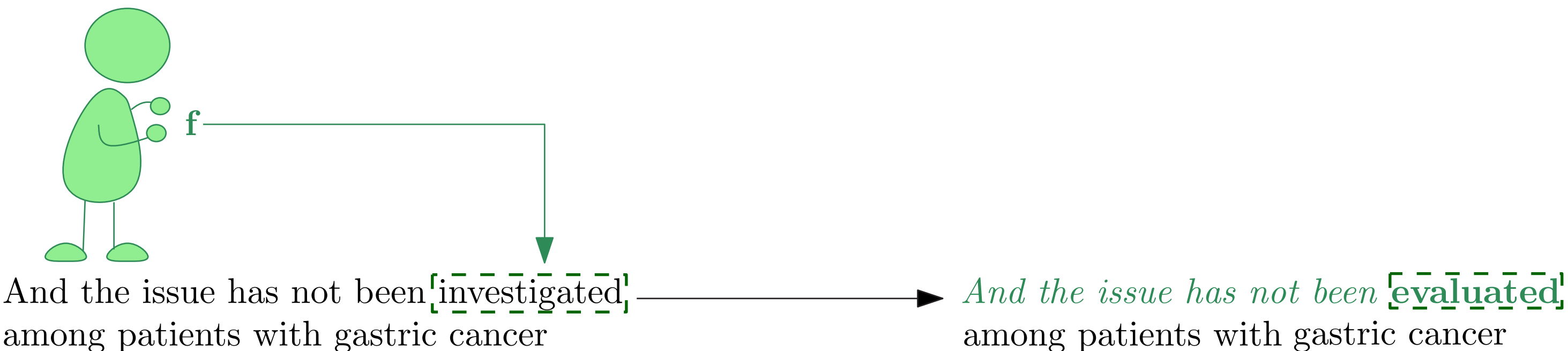
$$\tilde{y} = \arg \max_y Pr(y) \cdot Pr(x | y)$$

x: source y: translation
Pr(y): language model Pr(x | y): translation model

PREFIX-BASED IMT

Source: Et la question n' a pas encore été évaluée chez les patients atteints de cancer gastrique

Target translation: And the issue has not been evaluated in gastric cancer patients



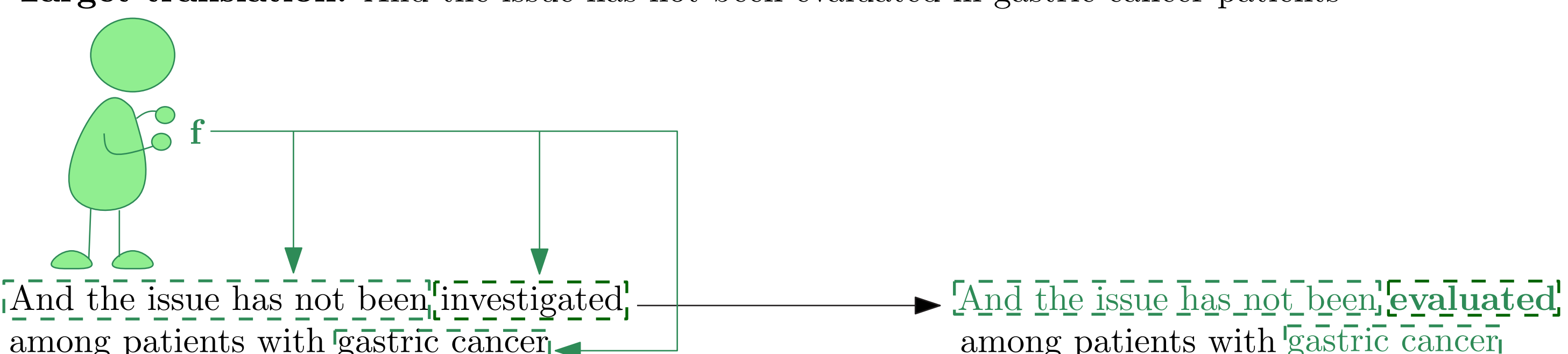
Prefix-based IMT session:

0	S	And the issue has not been investigated among patients with gastric cancer
1	U	And the issue has not been evaluated among patients with gastric cancer
	S	<i>And the issue has not been evaluated</i> with patients
2	U	<i>And the issue has not been evaluated</i> in patients
	S	<i>And the issue has not been evaluated in</i> patients
3	U	<i>And the issue has not been evaluated in</i> gastric
	S	<i>And the issue has not been evaluated in</i> gastric cancer patients
E	U	And the issue has not been evaluated in gastric cancer patients

SEGMENT-BASED IMT

Source: Et la question n' a pas encore été évaluée chez les patients atteints de cancer gastrique

Target translation: And the issue has not been evaluated in gastric cancer patients



Segment-based IMT session:

0	S	And the issue has not been investigated among patients with gastric cancer
1	U	<u>And the issue has not been</u> investigated among patients with <u>gastric cancer</u>
	S	<u>And the issue has not been</u> evaluated among patients with <u>gastric cancer</u>
	S	<u>And the issue has not been</u> <u>evaluated</u> in <u>gastric cancer</u> patients
E	U	And the issue has not been evaluated in gastric cancer patients

STATISTICAL FRAMEWORK

SMT fundamental equation:

$$\tilde{y} = \arg \max_y Pr(y | x)$$

Segment-based IMT:

$$\tilde{h}_1, \dots, \tilde{h}_N = \arg \max_{h_1, \dots, h_N} Pr(\hat{f}_1 h_1, \dots, \hat{f}_N h_N | x)$$

$\hat{f}_1, \dots, \hat{f}_N$: sequence of N segments validated by the user (feedback signal).

$\tilde{h}_1, \dots, \tilde{h}_N$: sequence of new translation segments.

$$\tilde{y} = \hat{f}_1 \tilde{h}_1, \dots, \hat{f}_N \tilde{h}_N$$

Search in the space of the translation, constrained by the sequence of segments $\hat{f}_1, \dots, \hat{f}_N$.

EXPERIMENTS

- Simulated environment.
- Corpora from different domains.
- Comparison of prefix-based against segment-based IMT.
- Measuring human effort (WSR and MAR)

RESULTS

Corpus	Language	BLEU	Prefix-Based		Segment-Based	
			WSR (%)	MAR (%)	WSR (%)	MAR (%)
EMEA	Fr-En	31.3	57.8	12.4	34.4	18.8
	En-Fr	30.2	58.4	12.5	40.4	16.3
EU	Es-En	48.2	45.6	10.2	28.3	15.0
	En-Es	48.7	44.6	9.7	29.8	13.5
TED	Zh-En	11.7	83.1	22.4	54.1	28.3
	En-Zh	8.7	86.3	55.7	59.2	72.4
Xerox	Es-En	54.5	35.8	10.5	23.2	16.9
	En-Es	62.2	28.3	7.9	22.1	12.5

- Substantial reduction of the typing effort (up to 29 points of WSR).
- Slight increase in the number of mouse actions (from 4 up to 6.5 points of MAR).

CONCLUSIONS

- New IMT approach that breaks down the prefix constraint.
- The user can select all correct words from each translation hypothesis.
- Human effort effectively reduced in a simulated environment.
- Future work: experiments with real users.

ACKNOWLEDGEMENTS

The research leading to these results has received funding from the Ministerio de Economía y Sostenibilidad (MINECO) under project SmartWays (grant agreement RTC-2014-1466-4), and Generalitat Valenciana under project ALMAMATER (grant agreement PROMETEOII/2014/030).