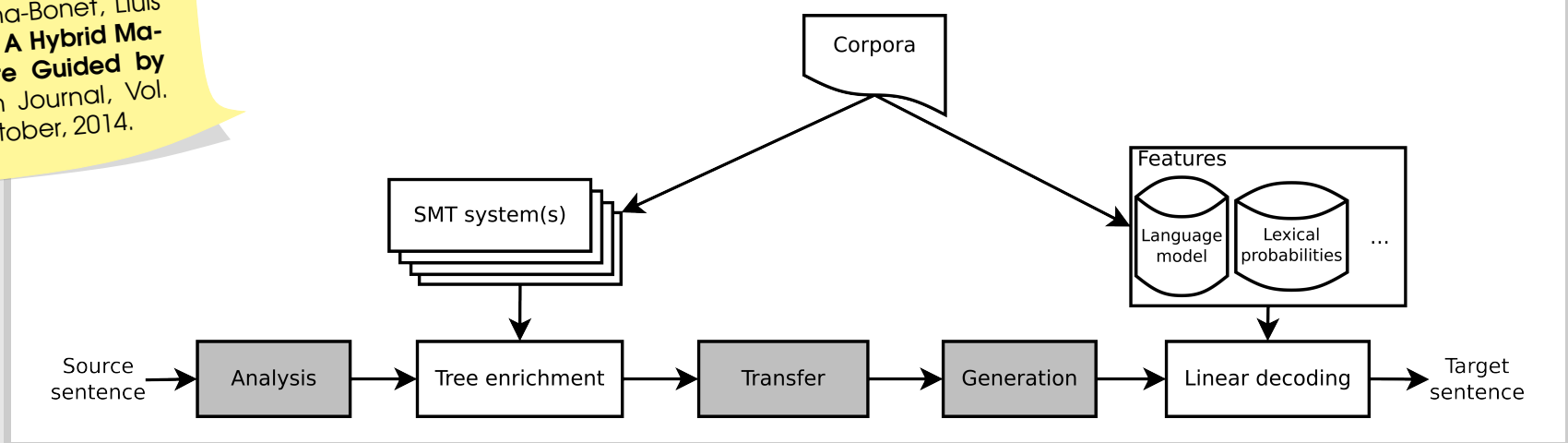


# Journey through Natural Language Processing

## Hybrid Machine Translation

### SMatxinT, a HMT Architecture Guided by Syntax

Gorka Labaka, Cristina España-Bonet, Lluís Màrquez and Kepa Sarasola. **A Hybrid Machine Translation Architecture Guided by Syntax**. Machine Translation Journal, Vol. 28, issue 2, pages 91–125, October, 2014.



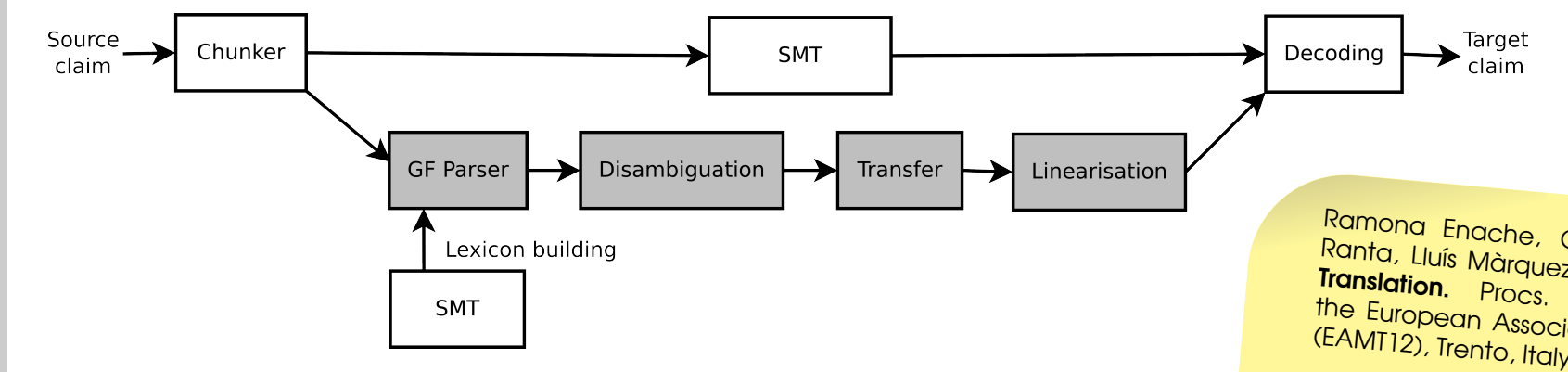
	TER	BLEU	NIST	GTM-2	MTR-st	RG-S*	ULC
Matxin	73.57	14.29	6.05	22.62	20.27	15.90	41.02
SMT <sub>b</sub>	68.49	15.93	6.45	23.64	21.59	16.39	44.56
SMT <sub>g</sub>	68.93	15.21	6.43	23.44	21.84	17.66	44.20
SMatxinT (m)	<b>66.70</b>	<b>17.14</b>	<b>6.72</b>	<b>24.58</b>	<b>22.52</b>	<b>18.66</b>	<b>48.00</b>
SMatxinT (r)	<b>66.63</b>	<b>17.18</b>	<b>6.73</b>	<b>24.59</b>	<b>22.51</b>	<b>18.52</b>	<b>48.03</b>

- Candidate chunk translations are calculated with an SMT system and used to enrich the RBMT tree-based source representation with more alternatives

- The most probable combination among the available fragments is obtained with a monotone statistical decoding (m) following the order provided by the RBMT system

- The evaluation on news (out-of-domain) shows a preference for SMatxinT. There is no necessity for reordering the RBMT choice (r)

### MOLTO, a HMT for Patent Translation



Ramona Enache, Cristina España-Bonet, Aarne Ranta, Lluís Màrquez. **A Hybrid System for Patent Translation**. Procs. 16th Annual Conference of the European Association for Machine Translation (EAMT12), Trento, Italy, May 8-30, 2012.

- A GF translator (RBMT) is built for the specific domain which, in turn, uses an in-domain SMT system to build its lexicon

- Another SMT system is on top of the GF to translate those phrases not covered by the grammar. Hard integration (HI) vs. Soft integration (SI)

- Evaluations consistently show a preference for the SI hybrid system in front of the two individual translators

	TER	BLEU	NIST	GTM-2	MTR-p	RG-S*	ULC
GF	58.90	26.56	5.57	22.74	38.76	29.00	16.17
SMT	25.32	63.18	9.99	44.58	71.64	72.65	67.14
HI	31.24	55.88	9.24	38.81	67.30	67.80	58.84
SI1.0	25.10	63.56	10.02	<b>44.86</b>	<b>71.96</b>	72.89	67.56
SI0.5	<b>25.02</b>	<b>63.60</b>	<b>10.03</b>	44.84	71.94	<b>72.93</b>	<b>67.60</b>
SI0.0	25.36	63.15	9.99	44.54	71.60	72.66	67.11

## Document-Level Machine Translation

- Using a document-level (DL) decoder and considering a feature function that rewards coherent translations
- Capturing semantic information in corpora with bilingual word vector models that can be integrated as semantic space language models (biSSM) at translation time

	TER	BLEU	NIST	MTR-p	RG-S*	SP-Op	ULC
DL	72.83	28.33	7.46	23.22	30.36	19.38	77.14
DL + monoSSM	72.61	28.48	7.52	23.28	30.33	19.61	77.49
DL + biSSM	<b>72.56</b>	<b>28.58</b>	<b>7.66</b>	<b>23.31</b>	<b>30.38</b>	<b>19.78</b>	<b>77.89</b>

En2Es on News

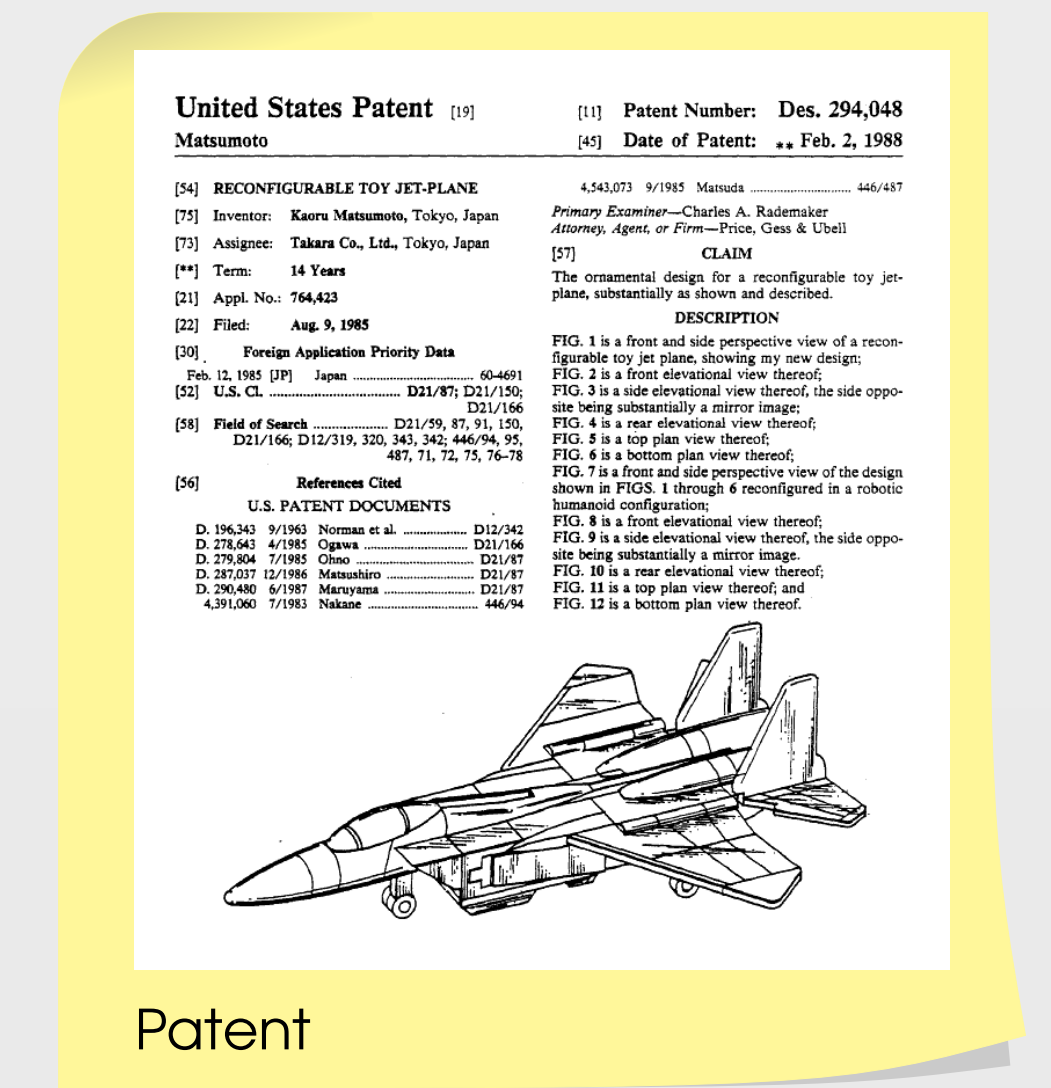
Eva Martínez García, Cristina España-Bonet, Lluís Màrquez. **The UPC TweetMT participation: Translating Formal Tweets using Context Information**. Procs. 'XXI Congreso de la Sociedad Española de Procesamiento de lenguaje natural', Alacant, Spain, September 2015.

	TER	BLEU	NIST	GTM-2	MTR-e	RG-S*	OI	ULC
SMT	13.56	78.06	12.04	<b>73.02</b>	54.09	82.21	83.31	66.62
DL	<b>13.44</b>	<b>78.19</b>	<b>12.07</b>	72.97	<b>54.15</b>	<b>82.45</b>	<b>83.48</b>	<b>67.11</b>

Es2Ca on Tweets



### Piece of news



### Patent



### Wikipedia article



### Tweets

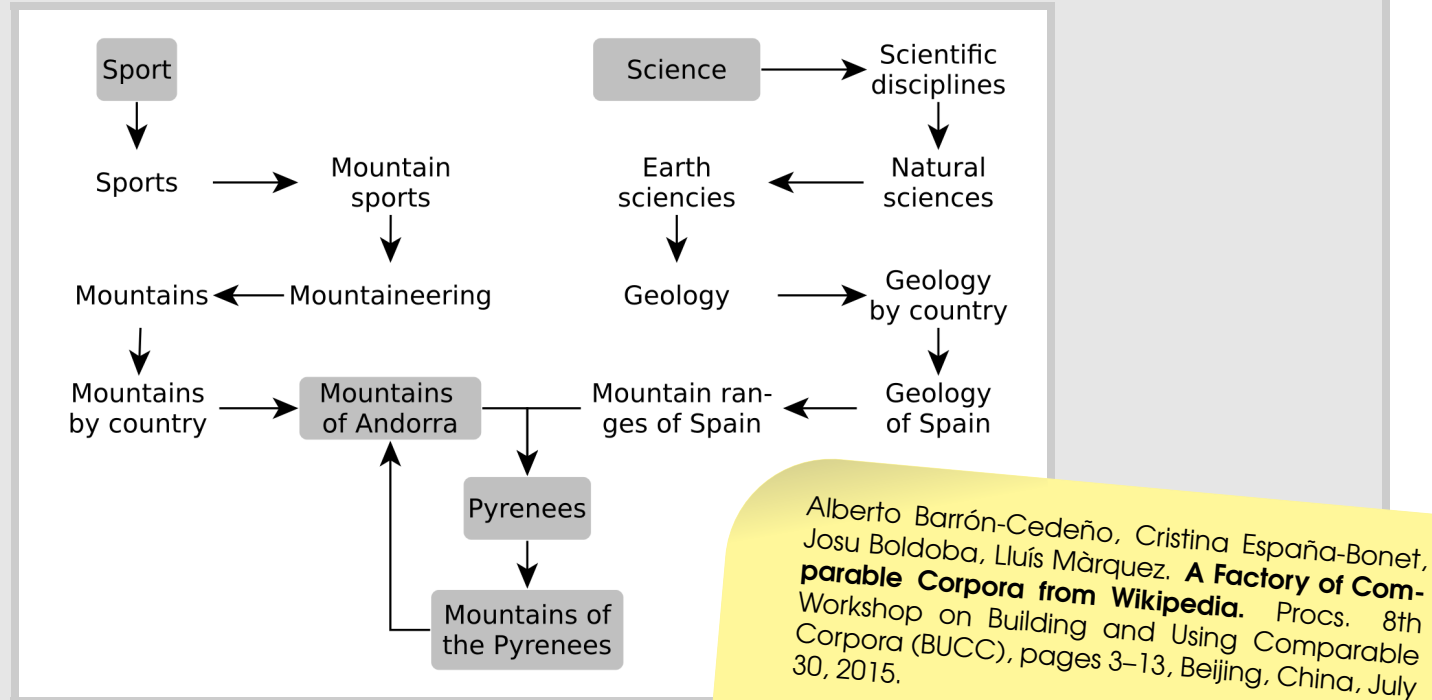
## Cristina España-Bonet

Universitat Politècnica de Catalunya – BarcelonaTech

## Comparable Corpora from Wikipedia

### Domain Corpora

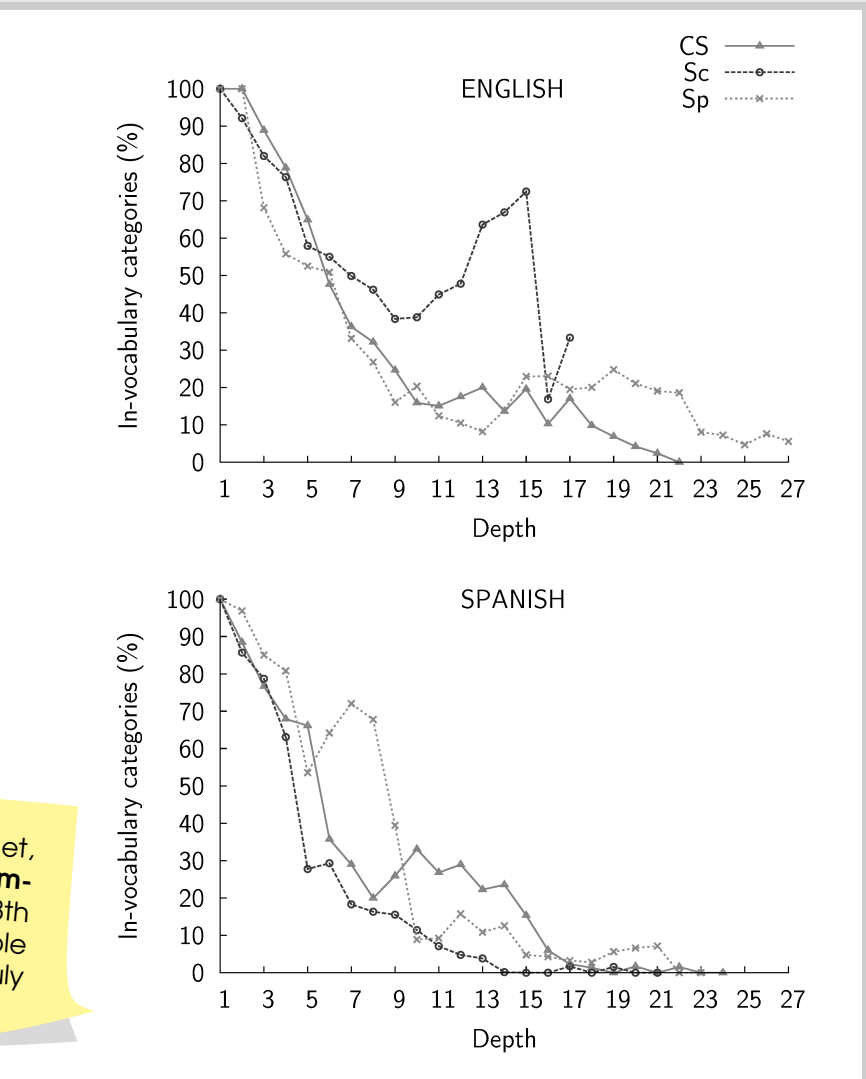
Given a domain, explore Wikipedia's category graph as a tree —breadth-first search departing from the main category visiting nodes only once



- Define an in-domain vocabulary (MFTs in root articles)
- Stopping criterion: levels with a minimum percentage of categories with in-domain vocabulary in the title

### Comparable Corpora

- Articles are related across languages via *interlanguage links*
- In-domain CC: monolingual extracted, union or intersection
- Usefulness of the corpora confirmed for training SMT systems



	Articles		Depth	
	50% en-es	60% en-es	50% en es	60% en es
CS	18,168	8,251	6	5
Sc	161,130	21,459	6	4
Sp	72,315	1,980	8	3

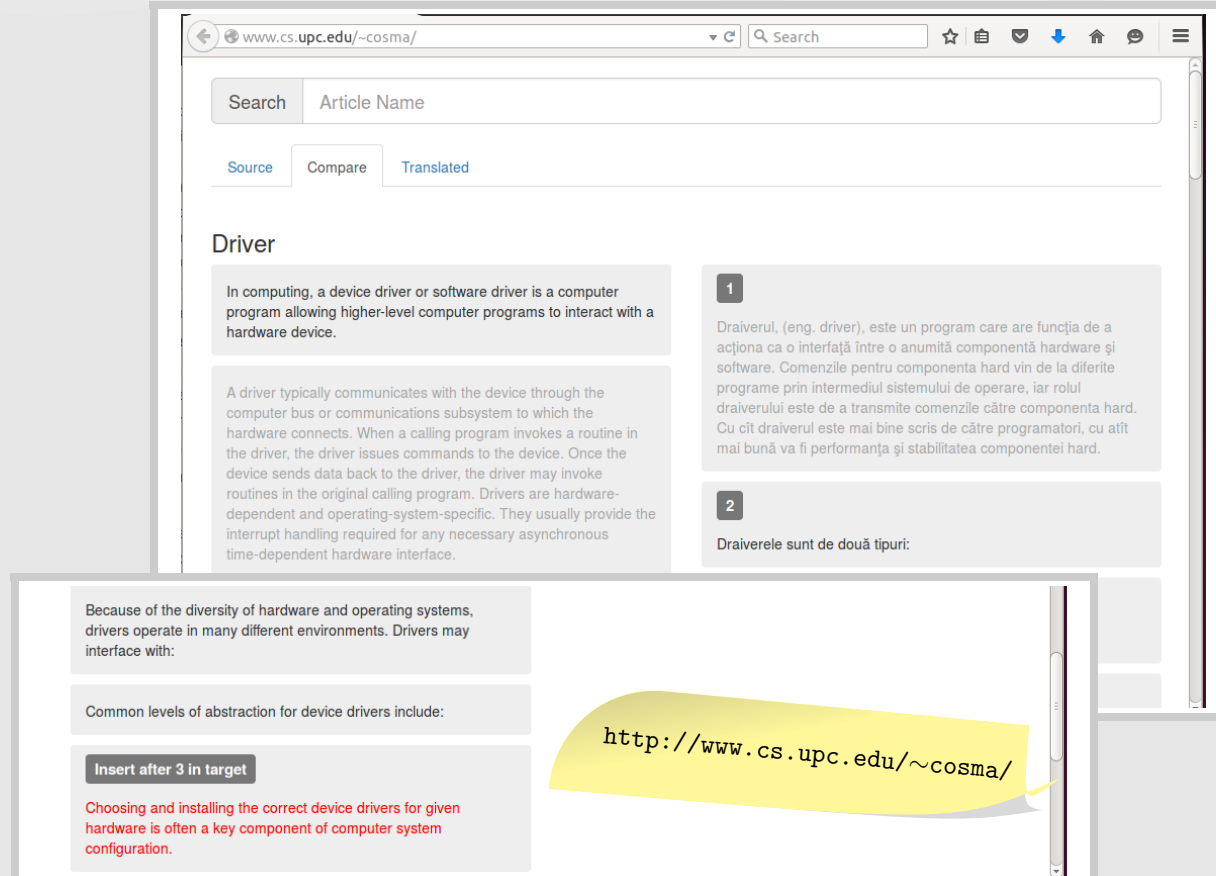
CS: Computer Science, Sc: Science, Sp: Sports

### Wikiparable, CC Extraction (ongoing work)

	Graph-based			IR-based		
	#Arts/Cat	Depth	$\rho$ Spear	#Arts/Cat	$\rho$ Spear	
1 English	50,514±121,881	5.9±2.8	0.38±0.20	61,239±73,248	0.27±0.15	
3 French	8,278±26,483	4.3±1.9	0.39±0.23	18,158±17,871	0.15±0.13	
4 Spanish	6,638±17,050	4.4±2.1	0.38±0.22	21,490±19,605	0.18±0.14	
2 German	2,752±9,573	3.4±1.9	0.42±0.25	12,887±18,876	0.13±0.13	
6 Arabic	2,999±9,546	3.6±2.3	0.45±0.28	4,622±4,082	0.12±0.13	
7 Romanian	1,398±8,683	3.4±1.8	0.42±0.26	1,750±1,839	0.06±0.13	
5 Catalan	1,140±4,693	3.3±1.9	0.44±0.22	5,959±5,058	0.14±0.13	
8 Basque	440±1,654	3.1±1.5	0.48±0.25	1,819±2,732	0.13±0.13	
9 Greek	356±1,982	2.8±1.6	0.55±0.22	1,604±2,378	0.16±0.16	
10 Occitan	104±598	2.4±1.3	0.51±0.29	419±2,040	0.14±0.17	

Mean values over 743 categories

Results for the most restrictive model for each architecture



## Wikipedia Enrichment

Prototype for including relevant information into a Wikipedia article in language  $L_1$  from the same article in  $L_2$

1. Paragraph alignment (similarity measures)
2. Relevant paragraph identification (from unaligned paragraphs)
3. Determination of the insertion position into the target
4. Preliminary translation of the relevant paragraphs by a translation engine