

# WP5

## Statistical and Robust Translation

Lluís Màrquez  
Cristina España-Bonet

Universitat Politècnica de Catalunya, TALP Research Center

– 1st year review –

Luxembourg, March 15th, 2011

- 1 General view
- 2 Ongoing work
- 3 Future work
- 4 Dissemination

# General view

## *Goal*

Extension of the grammar-based translation methods to widen their coverage and quality in unconstrained text translation.

Extension of the grammar-based translation methods to widen their coverage and quality in unconstrained text translation.

Especially *related to*:

**WP2** Grammar-based translation method

**WP7** Quasi-unconstrained domain, patents

**WP9** Evaluation

UPC

38

SMT technology, hybrid models, corpora processing, evaluation

**UPC**

**38**

SMT technology, hybrid models, corpora processing, evaluation

**UGOT**

**9**

Probabilistic extension of GF, synthetic corpora for SMT

**UPC**

**38**

SMT technology, hybrid models, corpora processing, evaluation

**UGOT**

**9**

Probabilistic extension of GF, synthetic corpora for SMT

**UHEL**

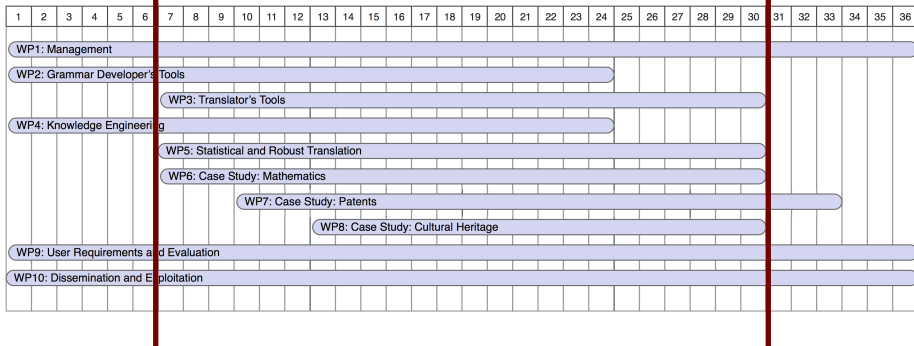
**6**

Usability and evaluation of the combined system

# General view

## Timeline

6 < month < 31





**Month 18** — Month 24 — Month 30

### **MS5**

First prototypes of the *baseline* combination models.

### **D51**

Description of the final collection of corpora.

Month 18 — **Month 24** — Month 30

### **MS7**

First prototypes of hybrid combination models.

### **D52**

Description and evaluation of the combination prototypes.

Month 18 — Month 24 — **Month 30**

### **MS8**

Translation tool complete.

### **D53**

WP5 final report: statistical and robust MT.

- 1 General view
- 2 Ongoing work**
  - Scheduled plan
  - Baselines
  - Hybrid systems
- 3 Future work
- 4 Dissemination

- Compilation and annotation of corpora from the patents domain.
- Training and adaptation of the base SMT systems.
- Statistical extension of the patents GF grammar.
- Evaluation and comparison of GF, SMT and cascade systems (baselines) in real domain data.
- First experiments with the combination approaches.

WP5 is tightly connected to **WP7** (Case of study: Patents).

### **Consequences:**

- An obvious delay in corpora compilation and annotation.
- Change of approach: from optimising base systems to dig into the hybrid system.



Mainly, just a change of order in tasks.

### **SMT baseline, Standard In-Domain System**

- **Corpus:** WP7 selected corpus
- **Language model:** 5-gram interpolated Kneser-Ney discounting, SRILM Toolkit
- **Alignments:** GIZA++ Toolkit
- **Translation model:** Moses package
- **Weights optimization:** MERT against BLEU
- **Decoder:** Moses

# Ongoing work

*SMT baseline, evaluation*

## BLEU

	EN2DE	DE2EN	EN2FR	FR2EN	DE2FR	FR2DE
<b>Bing</b>	0.33	0.43	0.43	0.45	0.20	0.24
<b>Google</b>	0.45	0.58	0.53	0.62	0.43	0.39
<b>Domain</b>	<b>0.58</b>	<b>0.65</b>	<b>0.62</b>	<b>0.70</b>	<b>0.56</b>	<b>0.53</b>



# Ongoing work

*SMT baseline, deep evaluation*

METRIC	DE2EN			EN2DE		
	Bing	Google	Domain	Bing	Google	Domain
1-WER	0.52	0.64	<b>0.72</b>	0.42	0.51	<b>0.69</b>
1-PER	0.66	0.76	<b>0.82</b>	0.56	0.64	<b>0.77</b>
1-TER	0.59	0.67	<b>0.76</b>	0.45	0.53	<b>0.71</b>
BLEU	0.43	0.58	<b>0.65</b>	0.33	0.45	<b>0.58</b>
NIST	8.25	9.67	<b>10.12</b>	6.53	8.05	<b>9.40</b>
ROUGE-W	0.40	0.48	<b>0.52</b>	0.34	0.41	<b>0.48</b>
GTM-2	0.30	0.40	<b>0.47</b>	0.25	0.32	<b>0.43</b>
METEOR-pa	0.60	0.69	<b>0.74</b>	0.36	0.45	<b>0.57</b>
<b>ULC</b>	0.09	0.29	<b>0.49</b>	0.03	0.19	<b>0.43</b>

# Ongoing work

*Two hybridisation approaches: Who leads?*

## 1. Integration led by **SMT**

Make available GF translations to a SMT system.

## 2. Integration led by **GF**

Complement with SMT options the GF translation structure.

### 1. Integration led by **SMT**

Make available GF translations to a SMT system.

- If **GF is able to generate Giza-like alignments**, phrases can be extracted in the SMT way and we can combine translation tables.

### 2. Integration led by **GF**

Complement with SMT options the GF translation structure.

- GF needs **robust and probabilistic parsing** for out of coverage content if has to be applied to open-domain text.

### **From many-to-many to one-to-many**

You want\_to\_go to the\_nearest park  
(0) (1) (2) (3) (4)

Quieres ir al parque mas cercano  
(0) (1)(2) (3) (4) (5)

1-0 1-1 2-2 3-4 3-5 4-3

(alignments from Phrasebook grammar)

### **Phrasebook grammar** (toy example)

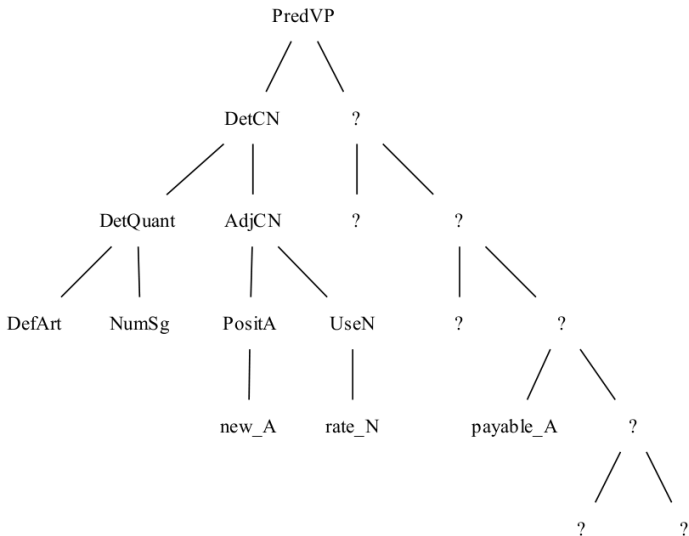
- Syntetic corpus generation
- Parallel corpus with 200 sentences
- Tiny for SMT
- Null intersection with SMT corpora

### **Patents grammar**

- Needed for real experiments (planned within next 6 months)

# Ongoing work

## *Hybridisation 2: robust parsing*



**Chunk parsing.** Detect:

- Basic noun phrases i.e. without PP attachment
- Verb phrases without the object
- Prepositions - mark the PP attachments

**Experiment** on parsing noun phrases from PennTreebank:

- 75% of the phrases were parsed
- Must improve NE and dates

- 1 General view
- 2 Ongoing work
- 3 Future work**
  - Related to the baselines
  - Related to the hybridisation
- 4 Dissemination



# Future work

## *Related to the baselines*

- Estimate a **GF baseline** on the test sets defined in WP7.
- **Naïve combination** of GF and SMT as a hybrid baseline.
- **Evaluation** of both systems and comparison with the SMT baseline.

- **Hard integration GF+SMT**  
Force fixed GF translations within a SMT system.
- Application of the soft integration **GF+SMT led by SMT** to the patents case, and extension with probabilistic estimations for GF phrases.
- Improvement of the robust parser and implementation of a soft integration **led by GF**.
- A first automatic **evaluation** of the resulting systems.

- 1 General view
- 2 Ongoing work
- 3 Future work
- 4 Dissemination**

### Talks

- *Online Parsing, Type Checking and Advanced Editor for Controlled Languages in GF.* Krasimir Angelov  
MOLTO's second meeting. 8 September, 2010, Varna.
- *Soft integration SMT/GF*  
Cristina España-Bonet and Lluís Màrquez  
MOLTO's internal workshop. 1-5 November, 2010, Chalmers University of Technology, Goteborg.
- *A TAG formalism for Parsing and Translation*  
Xavier Carreras  
MOLTO's internal workshop. 1-5 November, 2010, Chalmers University of Technology, Goteborg.

## Related reports

- *SMatxinT, the Spanish-to-Basque hybrid translator*  
Cristina España-Bonet, Gorka Labaka, Lluís Màrquez and  
Kepa Serasola  
Internal Report.

# WP5

## Statistical and Robust Translation

Lluís Màrquez  
Cristina España-Bonet

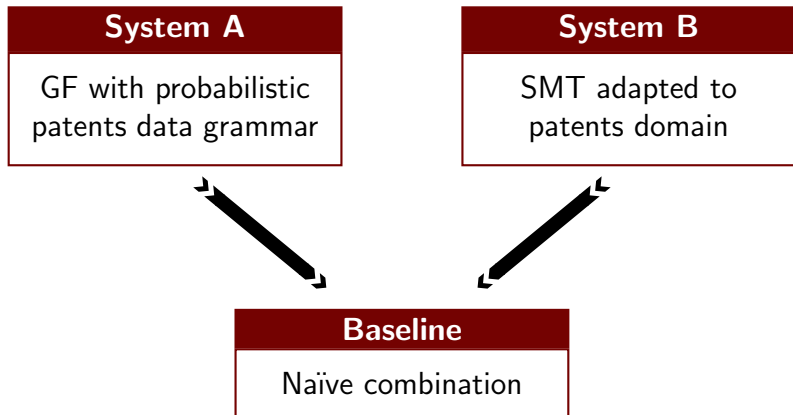
Universitat Politècnica de Catalunya, TALP Research Center

– 1st year review –

Luxembourg, March 15th, 2011

# Extra slides

*Baseline systems (Ongoing work: System B)*



# Extra slides

*SMT baseline analysis. English-German Translations, scores*

METRIC	DE2EN			EN2DE		
	Bing	Google	Domain	Bing	Google	Domain
1-WER	0.52	0.64	<b>0.72</b>	0.42	0.51	<b>0.69</b>
1-PER	0.66	0.76	<b>0.82</b>	0.56	0.64	<b>0.77</b>
1-TER	0.59	0.67	<b>0.76</b>	0.45	0.53	<b>0.71</b>
BLEU	0.43	0.58	<b>0.65</b>	0.33	0.45	<b>0.58</b>
NIST	8.25	9.67	<b>10.12</b>	6.53	8.05	<b>9.40</b>
ROUGE-W	0.40	0.48	<b>0.52</b>	0.34	0.41	<b>0.48</b>
GTM-2	0.30	0.40	<b>0.47</b>	0.25	0.32	<b>0.43</b>
METEOR-pa	0.60	0.69	<b>0.74</b>	0.36	0.45	<b>0.57</b>
ULC	0.09	0.29	<b>0.41</b>	0.03	0.19	<b>0.43</b>



Why such good scores?

---

<b>DE</b>	Verwendung nach Anspruch 23 , worin das molare Verhältnis von Arginin zu Ibuprofen 0,60 : 1 beträgt .
<b>EN</b>	The use of claim 23 , wherein the molar ratio of arginine to ibuprofen is 0.60 : 1 .

---

Why such good scores?

---

**DE**      Verwendung nach Anspruch 23 , worin das molare Verhältnis von Arginin zu Ibuprofen 0,60 : 1 beträgt .

**EN**      **The use** of claim 23 , wherein the molar ratio of arginine to ibuprofen is 0.60 : 1 .

---

**Domain**    The use of claim 23 , wherein the molar ratio of arginine to ibuprofen is 0.60 : 1 .

**Google**    The **method** of claim 23 , wherein the molar ratio of arginine to ibuprofen 0.60 : 1 **is** .

**Bing**      ~~The~~ Use of claim 23 , wherein the molar ratio of arginine to ibuprofen is 0.60 : 1 .

---

What's wrong?

---

<b>DE</b>	(±)-N-(3-Aminopropyl)-N,N-dimethyl-2,3-bis(syn-9-tetradecenyl-oxyl)-1-propanaminiumbromid
<b>EN</b>	(±)-N-(3-aminopropyl)-N,N-dimethyl-2,3-bis(syn-9-tetradecenyl-oxyl)-1-propanaminium bromide

---

What's wrong?

---

<b>DE</b>	(±)-N-(3-Aminopropyl)-N,N-dimethyl-2,3-bis(syn-9-tetradecenyloxy)-1-propanaminiumbromid
<b>EN</b>	(±)-N-(3-aminopropyl)-N,N-dimethyl-2,3-bis(syn-9-tetradecenyloxy)-1-propanaminium bromide

---

<b>Domain</b>	(±)-N-(3-Aminopropyl)-N,N-dimethyl-2,3-bis(syn-9-tetradecenyloxy)-1-propanaminiumbromid
<b>Google</b>	(±)-N-(3-aminopropyl)-N , N-dimethyl-2 , 3-bis (syn-9-tetradecenyloxy) is 1- propanaminiumbromid
<b>Bing</b>	(±)-N-(3-Aminopropyl)-N,N-dimethyl-2,3-bis(syn-9-tetradecenyloxy)-1-propanaminiumbromid

---

# Extra slides

*SMT baseline analysis. English-French Translations, scores*

METRIC	FR2EN			EN2FR		
	Bing	Google	Domain	Bing	Google	Domain
1-WER	0.54	0.66	<b>0.78</b>	0.57	0.63	<b>0.73</b>
1-PER	0.71	0.78	<b>0.86</b>	0.68	0.75	<b>0.82</b>
1-TER	0.59	0.70	<b>0.80</b>	0.60	0.66	<b>0.74</b>
BLEU	0.45	0.62	<b>0.70</b>	0.43	0.53	<b>0.62</b>
NIST	8.52	10.01	<b>10.86</b>	8.39	9.21	<b>9.96</b>
ROUGE-W	0.41	0.50	<b>0.54</b>	0.39	0.45	<b>0.49</b>
GTM-2	0.32	0.43	<b>0.53</b>	0.31	0.36	<b>0.45</b>
METEOR-pa	0.61	0.72	<b>0.77</b>	0.57	0.65	<b>0.71</b>
ULC	0.07	0.28	<b>0.44</b>	0.10	0.23	<b>0.39</b>

# Extra slides

*SMT baseline analysis. German-French Translations, scores*

METRIC	DE2FR			FR2DE		
	Bing	Google	Domain	Bing	Google	Domain
1-WER	0.42	0.52	<b>0.76</b>	0.30	0.43	<b>0.65</b>
1-PER	0.58	0.68	<b>0.77</b>	0.46	0.59	<b>0.74</b>
1-TER	0.47	0.56	<b>0.68</b>	0.32	0.46	<b>0.66</b>
BLEU	0.29	0.43	<b>0.56</b>	0.24	0.39	<b>0.53</b>
NIST	6.72	8.21	<b>9.10</b>	5.35	7.30	<b>8.88</b>
ROUGE-W	0.31	0.38	<b>0.45</b>	0.29	0.37	<b>0.44</b>
GTM-2	0.24	0.30	<b>0.41</b>	0.21	0.28	<b>0.41</b>
METEOR-pa	0.45	0.56	<b>0.64</b>	0.26	0.39	<b>0.51</b>
ULC	0.03	0.22	<b>0.41</b>	-0.03	0.19	<b>0.44</b>

### **Google**

Few OOVs but tokenization problems with compounds.

### **Bing**

Lack of specific vocabulary.

### **In-domain SMT**

Try to solve the problems of the general systems, but still:

- Improve compound detector.
- Fix structures are translated different depending on the vocabulary.

### GF System

- Composition of **parsing** and **linearisation** via an **abstract syntax** or interlingua

### Patents grammar

- **General** structure grammar
- **Compounds** grammar



**Statistical MT** can alleviate some of the **RBMT** flaws

**Rule-based MT** can alleviate some of the **SMT** flaws

### Rule-based MT can alleviate some of the **SMT** flaws

#### Missing constituents (verb)

---

<b>DE</b>	Verwendung nach Anspruch 2, wobei die Menge von Cumarin oder 7-Hydroxycumarin im Medikament 45 mg pro Medikamenten-Einheit <b>beträgt</b> .
<b>EN</b>	Use according to claim 2 wherein the amount of coumarin or 7-hydroxycoumarin in the medicament <b>is</b> 45 mg pro drug unit.

---

<b>SMT</b>	The use according to claim 2, wherein the amount of coumarine or 7-Hydroxycumarin in the medicament $\phi$ 45 mg per Medikamenten-Einheit.
------------	--

---

**Rule-based MT** can alleviate some of the **SMT** flaws

### Reordering problems (verbs & conjunctions)

---

<b>DE</b>	Verfahren nach Anspruch 20 oder 21, wobei das auf Platin basierende Analogon Cisplatin oder Carboplatin <b>ist</b> .
<b>EN</b>	The method of claim 20 or 21, wherein the platin-based analogue <b>is</b> cisplatin OR carboplatin.

---

<b>SMT</b>	A method according to claim 20 or 21, wherein the platinum based on analog cisplatin OR <b>is</b> carboplatin.
------------	--

---

### **0. Hard** integration

Force fixed GF translations within a SMT system.

### **1. Soft** integration led by **SMT**

Make available GF translations to a SMT system.

### **2. Soft** integration led by **GF**

Complement with SMT options the GF translation structure.

### **SMT leads translation, RBMT complements**

Complement the SMT translation table with RBMT options.

- **GF environment**

GF alignments for SMT, therefore **language-independent** approach.

(soon applied to WP7 languages)

### **GF alignments**

- Based on the relation between the concrete syntaxes and the abstract syntax.
- Many-to-many.
- Semantic wrt. abstract syntax.

### **SMT alignments**

- Based on corpus occurrences.
- One-to-many.

**GF** scored partial output as **new features** in SMT decoding.

$$\log P(e|f) \sim \lambda_{lm} \log P(e) + \lambda_g \log P(f|e) + \lambda_d \log P(e|f) \\ + \lambda_{di} \log P_{di}(e, f) + \lambda_w \log w(e) + \lambda_{GF} \log P_{GF}(e|f)$$

quite a challenge ||| todo un reto ||| 0.333 0.002 0.5 0.002 2.718  **$\log P_{GF}(e|f)$**

Requirements:

- GF predictions have to be probabilistic.
- Phrase pairs without prediction must be complemented.



### **GF leads translation, SMT decodes**

Complement the GF translation structure with SMT options.

- **GF**

Nowadays, there is no GF grammar for SMT corpora domains and no SMT corpora for GF grammar domains.

SMatxinT: Proof of concept.

- The GF system must parse and translate the input sentence.
- Phrases and segmentation are those given by the GF system.
- Each segment (and up) is sent to a generic SMT to provide more partial translations.
- A Moses-like decoder is fed with the resulting phrases to search for the highest scored translation.
- This statistical decoder performs no reordering and uses very simple features.

### **SMatxinT vs. MOLTO**

#### **General translator vs. in-domain translator**

- With SMatxinT, results are better for **out-of-domain** tests, where the difference between SMT and RBMT systems is less important, but systems (specially SMT) have a lower quality.
- With MOLTO, both systems will be **in-domain**, so they are expected to be high quality. Improvements here will be over already good translations.