

# Data Issues of the Multilingual Translation Matrix



Daniel Zeman

Charles University in Prague, Institute of Formal and Applied Linguistics  
Univerzita Karlova v Praze, Ústav formální a aplikované lingvistiky  
Malostranské náměstí 25, CZ-11800 Praha  
zeman@ufal.mff.cuni.cz



## PARALLEL DATA

Corpus	SentPairs	Tokens Ing1	Tokens Ing2
cs-en	782,756	17,997,673	20,964,639
de-en	2,079,049	55,143,719	57,741,141
es-en	2,123,036	61,784,972	59,217,471
fr-en	2,144,820	69,568,241	59,939,548
de-cs	652,193	17,422,620	15,383,601
es-cs	692,118	20,189,811	16,324,910
fr-cs	686,300	22,220,780	16,190,365
un.es-en	11,196,913	368,154,702	328,824,317
un.fr-en	12,886,831	449,279,647	372,627,886

**Table 1:** Number of sentence pairs and tokens for every language pair in the parallel training corpus (EuroParl + News Commentary v7; last two lines are the United Nations corpora). The xx-cs corpora were obtained as intersections of xx-en and cs-en.

## MONOLINGUAL DATA

Corpus	Segments	Tokens
newsc+euro.cs	819,434	18,491,692
newsc+euro.de	2,360,811	58,683,607
newsc+euro.en	2,430,718	65,934,441
newsc+euro.es	2,307,429	66,072,443
newsc+euro.fr	2,361,764	74,083,166
news.all.cs	14,552,899	244,728,011
news.all.de	24,446,319	462,924,303
news.all.en	42,161,804	1,039,806,242
news.all.es	8,627,438	249,022,213
news.all.fr	16,708,622	438,489,352
gigaword.en	70,592,779	2,546,581,646
gigaword.es	31,304,148	1,064,660,498
gigaword.fr	21,674,453	963,571,174

**Table 2:** Number of segments (paragraphs or sentences) and tokens for every monolingual corpus. Newsc = News Commentary; news.all = crawled news from all the years.

## BLEU

Direction	Baseline	news.all	gigaword
en-cs	0.1196	<b>0.1434</b>	
en-de	0.1426	<b>0.1629</b>	
en-es	0.2778	<b>0.3136</b>	<b>0.3136</b>
en-fr	0.2599	<b>0.2897</b>	0.2874
cs-en	0.1796	<b>0.2031</b>	0.2013
de-en	0.1877	0.2136	<b>0.2144</b>
es-en	0.2219	<b>0.2428</b>	0.2390
fr-en	0.2459	<b>0.2764</b>	0.2756
cs-de	0.1365	<b>0.1550</b>	
cs-es	0.1952	<b>0.2211</b>	0.2184
cs-fr	0.1953	<b>0.2167</b>	0.2147
de-cs	0.1212	<b>0.1400</b>	
es-cs	0.1281	<b>0.1489</b>	
fr-cs	0.1253	<b>0.1442</b>	

**Table 3:** BLEU scores. Baseline: LM on parallel corpus only. Other columns add news.all or gigaword respectively.

## UN CORPUS

Direction	Parallel	Mono	BLEU
en-es	news-euro-un	news.all	0.3194
en-es	news-euro	news.all	0.3136
en-es	un	un	0.2694
en-fr	news-euro	news.all	0.2897
en-fr	un	un	0.2541
es-en	un	un	<b>0.2688</b>
es-en	news-euro	news.all	0.2428
fr-en	news-euro	news.all	0.2764
fr-en	un	un	0.2392

**Table 4:** BLEU scores with and without the UN corpus.

## HUMAN

Direction	BLEU	Hwin
en-fr	0.2897	0.549350
en-de	0.1629	0.545068
fr-en	0.2764	0.534239
en-cs	0.1434	0.533774
en-es	0.3136	0.525881
de-en	0.2144	0.499220
cs-en	0.2031	0.463432
es-en	0.2428	0.259303

**Table 5:** Human comparison with other systems (percentage of times when this system produced better translation than another system, ignoring ties).

## THE SYSTEM

- Tokenization, in-house quotation marks normalization, supervised truecasing (based on lemmas ← TreeTagger or Morče)
- Word alignment: Giza++ operating on lemmas
- Hexagram language model
- Moses decoder, no lexical reordering model, no factored model
- All BLEU scores computed by the system on tokenized test data, truecased (except for the sentence-initial letter)

### Supervised truecasing

Training data with lemma factor:  
A token is lowercased unless its lemma is uppercased.

Changes sentence-initial tokens, words in English headings, highlighting uppercase etc. *Errors on common-proper ambiguities?*

mušaraf	+ův	po+ sled +ní	výstup	?
musharraf	s	last	act	?

udusit	ne+	závislé	soud + nictví	a	svobod +ná	média	.
to stifle	in+	dependent	judiciary	and	free	media	.

nic ne+	může	být	dále	od	pravd +y	.
nothing	could	be	further	from	the truth	.

mnoho	pákistán	+ců	rozčar +ova	+ných	pákistán	+skou	po+ liticko +u	garnit +u + rou
many	pakistan	+is	dis+ illusion	+ed	with pakistan	s	political	class

## Unsupervised morphemic segmentation

en-cs: Swapped words

en-cs: Wrong case

cs-es: Bad word order

Direction	Translation
en	The story of Libya's liberation, or rebellion, already has its defeated.
en-cs	Příběh Libye osvobození, nebo odboj, už svou porážku.
en-de	Die Geschichte von Libyens Befreiung oder Rebellion, hat bereits seine Niederlage.
en-es	La historia de la liberación de Libia, o la rebelión, ya tiene su derrotado.
en-fr	L'histoire de la libération de la Libye, ou la rébellion, a déjà son vaincu.
cs	Příběh libyjského osvobození, anebo libyjské rebélie, již má své poražené.
cs-de	Die Geschichte des libyschen Befreiung, oder der libyschen rebélie, hat bereits seine Verlierer.
cs-en	The story of the libyan liberation or libyan rebélie, already has its losers.
cs-es	La historia de la liberación de Libia, Libia o rebélie, ya tiene perdedores.
cs-fr	L'histoire de la libération de la Libye, ou rebélie, possède déjà des perdants.
de	Die Geschichte des libyschen Befreiungskampfes oder libyschen Rebellion kennt schon ihre Verlierer.
de-cs	Historie libyjského Befreiungskampfes nebo libyjského povstání už zná své poražené.
de-en	The history of the libyan liberation or libyan rebellion knows about their losers.
es	La historia de la liberación libia, o de la rebelión libia, ya tiene sus perdedores.
es-cs	Historie osvobození Libye, nebo libyjské povstání, protože má své poražené.
es-en	The history of the liberation of Libya, or of the rebellion in Libya, already has sus perdedores.
es-en.un	The history of liberation libyan, or of the rebellion libyan, already has its losers.
fr	La libération ou rebellion libyenne a déjà ses vaincus.
fr-cs	Propuštění nebo rebellion libyjské již své poražené.
fr-en	The release or libyan rebellion already has its losers.

en-cs & en-de: "defeated" mistranslated as "defeat"

cs: "rebélie" is very rare, synonyms exist

en-cs & fr-cs: missing verb "has" (often dropped when translating perfect tense)

fr-en & fr-cs: "libération" wrongly translated as "release" instead of "liberation"

es-en: The UN corpus has larger vocabulary ("perdedores") but it also prefers strangely ordered N-Adj phrases.

**Example 1:** Reference translations and system hypotheses for one sentence. Original language is Czech.