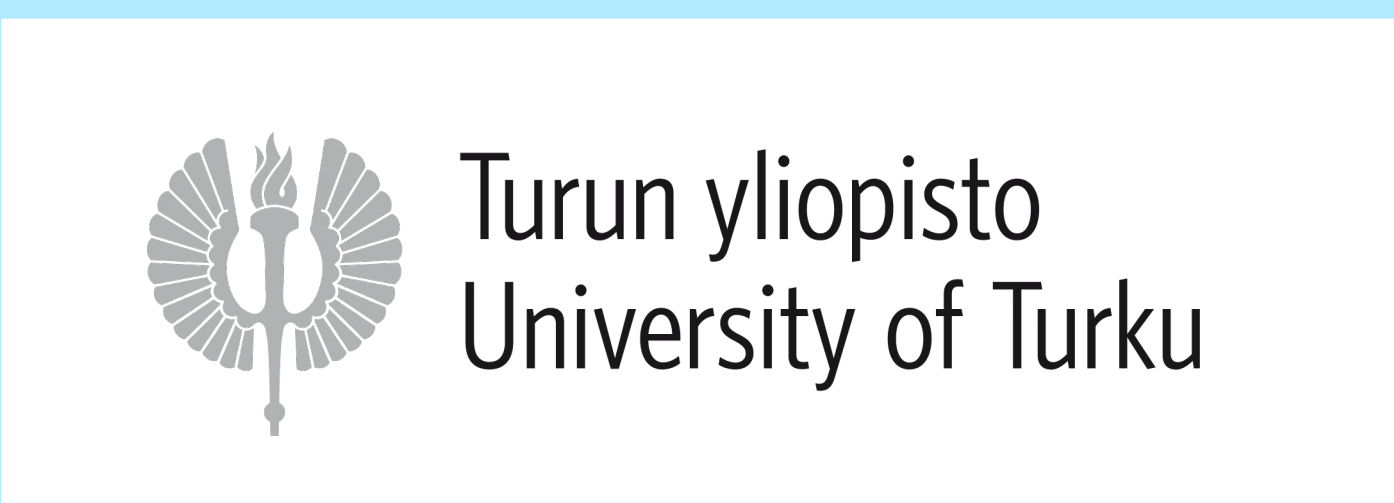


# Quality Assessment of the Reuters Vol 2 Multilingual Corpus

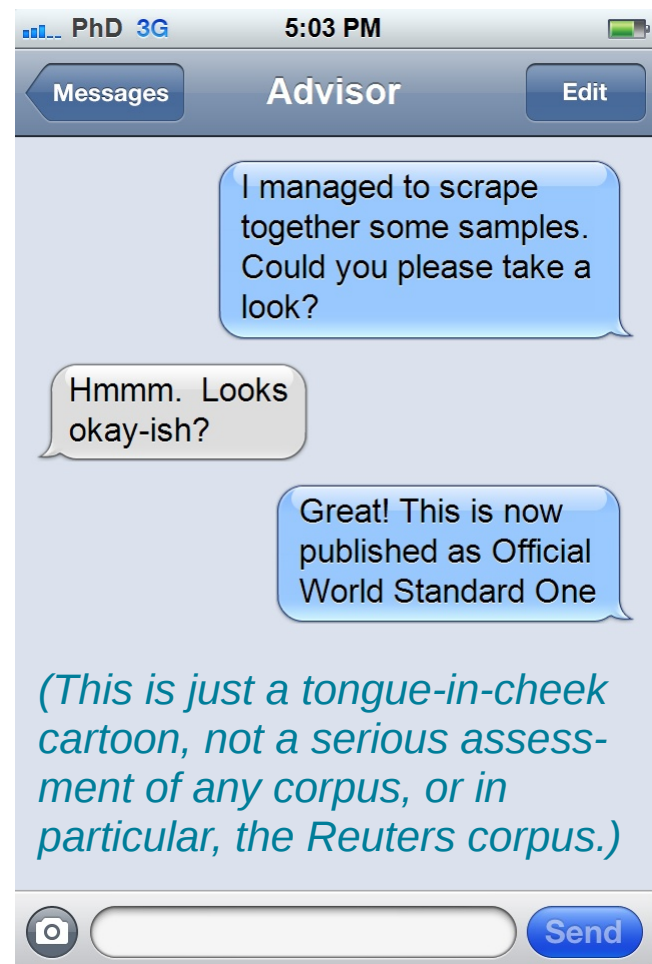
Robin Eriksson • University of Turku



# The Problem

Like every other deliverable, corpora are produced under limited timetable and budget.

Frequently, quality requirements are implicit, informal, and / or nonexistent.



We need better metrics, and tools for assessing and – ideally – improving corpora.

# The Corpus

The Reuters Vol. 2 Multilingual Corpus contains newsprint from 1996-1997 from 13 local Reuters offices.

<i>Location</i>	<i>Articles</i>
Denmark	11,185
France	85,393
(West) Germany	116,212
Italy	28,406
Japan	65,499
Latin America *	79,818
The Netherlands	17,940
Norway	9,409
Portugal	8,841
Russia	17,487
Spain *	18,655
Sweden	15,732
Taiwan	28,964

With two of the offices (\*) publishing in Spanish, but with the (unadvertised) inclusion of English, the corpus contains material in 13 languages.

This multilingual corpus contains material from the same timeframe as the monolingual English Vol 1 corpus.

Each newswire article is exported as a separate XML file with metadata about publication date, categories, etc.

# The Methodology

We specify a set of constraints which the samples in the corpus must obey. Violations are excluded or at least investigated.

- *Hard* constraints must never be violated.
- *Soft* constraints may be within normal variation.

By iteratively applying constraints,  
we shave off undesired artefacts  
from the clean core of the corpus.

Ultimately, we hope to have an ideal catalog of constraints which excludes all invalid samples.



# The Constraints

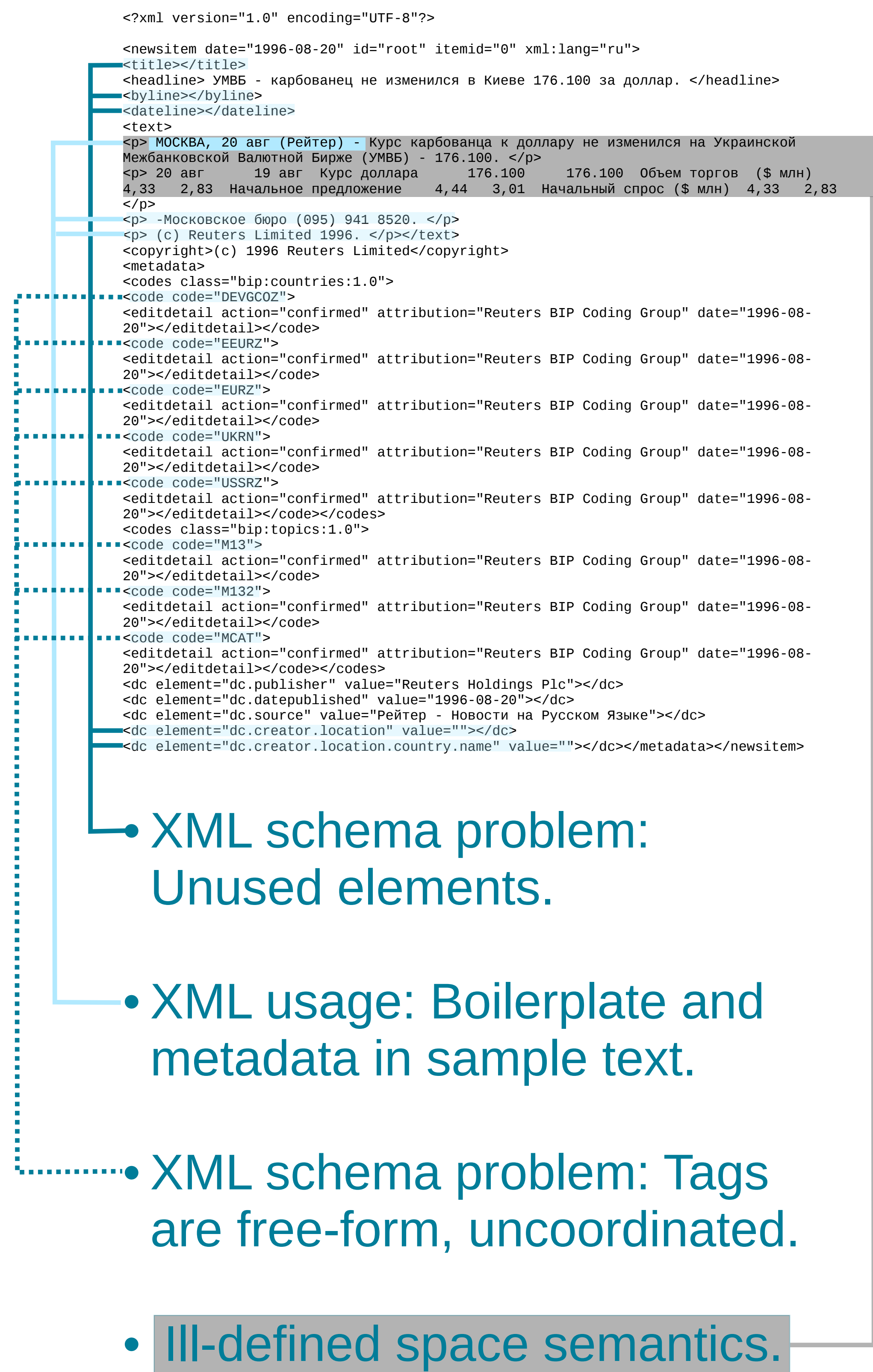
For this initial case study, a fairly simple and straightforward set of demo constraints was investigated.

- ***No duplicates.***  
In fact, lots of identical samples.
- ***Correct and consistent markup.***  
Are samples represented so that we can understand and manipulate them properly?
- ***Within entropy range.***  
Text in a particular language fits the expected byte entropy distribution. Aberrations may be corrupted or in the wrong language.

- **Consistent category tagging.** Investigation of duplicates and near-duplicates suggests 1% error rate. This is probably too low.

Also, language only bulk tagged.

# The Results



## Entropy investigations unearthed a number of corruption problems.

[illegible]

```

xml/339944.xml.<p> Weiter im Aufwind befand sich das britische Pfund, das auf 2,4764/4796
(2/4764/4791) $f kletterte. Spekulationen auf einen
Zinssatzanstieg von90x06DE07FK48051326B8Z4P. 27830607USCHLAND/BUNDESWEHREINBUNDESWEHR-Soldaten
steht Hiit. Die Bundeswehr hat am 27. April 2011 eine
einer Ausbildung für den Bosnien-Einsatz BIR24P. 27830607USCHLAND/BUNDESWEHREINBUNDESWEHR-Soldaten
xml/199997.xml.<p> Les députés UDF ont également décidé de constituer un groupe de travail
qui présentera des propositions en matière de COR.
27830607USCHLAND/BUNDESWEHREINBUNDESWEHR-Soldaten
Net cash flow 183.93 vs 189.09 Earnings per share (pesetas) 289.92 vs 307.16 Cash flow per
share (pesetas) 61.89 vs 62.99 Shareholders&sapos; Shares 717.82 vs 658.04 Provisions
27830607USCHLAND/BUNDESWEHREINBUNDESWEHR-Soldaten
53.95 95.82 Debt (medium+long-term) 280.70 vs 221.88 Debt (short-term) 163.67 vs
92.73.
xml/371118.xml.<p> 5F8F9P809081108P06E4ACD07A14D10800E2400000000 </p>
xml/35447.xml.<p> Получит комментарий непосредственно в PAO Газпром пока не
27830607USCHLAND/BUNDESWEHREINBUNDESWEHR-Soldaten
xml/336399.xml.<p> LUOWIGSBURG - Wustenrot Holding GmbH, PK mit erstem Übernahm
DE C19X5064875H1E0806424152783730BUuInflationInflation&sapos;bank sees inflation staying below
27830607USCHLAND/BUNDESWEHREINBUNDESWEHR-Soldaten
expected German inflation to remain below two percent for the next two years, with high

```

Near-duplicate investigation showed large number of similar samples, but without context, we cannot tell whether they are erroneous or just regular newsroom republication. More investigations needed!



Poster presented at  
the LREC 2016 in  
Portorož, Slovenia  
May 23-28, 2016.

On-line repository:  
<http://github.com/rcv2/>