# Enhancing Interoperability for Under-Resourced Languages
## A Case Study on Linking Lithuanian-English Data in the Cybersecurity Domain

Christian Chiarcos[1]   Maxim Ionov[2]   Andrius Utka[3]   Sigita Rackevičienė[4]

[1]University of Augsburg, Germany   [2]University of Cologne, Germany
[3]Vytautas Magnus University, Lithuania  [4]Mykolas Romeris University, Lithuania

---

**Languages**: Lithuanian - English

**Data**: cybersecurity domain
1) lexicon (terminology: TBX)
2) corpora (parallel: TMX / annotated: CoNLL)

**Challenge**: Publish that such that
3) anyone can easily re-use our data
4) we can integrate lexical data, linguistic annotations and parallel corpus
5) we access / query / interlink / process all data with off-the-shelf technology

**UniDive**: WG2 (mostly)
**Status:** On-going

### We have a solution that works nicely :)
Yet, aspects of the data modelling are still experimental and need to be supported by additional use cases

---

## TBX: Lithuanian-English Cybersecurity Termbase

- TermBase eXchange Format (TBX v.2)
- 233 cybersecurity concepts with Lithuanian and English designations, definitions and context examples:

  - Concept id
  - Subject field (subdomain)
  - Term group of synonymous terms
  - Description group with definitions
  - Description group with examples

Data silo, using a domain-specific, proprietary XML standard



---

## TMX: Parallel Corpus

- Translation Memory eXchange format

TMX 1.4b
open (but dated) standard in the translation industry



## Vert: Monolingual Annotations

- CoNLL-style format plus SGML markup
- morphosyntax for both languages



---

## TMX -> Web Annotation

We are not aware of any precedent for encoding parallel corpora in RDF
**Web Annotation** is a W3C standard for annotating textual, multimodal and other content on the web with labels or structured information
- We model the translation relation as an annotation that points to two or more "targets", identified by **XPath Selectors** pointing to the original file
- Additionally, we use the NIF vocabulary to include the **full text** as nif:Context

## Vert ->-CoNLL-RDF / NIF

- **CoNLL-RDF** is an established vocabulary and a library for encoding/wrapping TSV data in RDF
- Based on NIF, so, it can be directly **linked to TMX** contexts

---

## TBX -> OntoLex

- **OntoLex**: widely for lexical information in the web of data
  **=>** encode or wrap lexical information in RDF
  **=>** can be flexibly merged with other RDF/RDF-wrapped data, e.g., using the query language SPARQL
- We provide a custom TBX converter converter inspired by Cimiano et al. (2015)

Except for solutions currently pursued in the context of the Linguistic Linked Open Data (LLOD) community, we are not aware of any established technology that allows to seamlessly integrate parallel corpus data, lexical data and morphosyntactic annotations

Conceptually, this requires
- directed graphs => here: **RDF data model**
- community standards for how to encode linguistic information in graphs => here: **LLOD community standards**
- a standard format, designated databases and protocols => here: **RDF/SPARQL** ecosystem

### See handout for how to perform cross-lingual queries across corpus and terminology ;)

---

christian.chiarcos@uni-a.de
mionov@uni-koeln.de
andrius.utka@vdu.lt
sigita.rackeviciene@mruni.eu