



Extracting Terminological Concept Systems from Natural Language Text (Text2TCS)



Text2TCS

Terminological Concept Systems (TCS)

- Concept: abstractions of a set of physical or abstract entities
- Term: their designation by linguistic means
- Relation: hierarchical and semantic relations between concepts that structure and organize the TCS

TCS From Text

Text Coronavirus disease 2019 (COVID-19) is a disease that has been spreading through Europe.



Terms <term>*Coronavirus disease 2019*</term>
<term>*COVID-19*</term> <term>*disease*</term>
<term>*spreading*</term> <term>*Europe*</term>

TCS From Text

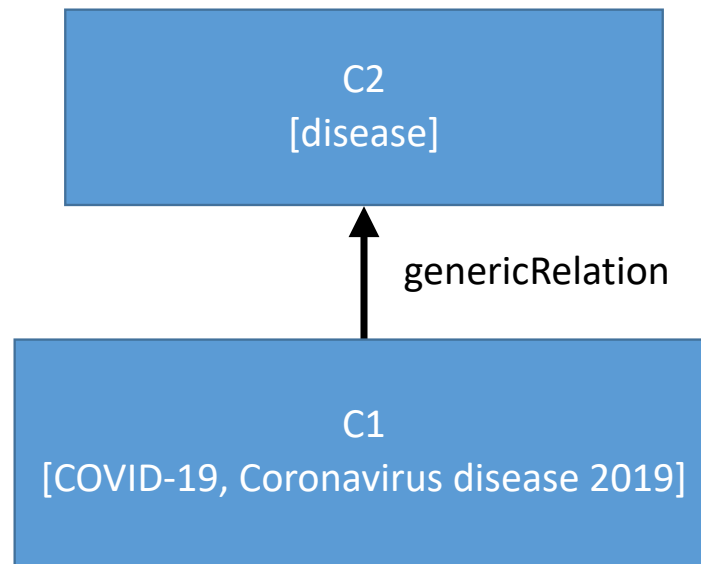
C1
[COVID-19]

TCS From Text

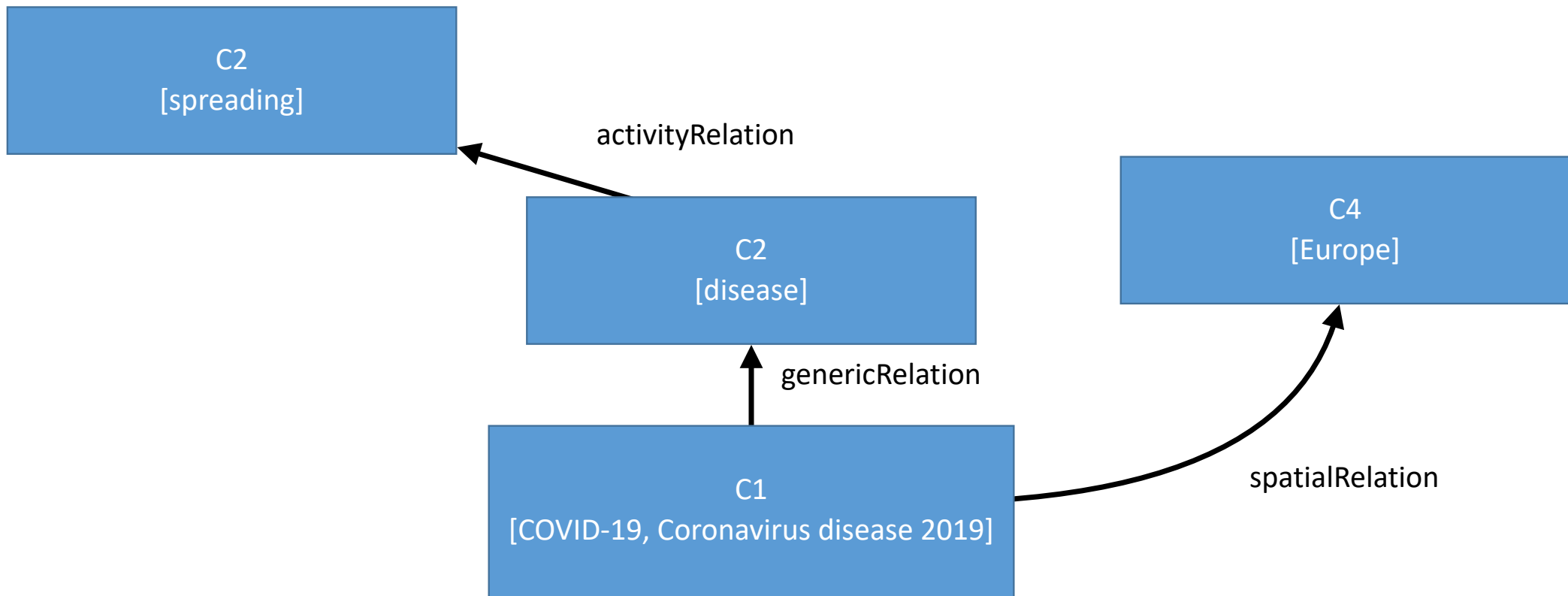
C1

[COVID-19, Coronavirus disease 2019]

TCS From Text



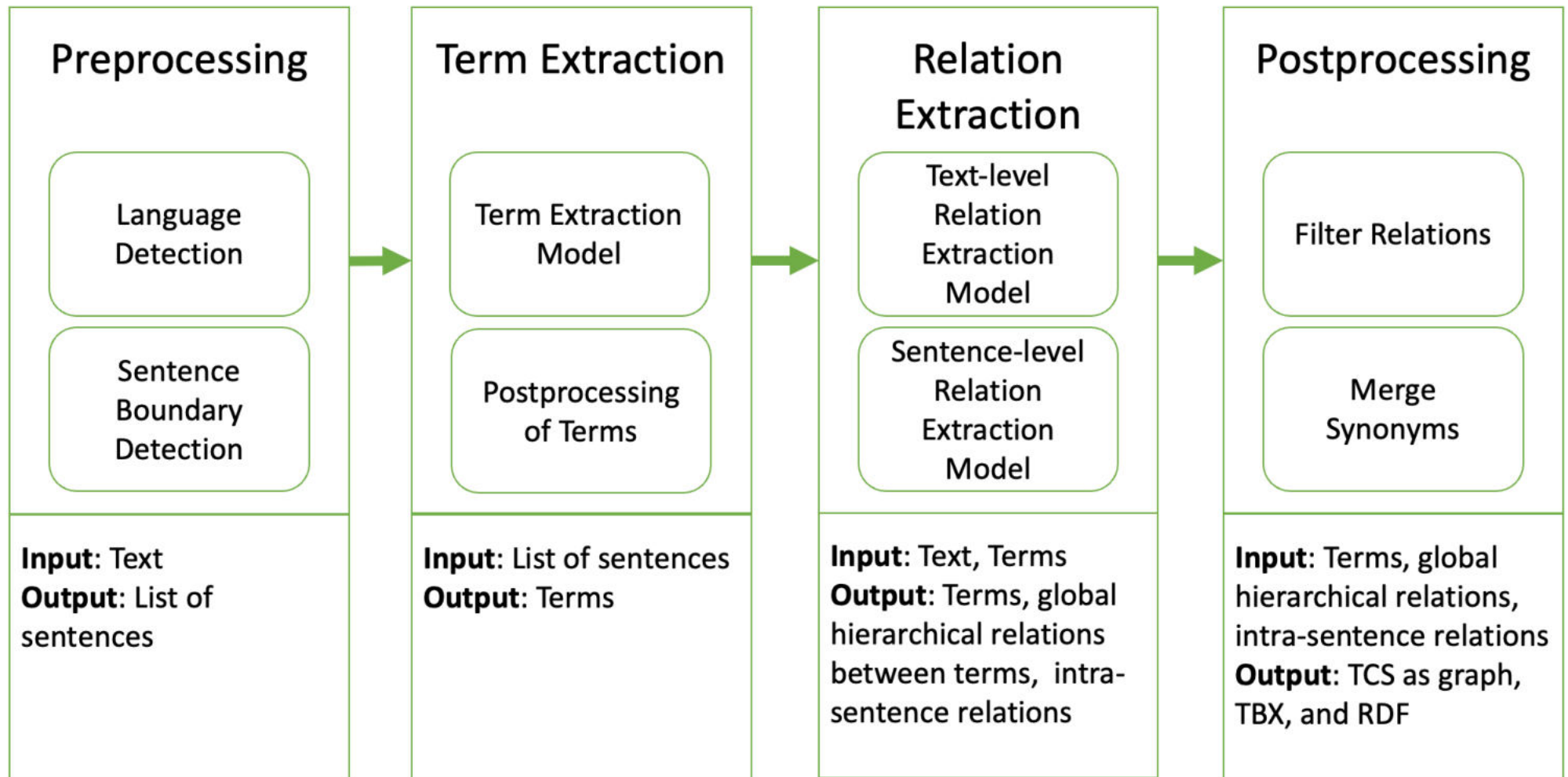
TCS From Text



Motivation

- Specialized communication needs to be able to rely on a consistent usage of terminology by different parties
 - A TCS alleviates that problem, but manually curating a TCS is cumbersome
- **Automated approaches are needed**

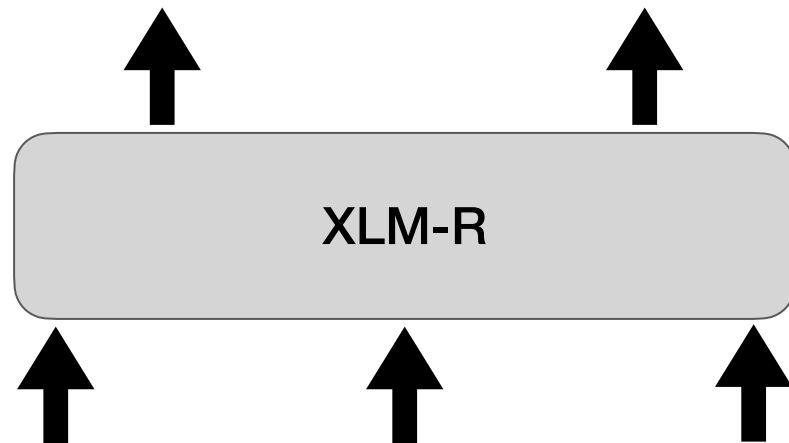
Text2TCS — the Pipeline



Method - Multilingual BERT models (XLM-R)

Term Extraction

['n', 'B-T', 'n', 'n', 'n', 'B-T', 'T']

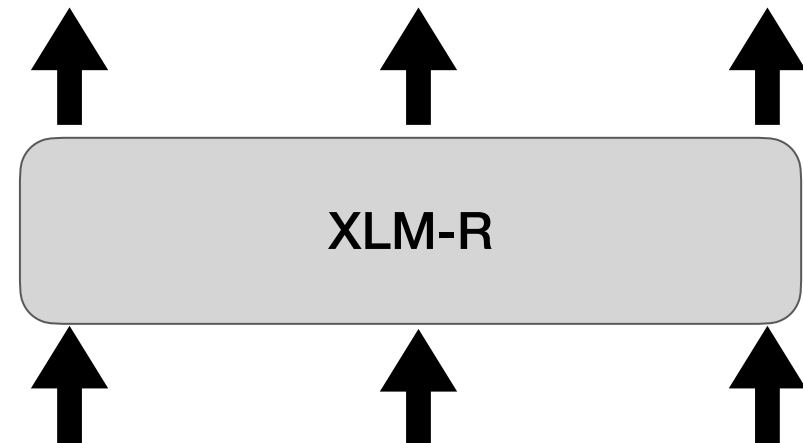


The cough was caused by COVID-19 disease.

Lang C. et al (2021) Transforming Term Extraction: Transformer-Based Approaches to Multilingual Term Extraction Across Domains. Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021.

Relation Extraction

causalRelation(COVID-19, cough)



cough. COVID-19. The cough was caused by COVID-19.

Wachowiak L. et al. (2021) Towards Learning Terminological Concept Systems from Multilingual Natural Language Text. in LDK 2021

Text2TCS on ELG

<https://live.european-language-grid.eu/catalogue/tool-service/8122>



RELEASE 2

My grid 

Dagmar Gromann 

[Technologies](#) [Resources](#) [Community](#) [Events](#) [Documentation](#) [About ELG](#)

[Go to catalogue](#)

COVID-19 is an infectious disease caused by the SARS-CoV-2 virus.

Features

Name	Value
------	-------

Graph Link	https://live.european-language-grid.eu/temp-storage/retrieve/01grt8z9-27759128ugv8d17c8w3anci6m2ait 
------------	---

TBX Link	https://live.european-language-grid.eu/temp-storage/retrieve/01grt8z9-jfkfp6bf2d04qvtw219iwhq4irtxr 
----------	---

BACK

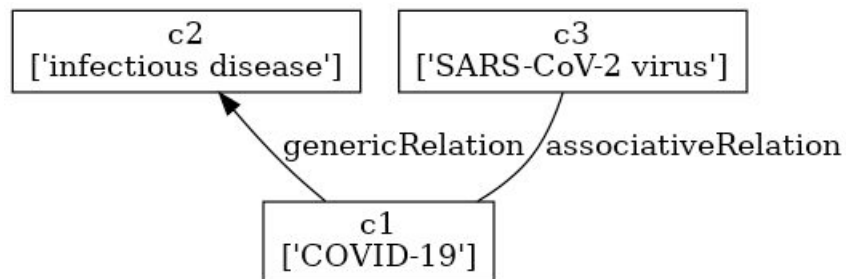
Annotations

- c1: COVID-19
- c2: infectious disease
- c3: SARS-CoV-2 virus

Text2TCS on ELG

<https://live.european-language-grid.eu/catalogue/tool-service/8122>

Concept Map



TBX/XML

```
<?xml version='1.0' encoding='utf8'?>
<?xml-model href="https://raw.githubusercontent.com/LTAC-Global/TBX-Core_dialect/master/Schemas/TBX
schematypens="http://relaxng.org/ns/structure/1.0"?>
<?xml-model href="https://raw.githubusercontent.com/LTAC-Global/TBX-Core_dialect/master/Schemas/TBX
/dsd/1/schematron"?>
<tbx type="TBX-Core" style="dca" xml:lang="en" xmlns="urn:iso:std:iso:30042:ed-2">
  <tbxHeader>
    <fileDesc>
      <sourceDesc>
        <p>TBX file automatically generated by Text2TCS (https://text2tcs.univie.ac.at/)</p>
      </sourceDesc>
    </fileDesc>
    <encodingDesc>
      <p type="XCSURI">TBXXCSV02.xcs</p>
    </encodingDesc>
  </tbxHeader>
  <text>
    <body>
      <conceptEntry id="c1">
        <transacGrp>
          <transac type="transactionType">origination</transac>
          <transacNote type="responsibility">Text2TCS</transacNote>
          <date>21-11-11_17h-34m</date>
        </transacGrp>
        <langSec xml:lang="en">
          <termSec id="c1-en-t0">
            <term>COVID-19</term>
          </termSec>
        </langSec>
        <descripGrp>
          <descrip type="genericRelation">c2</descrip>
        </descripGrp>
        <descripGrp>
          <descrip type="associativeRelation">c3</descrip>
        </descripGrp>
      </conceptEntry>
      <conceptEntry id="c2">
        <transacGrp>
          <transac type="transactionType">origination</transac>
          <transacNote type="responsibility">Text2TCS</transacNote>
          <date>21-11-11_17h-34m</date>
        </transacGrp>
        <langSec xml:lang="en">
          <termSec id="c2-en-t0">
            <term>infectious disease</term>
          </termSec>
        </langSec>
      </conceptEntry>
      <conceptEntry id="c3">
        <transacGrp>
```

Major findings

- Fine-tuning pre-trained neural language models works well on this TCS extraction task
- Works also well with Neural Machine Translation models - even for joint term and relation extraction (but too slow and big for ELG)
- Coordination between sentence- and text-level output challenging - which model to trust more?
- Pre-annotated data still needed
- All depends on good team work...

TEAM



DAGMAR GROMANN
Project
leader



LENNART WACHOWIAK
Machine
learning and IT



Text2TCS



CHRISTIAN LANG
Translation
and IT

BARBARA HEINISCH
Terminology
and usability





universität
wien



EUROPEAN
LANGUAGE
GRID



Text2TCS

Thank you for listening!

More information on:

<https://text2tcs.univie.ac.at/>

Text2TCS on ELG:

<https://live.european-language-grid.eu/catalogue/tool-service/8122>

Any feedback very welcome:

dagmar.gromann@univie.ac.at