

ELRC Action: Covering Confidentiality, Correctness and Cross-linguality

Tom Vanallemeersch[†], Arne Defauw[†], Sara Szoc[†], Alina Kramchaninova[†], Joachim Van den Bogaert[†], Andrea Lösch[‡]



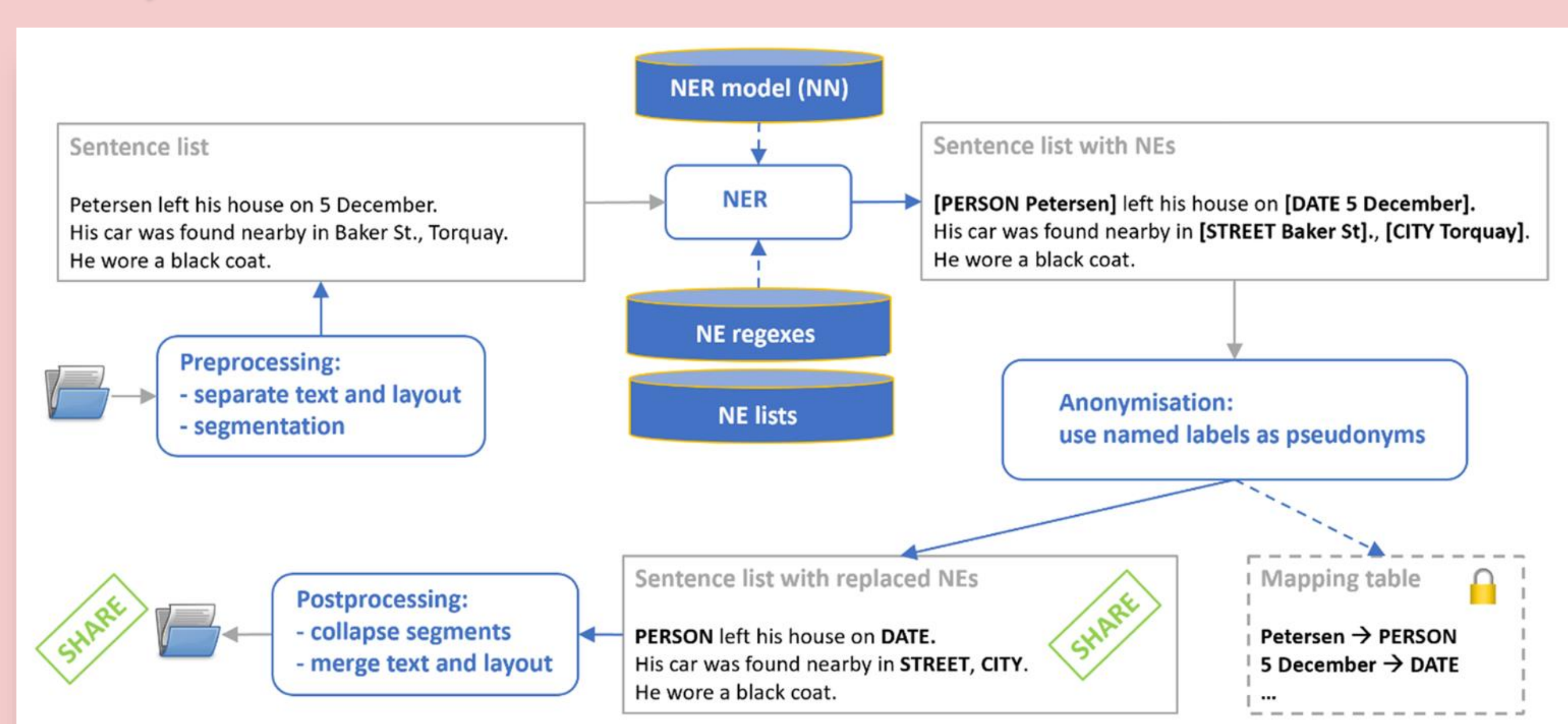
ELRC action:

- Aim: minimising language barriers across European Union
- Two language technology assessments involving consultation round:
 - Automated anonymisation
 - Multilingual Fake News Processing

Automated Anonymisation

- Anonymisation consists of **removing personal identifiable information (PII)**.
- It is important when sharing language data without violating the **General Data Protection Regulation (GDPR)**.
- It involves two steps:
 - Detecting **what** should be anonymised via a **named-entity recognition (NER)** system.
 - Determining **how** it should be anonymised.

- ⇒ Consultation round with stakeholders
- ⇒ Identification of scenario for workflow using NE labels:



- ⇒ Development of proof-of-concept software:

Anonymisation specification: tools illustrating scenarios

TM-Anonymizer monolingual TM-Anonymizer bilingual COMPRISE Text Transformer Biroamer

Language: en

File to anonymise (docx): No file selected.

User-defined NEs: No file selected.

Available tools:

- TM-Anonymizer¹
- COMPRISE Text Transformer²
- Biroamer³

Multilingual Fake News Processing

US BACON Reserves Hit 50 year low



Coronavirus Bioweapon - How China Stole Coronavirus from Canada and Weaponized it



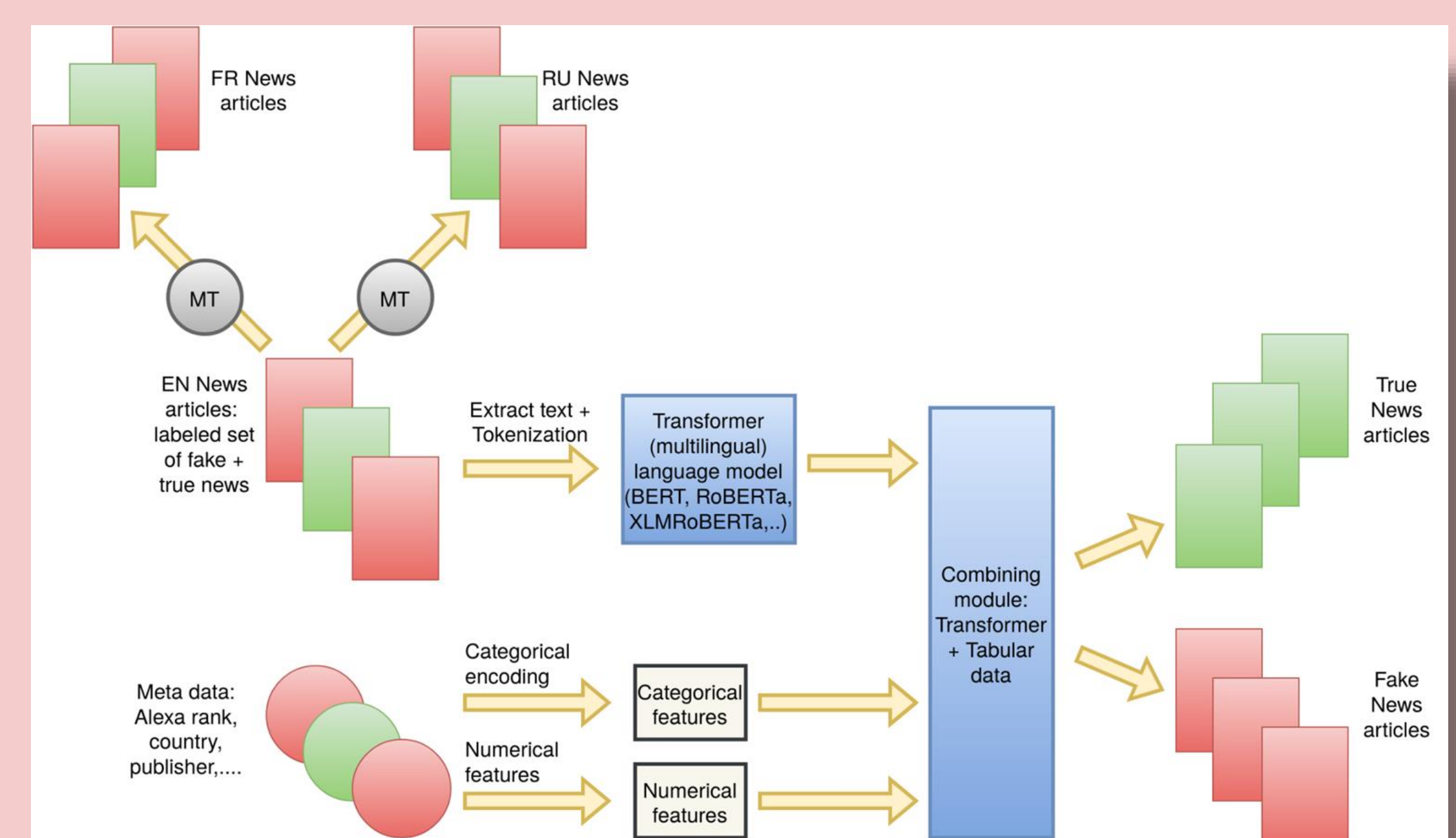
- Disinformation (fake news) on various topics is spreading quickly.
- Need for tools for automatic detection of fake news is gaining urgency.
- Disinformation is a global phenomenon: it is important to explore multilingual techniques for fake news detection.

- ⇒ Consultation round with stakeholders

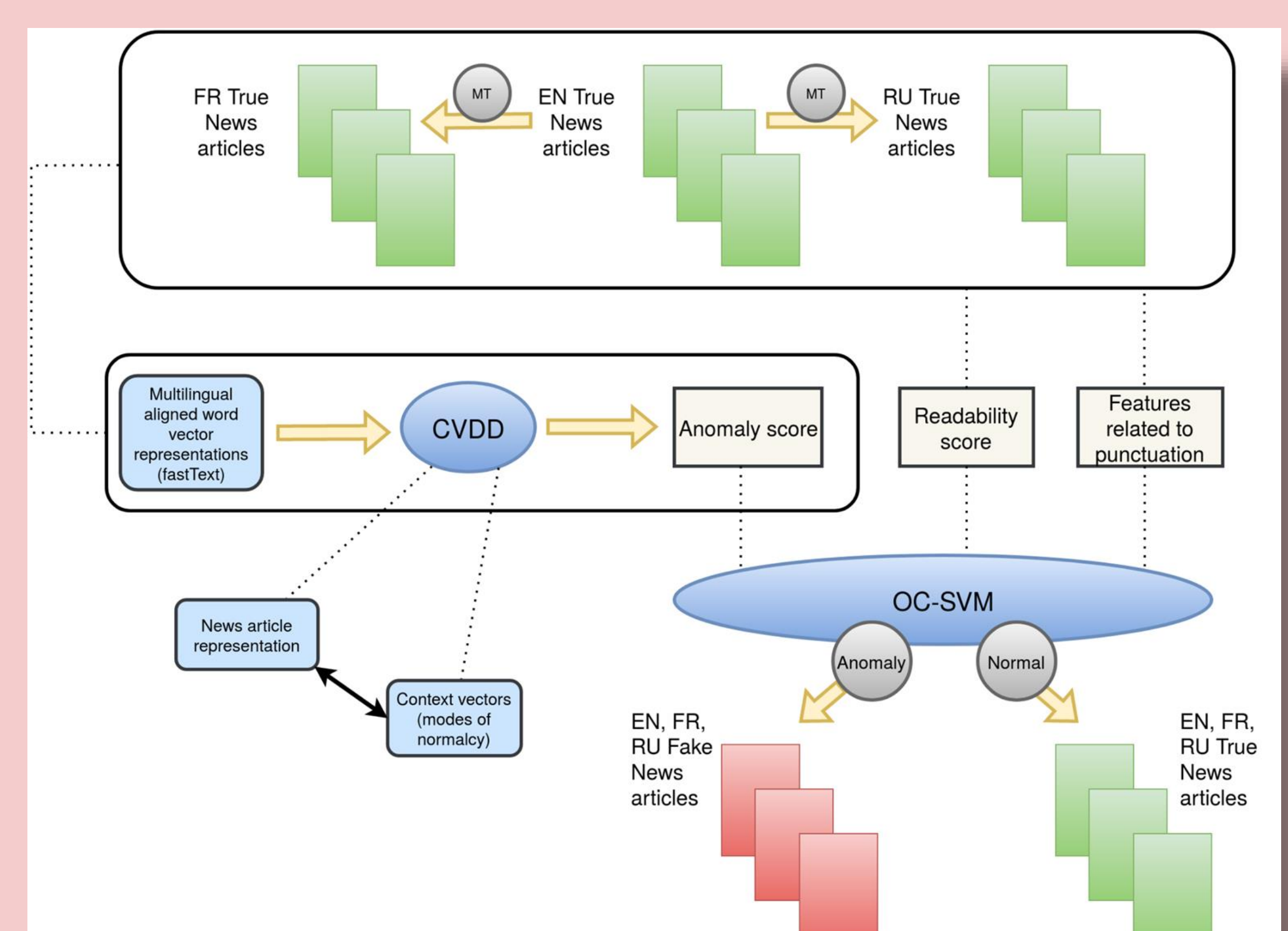
- ⇒ Main findings:

- Lack of multilingual resources for detection of fake news
- Quick evolution of topics

Supervised approach for multilingual fake news detection in Russian (RU) and French (FR) via use of machine translation (MT) and pretrained multilingual language models:⁴



Unsupervised approach by merely training on true news articles using Context Vector Data Description (CVDD) algorithm⁵ + aligned vector representations.⁶



References:

1. Kamran, A., et al. (2020). CEF Data Marketplace: Powering a long-term supply of language data. EAMT2020.
2. Adelani, D.I., et al. (2020). Privacy guarantees for de-identifying text transformations. Interspeech 2020.
3. Devlin J., et al. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. NAACL2019.
4. Bañón, M., et al. (2020). ParaCrawl: Web-scale acquisition of parallel corpora. ACL2020.
5. Ruff, L. et al. (2019). Self-attentive, multi-context, one-class classification for unsupervised anomaly detection on text. ACL2019.
6. Joulin, A. et al. (2018). Loss in translation: learning bilingual word mapping with a retrieval criterion. EMNLP2018.

[†]CrossLang
Kerkstraat 106, 9050 Ghent, Belgium, {firstname.lastname}@crosslang.com

[‡]DFKI
Stuhlsatzenhausweg 3, D-66123 Saarbrücken, Germany, andrea.loesch@dfki.de