



# Ontology-based translation memory maintenance

---

Andraž Repar<sup>1,2</sup>, Senja Pollak<sup>3</sup>

<sup>1</sup> International Postgraduate School, Ljubljana, Slovenia

<sup>2</sup> Iolar, Ljubljana, Slovenia

<sup>3</sup> Jozef Stefan Institute, Ljubljana, Slovenia

# Translation industry

According to SDL, *'Inconsistencies in the use of terminology' is the top cause of translation rework*

Iolar started developing its own terminology management system in cooperation with the JSI Institute:

- Terminology management
- (Bilingual) terminology extraction
- Definition extraction
- **Domain extraction**
- Good example extraction (Repar, Pollak, 2017)

# Outline of the talk

- Introduction - What is a translation memory
- Text mining approach
  - Dataset
  - Experiments
- Evaluation and results
- Conclusion and future work

# What is a translation memory

- Database of translations
- Standard in the translation industry
- Workflow: segment-translate-save

Why?

- Leveraging translation memories
- Reuse for e.g. machine translation

Type	Segments	Words	Characters	Percent	Placeables	Tags
<b>PerfectMatch</b>	0	0	0	0.00%	0	0
<b>Context Match</b>	0	0	0	0.00%	0	0
<b>Repetitions</b>	1	6	19	33.33%	0	0
<b>Cross-file Repetitions</b>	0	0	0	0.00%	0	0
<b>100%</b>	0	0	0	0.00%	0	0
<b>95% - 99%</b>	0	0	0	0.00%	0	0
<b>85% - 94%</b>	0	0	0	0.00%	0	0
<b>75% - 84%</b>	0	0	0	0.00%	0	0
<b>50% - 74%</b>	0	0	0	0.00%	0	0
<b>Internal:</b>						
<b>95% - 99%</b>	0	0	0	0.00%	0	0
<b>85% - 94%</b>	1	6	22	33.33%	0	0
<b>75% - 84%</b>	0	0	0	0.00%	0	0
<b>50% - 74%</b>	0	0	0	0.00%	0	0
<b>New</b>	1	6	19	33.33%	0	0
<b>Total</b>	<b>3</b>	<b>18</b>	<b>60</b>	<b>100%</b>	<b>0</b>	<b>0</b>

# Translation memory - an example

## Translation Memory eXchange (TMX) file format

```
<tu creationdate="20140527T094016Z" creationid="AR_00" changedate="20140527T094016Z" changeid="AR_00" lastusedate="20140527T094016Z">
  <prop type="x-Origin">TM</prop>
  <prop type="x-OriginalFormat">TradosTranslatorsWorkbench</prop>
  <tuv xml:lang="en-GB">
    <seg>The Subfund remained invested almost completely in USD-denominated corporate bonds with a short to medium residual term.</seg>
  </tuv>
  <tuv xml:lang="sl-SI">
    <seg>Podsklad je skoraj v celoti ostal vložen v dolarske podjetniške obveznice s kratko do srednje dolgo zapadlostjo.</seg>
  </tuv>
</tu>
<tu creationdate="20140527T100115Z" creationid="AR_00" changedate="20140527T100115Z" changeid="AR_00" lastusedate="20140527T100115Z">
  <prop type="x-Origin">TM</prop>
  <prop type="x-OriginalFormat">TradosTranslatorsWorkbench</prop>
  <tuv xml:lang="en-GB">
    <seg>The focus was on issuers from the IT and consumer goods sectors.</seg>
  </tuv>
  <tuv xml:lang="sl-SI">
    <seg>Usmeril se je na izdajatelje iz sektorjev informacijskih tehnologij in potrošniških dobrin.</seg>
  </tuv>
</tu>
```

*METADATA*

*SOURCE SEGMENT*

*TARGET SEGMENT*

# Why the need for maintenance

- Old technology - since the beginning of 1990s
- (Normally) one translation memory per project/domain
- Lots of users over a long period of time - resulting in:
  - Wrong languages in the translation memory
  - Out-of-domain content in the translation memory
- Initial domain too broad for modern applications (e.g. machine translation, terminology extraction)

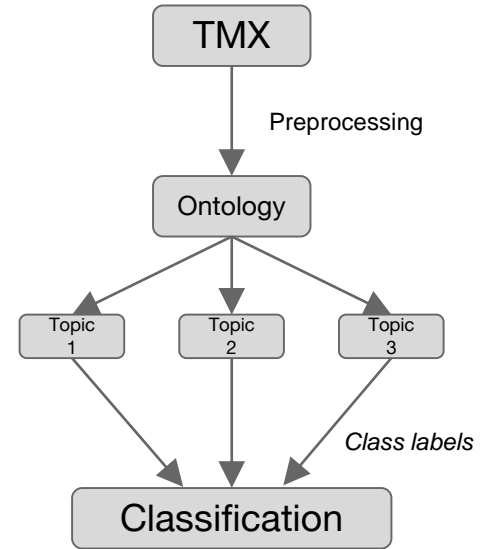
# Dataset

## MarLegFin translation memory

- Created sometime after 2000
- Almost 250,000 EN-SL segments
- Marketing, Legal, Finance
- However, other domains can also be found
- No way of manually extracting individual domains

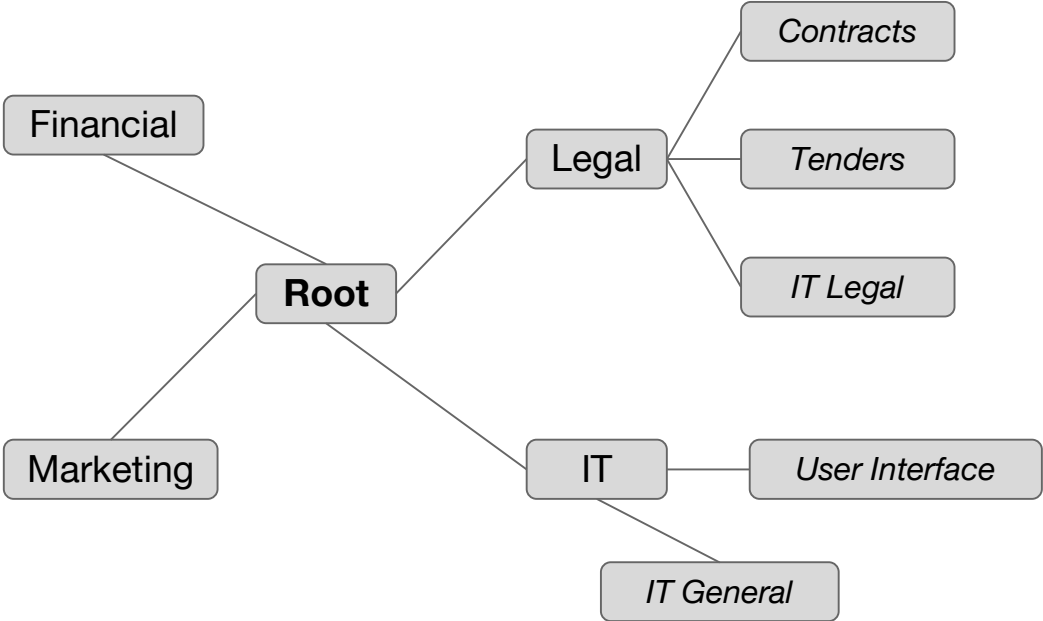
# Experiments

- Import data into OntoGen (Fortuna et. al., 2007)
  - One segment is one document ---- difficult to classify!
  - Removed all segments with less than 8 words
  - K-means clustering + manual grouping into topics
  - Manual evaluation of extracted topics
- Use the extracted topics as a shortcut for classification of new content
  - Topics as class labels, StringToWordVector
  - Test performance of Naive Bayes Multinomial, SVM and J48 in Weka (Hall et.





# Results





# Manual evaluation of OntoGen results

50 segments per topic

<i>Topic</i>	<i>Precision</i>
Financial	0.76
IT General	0.80
IT User Interface	0.86
Legal Contracts	0.80
Legal IT Legal	0.86
Legal Tenders	0.78
Marketing	0.80
<b>Average</b>	<b>0.81</b>

# Classifier performance - Main topics

Majority class accuracy = 0.406

	<b>J48</b>	<b>SMO</b>	<b>NB Multinomial</b>
<i>Accuracy</i>	0.597	0.619	0.671
<i>Precision</i>	0.615	0.608	0.678
<i>Recall</i>	0.597	0.619	0.671
<i>F-measure</i>	0.576	0.610	0.673

## Conclusion and future work

- OntoGen can be used for translation memory maintenance
- Classification performance was promising but not yet good enough

### Future work

- Increase frequency of domain terminology
- Segment grouping based on translation memory metadata