

# Looking for the Needle in a Haystack: Semi-automatic Creation of Latvian Multi-word Dictionary from Small Monolingual Corpora

Inguna Skadiņa

Institute of Mathematics and Computer Science

University of Latvia



NATIONAL  
DEVELOPMENT  
PLAN 2020



EUROPEAN UNION  
European Regional  
Development Fund

INVESTING IN YOUR FUTURE

# Multi-word expressions

Often called as “pain in neck”  
(Sag et al. 2002: 1) :

*raining cats and dogs*

*New York*

*tablet computer*

*part of speech*

*go on*

*make mistake*

Multi-word  
expressions  
(MWEs) are lexical  
items that

can be  
decomposed in  
**multiple lexemes**

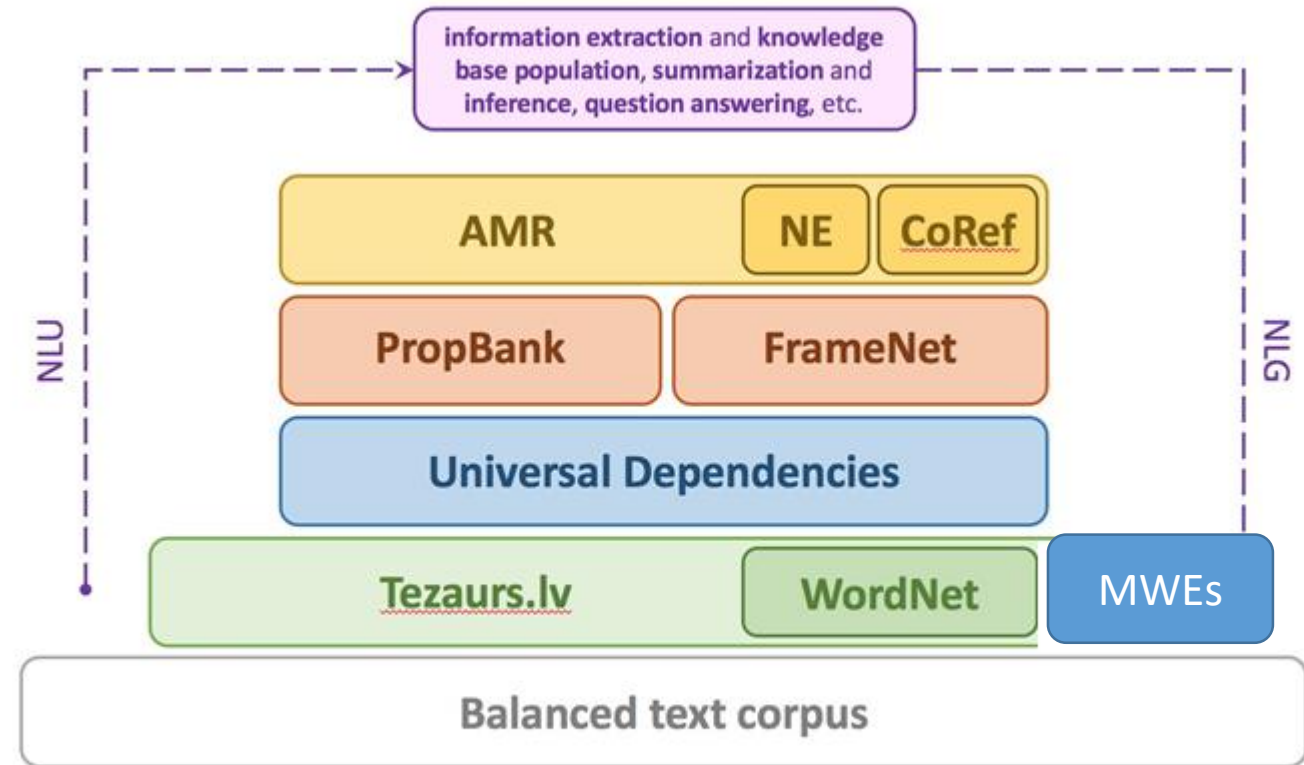


**display** lexical,  
syntactic, semantic,  
pragmatic and/or  
statistical  
**idiomaticity**

(Baldwin and Kim 2010: 269)

# Full Stack of Language Resources for Natural Language Understanding and Generation in Latvian

The project aims to create **multi-layered semantically annotated language resources** for Latvian, anchored in widely acknowledged multilingual representations (AMR, PropBank, FrameNet, Universal Dependencies, Grammatical Framework), that are required for the development of natural language understanding and generation applications.

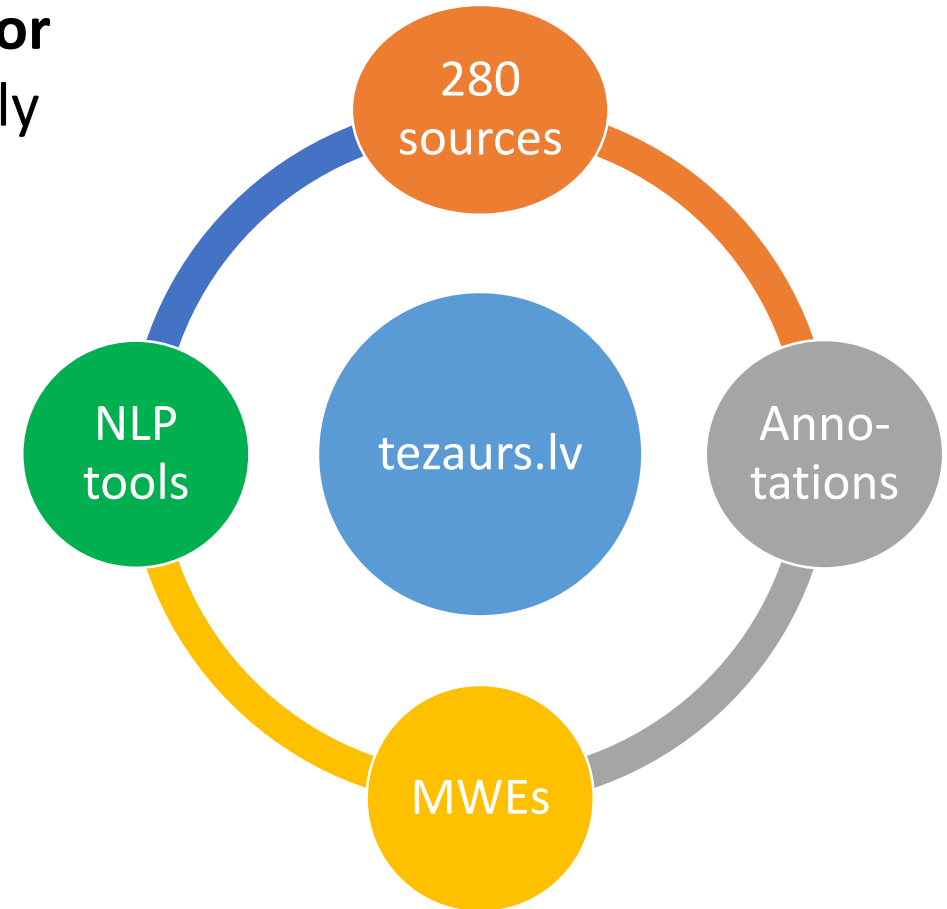


Full Stack of Language Resources for Natural Language Understanding and Generation in Latvian

# *Tezaurs.lv* - the largest open lexical database for Latvian

*Tezaurs.lv* aims to be the **central computational lexicon for Latvian**, bringing together all Latvian words and frequently used multi-word units and allowing for the integration of other LT resources and tools:

- currently contains 295 760 lexical entries
- compiled from more than 280 sources
- the dictionary is
  - enriched with phonetic, morphological, semantic and other annotations
  - enhanced with language processing tools allowing generation of inflectional forms and selection of corpus examples on the fly
- more than 2000 requests per day



Homonīmi

doma<sup>1</sup>  
doma<sup>2</sup>

**doma<sup>1</sup>** -as, s.

1. Domāšanas rezultāts (atzīņa, spriedums, ideja, pieņēmums).

// parasti vsk. Nodoms (ko darīt).

*Atmest domu* — atteikties no kāda nodoma.

// parasti vsk. Ideja, atzīņa (piemēram, mākslas darbā).

2. parasti dsk. Domāšanas process.

*legrimt (arī nogrimt) domās* — pilnīgi nodoties pārdomām.

3. parasti dsk. Uzskats, viedoklis.

*Domas dalās* — ir dažādi uzskati, viedokļi.

*Atsevišķas domas* — (a) Domas vai spriedums, kas nesaskan ar vairākuma domām. (b) Tiesas locekļa rakstveida paziņojums tai gadījumā, ja viņš nepievienojas tiesas lēmumam vai spriedumam.

Sub-senses

Glosses with synonymous cross-references (links)

Senses

Multi-word units

FRAZEOLŌGISMI: + ← Idioms (collapsed)

AVOTI: LLW ← Sources

MORFOLOĢIJA: lietvārds, sieviešu dzimte, 4. deklinācija —

	Vsk.	Dsk.
Nom.	<i>doma</i>	<i>domas</i>
Ģen.	<i>domas</i>	<i>domu</i>
Dat.	<i>domai</i>	<i>domām</i>
Akuz.	<i>domu</i>	<i>domas</i>
Lok.	<i>domā</i>	<i>domās</i>

Morphological description and the inflection table

KŌRPUSA PIEMERI: + ← Corpus examples (collapsed)

# The aim of this study

The aim is to extract lists of **good quality MWE candidates**, which

- can be delivered as open experimental MWE lexicon for Latvian and
- after manual inspection could be added to the largest Latvian open lexical database *tezaurs.lv*

# Strategies for MWE identification and extraction

MWE identification and extraction usually consists of two main steps:

- at first, morpho-syntactic patterns are applied to extract initial list of MWE candidates
- then, the initial list of MWE candidates is filtered by means of statistical measures

In this research:

- we start with application of statistical measures, allowing to identify wide range of MWE categories
- then apply linguistically motivated filters (patterns) to clean the list of initially extracted MWE candidates

# Limitation: rather small amount of data

Corpus	Size			
	Sentences (thousands)	Tokens (million)	Unique tokens (thousands)	Unique lemmas (thousands)
Balanced Corpus of the Modern Latvian language	148	5,54	408,01	111,59
Latvian-Lithuanian parallel corpus	223	3,24	307,53	87,88
Open Subtitles corpus	454	2,37	117,01	56,44

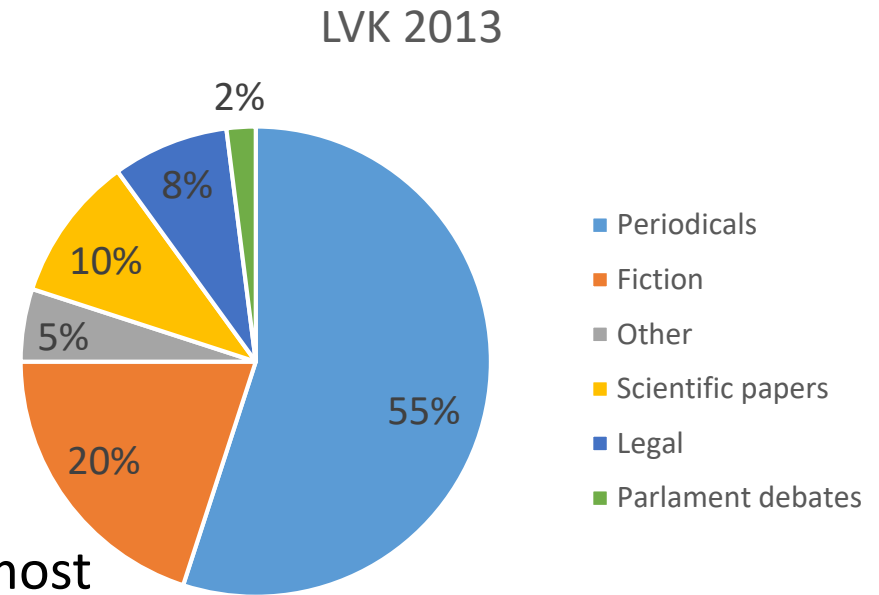


# Application of statistical measures

- **Frequency** - frequency usage of n-gramms
- **mi-score** (mutual information score): measures the strength of association (frequency of co-occurrence vs. separate occurrence). mi-score of 3 or higher is usually considered to be significant
- **t-score**: measures the confidence of association and could be applied also for low frequency words, the t-score of 2 or higher is considered to be statistically significant

# Application of statistical measures

- The Balanced Corpus of the Modern Latvian Language contains **only modern literal Latvian language texts**
- Sentences: 148 K, tokens 5,54 (million), unique tokens: 408,01 K, unique lemmas: 111,59 K
- Length of word sequences - up to 5
- Experiments:
  - at first mi-score and t-score were used individually and the most frequent MWE candidates were investigated
  - then we applied t-score as initial filter and afterwards sorted results by mi-score and frequency
- The threshold for t-score was set to 2.5, but for mi-score - 3



# Application of Statistical Measures

identified by mi-score		identified by t-score		identified by t-score, filtered by mi-score	
ordered by mi-score	ordered by frequency	ordered by t-score	ordered by frequency	ordered by mi-score	ordered by frequency
<b>very high (more than 70) mi-score, no MWE candidates</b>	<b>4 MWE candidates, short, but stable phrases, e.g., multi-word conjunctions</b>	<b>4 MWE candidates</b>	<b>4 MWE candidates, short, but stable phrases, e.g., multi-word conjunctions</b>	<b>4 MWE candidates, complex noun phrases</b>	<b>7 MWE candidates, verbal constructions or nouns followed by relative clause</b>
"(Caune, Rata, Grigule, Sviklis, Ugaine"	<b>kā arī (also)</b>	<b>kā arī (also)</b>	<b>kā arī (also)</b>	nolikums" (Latvijas Vēstnesis, 168 (3116), 22.10.2004.	<b>stājas spēkā (enter into force)</b>
"Pirts, baseini, vanna, solārijs, sports"	tas ir (it is)	(Ar grozījumiem, kas izdarīti ar (with amendments made by)	tas ir (it is)	<b>pensiju shēmas līdzekļu pārvaldītāju reģistrā (register of the pension scheme asset managers)</b>	<b>kas stājas spēkā (that enters into force)</b>
kurējās uguns vilinot kniņšļus gainājot	<b>stājas spēkā (enter into force)</b>	<b>likuma redakcijā, kas stājas spēkā (the law version that comes into force)</b>	kas ir (it is)	nolikums" (Latvijas Vēstnesis, 76 (3652), 11.05.2007.)	<b>likumu, kas stājas spēkā (the law that enters into force)</b>
aizā kurējās uguns vilinot kniņšļus	to, ka (the fact that)	<b>ne tikai (not only)</b>	<b>stājas spēkā (enter into force)</b>	nolikums" (Latvijas Vēstnesis, 124 (3072), 06.08.2004.)	<b>Ministru kabineta (Cabinet of Ministers)</b>

# Lemmatization

- When the MWE candidate list is ordered by mi-score, **four named entities** and **five terms** are among top 10 MWE candidates.
- When the list is ordered by frequency, **3 complex function words** and **2 frequent MWEs** are included.

Word sequences with highest mi-score (9 candidates)	Most frequent word sequences (5 cand.)
Arco Real Estate ' ' (company name)	kā arī (also)
pārvalde priekšnieks palīdzēja Linda Zubāne (assistant chief of administration Linda Zubāne)	pants punkts (article)
Černobiļa AES avārija sekas likvidēšana (Chernobyl nuclear plant disaster recovery)	, kas būt (which is)
šķirne ' Koričnoje Novoje ' (named entity)	tas , ka (the fact that)
jaukt dispersija kovariāta analīze iegūta (mixed variance covariance analysis provides)	kaut kas (something)
ar akūtu katarāli strutot endometriālu (with acute catarrhal stomach endometritis)	tas , kas (that/what)
ar hronisku katarāli strutot endometriālu (with chronic catarrhal stomach endometritis)	būt ļoti (to be very)
pārvalde priekšnieks palīdzēja Ieva Sietniece (assistant chief of administration Ieva Sietniece)	viens no (one of)
Valmiera / Rūjiena / Strenči-1 (list of names)	stāties spēkā (enter into force)
līcis piekraste krasts kā aizsargjosla (costal protection zone)	pēc tas (after)

# Filtering MWE Candidates

## Statistical filters (thresholds):

- high frequency (and mi-score between 4 to 11) is better signal that the string could be MWE, than high mi-score and low frequency (e.g. below 10)
- in case of t-score – high frequency together with high t-score is a signal for a good MWE candidate
- if t-score is used as initial filter and mi-score is used as the second filter,
  - most of MWE candidates will be frequent and have mi-score value be between 10 and 35,
  - or
  - will have high mi-score and low frequency.

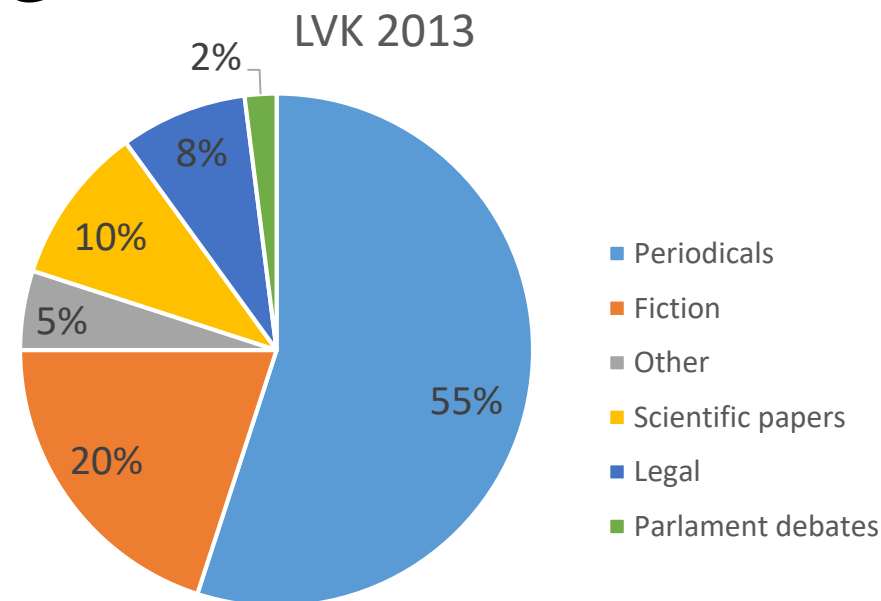
# Linguistic filters

- Simple regular expressions - filters out ungrammatical word sequences, e.g.:
  - *un (and), vai (or)* as the last word
  - *būt (to be)* at the beginning or end of string consisting of two words (e.g. *būt ļoti (to be very)*)
- Morpho-syntactic filters:
  - filters out word sequences that are ungrammatical  
*ir ļoti (is very)*, however *stājas spēkā (comes into force)*
  - extracts specific categories of MWEs
- Overlapping MWE candidates  
*stājas spēkā (comes into force), stājas spēkā ar (comes into force from), kas stājas spēkā (which comes into force)*

choice of the most appropriate MWE needs to be made by lexicographer.

# Results: Balanced Corpus of the Modern Latvian Language

- The corpus is balanced and contains **only modern literal Latvian language texts**, about 4,5 million words in total.
- Hypothesis: this corpus is a good source to identify different types of MWEs that occur in Latvian rather frequently.
- Findings: simple statistical measures applied on this corpus allow us mainly to identify good MWE candidates for the legal domain (e.g., *stājas spēkā* - *comes into force*)



# Limitation: 2-3 tokens

- When MWEs were identified by mi-score, all top 10 word sequences are MWEs. However, most of MWE candidates are named entities
- In case when MWEs were identified by t-score, 5 strings are named entities, 4 are terms and one (JP NVO RV) is string of characters.

mi-score (10 candidates)	t-score (9 candidates)
Legacy by Angosturs (named entity)	Arco Real Estate (company name)
Eastgate Properties Limited (company name)	Satja SAI Baba (company name)
Nike Riga Run (named entity – event)	<b>Pīrsons hī kvadrāts (Pearson's chi-square)</b>
ģenerāldirektors Jespers Koldings (general director Jesper Kolding)	<b>katarāli strutot endometrīts</b> (catarrulous endometritis)
Arco Real Estate (company name)	<b>JP NVO RV</b>
fon den Brinkena (name)	<b>amonijs nitrāts slāpekļis</b> (ammonium nitrate nitrogen)
Satja SAI Baba (company name)	Ge Money Bank (named entity – bank)
Satja Sai Baba(company name)	Parex Asset Management (named entity)
Latvian Art Theory (named entity)	New York Time (named entity)
<b>Pīrsons hī kvadrāts (Pearson's chi-square)</b>	<b>jaukt dispersija kovariāt</b> (mixed variance covariance)



# t-score as measure for term extraction

Frequency	Mi-score	MWE candidate
33	27.988560	ģenētiski modificēt kultūraugi ( <i>genetically modified crops</i> )
34	27.120896	ģenētiski modificēt mikroorganismi ( <i>genetically modified microorganism</i> )
42	26.717372	noziedzīgs nodarījums izdarīšana ( <i>committing criminal offences</i> )
112	26.346762	<b>LPP / LC</b> (name of party)
137	26.146896	ģenētiski modificēt organismi ( <i>genetically modified organism</i> )
168	25.801889	fondēta pensija shēma ( <i>funded pension schema</i> )
9	25.669791	konkurētspējīga priekšrocība pārvešana ( <i>competitive advantage transfer</i> )
8	25.222240	civila aizsardzība aizsargbūve ( <i>civil defence protection structure</i> )
37	25.169589	<b>° C temperatūra</b> ( <i>° C temperature</i> )
31	25.112838	infekcija slimība izraisītājs ( <i>infectious disease agent</i> )

The previous experiment shows that t-score allows better to identify terms that can be included into electronic dictionary. Therefore the threshold for t-score was raised up to 10: **8 terms, one named entity** (LPP/LC – name of party) and one sequence of words (*° C temperatūra*) was identified between top 10 candidates

# Extraction of verbal phrases

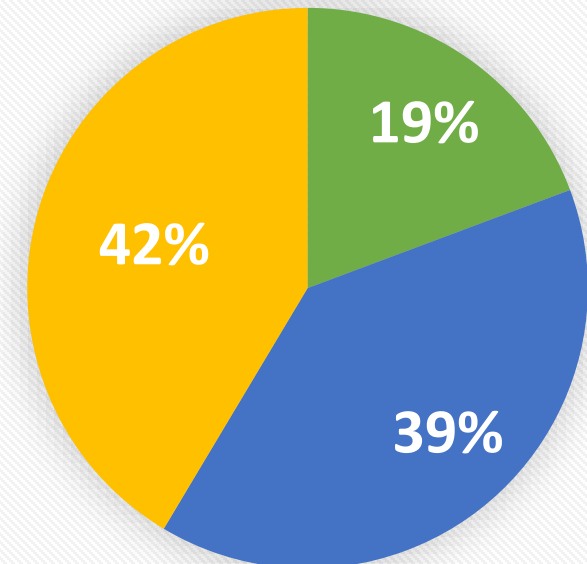
<b>Freq.</b>	<b>Mi-score</b>	<b>MWE candidate</b>
2006	44.740771	<b>stāties V spēks N</b> (come into force)
623	24.896466	<b>pieņemt V lēmums N</b> (to make decision / decide)
290	16.915062	<b>dot V iespēja N</b> (to enable)
247	15.360174	<b>tikt V gals N</b> (to manage)
218	14.595137	<b>veikt V pētījums N</b> (to do research)
141	11.759475	<b>sniegt V informācija N</b> (provide information)
130	11.393867	<b>pievērst V uzmanība N</b> (pay attention)
115	10.505124	<b>tiesības N saņemt V</b> (rights to receive)
104	10.187839	<b>ienākt V prāts N</b> (to come into one's head)
92	9.581201	<b>aizvērt V acs N</b> (to close eyes)

From top 10 MWE candidates, 7 could be accepted as MWEs: all are fixed in Latvian-English dictionary (Veisbergs 2005) and four of these MWEs are included in *tezaurus.lv*.

# Latvian-Lithuanian Corpus LiLa

- LiLa contains more than half of the texts that are originally written in Latvian (19,3%) or Lithuanian (39,3%)
- 41,4% of corpus are EU legal texts that were chosen due to lack of translations between Latvian and Lithuanian.
- The corpus contains 8,7 million tokens in total; 3,24 million tokens (no legal documents) were used in this experiment
- Texts originally written in Baltic languages represent domains:
  - **modern fiction (86%),**
  - periodicals (5,9%),
  - popular literature (5,6%)

■ LV-LT: 1695160  
■ LT-LV: 3448745  
■ EU documents: 3638145



# Latvian-Lithuanian Corpus

- Hypothesis: LiLa contains more frequently used fixed phrases and idiomatic expressions than the Balanced Corpus of Modern Latvian
- Results (ordered by mi-score): 7 MWE candidates are named entities, 1 fixed phrase, 1 part of longer phrase, 1 is a character string
- Results (ordered by frequency): 4 terms, 1 named entity, 1 complex function word, 4 other parts of longer phrases

mi-score (7 named entities)	Frequency (6 MWE candidates)
SIA "AD BALTIC" (company)	<b>apkure katls</b> (central heating boiler)
Ventspils peldbaseins relaksācija komplekss	<b>apkure iekārta</b> (heating system)
izstāde "Tech Industry" (event)	<b>Ventspils peldbaseins</b> (Ventspils swimming pool)
"Tech Industry" (event)	koksne granula (wooden pellet)
"AD BALTIC" (named entity)	katls iekārta (boiler equipment)
SEALEY POWER products (named entity)	granula apkure (pellet heating)
<b>W / m</b>	informācija par (information about)
<b>peldbaseins relaksācija komplekss</b> (swimming pool relaxation complex)	<b>sāls istaba</b> (salt room)
<b>Ventspils peldbaseins relaksācija</b> (Ventspils swimming pool relaxation)	<b>relaksācija komplekss</b> (relaxation complex)
REN TV Baltija (named entity)	<b>ne tikai</b> (not only)

# Open Subtitles Corpus

- Sentences: 454 K, tokens - 2,37 million, unique tokens - 117,01 K, unique lemmas - 56,44 K
- We used lemmatized corpus and applied t-score for identification and mi-score as the second filter. In addition, MWE candidates were filtered by frequency (in previous experiments only threshold 5 was used)
- The outcome of this experiment differs from the previous ones – besides named entities, different idiomatic expressions are identified

- We' il help you.

- Viktor, come on.

I owe you so much.

- Please let me help you.

- This is your friend talking.

We' re your family now.

Goodbye.

# Open Subtitles Corpus

- The outcome of this experiment differs from the previous ones – besides named entities, different **idiomatic expressions are identified**
- Similarly to the previous experiments, many (4 when frequency is at least 5 and 6 if frequency is at least 10/ 15) of extracted MWEs are named entities.
- However, different idiomatic expression are identified as well

Freq>=5	Freq>=10	Freq>=15
it ' s not going	it ' s not going	Bārts Šērmens (named entity)
viesnīca “ dižena Budapešta ” (hotel “Great Budapest”)	Bārts Šērmens (named entity)	“ Pearson Hardman
<b>dzīvot laimīgi līdz mūžs gals (live happily till end of life)</b>	dzeršana no zābaks (drink from boot)	“ Folsom foods ”
misis boss (named entity)	paskriet , paostīt , sarauties (run, sniff, cringe)	“ SouthJet ” 227
Bārts Šērmens (named entity)	Vašingtona māksla noziegums nodaļa (Washington Arts Crime Division)	“ Delta psi ” (named entity)
Rikijs Pontings (named entity)	<b>laimīgi līdz mūžs gals (happily till end of life)</b>	“ mežonīgs vepris ” (“wild hog” - name of bar)
dzeršana no zābaks (drink from boot)	“ SouthJet ” 227	pakārt viņš (hang him)
paskriet , paostīt , sarauties (run, sniff, cringe)	“ zēns ar ābols ” (“Boy with Apple” – painting)	<b>daudz laime dzimšana diena (happy birthday)</b>
<b>gulēt saldi (sleep well)</b>	“ Wayne enterprise ”	<b>ar tas nebūt nekāds sakars (nothing to do with it)</b>
kosmos kuģis (spaceship)	“ Pearson Hardman ”	dzeršana no zābaks (drink from boot)

# Conclusion

- In case of small amount of general domain (balanced) data automatic methods allow finding good MWE candidates – terms or named entities
- However, finding idiomatic expressions in small, general domain corpora is looking for the needle in haystack: only larger or more expressive corpus could help in identification process
- In case of small general parallel corpus:
  - the most reliable results are obtained for named entities
  - terms and complex function words could be also identified, but in this case more careful manual inspection is necessary
- If the aim of the MWE identification is to identify idiomatic expressions that recently have appeared in a language, then the corpus needs to represent more everyday language and have to be rather big, because idiomatic expressions are rare in balanced corpora that represent literary language and carefully edited texts.

# Thank you!

This work is supported by the European Regional Development Fund under the grant agreements No. 1.1.1.1/16/A/219 (*Full Stack of Language Resources for Natural Language Understanding and Generation in Latvian*)

NATIONAL  
DEVELOPMENT  
PLAN 2020



**EUROPEAN UNION**  
European Regional  
Development Fund

---

INVESTING IN YOUR FUTURE