

NEURAL MT AND OTHER LANGUAGE TECHNOLOGIES AT TILDE

Dr. Raivis SKADIŅŠ

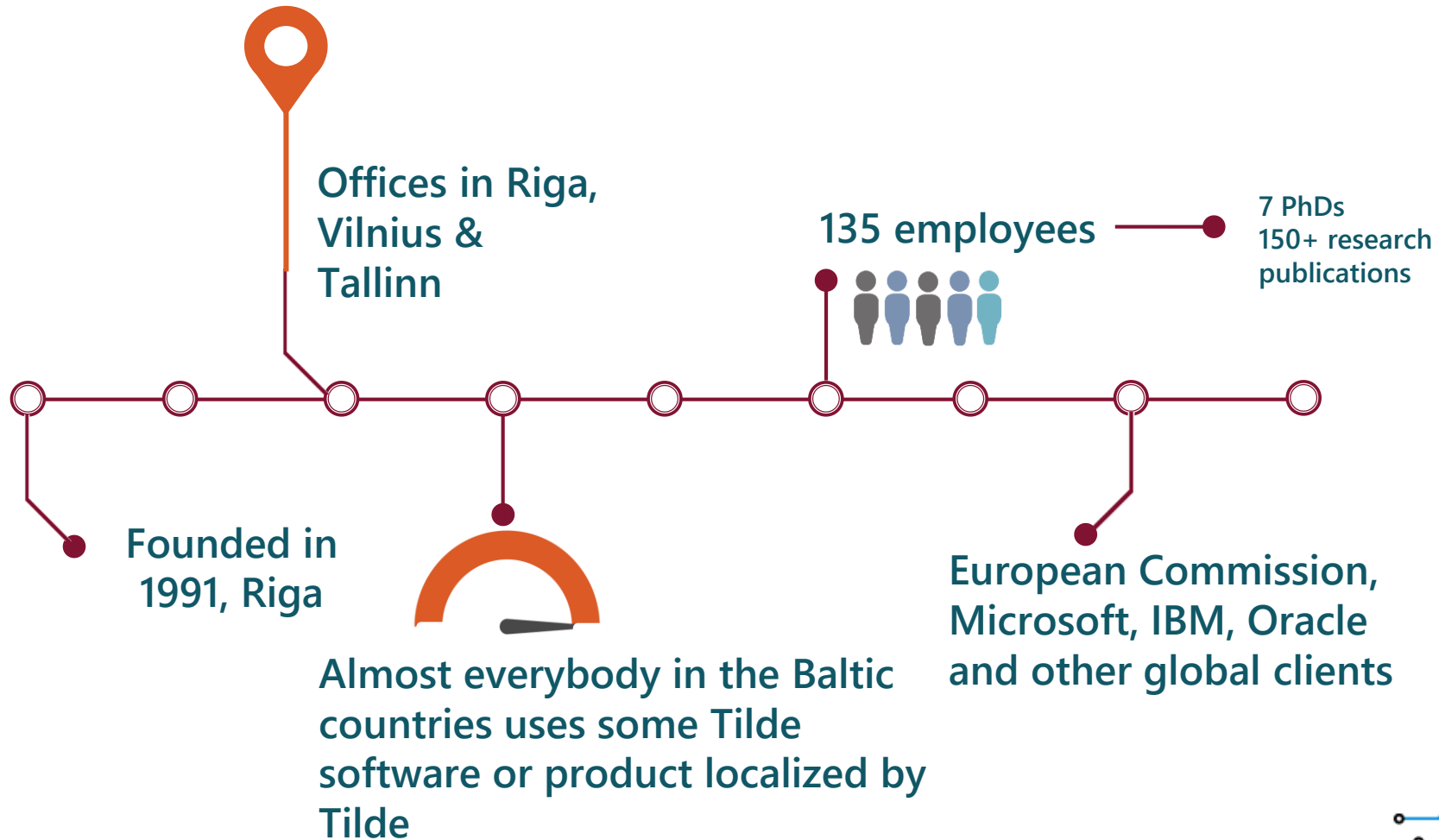
Tilde, Director of Research and Development

Fifth GF Summer School 2017, Riga, August 18, 2017

In my talk

- About Tilde and what we do
- Grammar Checking
- Neural Machine Translation





What we do



- All kinds of language technologies
 - spelling checkers
 - electronic dictionaries
 - terminology
 - encyclopedias
 - grammar checkers
 - machine translation
 - speech recognition and synthesis
 - virtual assistants and chatbots



What we do



- Wide range of clients
 - home and office users
 - localization companies
 - enterprise clients
 - governments
 - EU infrastructure projects
- Research projects



Grammar Checking



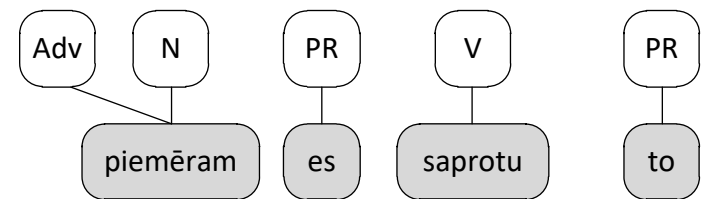
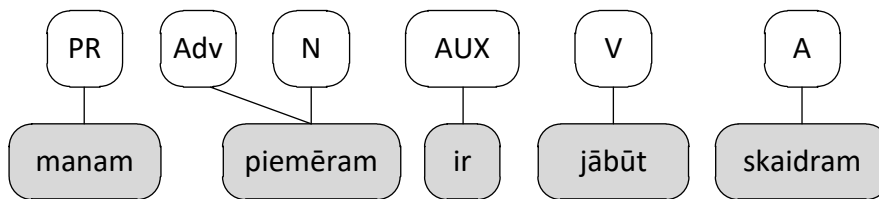
How we do it



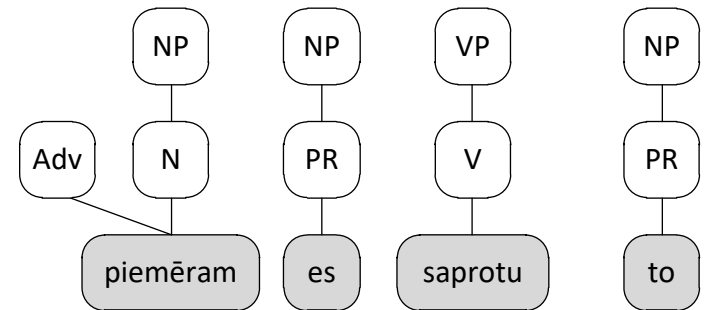
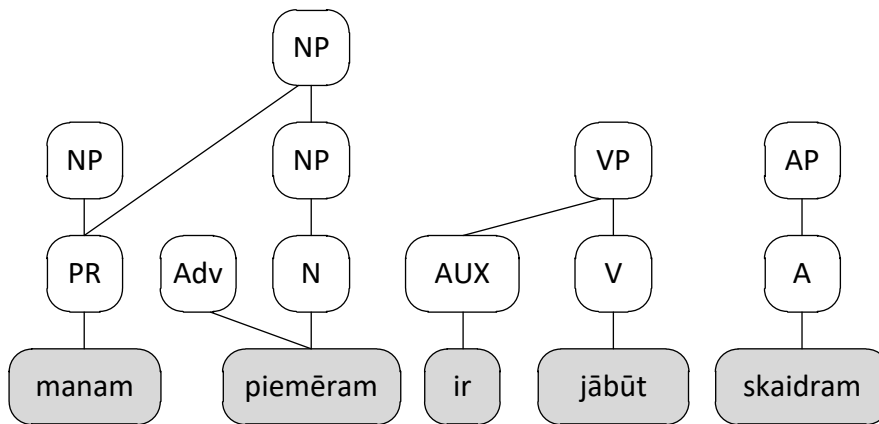
- If you can parse the sentence, then it is correct
- **But**, if you cannot parse it
 - It is wrong
 - Your grammar is incomplete
- Is it really so simple?
- Will any parser do?
- How to find the error? How to fix it?



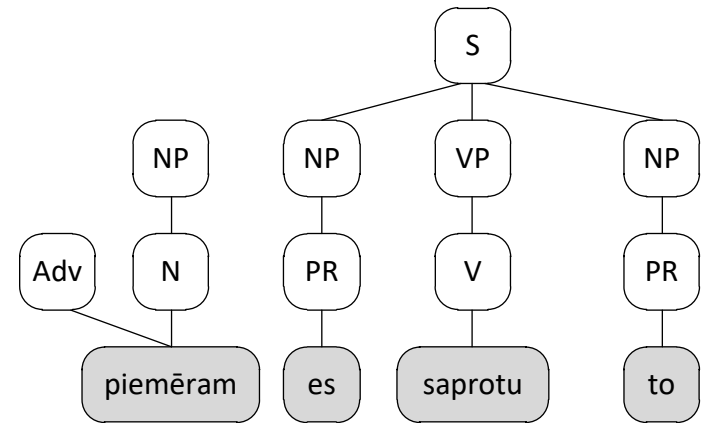
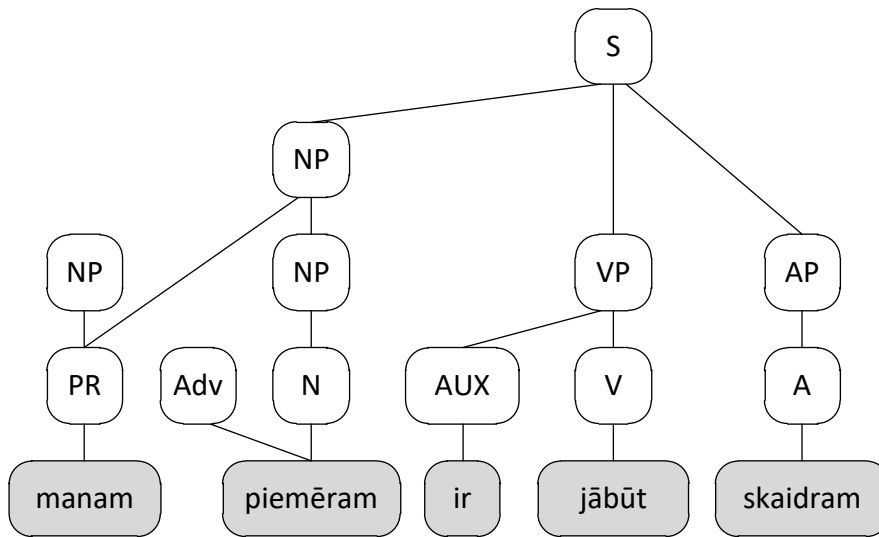
Example



Example



Example



Some examples of rules

NP -> attr:AP main:NP

Agree(attr:AP, main:NP, Case, Number, Gender)

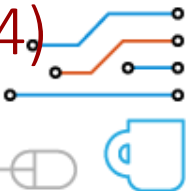
S -> subj:NP main:VP obj:NP

Agree(subj:NP, main:VP, Person)

subj:NP.Case == Nom

obj:NP.Case == Acc

- And there are hundreds of them; (Deksne et al., 2014)



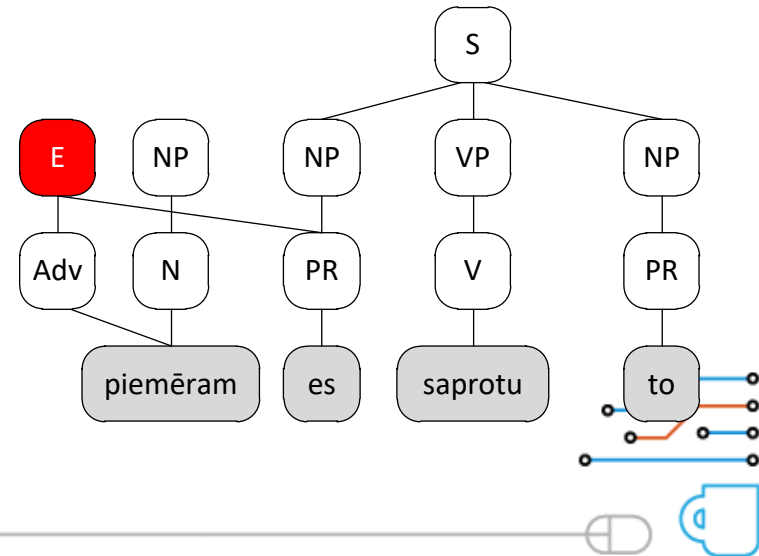
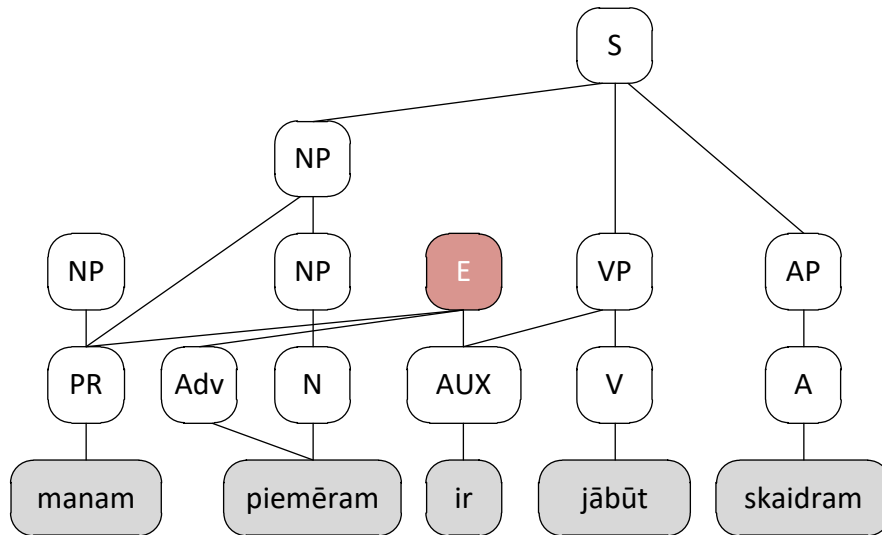
How to find the error?

- Two types of rules
 - Regular rules that describe syntax
 - Rules that describe errors
- We parse the sentence with both at the same time
- There is an error, if
 - an error rule has been applied
 - fragment where it has been applied cannot be parsed with regular rules

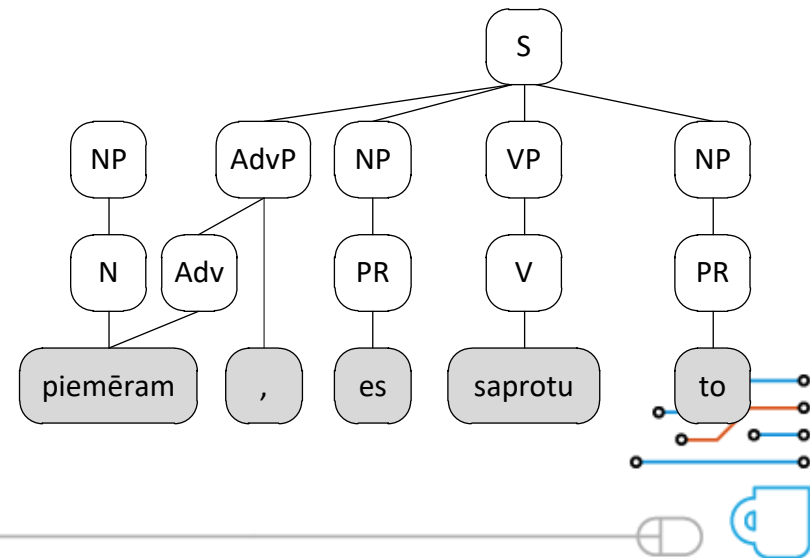
(Deksne & Skadiņš, 2011)



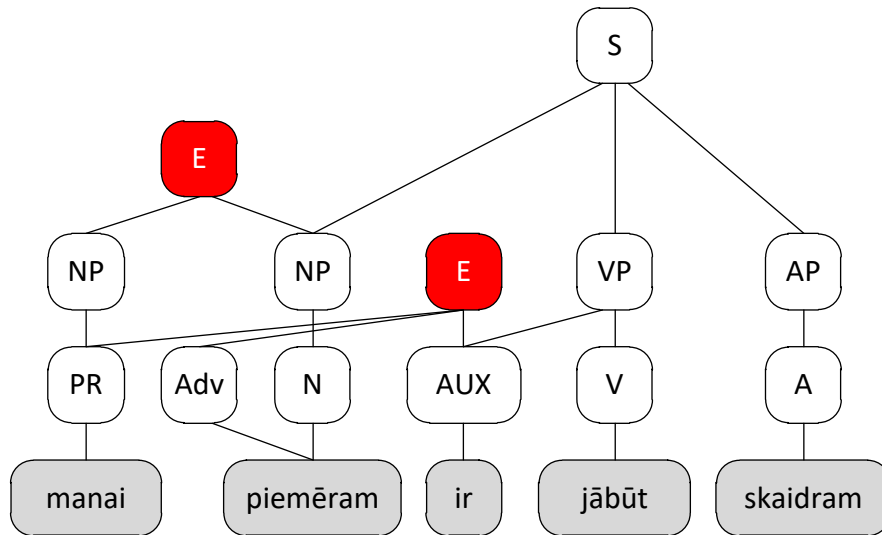
Example



Example



Example



Error rules



ERROR-1 -> attr:AP main:NP

Disagree(attr:AP,main:NP, Case, Number, Gender)

GRAMMCHECK MarkAll

attr:AP.Gender=main:NP.Gender

attr:AP.Number=main:NP.Number

SUGGEST(attr:AP+main:NP)



Error rules



ERROR-14 -> attr:N attr:G main:N

attr:N.Case==genitive

attr:N.Number==singular

attr:G.AdjEnd==definite

main:N.Number==plural

Agree(attr:G, main:N, Case, Number, Gender)

CapPattern fff

LEX Amerika savienots valsts



Rules



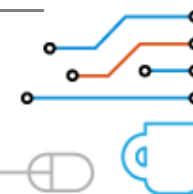
Rule type	Latvian	Lithuanian
Correct syntax rules	580	179
Error rules which depend on phrases described by correct syntax rules	263	72
Error rules which contain only terminal symbols	239	560
Total	1082	811



Evaluation



Corpus	Error type	Precision	Recall	F-measure
Lithuanian Balanced	all error types	0.898	0.412	0.564
	vocabulary errors	0.956	0.535	0.686
	incorrect usage of cases	0.734	0.259	0.383
Latvian Balanced	all error types	0.780	0.455	0.575
	punctuation in sub-clauses	0.757	0.643	0.695
	punctuation in participle clauses	0.617	0.671	0.643
Latvian Student papers (dev)	All error types	0.652	0.231	0.341
	punctuation in sub-clauses	0.706	0.586	0.641
	punctuation in participle clauses	0.656	0.560	0.604
Latvian Student papers (test)	all error types	0.753	0.203	0.320
	punctuation in sub-clauses	0.773	0.588	0.668
	punctuation in participle clauses	0.766	0.685	0.723



Machine Translation



**Rule-based
MT**

**Statistical
MT**

Neural MT



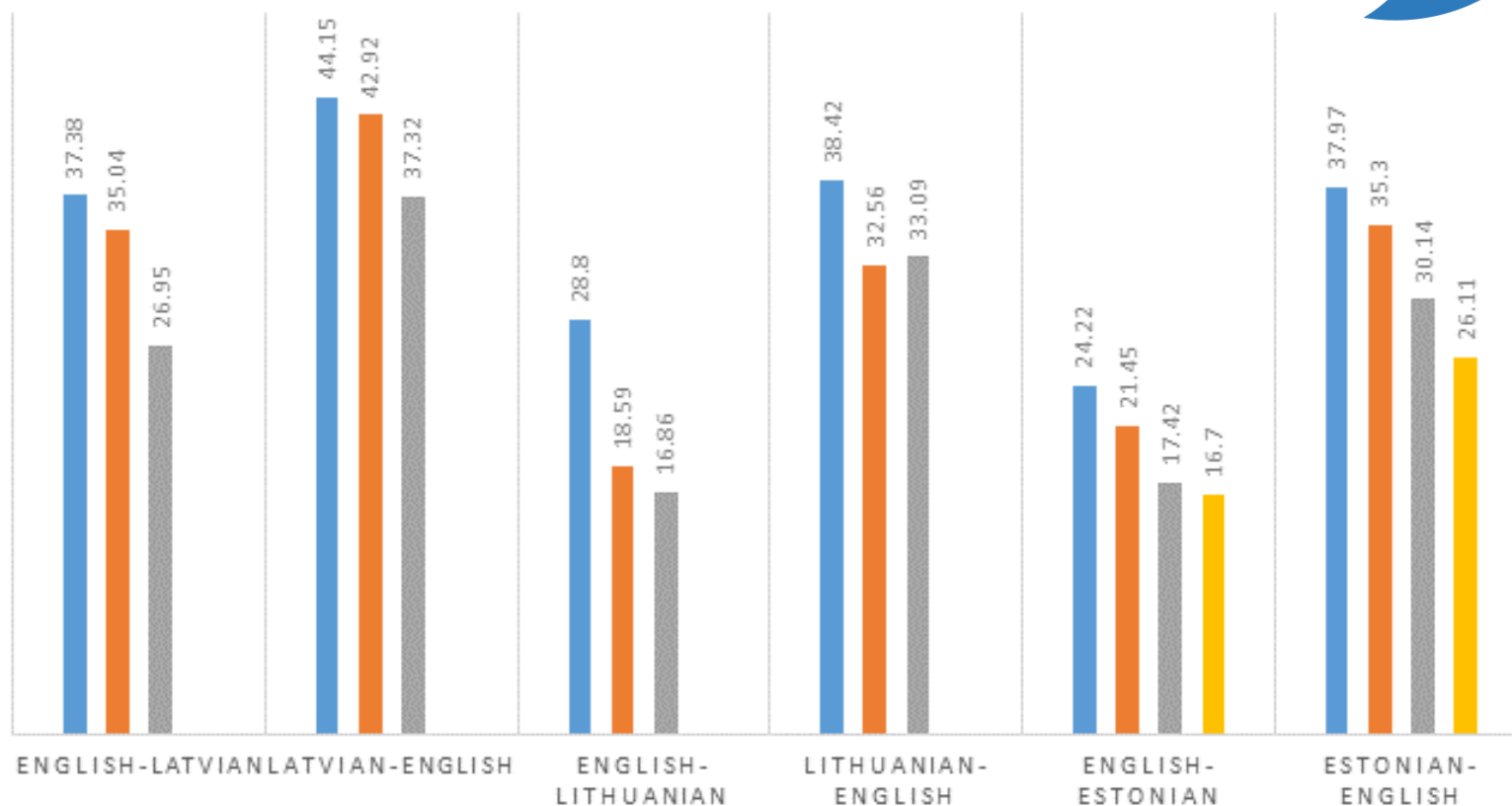
State-of-the-Art before neural MT?

Phrase-based statistical MT

Better
than
Google

BLEU SCORES OF GENERAL DOMAIN SMTS

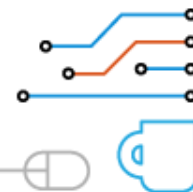
■ Tilde ■ Google ■ Microsoft ■ Tartu Uni



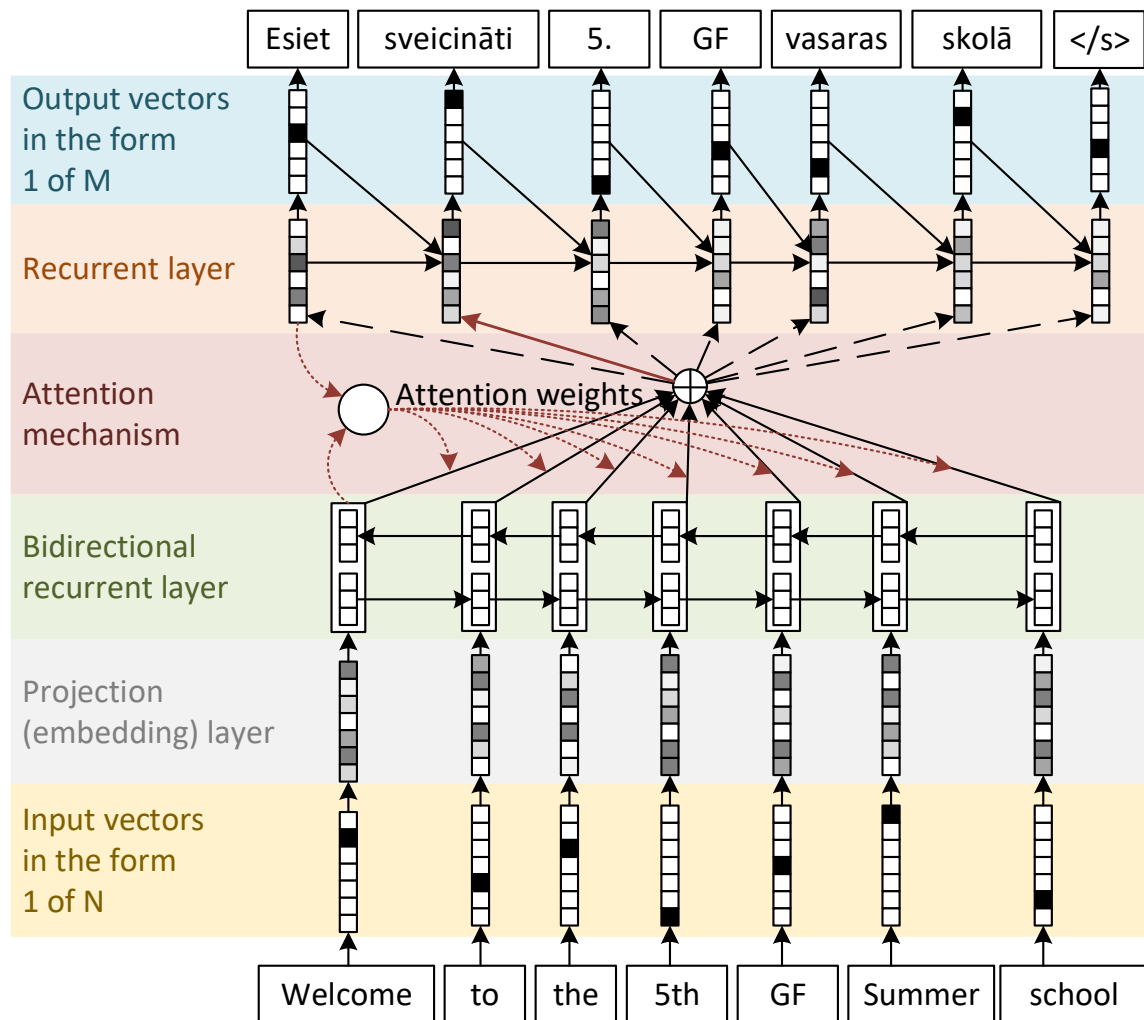
Dawn of the Neural MT



- New technology, 2015, 2016
- Very different architectures
- Many open questions
 - Is it good for Latvian and other under-resourced languages?
 - What is the quality?
 - Strengths and weaknesses?
 - Is it fast enough?
 - What infrastructure do we need?
 - etc.



Technology



- QT21 project
- Nematus and AmuNMT toolkits
- end-to-end NMT
- sub-word tokens (BPE)



Training data



Language pairs	Sentences in parallel corpus	Sentences in monolingual corpus
General domain		
en-et	21 900 622	48 567 363
et-en	21 900 794	217 724 716
ru-et	4 179 198	48 606 392
et-ru	4 179 153	138 001 100
en-lv	7 477 785	74 741 452
lv-en	7 476 956	95 259 699
<i>Pharmaceutical domain</i>		
en-lv	316 443	309 182



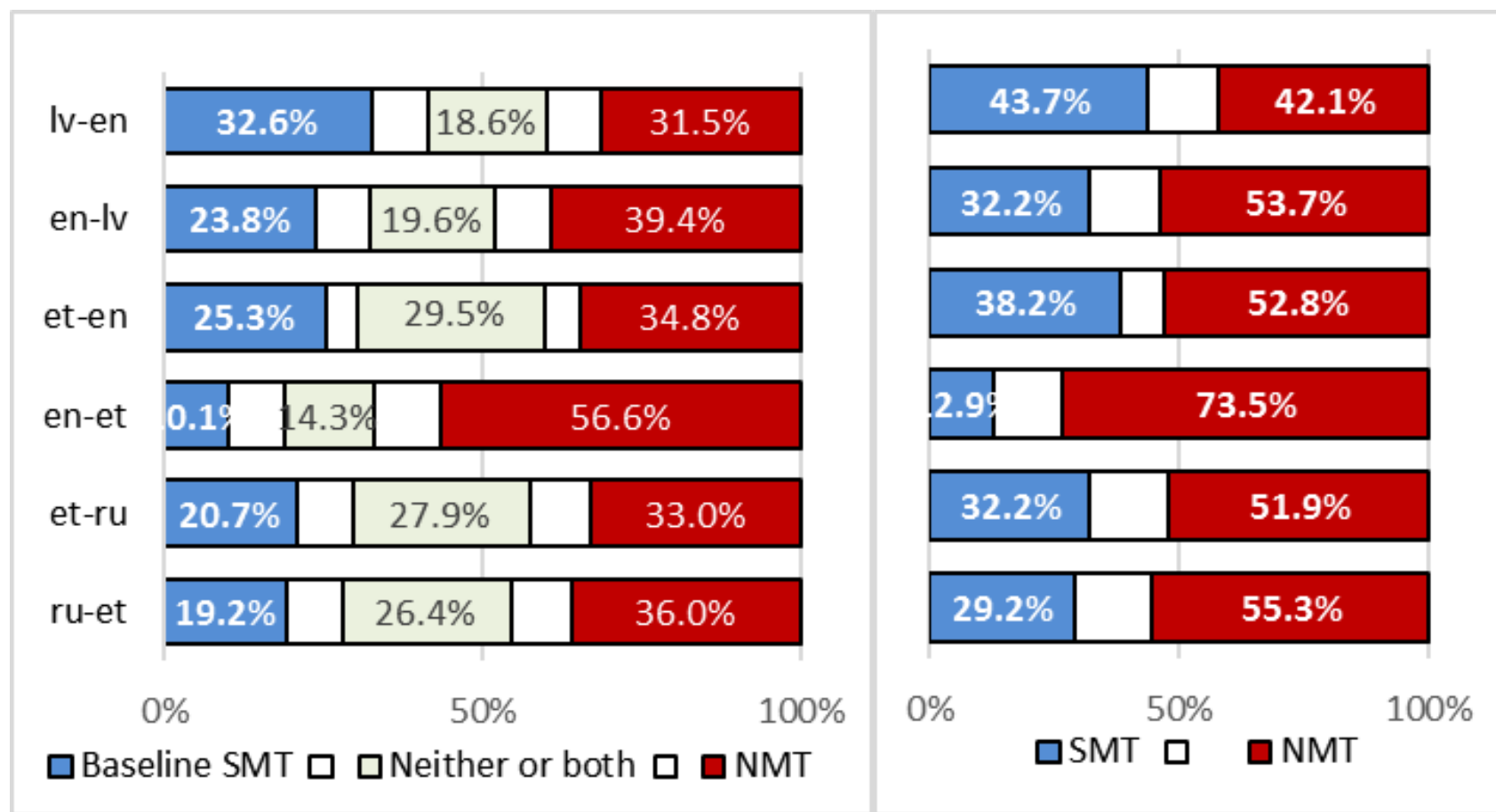
Automatic evaluation (BLEU)

Language pair	System	BLEU
en-et	Baseline SMT	22.53 (20.39-24.95)
	Google Translate (SMT)	19.80 (18.00-21.60)
	NMT	24.64 (22.76-26.54)
et-en	Baseline SMT	32.52 (30.55-34.53)
	Google Translate (SMT)	40.57 (38.48-42.84)
	NMT	31.74 (29.91-33.45)
ru-et	Baseline SMT	09.87 (08.73-11.01)
	Google Translate (SMT)	12.52 (11.03-14.01)
	NMT	09.02 (08.02-10.00)
et-ru	Baseline SMT	07.94 (07.07-08.82)
	Google Translate (SMT)	14.74 (13.18-16.15)
	NMT	09.39 (08.33-10.46)
en-lv	Baseline SMT	32.57 (29.96-35.33)
	translate.tilde.com (SMT)	37.54 (34.65-40.50)
	NMT	24.77 (22.94-26.72)
lv-en	Baseline SMT	28.79 (26.84-30.82)
	translate.tilde.com (SMT)	43.76 (41.25-46.45)
	NMT	29.62 (27.62-31.44)



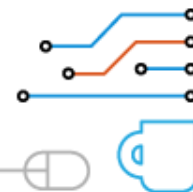
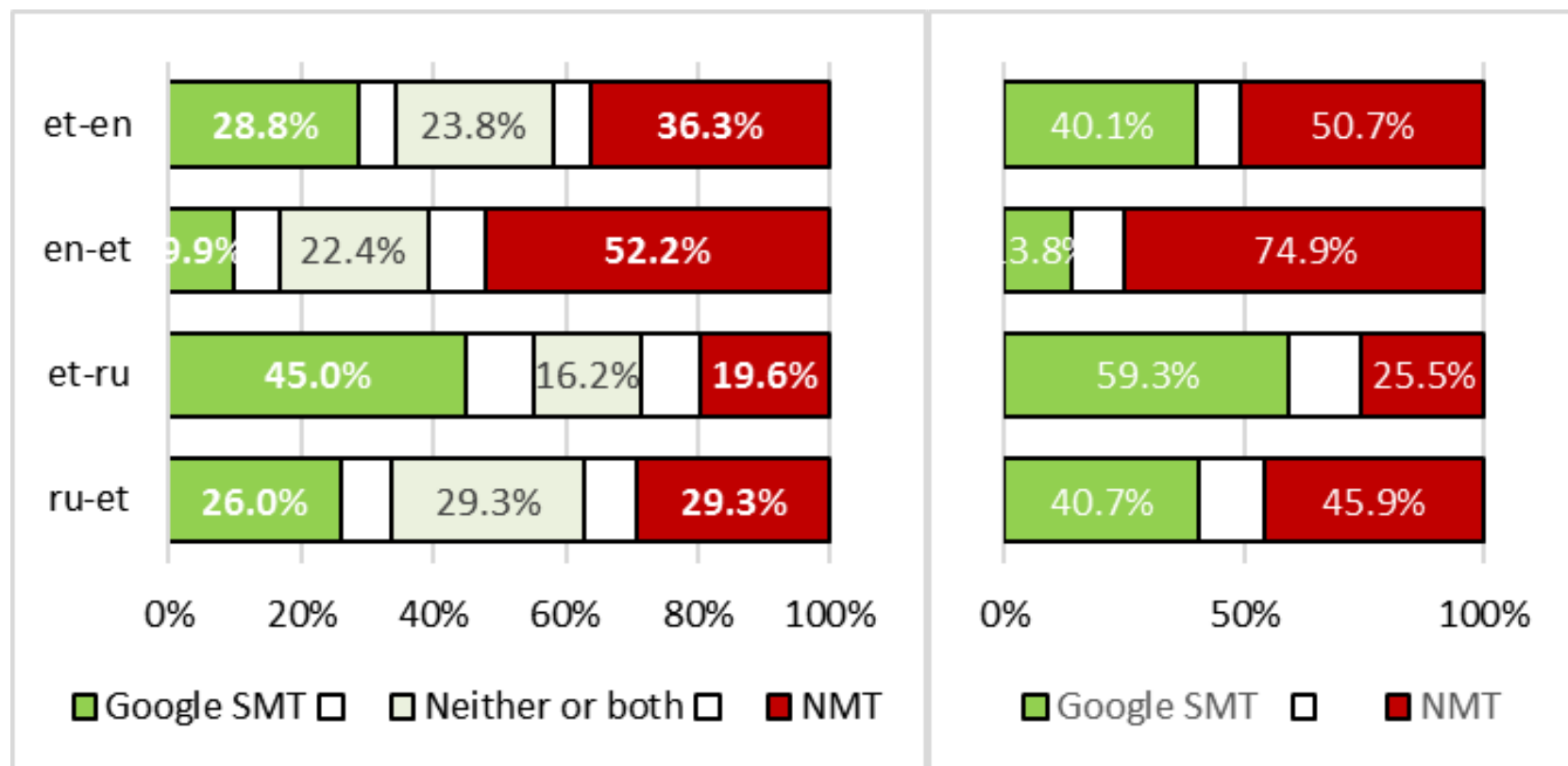
Human evaluation (system comparison)

SMT vs Neural MT



Human evaluation (system comparison)

Google Translate vs Neural MT





Tildes neironu mašīntul

translate.tilde.com/neural/lv

TILDE

Tildes neironu mašīntulkošana

EN LV LT ET

Tulkotājs

Lokalizācija un tulkošana

Tildes Birojs

Tildes Jumis

Runas atpazīnējs

Produkti

Pētniecība

Par Tildi

Pirkt

Tulkot tekstu

Tulkot dokumentu

Tulkot tīmekļa lapu

Atgriezties pie klasiskā Tildes tulkotāja

No: Angļu Latviešu Igaunu

Uz: Arābu Latviešu Igaunu

Tulkot

On December 19 2016, the Presidium of the Latvian Academy of Sciences (LAS) announced the most significant achievements of Latvian science in 2016. This year in addition to the top 11 achievements, 6 proposals were acknowledged with the Diploma of the President of the LAS.

Awards ceremony to honour the best achievements in scientific research of 2016 is due to take place on 14 February 2017, 16 PM at Riga Motor Museum, 6 S.Eizenšteina Str.

List of the top 11 + 6 achievements in science 2016 is available here.

Latvijas zinātņu akadēmijas Prezidijs (LAS) 2010. gada 19. decembrī paziņoja par Latvijas zinātnes visnozīmīgākajiem sasniegumiem 2016. gadā. Šogad papildus 11 sasniegumiem, 6 priekšlikumi tika apstiprināti AR LAS Priekšsēdētāja diplomu.

Apbalvošanas ceremonija notiek, lai godinātu labākos sasniegumus 2016. gadā zinātniskajā pētniecībā, kas notiks 2017. gada 14. februārī, 16 PM Rīgas automobilī, 6 S.Eizenšteina Str.

Saraksts ar 11 + 6 sasniegumiem zinātnē, 2016. gadā ir pieejams šeit.

Rediģēt

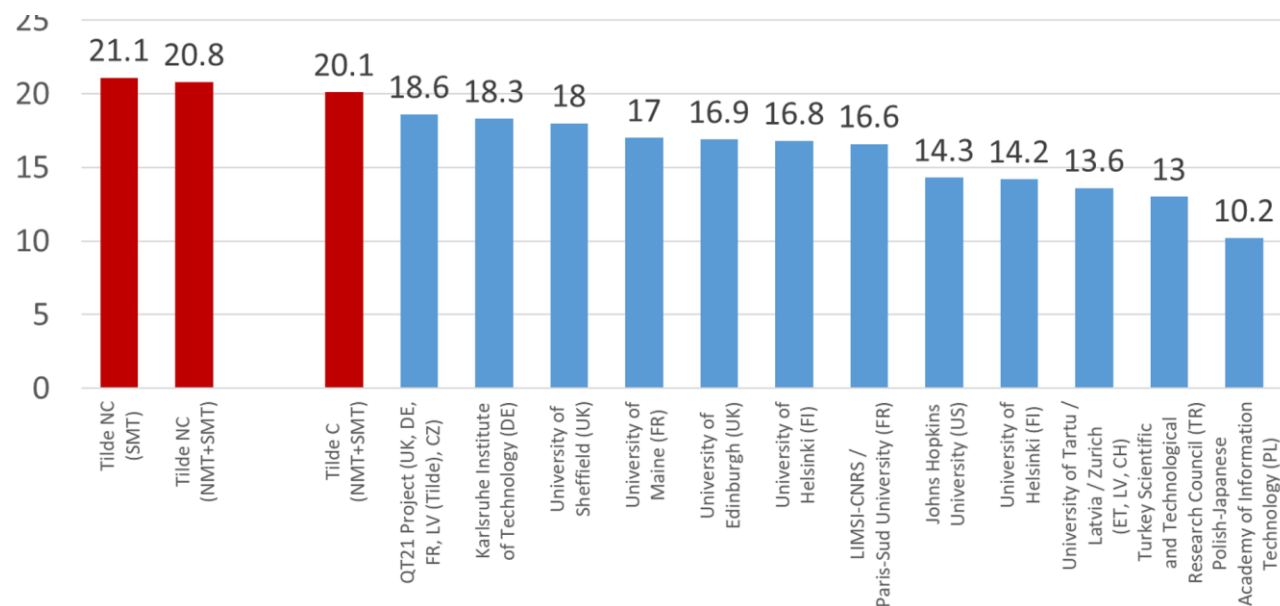
First Conclusions

- In most cases Neural MT outperforms Statistical MT in human evaluation. It is true also for under-resourced languages like Latvian and Estonian
- Fluency is much better, word agreement is better, translates even unseen words
but can hide semantic errors
- It is not a panacea, it is a field for new research and development



WMT 2017 Competition

- Yearly competition of MT researchers
- Latvian – first time this year
- Both human and automatic evaluation

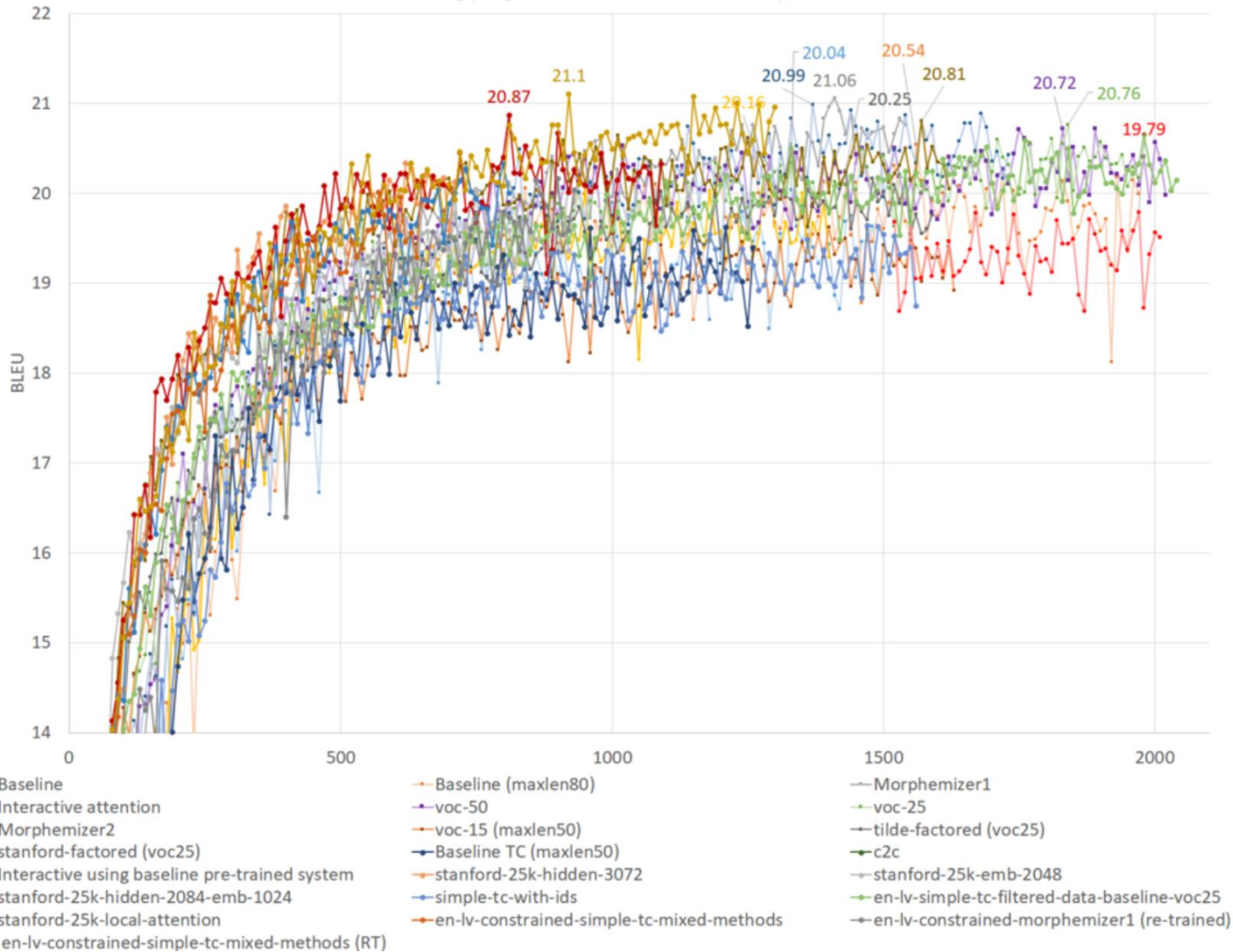


The winning system

- Nematus based NMT system
- Main improvements
 - data preprocessing and cleaning
 - special handling of numbers, ID etc. and rare words
 - hybrid with SMT
 - morphology aware sub-word units
 - factored NMT
 - back-translation of monolingual target language data
 - MLSTM recurrent neural network
- A lot of experiments with different configurations (~ 55 trained NMT systems)



Training progress for EN-LV constrained systems





Tilde's Machine Translation Systems for WMT 2017

Mārcis Pinnis, Rihards Krišlauks, Toms Miks, Daiga Dekšne, and Valters Šics

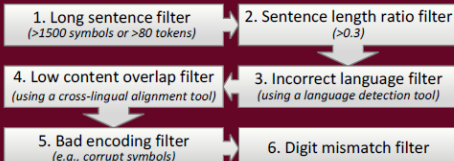
Tilde, Vienības gatve 75A, Rīga, Latvia

{marcis.pinnis,rihards.krislauks,daiga.deksne,toms.miks,valters.sics}@tilde.lv



Data Filtering

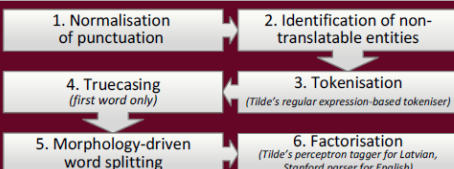
Due to noise in the training data, we performed data filtering.



Unique parallel/monolingual sentence counts

Scenario	Lang. pair	Before filtering	After filtering
Constrained	en-lv	1.92 M / 28.81 M	1.61 M / 27.75 M
	lv-en	1.92 M / 335.55 M	1.61 M / 330.23 M
Unconstrained	en-lv	15.78 M / 87.60 M	12.69 M / 81.68 M
	lv-en	15.78 M / 360.01 M	12.69 M / 351.99 M

Data Pre-processing



Synthetic Data

Unknown Words as Placeholders

To make NMT models more robust to rare and unknown phenomena, we supplement the training data with sentence pairs where one to three content words are replaced with unknown word (UNK) placeholders.

Back-translation of Monolingual Data

For domain adaptation, we use a synthetic parallel corpus that is acquired by back-translation of in-domain monolingual data from the target language using a reverse NMT system.

Scenario	Lang. pair	<UNK> placeholder sent.	Back-translated sent.	Total
Constrained	en-lv	1.48 M	3.09 M	6.19 M
	lv-en	1.48 M	3.09 M	6.19 M
Unconstrained	en-lv	11.66 M	21.69 M	46.04 M
	lv-en	11.66 M	21.36 M	45.71 M

Introduction

We present Tilde's WMT 2017 MT systems that were ranked as **the best performing systems by automatic evaluation.**

Machine Translation Systems

- **SMT Systems**—Moses phrase-based systems, fast-align word alignment, 7-gram translation models, 5-gram KenLM language models, trained on the Tilde MT platform.
- **NMT Systems**—Nematus NMT systems with **MLSTM recurrent units, morphology-driven word splitting**, vocabulary size of 25,000 for constrained systems and 50,000 for unconstrained systems, decoding beam size of 12, ensembles of 5 to 7 models, back-translated data used to train final systems.
- **NMT-SMT hybrid systems**—rare words (e.g., person names, location names, different scripts, etc.) are replaced with unknown word place-holders, sentences are translated with NMT systems, after which rare words are translated with SMT systems. In unconstrained systems, a named entity data base is used to improve person name translation quality.

Example of the NMT-SMT Hybrid Translation Process

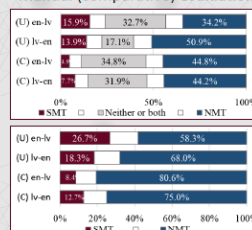
Translation step	Example sentence
Source text	šodien skatīties Ikauniecs-Admidinas startu Rio spēlēs.
Pre-processed text	šodien skat00 ieties I00 kaun00 iec00 es - Ad00 mi00 di00 nas start00 u Rio spēlē00 s .
Text with identified rare words	šodien skat00 ieties βIDβ - βIDβ start00 u Rio spēlē00 s .
NMT translation	watch the βIDβ - βIDβ start at the Rio Games today .
Moses XML with untranslated rare words	<nmt translation="watch the"> šodien skatīties </nmt> Ikauniecs <nmt translation="-"> -</nmt> Admidinas <nmt translation="start at the Rio Games today"> šodien startu Rio spēlēs</nmt> <nmt translation="-"> .</nmt>
Moses XML with identified untranslated person names	<nmt translation="watch the"> šodien skatīties </nmt> <ne translation="Ikauniecs" prob="1.0"> Ikauniecs </ne> <nmt translation="-"> -</nmt> <ne translation="Admidina Admidins" prob="0.95 0.05"> Admidinas </ne> <nmt translation="start at the Rio Games today"> šodien startu Rio spēlēs</nmt> <nmt translation="-"> .</nmt>
SMT translation	watch the Ikauniecs - Admidina start at the Rio Games today .
Post-processed translation	Watch the Ikauniecs-Admidina start at the Rio Games today.
NMT only transl. (for comparison)	Today, look at the start of the Isolence-Admidias in the Rio Games.

Evaluation

Automatic evaluation (submitted systems are marked in bold)

Scenario	Lang. pair	System	BLEU (CS)	BEER 2.0	Character
Constrained	en-lv	SMT	12.98±0.62	0.5086	0.6642
		NMT	†19.49±0.79	0.5478	0.5877
		Hybrid	†19.52±0.82	0.5482	0.5853
	lv-en	SMT	15.47±0.59	0.5219	0.6606
		NMT	†20.01±0.67	0.5494	0.6088
		Hybrid	†20.06±0.63	0.5496	0.6081
Unconstrained	en-lv	SMT	20.43±0.86	0.5491	0.6126
		NMT	20.04±0.78	0.5563	0.5832
		Hybrid	20.08±0.78	0.5567	0.5827
	lv-en	SMT	19.05±0.63	0.5515	0.6233
		NMT	†22.02±0.63	0.5677	0.5838
		Hybrid	†22.06±0.66	0.5683	0.5833

Manual (comparative) evaluation



Comparison of the best constrained system submissions

Lang. pair	System	BLEU (CS)	BEER 2.0	Character
en-lv	Tilde (hybrid)	†19.52±0.79	0.5482	0.5853
	QT21 combination	18.03±0.71	0.5403	0.6455
	KIT primary	17.72±0.69	0.5428	0.6051
lv-en	Tilde (hybrid)	†20.06±0.65	0.5496	0.6081
	UEDIN NMT	19.08±0.65	0.5462	0.6308
	JHU SMT	16.95±0.60	0.5281	0.6485

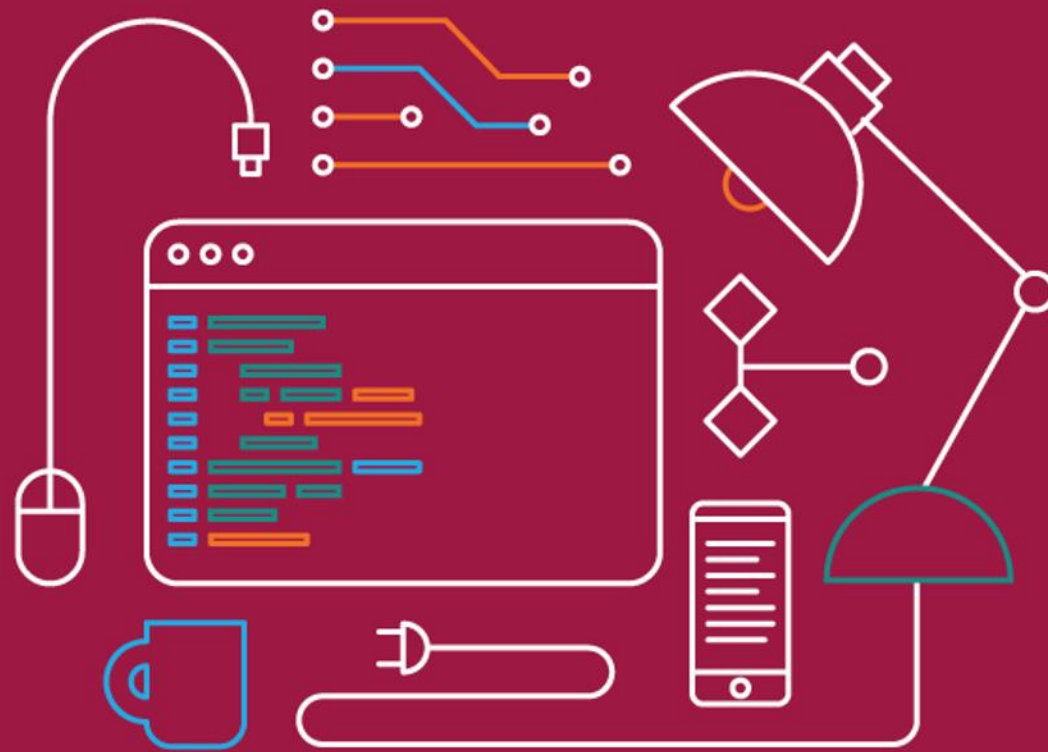
Get the poster in PDF:



The research has been supported by the European Regional Development Fund within the research project "Neural Network Modelling for Inflected Natural Languages" No. 1.1.1-1/16/A/215.

◦ (Pinnis et al., 2017)





THANK YOU!
QUESTIONS, DISCUSSIONS

References



- Deksne, D., & Skadiņš, R. (2011). CFG Based Grammar Checker for Latvian. In *Proceedings of the 18th Nordic Conference of Computational Linguistics NODALIDA 2011* (p. 275 278). Riga.
- Deksne, D., Skadiņa, I., & Skadiņš, R. (2014). Extended CFG Formalism for Grammar Checker and Parser Development. In A. Gelbukh (Ed.), *Computational Linguistics and Intelligent Text Processing, 15th International Conference, CICLing 2014, Proceedings, Part I* (pp. 237–249). Kathmandu, Nepal: Springer. <http://doi.org/10.1007/978-3-642-54906-9>
- Pinnis, M., Krišlauks, R., Miks, T., Deksne, D., Šics, V. (2017). Tilde's Machine Translation Systems for WMT 2017.

