



#### Multilingual and morphology-aware NMT

Dr. Raivis SKADIŅŠ

Tilde, Director of Research and Development

Second Finnish Workshop on Machine Translation, Helsinki, November 1, 2017



# In my talk

- How the WMT 2017 Shared Task was Won (EN-LV-EN)
- What we do with Estonian (Multilingual NMT, CNN etc.)





# How the WMT 2017 Shared Task was Won



-0



- WMT official training data for constrained systems
  - Parallel 4.51M
  - Mono: 27.75M (LV), 330.23M (EN)
- Data was noisy
  - we filtered out a lot
  - only 1.61M left
- For unconstrained systems
  - Parallel: 12.69M
  - Mono: 81.68M (LV), 351.99M (EN).
- For back-translation
  - random 1.61M sentences from the mono corpus



#### Preprocessing



- Normalization (-, ", «, ', etc.)
- Non-translatable tokens
  - emails → \$EMAIL\$
  - urls, filenames  $\rightarrow$  \$URL\$
  - XML tags → \$XML\$
- Tokenization
- Truecasing



0

#### 1000 1200 1400



# **Data filtering**

Original 4.5M

Baseline: 1.74M

• Filtered: 1.61M



# **Baseline NMT architecture**





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



#### **Recurrent Units**





• LSTM

Multiplicative
 LSTM

• MLSTM won





## **Multiplicative LSTM**



- Slightly modified LSTM
  - intermediate result  $m_t$  is used to calculate parameters for  $f_t, i_t, o_t$ LSTM uses the previous result  $h_{t-1}$



 $m_t = W_m \cdot h_{t-1} \cdot x$ 





#### **Vocabulary size**

We tried 25K,
 50K and 100K

but we should not





#### **Context length (maxlen)**

- We tried50 and 80
- Not that important



#### **Sub-word units**



- We used morphological analyzer to split off
  - prefixes
  - endings
  - compounds

Word	BPE	MWS
English		
legalization	leg@@ alization	legal@@ ization
legalize	leg@@ alize	legal@@ ize
legalized	leg@@ alized	legal@@ iz@@ ed
legalizes	legaliz@@ es	legal@@ izes
legalizing	legaliz@@ ing	legal@@ iz@@ ing
Latvian ("a	tbalss" is translated	d as "echo")
atbalss	at@@ balss	atbals@@ s
atbalsis	atbal@@ sis	atbals@@ is
atbalsi	atbal@@ si	atbals@@ i
atbalsīs	at@@ bals@@ īs	atbals@@ īs
atbalsīm	at@@ balsīm	atbals@@ īm



#### **Sub-word units**



 Morphology helps



0

#### **Factored models**



EN-LVStanford parser



### **Factored models**



Training progress for LV-EN constrained systems • LV-EN 27 Latvian POS tagger 23.97 23.48 23.46 22 BLEU 17 -lv-en-simple-tc-baseline-voc25 ----lv-en-simple-tc-morph-factored -nematus-lv-en-constrained-simple-tc-tagger5 ----nematus-lv-en-constrained-simple-tc-tagger7 12 200 400 1600 0 600 800 1000 1200 1400 1800

# NMT does not like rare words



- Rare named entities cause NMT to output unpredictable output
- We replaced all rare words with \$ID\$ in training data
- NMT learned to leave them untranslated
- We translated \$ID\$ with SMT or left untranslated



### **NMT-SMT Hybrid**



Impact on
 BLEU was
 minimal

Translation step	Example sentence				
Source text	Šodien skatieties <b>Ikaunieces-Admidiņas</b> startu Rio spēlēs.				
Pre-processed text	šodien skat@@ ieties <b>I@@ kaun@@ iec@@ es - Ad@@ mi@@ di@@ ņas</b> start@@ u Rio spēlē@@ s .				
Text with identified	šodien skat@@ ieties $eta  extsf{ID}eta$ - $eta  extsf{ID}eta$ start@@ u Rio spēlē@@ s .				
rare words					
NMT translation	watch the $eta {f ID}eta$ - $eta {f ID}eta$ start at the Rio Games today .				
Moses XML with untranslated rare words	<pre><nmt translation="watch the">šodien skatieties </nmt>Ikaunieces <nmt translation="-">-</nmt>Admidiñas <nmt translation="start at the Rio Games today">šodien startu Rio spēlēs</nmt><nmt translation=".">.</nmt></pre>				
Moses XML with identified untranslated person names	<pre><nmt translation="watch the">šodien skatieties </nmt><ne prob="1.0" translation="lkauniece">lkaunieces</ne> <nmt translation="-">-</nmt><ne prob="0.95  0.05" translation="Admidina  Admidins">Admidinas</ne> <nmt translation="start at the Rio Games today">šodien startu Rio spēlēs</nmt><nmt translation=".">.</nmt></pre>				
SMT translation	watch the ${f Ikauniece}$ - ${f Admidina}$ start at the Rio Games today .				
Post-processed	Watch the Ikauniece-Admidina start at the Rio Games today.				
translation					
NMT only transl. (for comparison)	Today, look at the start of the <b>Isolence-Admidias</b> in the Rio Games.				

# The final system



 We merged everything together



Training progress for EN-LV constrained systems



#### **Back-translation**

• It works

 Probably for everybody <sup>(C)</sup>





0



#### WMT 2017 EN-LV



#### **WMT 2017 LV-EN**







- ~55 NMT systems trained
- ~25 SMT systems trained



• (Pinnis et al., 2017)





• Some language pairs are really under resourced (ET-RU)

Inspired from Google Zero-Shot NMT





	Development				Test			
	RU-ET	ET-RU	EN-ET	ET-EN	RU-ET	ET-RU	EN-ET	ET-EN
SMT					12.52	14.74	22.53	32.52
MLSTM Shallow	17.51	18.46	23.79	34.45	11.1 <sup>1</sup>	12.3 <mark>2</mark>	26.14	36.78





	Development				Test			
	RU-ET	ET-RU	EN-ET	ET-EN	RU-ET	ET-RU	EN-ET	ET-EN
SMT					12.52	14.74	22.53	32.52
MLSTM Shallow	17.51	18.46	23.79	34.45	11.11	<mark>12.3</mark> 2	26.14	36.78
MLSTM Shallow Multiling.	<10	<10	<10	<10				





		Develo	opment		Test			
	RU-ET	ET-RU	EN-ET	ET-EN	RU-ET	ET-RU	EN-ET	ET-EN
SMT					12.52	14.74	22.53	32.52
MLSTM Shallow	17.51	18.46	23.79	34.45	11.11	<mark>12.3</mark> 2	26.14	36.78
MLSTM Shallow Multiling.	<10	<10	<10	<10				
GRU Shallow	13.7 <mark>0</mark>	<mark>13.7</mark> 1	17.95	27.84	10.6 <mark>6</mark>	<mark>11.1</mark> 7		
GRU Deep	17.03	17.42			<mark>10.3</mark> 3	<mark>12.3</mark> 6		





	Development				Test			
	RU-ET	ET-RU	EN-ET	ET-EN	RU-ET	ET-RU	EN-ET	ET-EN
SMT					12.52	14.74	22.53	32.52
MLSTM Shallow	17.51	18.46	23.79	34.45	11.11	12.32	26.14	36.78
MLSTM Shallow Multiling.	<10	<10	<10	<10				
GRU Shallow	<mark>13.7</mark> 0	<mark>13.7</mark> 1	17.95	27.84	10.6	<mark>11.1</mark> 7		
GRU Deep	17.03	17.42			<b>10.3</b> 3	12.36		
GRU Deep Multiling.	17.07	17.93	23.37	33.52	13.75	14.57	25.76	36.93



# **Fully Convolutional NMT**



• Latest experiments with Estonian-English

Method	BLEU
MLSTM-SU	36.78
GRU-DM	36.93
SMT	32.52
Fully convolutional	40.54



# **Acknowledgements & related papers**



- The research has been supported by the European Regional Development Fund within the research project "Neural Network Modelling for Inflected Natural Languages" No. 1.1.1/16/A/215.
- The research has been supported by The Ministry of Education and Research (<u>https://www.hm.ee/en</u>) and by The National Programme for Estonian Language Technology (<u>https://www.keeletehnoloogia.ee/en</u>).
- Pinnis, M., Krišlauks, R., Deksne, D., & Miks, T. (2017). Neural Machine Translation for Morphologically Rich Languages with Improved Sub-word Units and Synthetic Data. In Proceedings of the 20th International Conference of Text, Speech and Dialogue (TSD2017). Prague, Czechia.
- Pinnis, M., Krišlauks, R., Miks, T., Deksne, D., & Šics, V. (2017). Tilde's Machine Translation Systems for WMT 2017.
   In Proceedings of the Second Conference on Machine Translation (WMT 2017), Volume 2: Shared Task Papers (pp. 374–381). Copenhagen, Denmark: Association for Computational Linguistics. Retrieved from http://www.aclweb.org/anthology/W17-4737
- Pinnis, M., Krišlauks, R., Deksne, D., & Miks, T. (2017). Evaluation of Neural Machine Translation for Highly Inflected and Small Languages. In Proceedings of the 18th International Conference on Intelligent Text Processing and Computational Linguistics (CICLING 2017). Budapest, Hungary.