

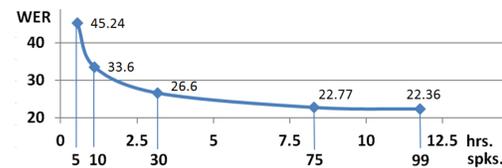
## 1. Overview

### Goal of this work

- Collect speech and text data in the Ukraine for the East Slavic language Ukrainian as a part of the *GlobalPhone* corpus (Schultz et al., 2013) with our Rapid Language Adaptation Toolkit (RLAT) (Vu et al., 2010)
- Develop a Ukrainian LVCSR system rapidly
- Use grapheme-to-phoneme models derived from existing dictionaries of other languages to reduce necessary manual effort for dictionary generation
- Apply state-of-the-art techniques for acoustic modeling and our day-wise text collection and language model interpolation strategy (Vu et al., 2010)

## 3. Baseline ASR System

- 51 phonemes
  - 38 basic phonemes set for Ukrainian (Buk, 2008) (Bilous, 2005)
  - 13 semi-palatalized consonants (Pylypenko, 2008) (Lytvynov, 2009) (Liudovyk, 2011)
- Pronunciation dictionary – manually generated with 882 simple “search and replace rules” from (Buk, 2008)
- Bootstrapping with RLAT using multilingual phone inventory *MM7* (Schultz and Waibel, 2001)
- Preprocessing: 143 MFCC (adjacent frames) → 42 (LDA)
- 11.75 hours from 99 speakers to train acoustic models
- Fully-continuous 3-state left-to-right HMM, Context-dependent AM: decision tree splitting stopped at 500 quinphones
- 3-gram language model: Vocabulary size=7.4k, PPL=594, OOV rate=3.6%
- Word error rate on development set: **22.36%**, on evaluation set: **18.64%**



## 4. Cross-lingual Dictionary Production

- Grapheme Mapping*: Mapping Ukrainian graphemes to the graphemes of the related language (*Rules before g2p*)
- Applying *g2p* model of the related language to the grapheme-mapped Ukrainian words
- Phoneme Mapping*: Mapping resulting phonemes of the related language to the Ukrainian phonemes (*Rules after g2p*)
- Optional: Post-processing rules to revise shortcomings (*Post-rules*)

Step	ru	bg	de	en	# Rules before g2p	# Rules after g2p	PER (%)	WER (%)	# Post-rules	PER (%)	WER (%)
1	биг	биг	bih	bih	43	56	12.4	22.80	57	1.7	<b>21.63</b>
2	ru_b ru_i ru_g	bg_b bg_i bg_g	de_b de_i	en_b en_ih	40	79	10.3	23.70	65	2.8	<b>22.09</b>
3	ua_b ua_i ua_h	ua_b ua_i ua_h	ua_b ua_i	ua_b ua_y	(68)*	66	32.7	27.10	39	28.6	26.36
4	ua_bj ua_i ua_h	ua_bj ua_i ua_h	ua_bj ua_i	ua_b ua_y	(68)*	63	46.8	34.86	21	36.6	34.02

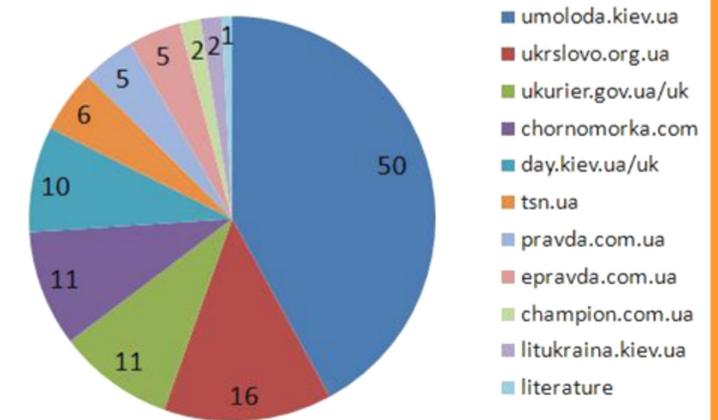
Cross-lingual pronunciation production for биг

Effort (# rules) and quality using cross-lingual rules

## 2. Ukrainian Resources

### Text Corpus

- Text crawled from online newspapers using RLAT
- Complemented fragments from Ukrainian literature and lyrics
- Applied language-dependent text normalization
- Selected prompts to record speech data for the training, development, and evaluation set
- Training set text used for the language model

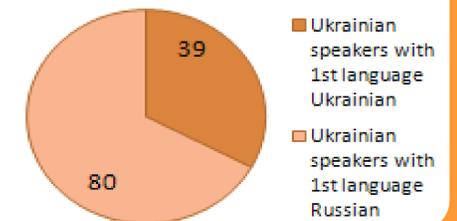


Text sources and the numbers of speakers reading prompts from them

### Speech Corpus

- Speech data collection in *GlobalPhone* style (Schultz, 2013), i.e. we asked Ukrainian speakers to read prompted sentences of newspaper articles

Set	Male	Female	#utterances	#tokens	Duration
train	38	61	11k	69k	11 h 45 mins
dev	4	6	1k	7k	1 h 14 mins
test	4	6	1k	7k	1 h 08 mins
Total	46	73	13k	83k	14 h 07 mins



## 5. System Optimization

### Acoustic Modeling of Semi-Palatalized Phonemes

- Information about semi-palatalized articulation is added as tag to pronunciation dictionary
- Tag added as question in clustering procedure
- Data decide during model clustering on the impact on the basic phoneme
- If no impact → share 1 common model; Otherwise → split models

### Language Model Improvement

- 5 day long “snapshot” crawls of 3 further Ukrainian online newspapers (texts with 94M running words)
- Interpolation of the individual LMs based on minimizing the PPL of the model on the dev set transcriptions
- 3-gram LM with a total of 40k words (PPL=373, OOV Rate=0.53%) performed best
- Word error rate on development set: **13.03%**, on evaluation set: **11.21%**

	WER (%) on dev	WER (%) on test
<b>Acoustic Modeling of Semi-Palatalized Phonemes</b>		
With semi-palatalized (baseline)	22.36	18.64
Without semi-palatalized	21.73	
Data-driven Semi-Palatalized Phone Modeling	<b>21.65</b>	
Grapheme-based	23.82	
<b>Language Model Improvement</b>		
5 day “snapshot” and vocabulary increase	<b>13.03</b>	<b>11.21</b>