

English and Czech of non-native speakers within NLP

Barbora Hladká

`hladka@ufal.mff.cuni.cz`

Charles University,
Faculty of Mathematics and Physics,
Institute of Formal and Applied Linguistics

**Think of learning a second language L2
by people with their native language L1.**

Many thanks to my dear collaborators

L2 = English

Martin Holub, Vincent Kríž, Tom Kocmi, Jindřich Libovický, Magda Ševčíková (CUNI), Pavel Ircing, Marie Kunešová, Jan Vaněk, Zbyněk Zajíc (ZCU)

- Center for large-scale multi-modal data interpretation (GAČR GAP103/12/G084, 2012-2019)

L2 = Czech

Jirka Hana, Hana Hanová, Martina Lupínková, Zdeňka Urešová (CUNI)

- Non-native Czech from the theoretical and computational perspective (GAČR 18-101857, 2016-2018)

The **Native Language Identification** task (NLI)
is to recognize the L1 of an L2 writer/speaker.

Topic

Most advertisements make products seem much better than they really are

Sample text

now a days the publicity is the best way to promoted a product and if you want to sale a product you should bring some information that makes , that the people who is seeing the advertisements make sure that the product very good and in the future this person could buy it .

Topic

Most advertisements make products seem much better than they really are

Sample text

now a days the publicity is the best way to promoted a product and if you want to sale a product you should bring some information that makes , that the people who is seeing the advertisements make sure that the product very good and in the future this person could buy it .

L1 = Spanish

In a nutshell

- 2005 – Very first experiments [Koppel et al., 2005]
- 2013 – **Shared task with Text**
CUNI [Hladká et al., 2013]
- 2015 – CUNI best model [Kříž et al., 2015]
- 2016 – Feature interpretation by CUNI
 - Initial experiments with Deep Learning by CUNI
 - **Shared task with Speech**
- 2017 – **Shared task with Text and Speech**
CUNI+ZCU

Settings

corpus	International Corpus of Learner English
L1s	Bulgarian, Czech, French, Russian, Spanish
data size	285 instances per L1
features	function words, character n -grams, error types, POS bigrams
# of feat.	1,035
approach	SVM
Acc*	80%

*Acc. – 10 cross-validation results

Shared Task 2013 – NLI with Text Data

TOEFL11 – a corpus of non-native English writing

- consists of essays on eight different topics
- written by non-native speakers of three proficiency levels (low/medium/high)
- 11 L1s: Arabic, Chinese, French, German, Hindi, Italian, Japanese, Korean, Spanish, Telugu, Turkish
- contains 1,100 essays per language with an average of 348 word tokens per essay

For more info see [Blanchard et al., 2013].

Shared Task 2013

Results

Rank.	System	# of feat.	Acc.	Approach
1	[Gebre et al., 2013]	73,626	84.6	SVM; tf-idf of uni-grams and bigrams of words
2	[Jarvis et al., 2013]	400K	84.5	SVM; {1,2,3}-grams of words, lemmas, POS tags, $df \geq 2$
22	[Hladká et al., 2013]	6,338	74.2	SVM; classifier combination; {1,2,3}-grams of lemmas, function words, POS tags, $df \geq 4$
29

CUNI best model

[Kríž et al., 2015]

Language modeling approach to feature extraction

- We built 11 special language models of English (M_i), each based on the texts with the same L1 language available in the training data.
- Then we compare M_i to a general language model of English (M_G).
- The cross-entropy of text t with empirical n-gram distribution p given a language model M with distribution q is

$$H(t, M) = - \sum_x p(x) \log q(x).$$

- **Normalized cross-entropy scores – used as features**

$$D_G(t, M_i) = H(t, M_i) - H(t, M_G) = - \sum_x p(x) \log \frac{q_i(x)}{q_G(x)},$$

where M_i are the special language models with distributions q_i , and M_G is the general language model with the distribution q_G .

CUNI best model

Error analysis

Aggregated confusion matrix

Sum of 10 confusion matrices obtained in 10-fold cross validation process

	ARA	DEU	FRA	HIN	ITA	JPN	KOR	SPA	TEL	TUR	ZHO
ARA	738	8	30	24	5	16	9	30	10	32	15
DEU	7	836	13	4	20	4	6	20	1	16	1
FRA	31	6	767	0	35	6	3	39	1	15	1
HIN	19	1	1	684	0	2	3	3	194	9	4
ITA	6	8	20	1	761	0	1	58	0	3	2
JPN	7	2	3	3	2	709	109	6	1	8	22
KOR	14	1	3	3	1	120	676	6	2	12	50
SPA	30	21	45	6	69	3	4	720	1	17	7
TEL	4	2	0	157	0	1	0	2	685	2	1
TUR	32	15	11	16	6	11	23	11	5	778	16
ZHO	12	0	7	2	1	28	66	5	0	8	781

State of the art as of Dec 2016

System	# of feat.	Acc.	Approach
[Bykh and Meurers, 2016]	16,841	– (85.4 on etest)	ensemble classifier; {1,2}-grams of lem- mas, words, POS tags, dependencies, suffixes, verb subcategorization patterns
[Gebre et al., 2013]	73,626	84.6	SVM; tf-idf of uni- grams and bigrams of words
[Jarvis et al., 2013]	400K	84.5	SVM; {1,2,3}-grams of words, lemmas, POS tags, $df \geq 2$
[Ionescu et al., 2014]	NA	84.1	Kernel-based learn- ing; {5-8}-grams of characters
[Križ et al., 2015]	55	82.4	SVM; language models using tokens, charac- ters, POS, suffixes

Traditional approaches to **vector representation of documents** – bag of words, bag of n-grams

- \oplus simple, robust
- \ominus word order is (almost) lost, short context
- \ominus ignoring semantics

Feature interpretation

Extracting n-grams

Extracting a compact and interpretable feature set

- intragroup homogeneity **and** intergroup heterogeneity
- L1 vs. other L1s **and** L1 vs. a different L1

Extracting n-grams

Inspiration by [Kyle et al., 2016]

- 1 **Get an initial feature set of 1-5-wordgrams**
- 2 **Get reduced list of n-wordgrams for each L1/L1-pair**
 - 1 Filter the initial feature set by document frequency in L1 set
 - 2 Compute *keyness value* for each n-gram using G-test
 - 3 Get positive (*overused*) and negative (*underused*) lists of n-grams sorted by keyness value
 - 4 Remove superstrings with lower keyness value
 - 5 Get reduced lists of n-grams
- 3 **Merge 11/55 reduced feature sets and get Reduced n-Grams**

Extracting n-grams

G-test

- log-likelihood statistical significance test

$$G = 2 \sum_i O_i * \ln\left(\frac{O_i}{E_i}\right)$$

- O_i observed frequencies, E_i estimated frequencies, $E_i = \frac{N_i \sum_i O_i}{\sum_i N_i}$
- the distribution of G is approximately a χ^2 distribution

- in our case $i = 1, 2$
 - positive (overused) n-grams $\frac{O_1}{N_1} > \frac{O_2}{N_2}$
 - negative (underused) n-grams $\frac{O_1}{N_1} < \frac{O_2}{N_2}$

Extracting n-grams

Example

ARA vs. other L1s – positive list of 104 n-grams

	ngram	G2	O1	E1	O2	E2	fr	df
1:	alot	764	242	30	127	339	369	234
2:	i	736	1149	485	4862	5526	6011	1939
3:	alot of	648	185	20	68	233	253	167
4:	and	389	7004	5544	61633	63093	68637	8803
5:	thier	330	216	53	439	602	655	382
...	
14:	dont	130	156	58	557	655	713	469
15:	guide	119	513	313	3361	3561	3874	718
17:	tour guide	103	393	232	2477	2638	2870	639
...	

Extracting n-grams

Example

ARA vs. other L1s – reduced positive list of 74 n-grams
(superstring threshold = 100%)

	ngram	G2	O1	E1	O2	E2	fr	df
1:	alot	764	242	30	127	339	369	234
2:	i	736	1149	485	4862	5526	6011	1939
3:	alot of	648	185	20	68	233	253	167
4:	and	389	7004	5544	61633	63093	68637	8803
5:	thier	330	216	53	439	602	655	382
...	
14:	dont	130	156	58	557	655	713	469
15:	guide	119	513	313	3361	3561	3874	718
17:	tour guide	103	393	232	2477	2638	2870	639
...	

Extracting n-grams

RnG interpretation

Manual analysis of (reduced) positive and negative lists of n-grams

A	B	C	D	E	F	G	H
		ALL F0 20/5	ARA	DEU	FRA	HIN	ITA
ARA	plus	spelling errors -- i; alot; thier;		spelling errors -- i; alot; thier;	spelling errors -- i; alot; thier	discourse markers -- because;	spelling errors -- alot; thier; discourse markers -- Also; function words -- and;
	minus	discourse markers -- However; function words -- But; punctuation -- period; comma; modal verbs -- can be;		discourse markers -- So; function words -- But; a lot of;	function words -- is; able; I think that; punctuation;	punctuation -- comma; function words -- of; has; But; the; they;their;	spelling errors -- alot of; function words -- I think that; punctuation -- semicolon;apostrophe;

Study by [Malmasi et al., 2015]

- 10 professors and researchers
- 5 L1s: Arabic, Chinese, German, Hindi, Spanish
- 30 essays from TOEFFL11 eval test set
 - 3 *easy* and 3 *hard* essays per L1
 - easy/hard wrt the Shared Task'13 predictions

Human NLI performance

Study by [Malmasi et al., 2015]

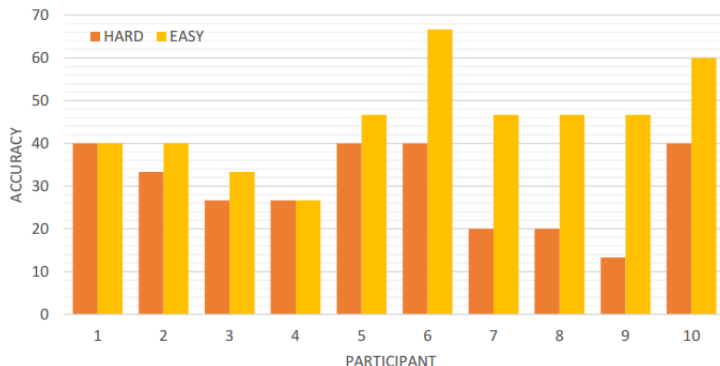


Figure 2: Prediction accuracy for each of our 10 participants under both easy and hard conditions.

Human NLI performance

Study by [Malmasi et al., 2015]

	Accuracy (%)		
	Easy	Hard	All
Random Baseline	20.0	20.0	20.0
NLI Plurality Vote	100.0	33.3	66.7
NLI Mean Probability	100.0	46.7	73.3
Top Human	66.7	40.0	53.3
Human Plurality Vote	73.3	40.0	56.7

Table 5: Comparing human participant performance against an NLI system on 30 selected texts.

Deep learning waves have lapped at the shores of computational linguistics for several years now, but 2015 seems like the year when the full force of the tsunami hit the major Natural Language Processing conferences. [Manning, 2015]

NLI experiments with DL by CUNI

Credits: Jindřich Libovický (2016)

Model	Acc
Paragraph vector model	
Distributed Bag-Of-Words [Le and Mikolov, 2014]	71.5%

Credits: Tom Kocmi (2017)

Model	Acc on dev test
Convolutional Neural Network (segment = sentence)	76%

ETS Corpus of Non-native Spoken English

- 11 L1 languages
 - ARA, CHI, FRE, GER, HIN, ITA, JAP, KOR, SPA, TEL, TUR
 - these L1s are identical to the TOEFL11 L1s
- 45s long recordings
- 5,132 instances in total
 - 3, 300 training, 965 dev test, 867 eval test

Shared Task 2016

Results

8 teams

- their best results on the development test set
- Unweighted Average Recall

team	UAR(%)
[Abad et al., 2016]	84.6
[Shivakumar et al., 2016]	78.6
[Gosztolya et al., 2016]	70.7
[Huckvale, 2016]	69.8
[Senoussaoui et al., 2016]	68.4
[Keren et al., 2016]	61.5
[Jiao et al., 2016]	52.2
[Rajpal et al., 2016]	39.8
baseline	45.1

Three subtasks

- 1 Essay Task: written responses only
- 2 Speech Task: spoken responses only
 - transcriptions and i-vectors provided for 45s audio files
 - the audio files are not available to the participants
- 3 Fusion Task: written and spoken resp. – **CUNI + ZCU**
 - Evaluation test set release on June 19, 2017
 - Result notification on June 26, 2017

Shared Task 2017

Data

Written responses

# of instances	ST 2013	ST 2017
9,900	training data	training data
1,100	dev test	training data
1,100	eval test	dev test
1,100	–	eval test

Spoken responses

ST 2016	data	ST 2017
# of instances	set	# of instances
3,300	training data	11,000
965	dev test	1,100
867	eval test	1,100

Shared Task 2017

Baselines on the dev test

	precision	recall	f1-score
essay	0.72	0.72	0.72
speech_transcriptions	0.52	0.52	0.52
speech_transcriptions+ivectors	0.76	0.76	0.76
fusion	0.75	0.75	0.75
fusion+ivectors	0.78	0.78	0.78

- SVM, word unigrams
- confusion matrices provided

Motivation

- Extracting features using NLP tools trained over grammatical text
- What impact do grammatical errors have on their performance?

Related work

- [Berzak et al., 2016b]: Treebank of Learner English, manual error annotation, manually annotated POS tags and UDs, POS tagging and parsing
- [Napoles et al., 2016]: NUCLE corpus, manual error annotation, dependency parsing
- [Dickinson and Ragheb, 2015]: Syntactically Annotating Learner Language of English (SALLE)
- [Cahill et al., 2014]: Self-training for parsing learner text

CzeSL – a corpus of non-native Czech writing

CzeSL-MAN \subseteq **CzeSL**

- 32 L1s
- contains 645 texts
- manual error annotation
- layer **T0** – anonymised transcripts of the original
- layer **T1** – errors in individual word forms fixed
- layer **T2** – all other types of errors fixed

For more info see [Rosen, 2016].

Topic My family

Student male, 16+, A2, staying in CR 0-1 year

Jmenujese Adam . Ja jsem Mongolska . Mongolska ma 21 kraji .
Moje rodina je hezka ještě velka . Mongolska je 3000 million
lidi . Ma tradiční píseňka , taneční . Mongolska tradiční písenka
je hezka . Ješte ma „ Morin khuur “ . Morin Khuur to je muzika
. Ten hezka tradiční pohádka , píseň . Mongolska má mnoho
tradiční svátík . Třeba Naadam , Tsagaarsur . Ješte mnoho
Velbloud , Kůn , Kravá , Koza , Ovce . Mongolsky lidi dobrý
. Mongolsko ma mnoho horý a nemam ocean . Mongolska
hlavní naměsto . Ulaanbaatar . ADAM , 18 Let
Bydlim v Čechagh už 6 měsíc .
1 . AHOJ

Sb-Pred-Obj annotation of CzeSL T0

Use the PDT guidelines to mark Sb, Pred, Pnom, Obj

← → /collections/collection_01/all.w.050 brat

1 Jmenujese Adam . Ja jsem Mongolska . Mongolska ma 21 kraji . Moje rodina je hezka jeste
velka . Mongolska je 3000 million lidi . Ma tradični píseňka , taneční . Mongolska tradični píseňka je
hezka . Jeste ma „Morin khuur“ . Morin Khuur to je muzika . Ten hezka tradični pohádka , píseň .
Mongolska má mnoho tradični svátík . Triba Naadam , Tsagaarsur . Jeste mnoho Velbloud , Kún , Kravá ,
Koza , Ovce . Mongolsky lidi dobrý . Mongolsko ma mnoho horý a nemam ocean . Mongolska hlavní
naměsto . Ulaanbaatar .

2 ADAM , 18 Let

3 Bydlim v Čechagh už 6 měsíc .

4 1 . AHOJ

Manual and automatic annotation of CzeSL T0

A	B	C	D	E	F	G	H
	T0 form	Lemma		Tag		Syntax	
	T0 form	Manual	Morphodita	Manual	Morphodita	Manual	UDPipe
1	Jmenujese	jmenovat_:T_:W	jmenujese	VB-S---1P-AA---	NNMS1----A----	Predicate	root
2							
3	Adam	Adam_:Y	Adam_:Y	NNMS1----A----	NNMS1----A----	Nominal	nsubj
4	.	.	.	Z:-----	Z:-----		
5							
6	Ja	já	já_x	PP-S1-1-----	PP-S1-1-----6	Subject	nmod
7	jsem	být	být	VB-S---1P-AA---	VB-S---1P-AA---	Predicate	cop
8							
9	Mongolska	Mongolsko_:G	Mongolsko_:G	NNFS1----A----	NNNS2----A----		root
10	.	.	.	Z:-----	Z:-----		
11							
12	Mongolska	Mongolsko_:G	Mongolsko_:G	NNFS1----A----	NNNS2----A----	Subject	root
13	ma	mit	ma-99_:B_:S	VB-S---3P-AA---	NNXXX----A---8	Predicate	cc
14	21	21	21	C=-----	C=-----	Object	nummod
15	kraji	kraj	kraj	NNIP2----A----	NNIP7----A----		conj
16	.	.	.	Z:-----	Z:-----		

Manual and automatic annotation of CzeSL T2

T2		
form	Morphodita	Morphodita
Jmenuji	jmenovat_:T_:W	VB-S---1P-AA--1
se	se_^(zvr._zájmeno/částice)	P7-X4-----
Adam	Adam_:Y	NNMS1----A---
.	.	Z:-----
Já	já	PP-S1--1-----
jsem	být	VB-S---1P-AA---
z	z-1	RR--2-----
Mongolska	Mongolsko_:G	NNNS2----A---
.	.	Z:-----
Mongolsko	Mongolsko_:G	NNNS1----A---
má	mít	VB-S---3P-AA---
	21	21 C=-----
krajů	kraj	NNIP2----A---
.	.	Z:-----

Manual and automatic annotation of CzeSL Problems

- Should lemmas be base forms from Interlanguage or the intended lemmas in Czech? **Examples:**
 - *Mongolska ma 21 kraji* - Lemma *Mongolska*? Or *Mongolsko*? What if they sometimes use *Mongolska* and sometimes *Mongolsko*.
 - *ja* for *já*
 - *kůn* for *kůň*
 - *naměsto* for *město*
- But there is little chance of reconstructing the IL lemma when other forms are used (*horý* - lemma *hora*? *horá*?) So probably std Czech lemma is safer.
- Sometimes it is hard to say, whether the student uses declension at all - i.e. whether say an adjective should be AAXXX or AAFS3

Manual and automatic annotation of CzeSL

Approaches?

- Different annotation instructions for different proficiency levels
- Use the STYX sentences [Hladká and Kučera, 2008] to re-train both the tagger and the parser
 - universal tagset
 - simpler sentence structures

Annotate according to the literal reading of the sentence rather than the corrections or other interpretations of potential intended meanings ... cases in which a literal reading is either impossible or contextually implausible are annotated according to the correct version ... [Berzak et al., 2016a]

Try to assume as little as possible about the intended meaning of the learner. You will have to use context to disambiguate at times, but in general, we are looking to annotate (morpho-) syntactic information and not semantic information. This is one of the most challenging principles to follow ... As much as possible, annotate the language “as is”. In other words, try not to think in terms of errors, but in terms of linguistic evidence ... [Dickinson and Ragheb, 2013]

References I



Abad, A., Ribeiro, E., Kepler, F., Astudillo, R., and Trancoso, I. (2016).

Exploiting phone log-likelihood ratio features for the detection of the native language of non-native english speakers.

In Interspeech 2016, pages 2413–2417.



Berzak, Y., Kenney, J., Spadine, C., Wang, J. X., Lam, L., Mori, K. S., Garza, S., and Katz, B. (2016a).

Trebank of learner english annotation manual - draft version 1.0, ud release 1.3 (may 2016).



Berzak, Y., Kenney, J., Spadine, C., Wang, J. X., Lam, L., Mori, K. S., Garza, S., and Katz, B. (2016b).

Universal dependencies for learner english.

CoRR, abs/1605.04278.

References II



Blanchard, D., Tetreault, J., Higgins, D., Cahill, A., and Chodorow, M. (2013).

TOEFL11: A Corpus of Non-Native English.

ETS Research Report Series, 2013(2):i–15.



Bykh, S. and Meurers, D. (2016).

Advancing linguistic features and insights by label-informed feature grouping: An exploration in the context of native language identification.

In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 739–749, Osaka, Japan. The COLING 2016 Organizing Committee.

References III

 Cahill, A., Gyawali, B., and Bruno, J. (2014).

Self-training for parsing learner text.

In Proceedings of the First Joint Workshop on Statistical Parsing of Morphologically Rich Languages and Syntactic Analysis of Non-Canonical Languages, pages 66–73, Dublin, Ireland. Dublin City University.

 Dickinson, M. and Ragheb, M. (2013).




Annotation for learner english guidelines, v. 0.1 (june 2013).

 Dickinson, M. and Ragheb, M. (2015).

On grammaticality in the syntactic annotation of learner language.

In Proceedings of The 9th Linguistic Annotation Workshop, pages 158–167, Denver, CO.

References IV

-  Gebre, B. G., Zampieri, M., Wittenburg, P., and Heskes, T. (2013). Improving Native Language Identification with TF-IDF Weighting. In *Proceedings of the Eighth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 216–223, Atlanta, Georgia. ACL.
-  Gosztolya, G., Grósz, T., Busa-Fekete, R., and Tóth, L. (2016). Determining native language and deception using phonetic features and classifier combination. In *Interspeech 2016*, pages 2418–2422.
-  Hladká, B., Holub, M., and Kríž, V. (2013). Feature Engineering in the NLI Shared Task 2013: Charles University Submission Report. In *Proceedings of the Eighth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 232–241, Atlanta, Georgia, USA. Microsoft Research, ACL.

References V



Hladká, B. and Kučera, O. (2008).

An annotated corpus outside its original context: A corpus-based exercise book.

In *ACL 2008: Proceedings of the Third Workshop on Innovative Use of NLP for Building Educational Applications*, pages 36–43, Columbus, OH, USA. Association for Computational Linguistics (ACL).



Huckvale, M. (2016).

Within-speaker features for native language recognition in the interspeech 2016 computational paralinguistics challenge.

In *Interspeech 2016*, pages 2403–2407.






Ionescu, R. T., Popescu, M., and Cahill, A. (2014).




Can characters reveal your native language? A language-independent approach to native language identification.

In *Proceedings of the 2014 Conference on EMNLP*, pages 1363–1373, Doha, Qatar. ACL.

References VI

-  Jarvis, S., Bestgen, Y., and Pepper, S. (2013). Maximizing Classification Accuracy in Native Language Identification. In *Proceedings of the Eighth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 111–118, Atlanta, Georgia. ACL.
-  Jiao, Y., Tu, M., Berisha, V., and Liss, J. (2016). Accent identification by combining deep neural networks and recurrent neural networks trained on long and short term features. In *Interspeech 2016*, pages 2388–2392.
-  Keren, G., Deng, J., Pohjalainen, J., and Schuller, B. (2016). Convolutional neural networks with data augmentation for classifying speakers' native language. In *Interspeech 2016*, pages 2393–2397.




References VII

-  Koppel, M., Schler, J., and Zigdon, K. (2005).
Determining an author's native language by mining a text for errors.
In *Proceedings of the Eleventh ACM SIGKDD International Conference on Knowledge Discovery in Data Mining*, KDD '05, pages 624–628, New York, NY, USA. ACM.
-  Kríž, V., Holub, M., and Pecina, P. (2015).
Feature Extraction for Native Language Identification Using Language Modeling.
In Angelova, G., Boncheva, K., and Mitkov, R., editors, *Proceedings of Recent Advances in Natural Language Processing*, pages 298–306, Hisarja, Bulgaria.
-  Kyle, K., Crossley, S. A., and Kim, Y. (2016).
Native language identification and writing proficiency.
International Journal of Learner Corpus Research, 1(2):187–209.

References VIII

-  Malmasi, S., Tetreault, J., and Dras, M. (2015).
Oracle and Human Baselines for Native Language Identification.
In Proceedings of the Tenth Workshop on Innovative Use of NLP for Building Educational Applications, pages 172–178, Denver, Colorado. Association for Computational Linguistics.
-  Manning, C. D. (2015).
Computational Linguistics and Deep Learning.
Computational Linguistics, 41(4):701–707.
-  Napoles, C., Cahill, A., and Madnani, N. (2016).
The effect of multiple grammatical errors on processing non-native writing.
In Proceedings of the 11th Workshop on Innovative Use of NLP for Building Educational Applications, pages 1–11, San Diego, CA. Association for Computational Linguistics.

References IX

-  Rajpal, A., Patel, T. B., Sailor, H. B., Madhavi, M. C., Patil, H. A., and Fujisaki, H. (2016).
Native language identification using spectral and source-based features.
In Interspeech 2016, pages 2383–2387.
-  Rosen, A. (2016).
Building and using corpora of non-native Czech.
In Brejová, B., editor, Proceedings of the 16th ITAT: Slovenskočeský NLP workshop (SloNLP 2016), volume 1649 of *CEUR Workshop Proceedings*, pages 80–87, Bratislava, Slovakia. Comenius University in Bratislava, Faculty of Mathematics, Physics and Informatics, CreateSpace Independent Publishing Platform.
-  Senoussaoui, M., Cardinal, P., Dehak, N., and Koerich, A. L. (2016).
Native language detection using the i-vector framework.
In Interspeech 2016, pages 2398–2402.

References X



Shivakumar, P. G., Chakravarthula, S. N., and Georgiou, P. (2016). Multimodal fusion of multirate acoustic, prosodic, and lexical speaker characteristics for native language identification. In *Interspeech 2016*, pages 2408–2412.