







## **A friend in need?** Language Technology for Second Language Learning

Elena Volodina

University of Gothenburg, Department of Swedish, Språkbanken

elena.volodina@svenska.gu.se









## Language Technology is ...

- ...a cross-disciplinary research area that covers development of computer programs for analysis, interpretation and generation of natural languages, etc... [*Wikipedia*]
- Other names: Computational Linguistics, Natural Language Processing (NLP), etc...









## Origins

1950's: Machine translation (MT), Russian-English, first attempts in USA











## Origins



1960's: MT failure => Computational linguistics as a new research field

\* grammar in both lan-ges

\* morphology (inflections)

\* syntax

\* semantics\* lexicon\* pragmatics



### **Current context**

Later: Language Technology/NLP is placed under Artificial Intelligence





### NLP + CALL = ICALL

ICALL

Natural Language Processing + technical competence (Computer Assisted) Language Learning + pedagogical competence

## **Focus on literacy**

- Dutch study:
  - → Average reading comprehension ~B1 level

Velleman, E., van der Geest, T.: Online test tool to determine the CEFR reading comprehension level of text. Procedia Computer Science 27 (2014)





Out ot 1.3 mln citizens 0.4 mln (potentially) need training in Estonian

Source: Statistikaamet: http://www.stat.ee/main-indicators.

## Swedish societal need

**Citizens with foreign background, 2002-2015** 



**2015**: out of **9,9 mln** citizens, **2,2 mln** have foreign backgrund, dvs **22,2 %** (Statistiska centralbyrån)

# What can we do?

#### cause versus symptoms



# ICALL tools for Second language (L2) learning





#### **Target group**

Grown-ups Children Analphabets Special needs etc





Data & Resources

Tools & Algorithms

Corpora Essays Word lists Grammar rules etc

POS-taggers Lemmatizers Sentence/text readability Error detectors Speech synthesis etc







#### Applicationdevelopment and maintenance

versus

#### Prototypedevelopment and evaluation (proof-of-concept)



## Lark Trills for Language Drills Text-to-speech technology for language learners

- Dictation and spelling exercise
- Focus on
  - evaluation of the quality of TTS
  - finding ways to give feedback on spelling errors

Elena Volodina, Dijana Pijetlovic. 2015. Lark Trills for Language Drills: Textto-speech technology for language learners. Proceedings of the 10th workshop on Building Educational Applications (BEA10), NAACL 2015, Denver, USA





#### for word & (inflected word) levels



### SPEED SPEIling Error Database

- For each correct item (base form + word class) we store:
  - session ID (no personal data, such as L1)
  - incorrect spelling(s)

## L2 spelling error database, SPEED

- <LexicalEntry uid="LexicalEntry-58d3459f-5acb-43f8-b60e-deb45a986c56">
Correct
<Sense id="speed--kelly-6950" uid="Sense-b1d45016-bdb5-4584-9ec5-11f780ecbf8a"/>
<word lang="swe" pos="AV" uid="word-4d9ed4cb-83b7-4293-a786-ff11f398e2d4</p>
Storresting sessionID="2013-05-13-22-27-28" time="22:58:25" uid="misspelling-3ba31d83-f1ad-4c99-bc33-2c6a3a3c7849
Storrevlig
- <modification uid="modification-4673c2b2-63a1-4c3c-b2a5-c6a80bc5dd20">
<feat att="updatedBy" val="laerka" uid="feat-ec50eb25-d870-4f53-b86c-ed620e3a332c"/>
<feat att="updatedBy" val="laerka" uid="feat-ec50eb25-d870-4f53-b86c-ed620e3a332c"/>
<feat att="modificationDateTime" val="2013-05-13T22:58:26.01+02:00" uid="feat-4032acab-afa3-42c3-b6b0-3f904e345b76"/>
<feat att="modificationAccepted" val="pending" uid="feat-22d76dd2-c26f-4ccd-b7a9-e6997082e0cd"/>

## **Error data**

Error types	Nr,%	Example (correct → *incorrect*)
Competence-based errors	55	
Consonant doubling	28	sto <b>pp</b> a → *sto <b>p</b> a*
Diacritics (å, ä, ö)	23	h <b>ö</b> gre → *h <b>o</b> gre*
Phonetic errors (e.g. voiced vs voiceless)	25	relevan <b>s</b> > *relevan <b>z</b> *
Consonant clusters (phoneme-grapheme mappings, incl. cases of homonyms)	20	<b>sk</b> ön → * <b>sj</b> ön*
Other (unclassifiable)	4	Israel → *visträv*
Performance-based errors	17	
Typos (neighbouring keys, addition, deletion, insertion, replacement)	17	förb <b>ä</b> ttra → *förb'ttra*
Across one word (phrases & sentences)	28	se en bild → *sen bild*



### SPEED SPEIling Error Database

Advantages of collecting a corpus by applying this method: participants are quickly attracted, while cost, time and effort of collecting a corpus are reduced

THIS is RESEARCH DATA!

And we need more of it!

## L2 infrastructure –

# a possible answer?

## What is infrastructure?



# An electronic research infrastructure

- (free accessible) data in electronic format
- technical platform for exploring data, including tools and algorithms for data analysis, and visualization
- a set of tools and technical solutions for new data collection and preparation, including data processing and annotation
- a network of experts in the relevant disciplines, incl. legal and ethical questions



# How can it help?



- Collect and annotate data (L2 essays, error logs, course books ...)
- Develop tools for analyzing L2 data (e.g essays, reading comprehension texts)
- Set up and maintain applications/databases
- Gain expert knowledge
  - to support research on L2 Swedish
  - to support course book writers, L2 teachers, L2 assessors, L2 students
  - to support instruction of future L2 teachers



## **Collecting data**



# **Two types of data**

- Produced **BY** L2 learners
  - → essays
  - $\rightarrow$  exercise logs
  - → errors
  - $\rightarrow$  (interviews)





- Produced by experts **FOR** L2 learners
  - → reading comprehension texts
  - → exercises
  - $\rightarrow$  (recordings of listening excerpts)

# **Challenges:** L2 learner-produced data

- Electronic L2 essays/logs are very difficult to collect
  - → NOT available online
  - → Need learner permits / copyright
  - $\rightarrow$  Need learner variables (gender, age, L1) / personal privacy act
  - → Sensitive in nature / anonymize
  - $\rightarrow$  Those who have it don't want or CAN'T share
- We need an infrastructure/environment for storing and collecting L2 data
  - $\rightarrow$  same variables for comparison
  - $\rightarrow$  same student same ID, etc

# **Challenges:** expert texts for L2 learners

- Electronic coursebooks
  - → NOT available online
  - → Aren't shared by publishing houses
  - → Need to be selected, bought, OCR-ed, proof-read, etc
  - → Can't be shared for copyright reasons
- We need extra information added (MANUALLY)
  - $\rightarrow$  text genres, topics, levels of difficulty
  - → exercise types, formats, target skills and competences

## **Curios "time & effort" fact:**

## **Data vs experiments**



# L2 essay pre-processing



## SweLL corpus

Sub-						Un-	
corpus	A1	A2	<b>B1</b>	B2	C1	known	Total
Tisus	-	-	-	27	78	-	105
Sw1203	-	-	33	45	11	1	90
SpIn	16	83	42	2	-	1	144
Total	16	83	75	74	89	2 (	339

## **SweLL corpus: topics**



Topics

#### Number of essays

## **SweLL corpus: L1s**



## SweLL corpus: age



## SweLL corpus: non-lemmatized items



Distribution of non-lemmatized tokens per level

## **COCTAILL corpus**





### COCTAILL quantitative explorations: target skills across levels



#### COCTAILL quantitative explorations: topics across levels







## COCTAILL

CEFR	Nr. of	Nr. of	Nr. of	Nr. of	Nr. of	Nr. of	Nr. of
level	books	authors	lessons	texts	tasks	sentences	tokens
						(texts)	(texts)
A1	4	10	37	101	160	1581	11132
A2	4	10	105	232	244	4217	37259
B1	4	12	83	345	389	6510	79402
B2	4	8	31	314	368	8527	101583
C1	2	2	22	115	333	5085	71991
Total	18 (12 titles)	42 (26 different names)	278	1106	1494	25920	301367



# **Exploiting data**







## SVALex L2 receptive vocabulary

Level	# items	# new items	# MWE	# doc.hapax	Doc.hapax examples	# EVP
A1	1,157	1,157	92	99	postnummer "zip code"	601
A2	3,327	2,432	300	635	jurist "lawyer"	925
B1	6,554	4,332	617	1,868	öga mot öga "face to face"	1,429
B2	8,728	4,553	880	3,051	snigelfart "snail speed"	1,711
C1	7,564	3,160	783	2,709	inom synhåll "within eyesight"	N/A
Total	15.681	15.681 1	.426	8.362		

http://vocabulary.englishprofile.org/staticfiles/about.html

user: englishprofile password: vocabulary





## SweLLex L2 productive vocabulary

Lev	#items	#new	#MWE	#hapax	Doc.hapax examples	#SVALex	#EVP
A1	398	398	15	0	-	1,157	601
A2	1,327	1,038	82	12	i kväll "tonight"	2,432	925
B1	2,380	1,542	206	36	fylla år "have birthday"	4,332	1,429
B2	2,396	959	264	58	fatta beslut "make a decision"	4,553	1,711
C1	3,566	1,545	430	152	sätta fingret "put a finger on sth"	3,160	N/A
C2	145	7	12	1	i bakhuvudet "in mind"	N/A	N/A

Total **5,475** 

http://vocabulary.englishprofile.org/staticfiles/about.html

user: englishprofile password: vocabulary

## Studera (Eng. "study")



http://cental.uclouvain.be/svalex/





• Aim: reuse corpora to support language learning



How can we automatically assess linguistic *complexity* (readability) on <u>text</u> or <u>sentence</u> level?

# Readability experiments

- Machine learning methods for automatic classification (WEKA, scikit-learn)
- Text- and sentence level
- Based on COCTAILL corpus
- 5 CEFR levels (A1-C1)
- 61 different linguistic features



#### From course books to automatic CEFR level assessment



## Features

Count	Lexical	Syntactic	Мо	rphological
Sentence length	A1 lemma IS	Avg DepArc length	Modal V to V	Verb IS
Avg token length	A2 lemma IS	DepArc Len $> 5$	Particle IS	V variation
Extra-long token	B1 lemma IS	Max length DepArc	3SG pronoun IS	Function W IS
Nr characters	B2 lemma IS	Right DepArc Ratio	Punctuation IS	Lex tkn to non-lex tkn
LIX	C1 lemma IS	Left DepArc Ratio	Subjunction IS	Lex tkn to Nr tkn
Bilog TTR	C2 lemma IS	Modifier variation	PR to N	Neuter N IS
Square root TTR	Difficult W IS	Pre-modifier IS	PR to PP	CJ + SJ IS
Semantic	Difficult N&V IS	Post-modifier IS	S-VB IS	Past PC to V
Avg senses per token	OOV IS	Subordinate IS	S-V to V	Present PC to V
N senses per N	No lemma IS	Relative clause IS	ADJ IS	Past V to V
	Avg. KELLY log freq	PP complement IS	ADJ variation	Present V to V
			ADV IS	Supine V to V
			ADV variation	Relative structure IS
			N IS	Nominal ratio
			N variation	N to V

## Online tool for text evaluation

#### https://spraakbanken.gu.se/larkalabb/texteval

Exercise Generator

Hit-Ex Learner Corpora Editor

Text evaluation

Svenska

\_\_\_\_RK

Language Acquisition Reusing Korp

Filmen hanlar om en pojke . Han heter NN och han gilla dansar ballet så mycket . När
han har idrott leklion , brukor han inte träna boxning så att träna ballet med många tjejer i
nästa klassrummet . Han <mark>tränar</mark> dansa mycket , var och när han kan . <mark>Hans</mark> lärare ger för
<mark>han</mark> ett <b>par skor</b> av ballet <mark>och han gömmer</mark> den mellan för <mark>två</mark> medrasser . Den <mark>mest</mark>
intressanta person i filmen är Billy . Han är en snäll pojke . Hans mamma dog , han bor
med hans pappa , bror och hans mormor . Han älskar hans mormor så mycket . Hon är
gammal och hon brukar göra konstiga saker . Billy sörjer när han saknar hans mamma .
Hans pappa är en <mark>gruva och han</mark> steijkar för att alla person måste jobba mycket hårt men
lön inte mycket pengar för familj , mat och deras liv . De jätte arg filmen utspelar @ ag i
England <mark>på 1984- talet</mark> . Liv av <mark>människor</mark> är fattiga . Man syns på kläderna i filmen . Man
trots att det måste vara på 1984- talet .



Learner essay Text readability

≡

English

Show all words of the following CEFR level(s)

A1
A2
B1
B2
C1

Additional options @

Mark all potentially incorrect words

Reset

Use Spellchecker

Edit text



David Alfter



Ildikó Pilán



Overall level: B1 Detailed evaluation Number of sentences Number of tokens

16 189

## Automatic !?!? annotation



→ Segmentation

- → Single-word errors (non-words vs real words)
- $\rightarrow$  Phrase-level errors
  - (grammar, combinability)
- → Word order errors

# L2 word-level normalization





#### • Levenstein distance (as is)

- Good for advanced levels (edit distance of 1)
- Fails at lower levels (with multiple edits)

	0		
À			
	1	5	
	-	1	and I

#### • LanguageTool + candidate ranking

- 73% correct variant selection
- Failed to identify 30% of spelling errors

## Levenstein distance



Level	Correct/total
A1	7/20
A2	13/20
B1	13/20
B2	15/20
C1	16/20

(1) substitution of one misspelled letter, e.g.: ursprang<sup>\*</sup>  $\rightarrow$  ursprung (origin);

(2) deletion of an extra letter, e.g.: sekriva\*  $\rightarrow$  skriva (to write), naman\*  $\rightarrow$  namn (name);

(3) insertion of one missing letter, i.e.
 sammanfata\* → sammanfatta
 (summarize).

# LanguageTool + ranking



	# tokens	% tokens
A1	204	9.7
A2	1118	6.0
<b>B1</b>	1650	5.5
B2	3526	10.8
C1	7511	12.4

scores for noun-verb and noun-adjective combinations included with a threshold of LMI  $\ge$  50

Number corrected tokens per level

## L2 "alternative" data

- Logs acc. to a defined research interest
- Steps:
  - Implement an activity for learners
  - Prepare database for storing (structured) data
  - Implement a way to browse logs, visualize statistics etc
  - If necessary add extra annotation steps (manual, automatic)

# Pilot 1 on L2 "alternative" data

- Identifying most predictive features for a language proficiency level (for diagnostic purposes)
  - Multi-word expressions
  - Syntactic properties (e.g. word order)
  - Knowledge of word morphology (e.g. inflections)



David Alfter

# L2 "alternative" data (logs)



### Exercise type evaluation

#### Bundled gaps (variant 1)

Which word fits into these gaps? Each gap contains the same word. Write the word.

Hennes	var på hans lår , gned in värme i
hans kalla ben .	

En annan taxi tar om skolbarnen .

Novelty hade flera trumf på

I första har hon spelat dragspel och fiol.

#### Evaluation

For which levels is this exercise type relevant?



#### Page 1 of 7

Continue

### https://spraakbanken.gu.se/larkalabb/exeval

# Pilot 2 on L2 "alternative" data

- Automatic assigning new words to a proficiency level
  - We predict the level automatically
  - Learners (of a known level) get the word in an exercise (or a series of exercises)
  - We see whether learners can cope with it



David Alfter

Ildikó Pilán

# L2 "alternative" data (logs)



https://spraakbanken.gu.se/larkalabb/wordguess-eesti Word guess



#### Hjälp:

slip; steal, sneak, creep



## Evaluation Reliability of tools



# The ultimate goal

#### L2 infrastructure activity development cycle



### L2 data for Lärka's research agenda











# Thank you!