RECLAS
Research Collegium for Language in Changing Society

# Towards learner corpora of English, Finnish and Swedish (+ some other languages) in Finland

Ari Huhta
Centre for Applied Language Studies
University of Jyväskylä, Finland
ari.huhta@jyu.fi

CLARIN L2 Workshop
Göteborg
December 6-8, 2017

# Goals

- Learn about annotation, error tagging, etc conventions (or decisions within CLARIN)
- Learn about automated tools for annotation etc
- Share information with other CLARIN members about the learning 'corpora' we have → possible collaboration in the future?
- Share with others our experiences about rating / CEFR-linking of the scripts / performances

# Outline / key points

- Research projects & language examinations or tests → learner performances → learner corpora
  - mostly at the University of Jyväskylä
  - mostly **writing** performances / scripts; transcribed & anonymised
  - mostly **unannotated** 'corpora'
  - variation & similarity in the writing tasks across 'corpora'
  - rated against **CEFR** scales (or a version of them) → linguistic characteristics of CEFR levels
  - mostly **teenaged** learners
  - mostly **Finnish-speaking** Finns or **immigrants**
  - other **information about the learners** who wrote the texts

- Fin-Clarin & SLATE collaboration
- Linking with the CEFR levels
- Future plans
- Questions and issues

# Where do the learner performances come from?

**Research projects** at the U. of Jyväskylä

- CEFLING 2007-09
- TOPLING 2010-13
- DIALUKI 2010-13 (and beyond)
- PhD student projects

**National Certificates** examination corpus (National Agency for Education & U. of Jyväskylä)

**National educational evaluation** study (Karvi – the Finnish Education Evaluation Centre) 2011-13 study)

RECLAS
Research Collegium for Language in Changing Society

# CEFLING - Combining Second Language Acquisition and Testing Approaches to Writing

*2007-2009 funded by the Academy of Finland*

- **writing**, 5 different tasks (messages, narrative, descriptive/argumentative)
  - English word derivation
- **cross-sectional**
- **7th**, 8th, 9th graders (aged 12-15)
- English as FL (Finnish L1), ca. 580 scripts
- Finnish as L2 (different L1s), , ca. 1 200 scripts
- Background questionnaire
- Part of the Finnish L2 corpus has been annotated

# CEFLING writing tasks 1
## (original task instruction was in Finnish)

### *Note to a friend*

You've set up a meeting with your English-speaking friend at a café. However, something has come up and you have other things to do. Send an email message to your friend.

- Explain why you can't come.
- Suggest a new time and place.

Remember to **begin** and **end** the message in appropriately. Write **in English / Finnish** in clear characters in the space below.

# CEFLING writing tasks 3

## *Message to an internet store*

Your parents have ordered a PC game for you from a British internet store. When you get the game you notice that it doesn't work properly. You get upset and decide to write an email message to the internet store. In the message, say

- who you are
- what your parents ordered
- why you're unhappy (mention at least two defects/problems)
- how you would like them to take care of  the matter
- give your contact information

Remember to **begin** and **end** the message appropriately. Write **in English / Finnish** in clear characters in the space below.

# CEFLING writing tasks 5

***Narrative***

Tell about the scariest / funniest / greatest experience in your life. Choose one.

- Tell what happened (what, where, when, and so on).
- Tell why the experience was scary / funny / great.

Write **in English / Finnish** in clear characters in the space below (continues on the reverse side).

# TOPLING - Paths in Second Language Acquisition

*2010-13 funded by the Academy of Finland*

- **writing**; similar tasks to CEFLING, 3 times over a period of 2 years
- **longitudinal** and other (e.g. keystroke logging)
- primary, lower secondary, upper secondary (also university)
- English as FL (Finnish L1), ca. 3 400 scripts
- Finnish as L2 (various L1s), ca. 2 550 scripts
- Swedish as L2/FL (Finnish L1), ca. 2 450 scripts

# DIALUKI – Diagnosing Reading and Writing in a Second or Foreign Language

*2010-13 funded by the Academy of Finland*

- **writing** 1-3 tasks, etc (reading, vocabulary, dictation; motivation; cognitive skills; background variables)
- **cross-sectional** (4th, 8th, 11th grade)
- English (Finnish L1), ca. 1 400 scripts (**longitudinal** several hundred)
- Finnish (Russian L1), ca. 400 scipts (longitudinal ca. 200)
- The longitudinal part for English still continues
- Automated analysis tools (Coh-Metrix & L2 Syntactic Complexity Analyzer) used in some sub-studies

# The Finnish National Certificates of Language Proficiency

National Agency for Education & U. of Jyväskylä

- High-stakes **language examination** in 9 languages, see http://ykitesti.solki.jyu.fi/en/

- Started in 1994, by now ca. 110 000 test takers

- A portion of speaking and writing performances are stored in the NC corpus, see http://yki-korpus.jyu.fi/

- Can be linked with test takers' test results and background variables

- Several hundred performances (in most languages)

# An example of a PhD student corpus at U. of Jyväskylä

- Ghulam Khushik's PhD study 2015 –
- **writing, cross-sectional** (& longitudinal)
- lower & upper secondary students
- English as FL (Sindhi L1, in Pakistan) & EFL (Finnish L1 from DIALUKI)
- 6 writing tasks from CEFLING and TOPLING
- Automated analysis tools: Coh-Metrix & L2 Syntactic Complexity Analyzer

→ Differences between CEFR levels; characteristics of levels A1 – B1; differences & similarities between Finnish vs Sindhi L1 speakers of English as FL

# Karvi (Finnish Education Evaluation Centre) corpora

- **National educational evaluation** study 2011-13
- 9th grade (15-year-olds)
- Finnish and Swedish-speaking schools
- English (3500 learners), French (2300), German (2700), Russian (1000), Swedish (1700)
- **writing** (all wrote 2 tasks)
- **speaking** (50-80% completed this)
- we are copying part of the writing and speaking performances for research purposes

# FIN-CLARIN & SLATE co-operation

- U. of Jyväskylä is a member of FIN-CLARIN
- SLATE (Second Language Acquisition and Testing in Europe)
  - network of researchers that combine language testing and SLA perspectives
  - Interested in investigating the **linguistic basis of the CEFR levels** (www.slate.eu.org & EuroSLA Monograph series #1; Bartning et al. 2010)
  - CEFLING, TOPLING & DIALUKI relate to the SLATE network

# www.slate.eu.org
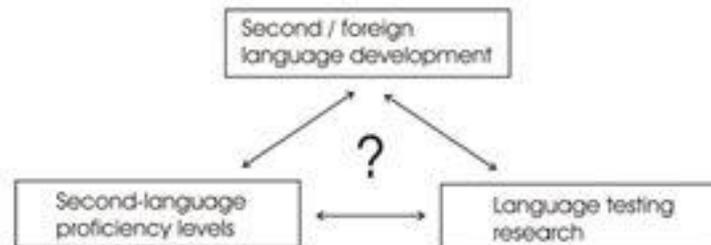
# Rating of performances against CEFR (and other) scales

- set of CEFR writing scales used in the CEFLING, TOPLING and DIALUKI projects (Alanen et al. 2010; Alanen et al. 2012; Huhta et al. 2014)
- One of purposes of these projects was to study the linguistic (lexical, grammatical & textual) features that characterise and/or differentiate CEFR levels
- This required that learners' writing performances could be placed at the CEFR levels → comparison of A1, A2, B1 etc in terms of their linguistic features
- Placement at the CEFR levels was done by rating the written scripts

RECLAS Society

| | OVERALL WRITTEN PRODUCTION | WRITTEN INTER-ACTION | CORRESPONDENCE & NOTES, MESSAGES, FORMS | CREATIVE WRITING & THEMATIC DEVELOPMENT & COHERENCE AND COHESION |
|---|---|---|---|---|
| **A1** | Can write simple isolated phrases and sentences. | Can ask for or pass on personal details in written form. | Can write a short simple postcard. Can write numbers and dates, own name, nationality, address, age, date of birth or arrival in the country, etc. such as on a hotel registration form. | Can write simple phrases and sentences about themselves and imaginary people, where they live and what they do. Can link words or groups of words with very basic linear connectors like 'and' or 'then'. |
| **A2** | Can write a series of simple phrases and sentences linked with simple connectors like 'and', 'but' and 'because'. | Can write short, simple formulaic notes relating to matters in areas of immediate need. | Can write very simple personal letters expressing thanks and apology. Can take a short, simple message provided he/she can ask for repetition and reformulation. Can write short, simple notes and messages relating to matters in areas of immediate need. | Can write about everyday aspects of his/her environment, e.g. people, places, a job or study experience in linked sentences. Can write very short, basic descriptions of events, past activities and personal experiences. Can write a series of simple phrases and sentences about their family, living conditions, educational background, present or most recent job. Can write short, simple imaginary biographies and simple poems about people. Can tell a story or describe something in a simple list of points. Can use the most frequently occurring connectors to link simple sentences in order to tell a story or describe something as a simple list of points. Can link groups of words with simple connectors like 'and', 'but' and 'because'. |

# After ratings: analysis of the quality of the ratings

- To investigate how consistently and similarly raters work (intra-rater & inter-rater consistency / reliability) & to increase the dependability of the placements of learners' scripts at the CEFR levels
  - double or multiple ratings required in most cases
- Many approaches to such analyses
  - proportions of agreement, correlations (several different), Cronbach's alpha, Kendal's tau (see e.g. Kaftandjieva & Takala 2000)
  - Multifaceted Rasch analysis (with Facets) (see Huhta et al. 2014)

# Facets analyses

- For analysing multiple ratings of learners' performances on one or multiple tasks based on one or two rating scales
- To enable a more reliable / accurate placement of the texts on the CEFR (or other) levels
- We can spot and remove (if we so decide):
  - (too) **inconsistent** raters
  - (too) unexpected / inconsistent ratings (invidual data points)
  - (too) **different** raters in terms of severity / leniency

# Facets results – an example of stage 1 of an analysis

# Future plans

- Continue using automated analyses of English learner texts

- Annotate at least some of the 'corpora': Explore automated annotation tools for English, Swedish & Finnish (possibly other languages, too)
  - Including error tagging
  - In compliance with standards created within CLARIN

- Continue collaboration with Karvi to make the writing and speaking performances collected by available for research

# Questions, issues

- Length of performances / scripts
  - very short texts are a problem for investigating individual texts or learners
- Certain kinds of rating issues (severity/leniency)
- Quality / accuracy of automated analyses of learner language
  - Effect of errors / non-standard language on the results?
  - Correction of learner texts: what and how much? (McNamara et al. 2014, p.155-6)
  - Jarvis: correction of misspellings increases accuracy of lexical analyses
  - Identification of sentence boundaries seems to be important (missing sentence final punctuation)

# Sample texts (DIALUKI, longitudinal, 6th grade)

**TASK**: A variant of the 'narrative' task used in CEFLING (telling about a memorable event from the previous year, e.g., during the summer holiday)

## E1016-102

Last year, when I was twelve years old, I visited my frend's house. I went swim and we fried barbegue. I don't know why I remembet this moment, it just was fun!

## E1026-102

In 2010 I'm going to Solkila with my dad and sister. We are going to watch Solki Rally. There is very awesome cars and many people. In Solkila were very hot day, about +°30. I wear a t-shirt, shorts and sunglasses. My dad was happy, because he pääsi watching cars. My sister was verybored and he doesn't want come to here.

# References

Alanen, R. Huhta, A. & Tarnanen M.. 2010. Designing and assessing L2 writing tasks across CEFR proficiency levels. In Bartning, Inge, Martin, Maisa & Vedder Ineke (eds.) *Communicative proficiency and linguistic development: intersections between SLA and language testing research.* EUROSLA Monograph Series, 1. 21-56. http://eurosla.org/monographs/EM01/EM01home.html

Alanen, R., Huhta, A., Jarvis, S., Martin, M. & Tarnanen, M.. 2012 Issues and challenges in combining SLA research and language testing. In Tsagari, Dina & Csepes, Ildiko (eds.) *Collaboration in Language Testing and Assessment*. (pp.15-30). Language Testing and Evaluation Series, Grotjahn, R. &. G. Sigott (general eds).  Frankfurt: Peter Lang. 15-30.

Huhta, A., Alanen, R., Tarnanen, M., Martin, M. & Hirvelä, T. 2014. Assessing learners' writing skills in a SLA study: Validating the rating process across tasks, scales and languages. *Language Testing 31*(3) 307–328.

Kaftandjieva, F. & Takala, S. 2002. Council of Europe scales of language proficiency: A validation study. In CEFR Case Studies, 106-129. https://www.coe.int/t/dg4/linguistic/Source/case_studies_CEF.doc

Khushik, G. & Huhta, A. (submitted). Syntactic complexity features across the Common European Reference Framework levels in the writing of Finnish and Pakistani learners of English.

McNamara, D., Graesser, A., McCarthy, P., & Cai, Z. 2014. *Automated evaluation of text and discourse with Coh-Metrix*. Cambridge: CUP.

# CEFLING writing tasks 2

***Message to your teacher***

You've been away from school for a week. Soon you'll have an English exam. Your teacher, Mary Brown, speaks only English. Send an email message to the teacher.

- Tell her why you've been away.
- Ask two things about the exam.
- Ask two things about the English lessons that were held during the week.

Remember to **begin** and **end** the message appropriately. Write **in English / Finnish** in clear characters in the space below.

# CEFLING writing tasks 4

***Opinion***

Choose one of the topics and write about what you think about the matter. Give reasons for your opinion.

1. Boys and girls should go to different classes at school.
2. No mobile phones at school!

Write **in English / Finnish** in clear characters in the space below (continues on the reverse side).