# Lärka: an online platform where language learning meets natural language processing

*Ildikó Pilán, David Alfter and Elena Volodina*
*Swedish Language Bank, University of Gothenburg, Sweden*

## Lärka: introduction

- A **free**ly available Intelligent Computer Assisted Language Learning (**ICALL**) **platform** developed at Språkbanken
- **Builds on** a variety of language technology **resources** (corpora, lexical resources, NLP tools)
- **Aim**: **support** language **teachers**, **learners** and **researchers**
- **Web-based,** functionalities also available as **web services**

## Multiple-choice exercises

- Exercises for **learners** of **Swedish** and students of **linguistics**

- 1 correct + **distractors** (incorrect options)



*item generation on-the-fly*

*seed sentences from corpora*

*automatic distractor suggestion*

*instant feedback*

## Word Guess

- **Gamified** lexical resourse re-use
- Hangman game mechanism
- **Clues**: definitions, translations



https://spraakbanken.gu.se/larkalabb/wordguess

## Listening exercise

- **Listen** to a word and **type** it
- **Clues**: sentences or first letter
- **Nuance** text-to-speech technology
- **Machine-learning** for the automatic classification of target **word complexity**



https://spraakbanken.gu.se/larkalabb/liwrix

# TextEval: text complexity analysis

- **Aim**: classification of second and foreign language (L2) texts into **proficiency (CEFR) levels**
- Texts written **by (*productive*)** or **for (*receptive*)** L2 learners can be analyzed for complexity
- **Machine-learning** based CEFR level assessment and **statistics**
- Per-**token color**-coded CEFR levels
  - **darker** shades: receptive vocabulary
  - **lighter** shades: productive vocabulary



*text to analyze*

*results*

## Feature set

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

- **Supervised** machine learning
- **Training data**:
  - receptive: L2 coursebook corpus (COCTAILL)
  - productive: L2 essay corpus (SweLL)
- **Accuracy**: 81% (receptive)

  72% (productive)

- From KELLY, an L2 frequency word **lists** with **mappings** to **CEFR** levels based on frequency bands
- **CEFR levels** of tokens, not words as **lexical** features (data sparsity)

*Incidence score*
(normalized category count)

$$\frac{1000}{N^{token}} \times N^{category}$$

$$\frac{N^{category}}{N^{noun} + N^{verb} + N^{adj} + N^{adv}}$$

**http://spraakbanken.gu.se/larka/texteval**

# HitEx: exercise item selection

## Overview

➤ **Aim**: identify sentences from corpora to be used as L2 exercise items

➤ **Why**: a. **authentic** language use
   b. available in **large quantities**

➤ **Challenges:**
   a. sentences might be too **difficult** or
   b. sometimes **out of context**.

➤ **HitEx**: a **hybrid** system for sentence selection with **25 criteria**
   - **machine learning** ➝ L2 complexity
   - **rules** ➝ all other criteria

➤ Criteria set as **filters** or **rankers**

## Architecture



## Common European Framework of Reference for Languages (CEFR)



http://gostudylink.net/en/blog/cefr-levels-explained

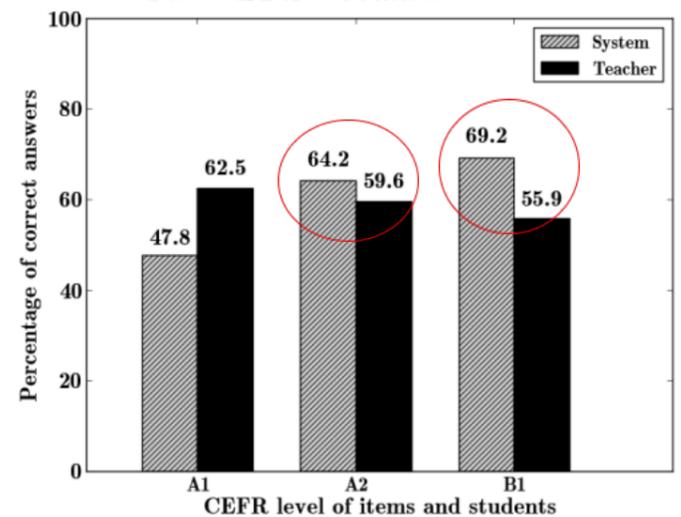## How well does it work?

**Ideal item difficulty (IID)**

$$IID = P_c + (100-P_c) / 2$$

*correct %*

IID = **0.645**
(= 64,5 % students answering correctly)



*A1: harder than it should be*
*A2, B1: appropriate difficulty*

## User interface



*computed on-the-fly*

*highly customizable*

*machine learning*

*detailed information per sentence*

**https://spraakbanken.gu.se/larkalabb/hitex**