

# Description of the SweLL corpus

---

## Contents

<b>General information</b>	1
<b>Metadata description for the SweLL corpus</b>	3
Personal metadata	3
Task metadata	5
Essay metadata	8
Schools metadata	8
<b>Manual coding / annotation in the SweLL corpus</b>	9
Pseudonymization codes	9
Correction annotation codes	11

---

## General information

**Total nr essays:** 525 (June, 2, 2020; **the collection is not finalized and will grow further**)

**NOTE!** In the SweLL corpus we use the term **CORRECTION ANNOTATION** instead of a more traditional **ERROR ANNOTATION**

**Description** of the project: <https://spraakbanken.gu.se/en/projects/swell>

The SweLL corpus is maintained at the University of Gothenburg, Språkbanken-Text

<<https://spraakbanken.gu.se>>

**Personal data management:** Essays were collected following a consent from the learners. The consent allows the use of essays for research by registered (approved) users. Handwritten essays were transcribed using secure encrypted environment (SweLL kiosk). All essays were manually **pseudonymized** (using SweLL kiosk) based on the Pseudonymization guidelines:

[https://spraakbanken.github.io/swell-project/Anonymization\\_guidelines](https://spraakbanken.github.io/swell-project/Anonymization_guidelines)

**Mode:** All essays were written as an exam, in a classroom. Most of the essays were written by hand and were transcribed later according to the Transcription Guidelines:

[https://spraakbanken.github.io/swell-project/Transcription\\_guidelines](https://spraakbanken.github.io/swell-project/Transcription_guidelines)

**Time constraints** and access to **allowed materials** varies between tasks (see details in Task Metadata for each particular task).

Pseudonymized essays were **normalized** and **corr-annotated** by L2 specialists.

To get information on **access** to the SweLL corpus, see the webpage:

<https://spraakbanken.gu.se/en/projects/swell/swell4users>

**To cite the corpus:**

Elena Volodina, Lena Granstedt, Arild Matsson, Beáta Megyesi, Ildikó Pilán, Julia Prentice, Dan Rosén, Lisa Rudebeck, Carl-Johan Schenström, Gunlög Sundberg and Mats Wirén (2019). The SweLL Language Learner Corpus: From Design to Annotation. *Northern European Journal of Language Technology*, Special Issue.

**Read the article:** <https://nejlt.ep.liu.se/article/view/1374>

Newer articles will be added on completion of the project

---

# Metadata description for the SweLL corpus

---

## Personal metadata

<i>General</i>	
Student ID	<b>370</b> unique students ( <b>numbers are not final!</b> ), e.g. <b>C16</b> . Letter prefix (for a school) + a running number
Birth year in 5-year intervals	1950-1954 -- 2000-2004
Gender	Kvinna, Man, Annat, Vill inte säga
Time in Sweden (sum in months)	0 - 315
Native language(s)	<b>72</b> unique languages in 101 unique combinations of 1-4 languages
<i>Education</i>	
Elementary education outside Sweden (nr months)	12 - 156
Elementary education in Sweden (nr months)	12 - 132
Introductory education in Sweden (nyanlända) (months)	1 - 36
Gymnasial education outside Sweden (nr months)	2 - 168

Gymnasial education in Sweden (nr months)	1 - 72
Professional education outside Sweden (nr months)	12 - 60
Professional education in Sweden (nr months)	3 - 72
University education outside Sweden (nr months)	12 - 228
University education in Sweden (nr months)	6 - 96
Professional degree	Free text (e.g. Ekonom)
Educational degree	Free text
Education: additional comments	Free text
<i>Language information</i>	
Education in L1 (hemspråk, education in Sweden)	Language name
Length of education in L1, nr months (hemspråk)	1 - 168
Swedish proficiency courses	Self-instruction, formal education
Length of Swedish proficiency courses, nr months	1 - 336
All known languages	List of languages

Other known language(s) except mother tongue(s)	List of languages
Best written language(s)	List of languages
Best spoken language(s)	Language name(s)
Language(s) used with the family	Language name(s)
Language(s) used with friends	Language name(s)
Metacomment	Free comment added by an assistant

## Task metadata

<i>Administrative Metadata</i>	
Task ID	Task ID in the corpus. A letter prefix (for a school) + <b>T</b> (ask) + a running number, e.g. <b>AT14</b> . Total of <b>38 unique tasks</b>
Semester (time span)	<b>VT-2018, HT-2018, VT-2019</b>
Task date	Year-week, e.g. <b>2018-W20</b>
Course type / school form	<ul style="list-style-type: none"> <li>● Ungdomsgymnasiet</li> <li>● Universitetet</li> <li>● Vuxenutbildningen</li> </ul>

Course level	<ul style="list-style-type: none"> <li>● Behörighetsgivande kurs</li> <li>● Förberedande kurs</li> <li>● Grundläggande SVA dk3</li> <li>● Inplaceringsprov SFI</li> <li>● SFI B / C / D</li> <li>● SVA 2 / 3</li> <li>● TISUS</li> </ul>
Grading scale	<ul style="list-style-type: none"> <li>● A-F</li> <li>● G/U</li> <li>● SFI inplacering</li> <li>● Uppgiften har inte betyg</li> </ul>
<i>Writing task details</i>	
Task type	<ul style="list-style-type: none"> <li>● Behörighetstest</li> <li>● Formativ skrivuppgift</li> <li>● Inplaceringsprov</li> <li>● Mitterminsprov</li> <li>● Slutprov</li> <li>● Test inför NP (Nationella Prov)</li> </ul>
Task - form / mode	Handskriven, Digital
Task duration (in minutes)	25--180
Text type / genre	<ul style="list-style-type: none"> <li>● Argumenterande</li> <li>● Berättande</li> <li>● Beskrivande</li> <li>● Förklarande</li> <li>● Informellt mejl</li> <li>● Instruerande</li> <li>● Resonerande</li> <li>● Utredande</li> <li>● Återgivande</li> </ul>
Task instructions	Free text alt. reference to an attachment
Allowed aids	Bilingual dictionary, Monolingual dictionary, Internet, etc.

Additional material	Free text comment
Additional comments	Free text comment
Additional comment on coursebooks used	Book title(s) / free text
Approximate level	<ul style="list-style-type: none"> <li>● Nybörjare</li> <li>● Fortsättning</li> <li>● Avancerad</li> </ul>
Tasks - subject / topic	<ul style="list-style-type: none"> <li>● 1. Din första arbetsplats alt. Kvinnors arbete; 2. Mejl till kusin på besök i Sverige</li> <li>● 1. När jag var liten och gick till skolan första gången; 2. Mejl till kusin på besök i Sverige</li> <li>● Argumenterande text om arbetsmoral</li> <li>● Argumenterande text/brev</li> <li>● Berätta hur du bor!</li> <li>● Berätta utifrån texten "Alma berättar"</li> <li>● Beskriv - En god relation</li> <li>● Beskriva - Min första kärlek</li> <li>● Brott - orsaker och konsekvenser</li> <li>● Demokratiska val - hur gammal ska man behöva vara för att rösta och varför?</li> <li>● Diskuterande text om pengars betydelse</li> <li>● En kulturupplevelse</li> <li>● En plats du tycker om</li> <li>● En viktig plats</li> <li>● Enkel utredande text om litterära teman</li> <li>● Familjen</li> <li>● Ge tips och råd</li> <li>● Ge tips och råd - en anställningsintervju</li> <li>● Insändare</li> <li>● Kommunikation och sociala medier</li> <li>● Mejl till en vän</li> <li>● Mina första intryck</li> <li>● Objektivt utredande uppgift</li> <li>● Om din bostad och om att bo</li> <li>● Referat av texten "Giftermål ett större steg än barn"</li> <li>● Skriv en insändare</li> <li>● Skriv ett brev</li> </ul>

	<ul style="list-style-type: none"> <li>• Skriv ett mail</li> <li>• Skriv om en känd person</li> <li>• Två sätt att uppfostra</li> <li>• Utredande text (pm) övning inför NP</li> <li>• Världens lyckligaste länder</li> <li>• Övnings-pm inför NP</li> </ul>
Task-url	In certain cases where handouts were used, attachments are available at the urls.

## Essay metadata

Grade	If available, according to the grading scale for a particular Task
-------	--

## Schools metadata

<i>School listing</i>		<i>Description</i>
A	Vuxenutbildningscentrum	Inplacering SFI-utbildning: A-D
B	Gymnasieskola	
C	Komvux/SFI	SFI A-D
E	Behörighetsgivande kurser	motsv gymnasiet
F	Behörighetsgivande kurser	motsv gymnasiet
G	SFI-provet	SFI A-D
H	TISUS-prov	motsv gymnasiet
J	Grundläggande vuxenutbildning	

# Manual coding / annotation in the SweLL corpus

---

## Pseudonymization codes

Pseudonymization guidelines: [https://spraakbanken.github.io/swell-project/Anonymization\\_guidelines](https://spraakbanken.github.io/swell-project/Anonymization_guidelines)

(Needs updating!)

Category type	Codes	Pseudonym / details
Names	firstname_male firstname_female firstname_unknown initials middlename surname	replaced by an equivalent
Geographical data	city (+ <i>foreign</i> for non-Swedish ones) area (+ <i>foreign</i> for non-Swe ones) country geo place region street_nr zip_code	Swedish names replaced with dummy names; All other names with equivalent ones (i.e. cities with other cities)
Institutions	school work other_institution	replaced with equivalents
Transportation	transport_name transport_nr	replace with X-linjen randomly
Age	age_digits age_string	replaced with a random number in the $\pm 2$ span from the actual number

Dates	date_digits day month_digit month_word year	11/11/1111 replace randomly 11/11 replace randomly ±2 randomly
Miscellaneous		
	account_nr	
	e-mail	email@dot.com
	extra	
	license_nr ( <i>e.g. cars</i> )	ABS 000
	other_nr_seq	
	phone_nr	0000-000000
	person_nr	123456-000
	url	url@com
	Zip-code	000 00
Sensitive	edu (education) prof (profession) fam (family members) sensitive (free text with hints on ethnic & sexual info, religious & political views)	Markup only, no replacement

## Correction annotation codes

Normalization guidelines: [https://spraakbanken.github.io/swell-project/Normalization\\_guidelines](https://spraakbanken.github.io/swell-project/Normalization_guidelines)

Correction annotation guidelines: [https://spraakbanken.github.io/swell-project/Correction-annotation\\_guidelines-and-codebook](https://spraakbanken.github.io/swell-project/Correction-annotation_guidelines-and-codebook)

<b>ORTHOGRAPHY</b>	
<i>O</i> (misspelling)	
<i>O-Cap</i> (capitalization)	
<i>O-Comp</i> (compounding)	
<b>LEXIS (derivational morphology included)</b>	
<i>L</i> (any other lexical problem)	
<i>L-Der</i> (wrong derivation mechanism used)	
<i>L-FL</i> (mix of foreign language in Swedish)	
<i>L-Ref</i> (reference error)	
<i>L-W</i> (wrong word choice)	
<b>MORPHOLOGY (inflectional) ≈ PHRASE LEVEL</b>	
<i>M-Adj/adv</i> (adjective instead of adverb)	
<i>M-Case</i> (case problem)	
<i>M-Def</i> (definiteness problem)	
<i>M-F</i> (problem on form / morphology level)	
<i>M-Gend</i> (problem with gender)	
<i>M-Num</i> (problem with number)	
<i>M-Other</i> (any other, not listed, problem on morphology level)	
<i>M-Verb</i> (any problem on the verb or verb phrase level)	
<b>SYNTAX ≈ CLAUSE LEVEL</b>	
<i>S-Adv</i> (word order, sentence adverbial placement)	
<i>S-Clause</i> (clause vs phrase level)	

<i>S-Comp</i> (phrase vs compound structure)	
<i>S-Ext</i> (extensive change)	
<i>S-FinV</i> (word order, verb placement)	
<i>S-M</i> (missing word)	
<i>S-MSubj</i> (missing subject)	
<i>S-Other</i> (any other problem/syntact. level)	
<i>S-Type</i> (change of construction type/phrase level)	
<i>S-WO</i> (any other word order problem)	
<b>PUNCTUATION</b>	
<i>P</i> (general miss on punctuation)	
<i>P-M</i> (missing)	
<i>P-R</i> (redundant)	
<i>P-Sent</i> (sentence segmentation)	
<i>P-W</i> (wrong)	
<b>OTHER</b>	
<i>C</i> (consistency change)	
<i>Cit-FL</i> (use of foreign language for citation)	
<i>Com!</i> (comment on an essay level)	
<i>OBS!</i> (comment on a token level)	
<i>X</i> (unintelligible)	
<i>Unid</i> (unidentified)	

---