## **Term Paper**

In this final project, I wrote a grapheme-to-phoneme (g2p) conversion program for Swedish and evaluated it against the Swedish pronunciation dictionary from WikiPron.<sup>1</sup>

Swedish has 9 vowels and 20 consonants in its alphabet; 26 of these letters are the same as those in English while the three others,  $\langle \mathring{a} \rangle$ ,  $\langle \ddot{a} \rangle$ , and  $\langle \ddot{o} \rangle$ , are additional, commonly used letters in the language. In my opinion, despite being relatively more phonetic than English, Swedish presents a difficult case for such tasks like g2p especially due to its large vowel phoneme inventory and specific segmental rules in addition to the tonal word accent system (i.e., pitch accent) that plays a very important role in pronunciation helping with situations such as differentiating *anden* 'the duck' (pronounced with accent 1 where the pitch on the stressed syllable (i.e.,  $\acute{a}n$ ) goes up and then down again in the unstressed syllable (i.e.,  $\acute{a}e$ ) from *anden* 'the spirit' (pronounced with accent 2 where the pitch rises in the stressed syllable (i.e.,  $\acute{a}e$ ) as in accent 1, and it also rises again on the following, (un)stressed syllable (i.e.,  $\acute{a}e$ ) as well) (Riad 2014). Stress placement is also a major factor in determining the pronunciation of any given word – in particular, the vowel length. However, according to Riad (2014), it is mainly morphology that controls where stress falls within words while there is not a strong phonological generalization for stress placement.

There is a wide range of variation when it comes to dialects in Swedish. To the best of my knowledge, all the rules I implemented are based on the information that is said to be valid for Central Sweden Swedish. In my g2p rules, I did not take the pitch accent into consideration. Another assumption I made is closely related to vowel pronunciation: Vowels in unstressed syllables are always short; however, since vowel length is such an important aspect of Swedish, I wanted to write rules to reflect short and long vowels. Since a vowel in a stressed syllable can be either short or long, the natural assumption I ended up making is that all the syllables are stressed, and it is only the phonetic environment that differentiates short vowels from their long versions and vice versa. I also made major simplifications and some minor adjustments in some rules to make them work with each other some of which I detail further below. Lastly, I left out some rules that were too complicated to handle and tried to keep the rules I used as general as possible without losing the overall "Swedishness".

The first major issue I encountered was the phoneme inventory of Swedish. Different sources used slightly different phonetic symbols especially for some vowel sounds (e.g., /e/ vs. /ɛ/ for <e>), and there is a vowel lowering rule that I decided to use, which added an initial 4 phonemes to the set of 18 phonemes that most sources used. I had to place the vowel lowering rule before the short- and long-vowel rules so that nothing is overridden by those two rules.

Another major issue was with some consonant rules, especially regarding the graphemes <c>, <k>, <ck>, <g>, <gn>, and <j>. Most rules ended up being overridden by some others, and I fixed the failures by reordering the rules and making small adjustments to some phonetic environments making them more specific (e.g., in *Rule #11 – <gn> as /\eta n/*, the preceding environment for the rule to apply is "elsewhere" according to the Swedish Alphabet Wikipedia

<sup>1</sup> The file is found here: https://github.com/kylebgorman/wikipron/blob/master/data/scrape/tsv/swe\_latn\_broad.tsv

page; however, I restricted it to only vowels since it would otherwise override *Rule #10* giving us an incorrect pronunciation of [ŋnɪst] for the word *gnist* 'telegraph'.)

Lastly, there are lots of exceptions to all these rules, some sounds, such as /fj/, are very difficult to handle due to being very common and being represented by many different spellings, and there is also optionality for some letter-pronunciation correspondences. For example, the consonant combination <ch> can be pronounced both as /fj/ and /c/. Initially, I thought I could create a union of possible outputs, but I got a "Multiple top writes" error, which led me to pick the most plausible phoneme based on my personal experiences with the language, namely /fj/.

Overall, I ended up creating 24 rules, which can be found <u>here</u>. Once I finished my rules, I wrote my tests in <u>g2p\_test.py</u> and all the word-pronunciation pairs I picked for the tests were successful.

To my surprise, it was the second half of the project – evaluating the grammar against the Swedish pronunciation dictionary from WikiPron – that I struggled the most with. I worked with the broad transcription data, which needed the following changes:

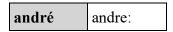
- Since I ignored accent in my grammar, I needed to remove the <sup>1</sup> and <sup>2</sup> used to mark accent in the file. For example:
  - o abbe ² a b ε
- There were too many repeated entries that I had to remove.
  - After cleaning the data, I deleted the TSV file that I had cloned from here previously, which had many repeated words and WikiPron links as pronunciations for some words, but now the data in the hyperlinked page seems to be devoid of all those problems that I had to clean by hand before making me unable to provide examples here.
- There were IPA symbols such as /r, \_, ɔ, œ/ and some more in the pronunciations from the WikiPron data that I did not include in my grammar. Therefore, I had to leave out words that had such symbols in their pronunciations. For example:
  - o adelssläkt Ď: dəls σslεk: t

The dictionary that I ended up with has 3328 word-pronunciation pairs, which is 542 pairs fewer than the original WikiPron file. There is one other change I had to make in the transcription file that I am going to mention below.

Since I was unable to figure out how to do the evaluation and then compute the Word-Error Rate (WER), it took me a long while to get back to the project, but I have finally finished the evaluation script, accessible as g2p\_eval.py, which gave a composition failure error on the first run. Since I did not want to end this whole project resorting to using try:/except: (which is still included in the program) and being unable to compute the WER, I tried the following:

- I created a TSV file out of the word-pronunciation pairs I used in my g2p\_test.py script, ran the g2p\_eval.py on it to receive a WER of 0.0, which gave me the idea that there may be something wrong with the WikiPron data that I cleaned up before.
- I realized the white spaces between the segments in the WikiPron data were at least one of the problems, so I removed them using regex.

- Since I still got the same Composition Failure error, I decided to try to evaluate portions of the WikiPron data; I got a WER of 92.00 from the first 50 words in the file, which was the first successful step!
  - O I first tried to divide the file into halves, but since I kept getting the same error with both halves and halves of halves and so on, which soon became too difficult to handle, I move on with portions of 20 pairs after the first 50, and found the problem:
    - The 76<sup>th</sup> word in the list has the letter <é>, which was not included in my SIGMA STAR, since it is not a part of the alphabet.<sup>2</sup>



• I added the letter to my grammar and ran the program without any failures this time.

While the accuracy, which I am aware was not asked to be reported in the project, is 0.22, the WER is 78.46 as seen below:

```
(base) Selins-Air-2:LING83800_Term_Project selinalkan$ ./g2p_eval.py --input sv_desegmented.tsv
Accuracy: 0.22
WER: 78.46
```

I have been studying Swedish for some years now, and not having any native speakers to interact with around me make it hard to check and/or confirm the pronunciations of words. Among the different resources which sometimes provide the chance to listen to pronunciations, Wiktionary has always been my first choice since it also mostly provides the IPA transcriptions. However, there are entries with no audio files to listen to and the IPA transcriptions are not always available and/or accurate, which causes problems for me and other learners to practice the language. Therefore, even though this is a very specific and practical issue, it still constitutes a good enough reason to automate and improve the g2p technology in Swedish as it is the case for most other languages, too. My work on Swedish has also proved how difficult it can be to capture all the pronunciation rules and numerous exceptional pronunciation patterns in languages.

Lastly, the graphemes I used are based on the <u>Swedish Alphabet Wikipedia page</u>. For the phonemes, I consulted three sources: the book titled <u>The Phonology of Swedish</u> by Tomas Riad, the <u>Swedish Phonology Wikipedia page</u>, and the <u>Swedish IPA Wikipedia page</u>. The phonological rules are based on the following: the <u>Swedish Alphabet Wikipedia page</u>, the <u>Swedish Phonology Wikipedia page</u>, <u>The Phonology of Swedish</u> by Tomas Riad, Swedish Phonology on <u>Glottopedia</u>, <u>this</u> blogpost on the Language Learners website.

<sup>&</sup>lt;sup>2</sup> In the *Usage Notes* for Swedish regarding <é>, this Wiktionary page also indicates the letter in question is not an independent letter in the Swedish alphabet.

## References

Help:IPA/Swedish. (2022, May 18). In Wikipedia.

https://en.wikipedia.org/w/index.php?title=Help:IPA/Swedish&oldid=1088478735

Nario, J. (2021, August 24). Tricks to master the Swedish accent: Common Swedish pronunciation problems students have & how to solve them. *Language Trainers*. <a href="https://www.languagetrainers.ca/blog/swedish-accent-pronunciation/#:~:text=In%20Swedish%2C%20speakers%20use%20single,to%20distinguish%20between%20two%20words.&text=Double%20consonants%20also%20affect%20pronunciation,(vägen%2C%20the%20road).

Riad, T. (2013). *The phonology of Swedish*. Oxford University Press. DOI: https://doi.org/10.1093/acprof:oso/9780199543571.001.0001

Swedish alphabet. (2022, May 5). In Wikipedia.

https://en.wikipedia.org/w/index.php?title=Swedish\_alphabet&oldid=1086274115

Swedish Phonology. (2018, March 2). In Glottopedia.

 $\frac{\text{http://www.glottopedia.org/index.php/Swedish\_Phonology\#:}\sim:\text{text=Swedish}\%20\text{makes}\%}{20\text{use}\%20\text{of}\%20\text{nine}\%20\text{short}\%20\text{vowels.\&text=In}\%20\text{many}\%20\text{cases}\%20\text{and,}\%5\text{B}\epsilon\%}{5\text{D}\%20\text{are}\%20\text{only}\%20\text{allophones.}}$ 

Swedish phonology. (2022, May 20). In Wikipedia.

https://en.wikipedia.org/w/index.php?title=Swedish\_phonology&oldid=1088877117 é. (2022, August 6). In *Wiktionary*. https://en.wiktionary.org/wiki/é#Usage\_notes\_2