

# CLUZH at SIGMORPHON 2020 Shared Task on Multilingual Grapheme-to-Phoneme Conversion

Peter Makarov

Simon Clematide

Institute of Computational Linguistics

University of Zurich, Switzerland

makarov@cl.uzh.ch

simon.clematide@cl.uzh.ch

## Abstract

This paper describes the submission by the team from the Institute of Computational Linguistics, Zurich University, to the Multilingual Grapheme-to-Phoneme Conversion (G2P) Task of the SIGMORPHON 2020 challenge. The submission adapts our system from the 2018 edition of the SIGMORPHON shared task. Our system is a neural transducer that operates over explicit edit actions and is trained with imitation learning. It is well-suited for morphological string transduction partly because it exploits the fact that the input and output character alphabets overlap. The challenge posed by G2P has been to adapt the model and the training procedure to work with disjoint alphabets. We adapt the model to use substitution edits and train it with a weighted finite-state transducer acting as the expert policy. An ensemble of such models produces competitive results on G2P. Our submission ranks second out of 23 submissions by a total of nine teams.

## 1 Introduction

G2P requires mapping a sequence of characters in some language into a sequence of International Phonetic Alphabet (IPA) symbols, which represent the pronunciation of this input character sequence in some abstract way (not necessarily phonemic, despite the name of the task) (Figure 1).

Multilingual G2P is Task I of this year’s SIGMORPHON challenge. It features fifteen languages from various phylogenetic families and written in different scripts. We refer the reader to Gorman et al. (2020) for an overview of the language data. Each language comes with 3,600 training and 450 development set examples. It is permitted to use external resources as well as to build a single multilingual model.

We participate in this shared task with an adaptation of our SIGMORPHON 2018 system (Makarov

*fathaigh*  $\mapsto$  /fa:/ (“giants”)  
Irish of Cois Fhairrge (de Bhaldraithe, 1953)

Figure 1: Example of G2P.

and Clematide, 2018b), which was particularly successful in type-level morphological inflection generation. Our system is a neural transducer that operates over explicit edit actions and is trained with imitation learning (Daumé III et al., 2009; Ross et al., 2011; Chang et al., 2015, IL). It has a number of useful inductive biases, one of which is the familiar bias towards copying the input (implemented as the traditional copy edit). This is particularly useful for morphological string transduction problems, which typically involve small and local edits and where most of the input is preserved in the output. This contrasts with models that rely purely on generating characters such as generic encoder-decoder models, which as a result suffer, particularly on smaller-sized datasets.

Copying requires that the input and output character alphabets overlap, preferably substantially. This also allows our IL training to leverage a simple-to-implement expert policy (which during training provides demonstrations to the learner of how to optimally solve the task). The optimal completion of the target given the prediction generated so far during training requires finding edits that would extend the prediction so that the Levenshtein distance (Levenshtein, 1966) between the target and the partial prediction + the future suffix is minimized. Unfortunately, this objective alone would not discriminate between multiple edit action sequences that relate the input and the partial prediction + the future suffix. To address this spurious ambiguity, our IL training adds edit sequence scores, computed using traditional costs,<sup>1</sup> into the

<sup>1</sup>Copy costs zero, all other edits cost one.

objective. This naturally encourages the system to copy, however this would fail on any editing problem with disjoint alphabets.

G2P poses an interesting challenge for a system like ours. On the one hand, G2P shares many similarities with morphological string transduction: The changes are mostly local, it would suffice to perform traditional left-to-right transduction, and a substantial part of the work is arguably applying equivalence rules (e.g. the German letter “g” most often converts to /g/, “a” to /a/ or /a:/), which is similar to copying. Yet, a general solution to G2P cannot rely on overlapping alphabets since many scripts do not share many symbols, if any at all, with IPA (e.g. Korean or Georgian).

Our solution adapts the model to use substitution edits and trains it with a weighted finite-state transducer acting as the expert policy.

## 2 Model description

The underlying model is a neural transducer introduced in [Aharoni and Goldberg \(2017\)](#). It defines a conditional distribution over traditional edits  $p_\theta(\mathbf{y}, \mathbf{a} \mid \mathbf{x}) = \prod_{j=1}^{|\mathbf{a}|} p_\theta(a_j \mid a_{<j}, \mathbf{x})$ , where  $\mathbf{x}$  is an input sequence of graphemes and  $\mathbf{a} = a_1 \dots a_{|\mathbf{a}|}$  is an edit action sequence. (The output sequence of IPA symbols  $\mathbf{y}$  is deterministically computed from  $\mathbf{x}$  and  $\mathbf{a}$ .) The model is equipped with a long short-term memory (LSTM) decoder and a bidirectional LSTM encoder ([Graves and Schmidhuber, 2005](#)). The challenge is training this model: Due to the recurrent decoder, it cannot be trained with exact marginal likelihood unlike the more familiar weighted finite-state transducer ([Mohri, 2004](#); [Eisner, 2002](#), WFST) or its neuralizations ([Yu et al., 2016](#)). For a more detailed description of the model, we refer the reader to [Makarov and Clemenide \(2018a\)](#).<sup>2</sup>

**IL training** [Makarov and Clemenide \(2018a\)](#) propose training the model using IL, a general model fitting framework for sequential problems over exponentially sized output spaces. IL has been applied successfully to natural language processing (NLP) problems, e.g. transition-based parsing ([Goldberg and Nivre, 2012](#)) and language generation ([Welleck et al., 2019](#)). IL relies on the availability of demonstrations of how the task can optimally

<sup>2</sup>The model uses shared input character / action embeddings of size 100 and one-layer LSTMs with hidden-state size 200.

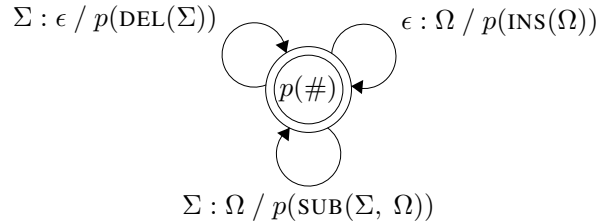


Figure 2: Stochastic edit distance ([Ristad and Yianilos, 1998](#)): A memoryless probabilistic FST.  $\Sigma$  and  $\Omega$  stand for any input and output symbol, respectively.

be solved given any configuration. Due to the nature of many NLP problems, such demonstrations can often be provided by a rule-based program (known as expert policy).

[Makarov and Clemenide \(2018a\)](#) use a combination of Levenshtein distance and edit sequence cost as the task objective ( $\beta \text{ED}(\hat{\mathbf{y}}, \mathbf{y}) + \text{ED}(\mathbf{x}, \hat{\mathbf{y}})$ ,  $\beta \geq 1$ ) and devise an expert policy for it. Given a target sequence  $\mathbf{y}$ , a partially completed prediction  $\hat{\mathbf{y}}_{1:n}$ , and the remaining input sequence  $\mathbf{x}_{k:l}$ , the expert needs to (1) identify the set of target suffixes  $\mathbf{y}_{j:m}$  that when appended to  $\hat{\mathbf{y}}_{1:n}$ , lead to a prediction with minimum Levenshtein distance from the target, and (2) check which of the edit sequences producing those suffixes have the lowest cost, i.e. minimum Levenshtein distance from the remaining input.

The second part is crucial for training accurate models especially in the limited resource setting, as it reduces spurious ambiguity arising under the first part of the objective alone. It is also the second part of the training objective that hinges on the overlap of the input and output alphabets, as this permits minimization using the edit distance dynamic program with traditional costs.

### 2.1 Adaptation to G2P

The adaptation is two-fold: First, we introduce substitution edits, which have previously not been employed to keep the total number of edit actions to a minimum. For each output character  $c$ , there is now a substitution action  $\text{SUBS}[c]$  which substitutes  $c$  for any input character  $x$ .

When the alphabets are disjoint, the completing edit sequences cannot be very informatively scored using traditional edit costs. For example, for the data sample  $\text{кит} \mapsto /k^i t/$  (Russian: “whale”), we would like the following most natural edit sequence to attain the lowest cost:  $\text{SUBS}[k], \text{INS}[i], \text{SUBS}[i], \text{SUBS}[t]$ . Yet, it is clear

that under traditional costs, this sequence attains the same cost as any other that consists of three substitutions and one insertion. Our solution to this is to learn costs from the training data to ensure an intuitive ranking of edit sequences.

**SED policy** Learning costs as well as computing string distance can be achieved with a very simple WFST: Stochastic Edit Distance (Ristad and Yianilos, 1998, SED), which is a probabilistic version of Levenshtein distance (Fig. 2). We use traditional multinomial parameterization.

Before starting training the neural transducer, we train a SED model using the Expectation–Maximization algorithm (Dempster et al., 1977). We use the following update in the M-step:  $\theta^{(t+1)} \propto \max(0, \tilde{\theta} + \alpha)$ , where  $\tilde{\theta}$  is the unnormalized weight computed in the E-step and  $0 < \alpha < 1$  is a sparse Dirichlet prior parameter associated with this edit. This corresponds to sparse regularization via Dirichlet prior (Johnson et al., 2007), which results in many edits having zero probability. We found this training to lead to more accurate SED models. Furthermore, it dramatically reduces the size of the edit action set that the neural transducer is defined over.

SED is integrated into the expert policy. During training, given a configuration consisting of a partial prediction, a remainder of the input, and the target, we query the expert policy for next optimal edits. We minimize the first part of the objective much like before, and we minimize the second part by decoding SED with the Viterbi algorithm.

Suppose we transduce the French word  $\mathbf{x} = \text{abject}$  (“vile”) into the target  $\mathbf{y} = \text{a b } \zeta \text{ } \epsilon \text{ k t}$ . Suppose also that the neural transducer currently attends to character  $x_4 = e$  and the prediction built so far during training is  $\hat{\mathbf{y}}_{1:7} = \text{a b } \zeta \text{ } e$  (note the error). We query the SED policy to get the optimal edit action whose likelihood we will maximize. First, much like before, we find that the following edits are optimal with respect to the first term of the training objective (call them *permissible*) as they do not increase the Levenshtein distance of the prediction from the target (assuming all subsequent edits are permissible too):  $\text{SUBS}[\epsilon]$ ,  $\text{INS}[\epsilon]$ ,  $\text{DEL}$ ,  $\text{SUBS}[\ ]$ ,  $\text{INS}[\ ]$ . (This can be verified by looking at the Levenshtein distance prefix matrix for strings  $\hat{\mathbf{y}}_{1:7}$  and  $\mathbf{y}$ .) Each such edit starts a suffix that completes the target, e.g. it is “ $\epsilon$  k t” for  $\text{SUBS}[\epsilon]$  and “k t” for  $\text{SUBS}[\ ]$ . Next, we use SED to rank the permissible edits by cost-to-go. For

each of the edits and their corresponding suffixes, the expert needs to execute the edit (e.g.  $\text{SUBS}[\epsilon]$  writes  $\epsilon$  and moves the attention to  $x_5 = e$ ) and then decode SED with Viterbi on the the remaining input and the suffix (both possibly modified by the edit). In this way, we obtain that  $\text{SUBS}[\ ]$  is the optimal action with the lowest cost-to-go (=negative sum of the log probabilities of the edit and of the Viterbi path) of 15.28 (vs 17.65 for  $\text{SUBS}[\epsilon]$ , 21.09 for  $\text{INS}[\epsilon]$ , 17.31 for  $\text{DEL}$ , and 17.31 for  $\text{INS}[\ ]$ ).<sup>3</sup>

**Exploration** This time, we also train the transducer with an aggressive exploration schedule:  $p_{\text{sampling}}(i) = \frac{1}{1+\exp(i)}$ , where  $i$  is the training epoch number. After a couple of training epochs, training configurations are generated entirely by executing edit actions sampled from the model.

### 3 Submission details

We train separate models for each language on the official training data and use the development set for model selection.<sup>4</sup> Our submission does not use any additional lexical resources.

For most of the models, we employ Unicode decomposition normalization (NFKD)<sup>5</sup> as a data preprocessing step. Importantly, this helps decomposing Unicode syllable blocks used e.g. in Hangul.

The size of the development set is rather small (450 examples), and having examined the data, we suspect that overly relying on the development set for model selection might hurt generalization. For example, the French development set contains three exceptions to the “ill”-/j/ equivalence; thus, a single model that achieves a high score on the development set might, in fact, be overfitting. To counter this, we build an eleven-model-strong majority-vote ensemble. Fortunately, training a neural transducer is fast as one epoch takes just about four minutes on average on a single CPU, due to the relatively small number of model parameters.

<sup>3</sup>This particular SED is trained on the French training data for 3 EM epochs with Dirichlet prior  $\alpha = 1e-05$  for all edits.

<sup>4</sup>We train the SED model for 20 epochs of EM with  $\alpha = 0.25$  for insertions and 0.5 for all other edits. We train the neural transducer for a maximum of 60 epochs with a patience of 12 epochs. We use mini-batches of size 5. We decode using beam search with beam width 4.

<sup>5</sup>Using NFKD instead of NFD was a bit unfortunate because some superscript diacritics get normalized to their regular size. Luckily, as pointed out to us by Kyle Gorman, there is a unique mapping from NFKD to NFC for the spaced output format of this task. See <http://www.unicode.org/reports/tr15/> for Unicode normalization forms.



ady	ʻ/ε/17	ə●/ε/9	ʃ/ʒ/8	ε/ə●/7	j●/ε/6	ε/ʻ/6	ε/ə●/5	ʒ/1/5	:/ε/5	a/ə/5
arm	ɔ/o/17	ε/ə●/12	○/●/12	ə●/ε/3	t/d/3	g/k <sup>h</sup> /2	ʃ <sup>h</sup> /ʒ/2	ε/j/2	χ/ʒ/2	tj/d●z/1
bul	r/r/26	o/ɔ/22	ə/a/14	a/ə/12	ç/ε/9	ε/ç/9	a/v/7	ʔ/1/5	v/ə/5	ε <sup>ʃ</sup> /5
dut	ə/ε/9	ε/j●/4	a:/a/4	e:/ə/4	ə/e/3	t/d/3	:/ε/3	o:/ɔ/2	ε/ε●/2	n/m/2
fre	ε/●ā/2	ε/●s/2	a/a/2	o/ɔ/2	w/ɔ/2	●j●ā/ε/1	ε/●k●s/1	ɔp/o/1	●ā/ε/1	e/ε●ʒ/1
geo	i/i/103	i/i/48	χ/x/5	ʏ/ʒ/4	ʒ/ʏ/3	x/χ/3	a/a/2	●s/ε/1		
gre	r/r/27	o/ɔ/19	r/r/15	e/ε/9	j/i/3	n●/ε/2	ç/i/2	m/m/2	●m●e/ε/1	ε/s●/1
hin	ε/ə●/10	ə●/ε/5	ε/●ə/2	ε:/ə/2	ε/●/2	α/a/2	i/i:1	●fi <sup>h</sup> /1	i/i/1	ə/j/1
hun	ʃ/ʒ/3	ε:/3	e:/i●n●t/1	ε/m●v●/1	m/e <sup>j</sup> /1	ts/x:/1	s:/ʃ●s/1	h●/ε/1	●h/○j/1	○/●/1
ice	:/ε/11	ε:/9	t●/ε/4	v/f/3	ε/ç/3	t/d/2	<sup>h</sup> /ε/2	h●/ε/2	γ●y/u:/1	c <sup>h</sup> /k/1
jpn	ε/ʒ/8	ε/ç/6	ε:/3	:/●ɸ <sup>B</sup> /2	:/●q/2	ɸ <sup>B</sup> /j●q/1	ε/●e/1	ɸ <sup>B</sup> /e/1	o:/ā/1	s●i/ū/1
kor	ε:/72	:/ε/18	ə:/Δ/11	z/ç/4	Δ/ə/4	d/t/4	g/k/3	ε/p●/3	d/t/2	l/n/2
lit	ε/ç/15	n/η/14	v/a:/12	:/ε/8	j/ε/7	o/ɔ/7	ε:/6	ε <sup>ʃ</sup> /5	ε/æ:/3	a:/v/3
rum	○/●/8	●/○/8	ç/j/6	r/r/5	j/●i/4	:/j/3	i/j/3	ε/e●/2	o/e/2	j/i/2
vie	ε:/2	+/ε/1	ε/e●/1	w●/ε/1	ʻ/ε/1	○m/ε/1	ə/e/1	ʔ/n/1	ʔ/ε/1	a/ɔ/1

Table 2: Ten most frequent errors per language. Notation: prediction / gold / error frequency. ● denotes whitespace. Computed using the UTF-8 aware version of the ISRI Analytic Tools for OCR Evaluation.<sup>8</sup>

**How good is SED policy?** Somewhat surprisingly, using SED as part of the expert policy results in competitive performance. Yet, SED is a very crude model (e.g. because of the lack of context, when used as a conditional model, SED assigns less probability to any edit sequence containing insertions than the same sequence but with all the insertions removed; this e.g. makes it unusable as a standalone model for G2P). On top of this, we also do not use learned roll-out, which would be recommended when training with a sub-optimal expert (Chang et al., 2015). We leave it for future work to examine whether the neural transducer’s performance on G2P would improve from replacing SED with a more powerful model.

## 5 Conclusion

This presents the approach taken by the CLUZH team to solving the SIGMORPHON 2020 Multilingual Grapheme-to-Morpheme Conversion challenge. Our submission is based on our successful SIGMORPHON 2018 system, which is a majority-vote ensemble of neural transducers trained with imitation learning. We adapt the 2018 system to work on transduction problems with disjoint input and output alphabets. We add substitution actions (not available in previous versions of the system) and employ a memoryless probabilistic finite-state transducer to define the expert policy for the imitation learning. We use majority-vote ensembling to counter the overfitting to the small development sets. These simple modifications result in a highly

competitive performance even without the use of any external resources or learning a single multilingual model. Our ensemble ranks second out of 23 submissions by a total of nine teams. Our error analysis indicates that addressing many of the errors requires additional information such as knowing the word’s lexical category, morphological segmentation, or etymology. We will make our code publicly available.

## Acknowledgment

We would like to thank the organizers for their great effort in these turbulent times. We thank Kyle Gorman for taking the time to help us with our Unicode normalization problem. This work has been supported by the Swiss National Science Foundation under grant CR-SII5\_173719.

## References

- Roei Aharoni and Yoav Goldberg. 2017. Morphological inflection generation with hard monotonic attention. In *ACL*.
- Tomás de Bhaldraithe. 1953. *Gaeilge Chois Fhairrge: An Deilbhocht*. Institiúid Ard-Léinn Bhaile Átha Cliath.
- Kai-Wei Chang, Akshay Krishnamurthy, Alekh Agarwal, Hal Daume III, and John Langford. 2015. Learning to search better than your teacher. In *ICML*.
- Hal Daumé III, John Langford, and Daniel Marcu. 2009. Search-based structured prediction. *Machine learning*.

<sup>8</sup><https://github.com/eddieantonio/ocreval>

- Arthur P Dempster, Nan M Laird, and Donald B Rubin. 1977. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1).
- Jason Eisner. 2002. Parameter estimation for probabilistic finite-state transducers. In *ACL*.
- Yoav Goldberg and Joakim Nivre. 2012. A dynamic oracle for arc-eager dependency parsing. In *COLING*.
- Kyle Gorman, Lucas F.E. Ashby, Aaron Goyzueta, Arya D. McCarthy, Shijie Wu, and Daniel You. 2020. The SIGMORPHON 2020 shared task on multilingual grapheme-to-phoneme conversion. In *Proceedings of the 17th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology*.
- Alex Graves and Jürgen Schmidhuber. 2005. Frame-wise phoneme classification with bidirectional LSTM and other neural network architectures. *Neural Networks*, 18(5).
- Mark Johnson, Thomas L Griffiths, and Sharon Goldwater. 2007. Bayesian inference for PCFGs via Markov Chain Monte Carlo. In *NAACL-HLT*.
- Vladimir I Levenshtein. 1966. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet physics doklady*, 10(8).
- Peter Makarov and Simon Clematide. 2018a. Imitation learning for neural morphological string transduction. In *EMNLP*.
- Peter Makarov and Simon Clematide. 2018b. UZH at CoNLL-SIGMORPHON 2018 shared task on universal morphological reinflection. *Proceedings of the CoNLL SIGMORPHON 2018 Shared Task: Universal Morphological Reinflection*.
- Mehryar Mohri. 2004. Weighted finite-state transducer algorithms. An overview. In *Formal Languages and Applications*, volume 148 of *Studies in Fuzziness and Soft Computing*. Springer Berlin Heidelberg.
- Eric Sven Ristad and Peter N Yianilos. 1998. Learning string-edit distance. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. 2011. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*.
- Ho-Min Sohn. 2001. *The Korean language*. Cambridge University Press.
- Sean Welleck, Kianté Brantley, Hal Daumé III, and Kyunghyun Cho. 2019. Non-monotonic sequential text generation. In *ICML*.
- Lei Yu, Jan Buys, and Phil Blunsom. 2016. Online segment to segment neural transduction. In *EMNLP*.