

IMPLEMENTING LOGIC-BASED APPROACHES TO PHONOLOGY

Hossep Dolatian

Stony Brook University

September 22, 2022

TABLE OF CONTENTS

- INTRODUCTION
- EXAMPLE OF LOGICAL TRANSDUCTION
- FORMAT FOR LOGICAL TRANSDUCTIONS
- COMPILING A LOGICAL TRANSDUCTION
- WHAT'S NEXT AND OPEN PROBLEMS
- CONCLUSION

EXPLAINING THE TITLE

- Title: *Implementing logic-based approaches to phonology*
- What we'll do:
 1. See what logical approaches *look* like
 2. Provide software support for running them
 3. Discuss general problems for future research and expansion
- URL to software and tutorial: tinyurl.com/hossepBMRS

1. Phonologists like finding patterns and making rules for them:
 - ▶ Dataset: { ba, pa, aba, *apa }
 - ▶ Rule: $p \rightarrow b / a_a$
2. After you make a rule on pen-and-paper, it's good to double-check the rule to make sure it's correct
3. How do you do that?
 - ▶ Manually, easy but risky
 - ▶ Computationally, more effortful but more reliable

DEBUGGING PEN-AND-PAPER LINGUISTICS

- To computationally debug a phonological analysis, you do the following:
 1. Convert pen-and-paper analysis (A1) to an explicit format (F)
 2. Compile the format F into a machine M
 3. Use the machine to double-check your analysis

DEBUGGING PEN-AND-PAPER LINGUISTICS

- To computationally debug a phonological analysis, you do the following:
 1. Convert pen-and-paper analysis (A1) to an explicit format (F)
 2. Compile the format F into a machine M
 3. Use the machine to double-check your analysis
- A conventional-ish way to do that is using finite-state machines:
 1. Manually convert your rule ($p \rightarrow b/a_a$) into a hand-drawn FST
 2. Compile your FST on Pynini, HFST, or whatever
 3. Run the FST on your data
- For FST users, there's a LOT of work on all these areas:
 - ▶ Conventional formats: diagrams, transition lists, ATT formats
 - ▶ Reliable compilers: Pynini, HFST, a myriad of packages
 - ▶ Some of these compilers can do steps (1-3) simultaneously with little human intervention

DEBUGGING PEN-AND-PAPER LINGUISTICS

- To computationally debug a phonological analysis, you do the following:
 1. Convert pen-and-paper analysis (A1) to an explicit format (F)
 2. Compile the format F into a machine M
 3. Use the machine to double-check your analysis
- A conventional-ish way to do that is using finite-state machines:
 1. Manually convert your rule ($p \rightarrow b/a_a$) into a hand-drawn FST
 2. Compile your FST on Pynini, HFST, or whatever
 3. Run the FST on your data
- For FST users, there's a LOT of work on all these areas:
 - ▶ Conventional formats: diagrams, transition lists, ATT formats
 - ▶ Reliable compilers: Pynini, HFST, a myriad of packages
 - ▶ Some of these compilers can do steps (1-3) simultaneously with little human intervention
- We'll explore doing debugging and conversion with logic-based approaches

TABLE OF CONTENTS

- INTRODUCTION
- **EXAMPLE OF LOGICAL TRANSDUCTION**
- FORMAT FOR LOGICAL TRANSDUCTIONS
- COMPILING A LOGICAL TRANSDUCTION
- WHAT'S NEXT AND OPEN PROBLEMS
- CONCLUSION

EXAMPLE LOGICAL TRANSDUCTION

- FSTs are a common computational framework for formalizing rules
- FSTs are equivalent to logical transductions, which are also explicit
- Example transduction for $p \rightarrow b/a_a$

Alphabet	Σ	=	$\{a,p,b\}$
Functions	F	=	$\text{succ}(x), \text{pred}(x)$
Predicates	$\mathbf{a_a}(x)$	$\stackrel{\text{def}}{=}$	$a(\text{pred}(x)) \wedge a(\text{succ}(x))$
Output formula	$\phi a(x^1)$	$\stackrel{\text{def}}{=}$	$a(x)$
	$\phi b(x^1)$	$\stackrel{\text{def}}{=}$	$b(x) \vee (p(x) \wedge \mathbf{a_a}(x))$
	$\phi p(x^1)$	$\stackrel{\text{def}}{=}$	$p(x) \wedge \neg \mathbf{a_a}(x)$

EXAMPLE LOGICAL TRANSDUCTION

- For this example transduction, an input string /papab/ becomes [pabab]

Input: $p_0 \xrightarrow{\triangleleft} a_0 \xrightarrow{\triangleleft} p_0 \xrightarrow{\triangleleft} a_0 \xrightarrow{\triangleleft} b_0$

Output: $p_1 \xrightarrow{\triangleleft} a_1 \xrightarrow{\triangleleft} b_1 \xrightarrow{\triangleleft} a_1 \xrightarrow{\triangleleft} b_1$

- Note the following details:
 - ▶ Input elements belong to a set of nodes indexed by 0
 - ▶ Output elements belong to a set of nodes indexed by 1
 - ▶ Successor relations are unchanged and order-preserving (postpone discussion)

FORMAT FOR LOGICAL TRANSDUCTIONS

- We converted our pen-and-paper rule to a small set of logical statements, but can we check that it's correct?

Alphabet	Σ	=	$\{a,p,b\}$
Functions	F	=	$\text{succ}(x), \text{pred}(x)$
Predicates	$\mathbf{a_a}(x)$	$\stackrel{\text{def}}{=}$	$a(\text{pred}(x)) \wedge a(\text{succ}(x))$
Output formula	$\phi a(x^1)$	$\stackrel{\text{def}}{=}$	$a(x)$
	$\phi b(x^1)$	$\stackrel{\text{def}}{=}$	$b(x) \vee (p(x) \wedge \mathbf{a_a}(x))$
	$\phi p(x^1)$	$\stackrel{\text{def}}{=}$	$p(x) \wedge \neg \mathbf{a_a}(x)$

- First step: make sure this format is explicit enough
 - But it's not because there's a lot of details in fonts/colors/superscripts
- Need to turn this format into a more explicit one

TABLE OF CONTENTS

- INTRODUCTION
- EXAMPLE OF LOGICAL TRANSDUCTION
- **FORMAT FOR LOGICAL TRANSDUCTIONS**
- COMPILING A LOGICAL TRANSDUCTION
- WHAT'S NEXT AND OPEN PROBLEMS
- CONCLUSION

FORMAT FOR LOGICAL TRANSDUCTIONS

- For FST formats, some formats are very explicit and machine-readable, e.g, an FST is a sequence of transitions:
(input state, input label, output state, output label)
- Can essentially write all of an FST's information into a simple text file

FORMAT FOR LOGICAL TRANSDUCTIONS

- For FST formats, some formats are very explicit and machine-readable, e.g., an FST is a sequence of transitions:
(input state, input label, output state, output label)
- Can essentially write all of an FST's information into a simple text file
- For logical transductions to our knowledge, there isn't a pre-existing format that's machine-readable and 100% explicit
- Thankfully, Eric Meinhardt has been making one using YAML syntax with inspiration from the BMRS notation (Bhaskar et al., 2020)

FORMAT FOR LOGICAL TRANSDUCTIONS

- Compare our logical statements to the text file **BetweenAVoicing.txt**

Alphabet	Σ	=	$\{a,p,b\}$
Functions	F	=	$succ(x), pred(x)$
Predicates	$\mathbf{a_a}(x)$	$\stackrel{\text{def}}{=}$	$a(pred(x)) \wedge a(succ(x))$
Output formula	$\phi a(x^1)$	$\stackrel{\text{def}}{=}$	$a(x)$
	$\phi b(x^1)$	$\stackrel{\text{def}}{=}$	$b(x) \vee (p(x) \wedge \mathbf{a_a}(x))$
	$\phi p(x^1)$	$\stackrel{\text{def}}{=}$	$p(x) \wedge \neg \mathbf{a_a}(x)$

FORMAT FOR LOGICAL TRANSDUCTIONS

- Compare our logical statements to the text file **BetweenAVoicing.txt**

Alphabet	Σ	=	$\{a,p,b\}$
Functions	F	=	$\text{succ}(x), \text{pred}(x)$
Predicates	$\mathbf{a_a}(x)$	$\stackrel{\text{def}}{=}$	$a(\text{pred}(x)) \wedge a(\text{succ}(x))$
Output formula	$\phi a(x^1)$	$\stackrel{\text{def}}{=}$	$a(x)$
	$\phi b(x^1)$	$\stackrel{\text{def}}{=}$	$b(x) \vee (p(x) \wedge \mathbf{a_a}(x))$
	$\phi p(x^1)$	$\stackrel{\text{def}}{=}$	$p(x) \wedge \neg \mathbf{a_a}(x)$

- Some differences:
 - ▶ Specify alphabets, copy set size, and word boundaries
 - ▶ Specify predicates vs. output functions
 - ▶ Clarify if a node (and its labels) belongs to input vs. a copy via underscore
 - ▶ No support for negation or Boolean operations
 - ▶ Logic uses just **IF CONDITION RESULT ELSE** statements
- such formats simulate a BMRS system = rational function

TABLE OF CONTENTS

- INTRODUCTION
- EXAMPLE OF LOGICAL TRANSDUCTION
- FORMAT FOR LOGICAL TRANSDUCTIONS
- **COMPILING A LOGICAL TRANSDUCTION**
- WHAT'S NEXT AND OPEN PROBLEMS
- CONCLUSION

COMPILING A LOGICAL TRANSDUCTION

- Given an explicit format for a logical transduction, next step is to compile it
- To my knowledge, there isn't a pre-existing compiler for BMRS-style logical transductions
- I made my own with Python that...
 1. Takes a text file with YAML syntax
 2. Converts the file to a Pythonic translation
 3. Uses dynamic programming (memoization) to run the translation on an input word

- Everything can work via the command line
 1. Download the repo
 2. On the commandline or terminal, type...

```
python3 bmrs.py run BetweenAVoicing.txt 'papab'
```
- The code will convert the **BetweenAVoicing.txt** file to a Python file **BetweenAVoicing.py** with all the logical info
- The file is run on the input word /papab/ to generate [pabab]

- Other commands you can run are...
 1. Convert the YAML file to Python
python3 bmrs.py convert BetweenAVoicing.txt
or
python3 bmrs.py convert BetweenAVoicing.txt BetweenAVoicing.py
 2. Run the compiler for the word /papab/ using either the text file or Python file
python3 bmrs.py run BetweenAVoicing.txt 'papab'
or
python3 bmrs.py run BetweenAVoicing.py 'papab'
- Conversion errors and running errors are logged
- You can see the input-output graphs as a TSV file (nicer to look on Excel)

TABLE OF CONTENTS

- INTRODUCTION
- EXAMPLE OF LOGICAL TRANSDUCTION
- FORMAT FOR LOGICAL TRANSDUCTIONS
- COMPILING A LOGICAL TRANSDUCTION
- **WHAT'S NEXT AND OPEN PROBLEMS**
- CONCLUSION

WHAT YOU CAN DO

- The YAML format can encode any BMRS-style system that generates a rational function
- The Python compiler can thus also run any rational function
- If you want to debug your pen-and-paper analysis or your logical statements, you can convert your analysis to a text file and then run it with the Python software
- Bugs will likely be found
- Repo includes example YAML files for inserting nodes, deleting nodes, and repeating nodes

WHAT YOU CAN'T DO (EVERYWHERE YET)

- **Features on YAML:**

- ▶ As of now, the YAML format cannot directly encode features or natural classes.
- ▶ If you want to make a transduction for intervocalic voicing, you'll need to create long predicates that pick out every vowel and consonant
- ▶ Hopefully this will be fixed soon.

WHAT YOU CAN'T DO (EVERYWHERE YET)

- **Features on YAML:**

- ▶ As of now, the YAML format cannot directly encode features or natural classes.
- ▶ If you want to make a transduction for intervocalic voicing, you'll need to create long predicates that pick out every vowel and consonant
- ▶ Hopefully this will be fixed soon.

- **Features on Python**

- ▶ The Python translation and compiler however CAN do this.
- ▶ So you'll need to manually create your own Python file and use features there.
- ▶ Check out the file **IntervocalicVoicing.py** to illustrate

WHAT YOU CAN'T DO (EVERYWHERE YET)

- **More expressive functions on YAML:**
 - ▶ The YAML format can encode any rational function
 - ▶ But it can't encode anything more expressive

WHAT YOU CAN'T DO (EVERYWHERE YET)

- **More expressive functions on YAML:**

- ▶ The YAML format can encode any rational function
- ▶ But it can't encode anything more expressive

- **More expressive functions on Python**

- ▶ The Python translation and compiler however can encode some non-rational but regular functions
- ▶ You can write functions that concatenate two rational functions: $X \rightarrow XX$ (=the **Copying.py** file)
- ▶ You can write a non-rational function with user-specified succ/pred relations like $X \rightarrow XX^r$ can encode a non-rational function with
- ▶ But you can't write a regular function that needs quantifiers to find succ/pred values. For example, a function $u1-u2-...-un \rightarrow un...u2u1$ where ux is a substring

WHAT YOU CAN'T DO (EVERYWHERE YET)

- **More expressive functions on YAML:**

- ▶ The YAML format can encode any rational function
- ▶ But it can't encode anything more expressive

- **More expressive functions on Python**

- ▶ The Python translation and compiler however can encode some non-rational but regular functions
- ▶ You can write functions that concatenate two rational functions: $X \rightarrow XX$ (=the **Copying.py** file)
- ▶ You can write a non-rational function with user-specified succ/pred relations like $X \rightarrow XX^r$ can encode a non-rational function with
- ▶ But you can't write a regular function that needs quantifiers to find succ/pred values. For example, a function $u1-u2-...-un \rightarrow un...u2u1$ where ux is a substring

- The main problem is that the BMRS format requires that functions are order-preserving. The YAML file follows this. But I've been thinking of ways to enrich BMRS notation to allow some types of non-order-preserving functions

- There's been a lot of new datasets that have been converted to logical transductions and BMRS notation (Strother-Garcia, 2019; Dolatian, 2020; Chandlee and Jardine, 2021)
 - convert all them to YAML to help find bugs in the YAML format
 - compile into Python to ensure correctness

OPEN PROBLEMS

- There's been a lot of new datasets that have been converted to logical transductions and BMRS notation (Strother-Garcia, 2019; Dolatian, 2020; Chandlee and Jardine, 2021)
 - convert all them to YAML to help find bugs in the YAML format
 - compile into Python to ensure correctness
- BMRS allows operations such as composition and others (Oakden, 2021)
 - create algorithms for doing these operations over YAML formats
 - create algorithms to detect interaction across rules in YAML

OPEN PROBLEMS

- BRMS and its implementations (YAML format, Python compiler) currently only handle strings...
- convert model-theoretic notation for other representations (trees, tiers, gestures, autosegmental graphs) to YAML-friendly syntax
- generalize the Python compiler
- BMRS and its implementations (YAML format, Python compiler) currently are just transducers.
- would be nice to also handle constraint-checking as an automaton.

OPEN PROBLEMS

- BRMS and its implementations (YAML format, Python compiler) currently only handle strings...
- convert model-theoretic notation for other representations (trees, tiers, gestures, autosegmental graphs) to YAML-friendly syntax
- generalize the Python compiler
- BMRS and its implementations (YAML format, Python compiler) currently are just transducers.
- would be nice to also handle constraint-checking as an automaton.
- Spoiler: I'll probably never do these (by myself), but feel free to pitch an idea for collaboration :)

TABLE OF CONTENTS

- INTRODUCTION
- EXAMPLE OF LOGICAL TRANSDUCTION
- FORMAT FOR LOGICAL TRANSDUCTIONS
- COMPILING A LOGICAL TRANSDUCTION
- WHAT'S NEXT AND OPEN PROBLEMS
- CONCLUSION

CONCLUSION

- It's good practice to double-check your analysis via computational implementations and simulations
- Can now do that with logical transductions
- Implementations open doors for yourself to contribute in this field, whether as a linguist or computer scientist
- Honestly, don't be shy if you're interested in one of the open problems :)

REFERENCES

- Bhaskar, S., J. Chandlee, A. Jardine, and C. Oakden (2020). Boolean monadic recursive schemes as a logical characterization of the subsequential functions. In A. Leporati, C. Martín-Vide, D. Shapira, and C. Zandron (Eds.), Proceedings of the 14th International Conference on Language and Automata Theory and Applications (LATA 2020).
- Chandlee, J. and A. Jardine (2021). Computational universals in linguistic theory: Using recursive programs for phonological analysis. Language 97(3), 485–519.
- Dolatian, H. (2020). Computational locality of cyclic phonology in Armenian. Ph. D. thesis, Stony Brook University.
- Oakden, C. D. (2021). Modeling phonological interactions using recursive schemes. Ph. D. thesis, Rutgers The State University of New Jersey, School of Graduate Studies.
- Strother-Garcia, K. (2019). Using model theory in phonology: a novel characterization of syllable structure and syllabification. Ph. D. thesis, University of Delaware.