

# Analysis of historical sound change via the *Index Diachronica*

---

Wilson Biggs

April 27, 2023

# Introduction

---

# What is the *Index Diachronica*?

- A hand-compiled reference index of historical sound changes, intended for use in creating constructed language families
- Originally a collection of posts in a thread on the *Zompist Bulletin Board*; compiled into one index by Galen Buttitta
- “The Index is not an academic source, nor is it perfect – far from it.”
  - Notation is only *mostly* consistent
  - A fair number of errors and dubious inclusions
  - Can still draw insights, but important to keep the caveats in mind
- Adapted from PDF to a more convenient and easily searchable HTML version, the *Searchable Index Diachronica*, by chri d. d.

# What is the *Index Diachronica*?

## **17.7.2.1.1 Anglo-Frisian to Old English**

*Pogostick Man*, from Wikipedia contributors (2011), “Phonological history of English”. *Wikipedia, the Free Encyclopedia*. <[http://en.wikipedia.org/w/index.php?title=Phonological\\_history\\_of\\_English&oldid=453796112](http://en.wikipedia.org/w/index.php?title=Phonological_history_of_English&oldid=453796112)>

ã: → ð:

V[+nas] → V[-nas]

{i,u} → Ø/ \_# ! V[-long]C\_#

k ʏ g → tʃ j dʒ / “in certain complex circumstances”

# What were my goals?

- What are the most common sound changes?
- How are sound changes influenced by their environments?
  - More specifically: how do neighboring consonants affect vowel changes?

## Parsing the data

---

# Parsing the data

- The majority of my work
- Used a separate parsing script
- Lots of trial and error, running the parsing script over and over and fixing issues until nothing broke
- Had to manually edit the data a bit to remove things that would have made it impossible to parse
- Lots of complexity – my explanation of the steps will be leaving out some of the edge cases I had to handle



# Pulling the content from HTML

The first step was to pull sound changes from the HTML and clean up any extra HTML tags. This was done using **Beautiful Soup**. The data I wanted was the name of the branch, the branches' and sounds' IDs, and the full text of the sound change rule.

```
index-diachronia-analysis [WSL: Ubuntu] - sid-tidy-with-  
edits.html  
  
8759 <section class="showtarget" id="Old-English">  
8760 <h2>17.7.2.1.1 Anglo-Frisian to Old English</h2>  
8761 <p><i>Pogostick Man</i>, from Wikipedia contributors  
(2011), "Phonological history of English". <i>Wikipedia,  
the Free Encyclopedia</i> [ ... ]</p>  
8762 <p class="schg" id="Old-English-ð">ð: ð ð: </p>  
8763 <p class="schg" id="Old-English-V+nas">V[+nas] ð V[-n  
as]</p>  
8764 <p class="schg" id="Old-English-i,u">{i,u} ð Ø / _# !  
V[-long]C_#</p>  
8765 <p class="schg" id="Old-English-k-y-q">k y q ð tʃ d dʒ  
/ "in certain complex circumstances"</p>  
8766 </section>
```



# Pulling out the environment

$\{i,u\} \rightarrow \emptyset / \_ \# ! V[-\text{long}]C\_ \#$



$\_ \# ! V[-\text{long}]C\_ \#$

index-diachronica-analysis [WSL: Ubuntu] - data\_parsing\_script.py

```
151 # Split the rule up. This will inevitably include stuff I don't want but we can work out how to remove that stuff later
152 # First split by the environment separator
153 env_split = rule_string.split(" / ", 1)
154
155 environment = ''
156
157 if len(env_split) > 1:
158     environment = env_split[1]
159 else:
160     # If no environment, but rule ends with some text in parentheses or quotes, consider that the environment
161     parens_match = re.search(r'(.+[^+]) (\(.(+)\)|"(.+)"$)', rule_string)
162     if parens_match:
163         env_split[0] = parens_match.group(1)
164         environment = parens_match.group(2)
165
166 split_rules: list[Tuple[str, list[str], str]] = []
```

# Parsing “optionals”

$e(:)j \rightarrow i$



$ej \rightarrow i$

$e:j \rightarrow i$

```
index-diachronica-analysis [WSL: Ubuntu] - data_parsing_script.py
168 # If there are any optional bits, run the split with all possible combinations of with and without them
169 optionals = [(match.start(), match.end(), match.group(0)) for match in re.finditer(r'(\.?*\?)', env_split[0])]
170 if (optionals):
171     combinations = list(itertools.chain.from_iterable(itertools.combinations(optionals, l) for l in range(len(optionals) + 1)))
172     for combo in combinations:
173         combo_string = env_split[0]
174         # print(combo)
175         combo_string = remove_combos(combo_string, combo)
176         combo_string = combo_string.replace('(','').replace(')','')
177         # print(combo_string)
178         split_rules += parse_rule_steps(combo_string)
179 else:
180     split_rules += parse_rule_steps(env_split[0])
181
182 # only uniques
183 for split_rule in [sr for i, sr in enumerate(split_rules) if sr not in split_rules[:i]]:
184     rule = Rule()
185
186     rule.id = rule_id
187     rule.branch_id = branch.id if branch else None
188     rule.branch_index = branch.index if branch else None
189     rule.original_text = decoded
190     rule.environment = environment
191     rule.from_sound, rule.intermediate_steps, rule.to_sound = split_rule
192
193     rules.append(rule)
```

# Parsing steps and multi-sound rules

a: ɔ: → ε: o: → e: ow → ej əw



from\_sound = "a:", intermediate\_steps=["ε:", "e:"], to\_sound = "ej"

from\_sound = "ɔ:", intermediate\_steps=["o:", "ow"], to\_sound = "əw"

```
index character analysis (RWS, Ubuntu) : data.parsing_script.py
17 def parse_rule_steps(steps): return list(Tuple(str, list(str), str))
18 """Parse out the steps of a rule"""
19 rules: list(Tuple(str, list(str), str)) = []
20
21 steps = re.sub(r'[- '], '', steps)
22
23 # split by the rule step separator
24 rule_split = (s.strip() for s in steps.split("-"))
25
26 if len(rule_split) < 2:
27     print("Too few steps in rule: {}".format(steps))
28     return []
29
30 if len(rule_split) > 2:
31     from_sounds = split_sounds(rule_split[0])
32     intermediate_steps = [split_sounds(step) for step in rule_split[1:-1]]
33     to_sounds = split_sounds(rule_split[-1])
34     if not (
35         len(from_sounds) == len(intermediate_steps[0]) == len(to_sounds)
36         and all(len(i) == len(intermediate_steps[0]) for i in intermediate_steps)
37     ):
38         print("Warning: mismatched lengths for rule: {}".format(steps))
39         return []
40
41     for index, from_sound in enumerate(from_sounds):
42         unbracketed_from = handle_brackets(from_sound)
43         unbracketed_to = handle_brackets(to_sounds[index])
44         for sub_from in unbracketed_from:
45             for sub_to in unbracketed_to:
46                 rules.append((unbracketed_from, [index], unbracketed_to, sub_to))
47
48     else:
49         from_sounds = split_sounds(rule_split[0])
50         to_sounds = split_sounds(rule_split[-1])
51         if len(from_sounds) != len(to_sounds):
52             print("Warning: mismatched lengths for rule: {}".format(steps))
53             return []
54
55         for index, from_sound in enumerate(from_sounds):
56             unbracketed_from = handle_brackets(from_sound)
57             unbracketed_to = handle_brackets(to_sounds[index])
58             for sub_from in unbracketed_from:
59                 for sub_to in unbracketed_to:
60                     rules.append((unbracketed_from, [], sub_to))
61
62 return rules
```

# Parsing brackets

$\{i,u\} \rightarrow \emptyset$



from\_sound = "i", to\_sound = "∅"

from\_sound = "u", to\_sound = "∅"

```
index@lachryna-analysis [WSL: Ubuntu] - data_parsing_script.py
65 def handle_brackets(sound: str) → set[str]:
66     """Handles bracketed sounds"""
67     sounds: set[str] = set()
68
69     # Find innermost bracket and recurse
70     bracket_start: int = 0
71     bracket_end: int = 0
72
73     # Get's do bracket matching. Ign
74     for (index, character) in enumerate(sound):
75         if character == '{':
76             bracket_start = index
77         elif character == '}':
78             bracket_end = index
79             break
80
81     if (bracket_start == 0 and bracket_end == 0):
82         return set([sound])
83
84     if (bracket_end == 0):
85         bracket_end = len(sound)
86
87     # print(sound[bracket_start : bracket_end+1])
88
89     inner_sounds = sound[bracket_start + 1 : bracket_end].split(',')
90
91     for inner_sound in inner_sounds:
92         repl_sound = sound[:bracket_start] + inner_sound + sound[bracket_end + 1:]
93         sounds |= handle_brackets(repl_sound)
94
95     return sounds
```

What are the most common  
sound changes?

---

# Questions to ask first

- What does it mean for two sound changes to be ‘the same’?
  - What if they have the same ‘from’ and ‘to’ sound, but different environments or intermediate steps?
  - I looked at both options - considering these the same, and considering these different.
  - Not much new information was gained by including environments or intermediate steps – the results were just the ‘most common sound changes with no environment or intermediate steps’
- Lots of sound changes are copied between different daughter languages in a single branch. Let’s count those as just 1 sound change.

# The most common sound changes

Looking at only the **from\_sound** and **to\_sound**...

from_sound	to_sound	count
h	Ø	46
w	Ø	44
ʔ	Ø	40
k	Ø	37
j	Ø	34
e	i	34
a	e	32
i	e	30
u	o	29
o	u	29

- So many deletions! Makes sense, since what can be deleted is less limited than 'what can become /i/'
- /h/ being removed is most common, which makes some sense – i.e. British English “history”
- Nothing really surprising here

# The most common sound changes

What if we ignore deleted sounds?

from_sound	to_sound	count
e	i	34
a	e	32
i	e	30
u	o	29
o	u	29
ts	s	27
s	ʃ	26
g	k	24
k	g	23
a	o	23

- Still nothing too surprising here
- Mostly just simple vowel changes
- /ts/ → /s/ is the most interesting thing here



How neighboring consonants  
affect vowel changes

---

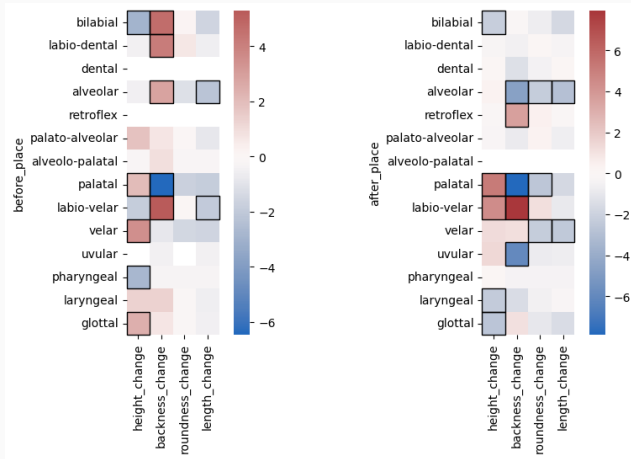
## More data parsing!

- Filter down my rules to just vowel changes
- Pull neighboring consonants from the environment and/or rule
  - Parse environments – consonants before/after the underscore
  - e.g. `n_mV#` → `before = "n"`, `after = "m"`
- Used **Gruut IPA** to split strings of IPA characters into individual phones (`/tʃu:z/` → `/tʃ/` + `/u:/` + `/z/`)
- Used **ipapy** to break down phones into their component features:
  - Consonants: voicing, place, manner, modifiers (palatalized, etc.)
  - Vowels: length, height, backness, roundness, modifiers (centralized, etc.)

# How I performed the analysis

- Assigned a number to the values of each vowel feature (0 = open, 3 = mid, 6 = close) and found the difference between start and end sound to get that feature's "change"
- Looked at preceding and following consonants separately
- Did statistical t-tests to determine if the 'average changes' for different consonant features were significantly different from each other, and in what direction they differed
- Visualized these results as a heatmap
- Outlined squares signify statistical significance
- Excluded duplicate sound changes shared by daughter languages within the same branch, like before

# Place of articulation



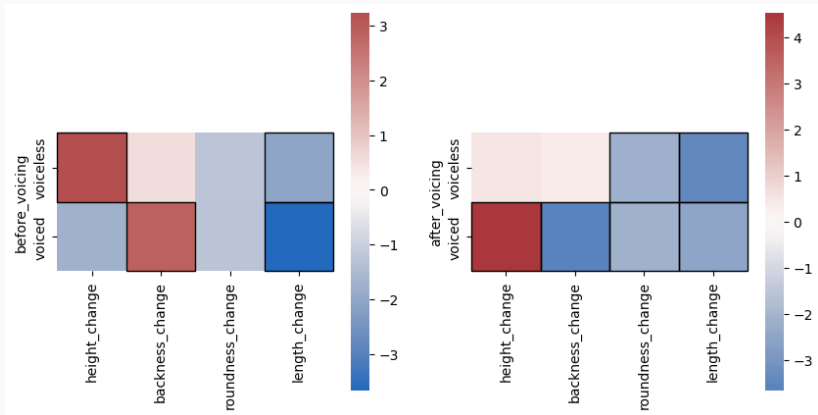
Height: negative = lowering (/u/ → /o/), positive = raising (/o/ → /u/)

Backness: negative = fronting (/u/ → /i/), positive = backing (/i/ → /u/)

Roundness: negative = unrounding, positive = rounding

Length: negative = shortening, positive = lengthening

# Voicing



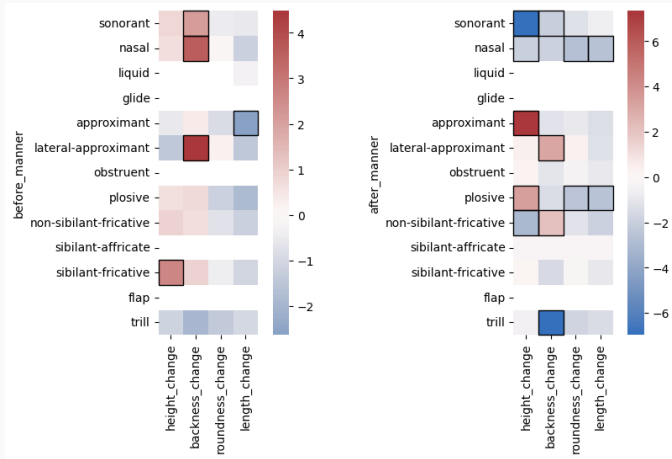
Height: negative = lowering (/u/ → /o/), positive = raising (/o/ → /u/)

Backness: negative = fronting (/u/ → /i/), positive = backing (/i/ → /u/)

Roundness: negative = unrounding, positive = rounding

Length: negative = shortening, positive = lengthening

# Manner



Height: negative = lowering (/u/ → /o/), positive = raising (/o/ → /u/)

Backness: negative = fronting (/u/ → /i/), positive = backing (/i/ → /u/)

Roundness: negative = unrounding, positive = rounding

Length: negative = shortening, positive = lengthening

# Takeaways

- Some results lined up with what I expected (palatals are associated with fronting, /w/ is associated with backing) but I couldn't explain most results
- Each combination of features was not evenly distributed, so they might be contaminating each other's results
- I tried using linear regression to try and untangle this, but those results were even messier - not included here

# Conclusion

---



## Future ideas

- Untangle the effects of these different variables to try and determine which features are actually having which effects
- Do the same analysis, but for consonant changes
- Fix some of the errors in my data – there is a big PDF of corrections to the *Index Diachronica* that someone has compiled
- Re-write the *Index* using more standard notation, like PhoMo (a sound-change notation intended to be parsed by software and used to automatically apply sound changes to words)
- Interactive tools – e.g. put in a word or IPA string, get the most likely ways that word could evolve

Questions?

Questions?