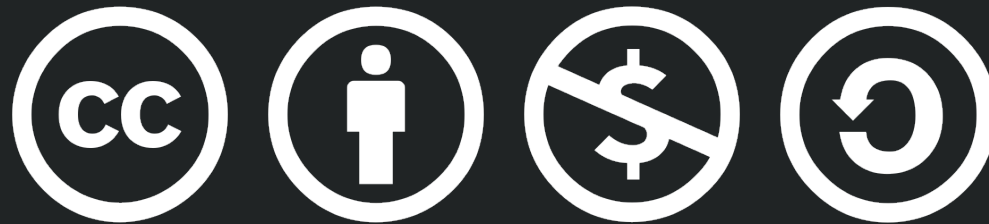


PyCon Canada 2016 Presentation

Quantifying the visual structure of written language

by Nick Anderegg



This work is licensed under the
Creative Commons Attribution-NonCommercial-ShareAlike 4.0
International License.

To view a copy of this license, visit
<http://creativecommons.org/licenses/by-nc-sa/4.0/>.

Quantifying the visual structure of written language

Nick Anderegg

Outline

- Research background
- My research
- Computational motivation
- Algorithm
- Take-aways

Background



Why study reading?

- Understand reading disabilities
- Improve natural language processing
- It's interesting

English reading

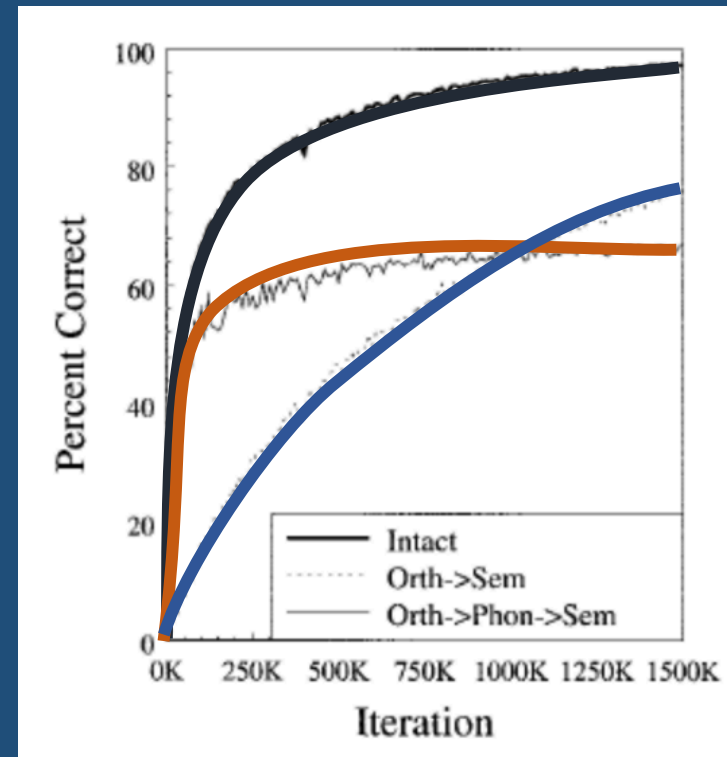
- Cooperative division of labor
- Depends on ease of computation on each route
- Neither route is perfect

Orthography

≈ The written form
of a language

Phonology

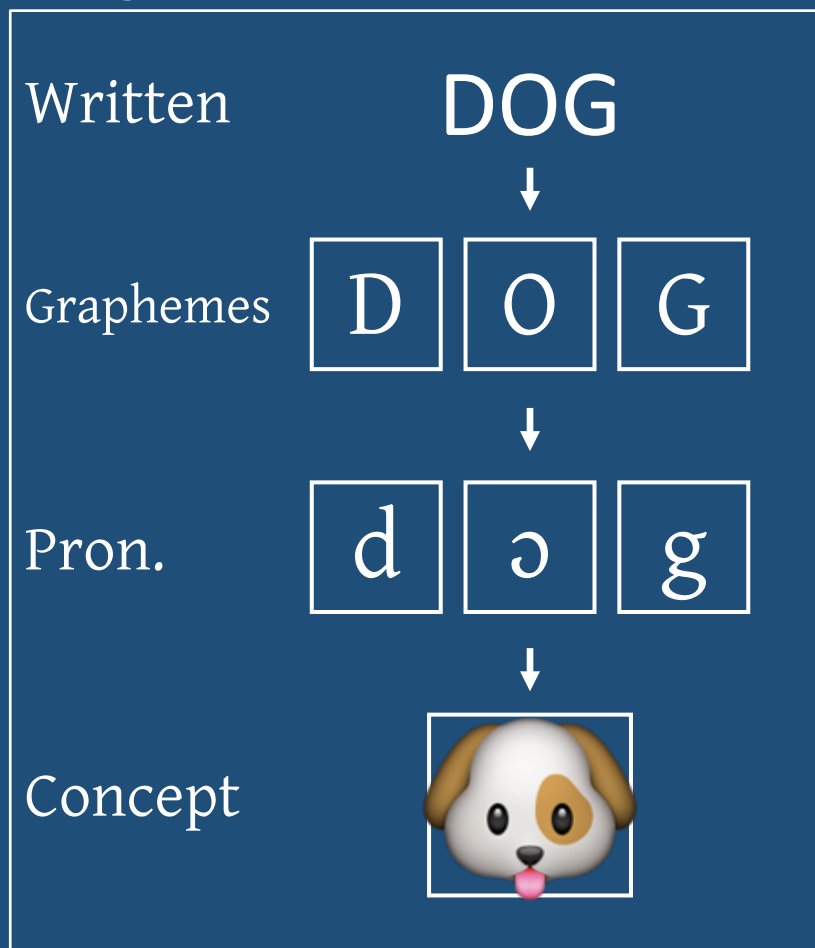
≈ The sound system
of a language



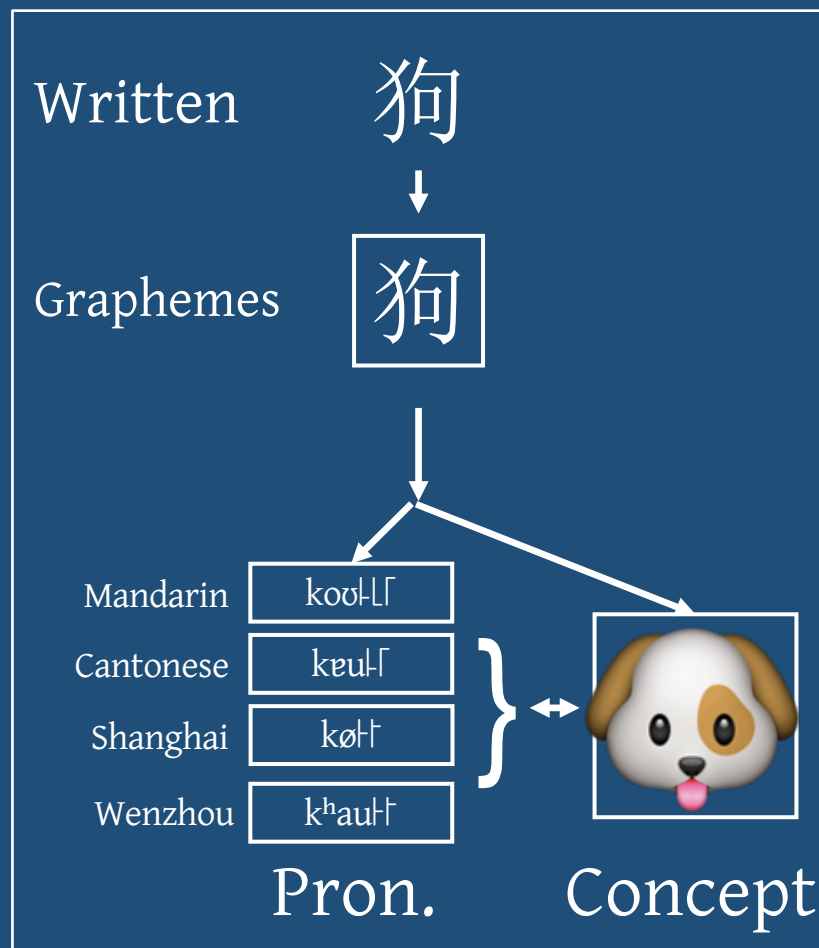
Harm and Seidenberg (2004)

Finding a universal model

English



Chinese



Encoded information: Phono-semantic compounds

- Two indicators/radicals
- Sub-character phonology
 - Imperfect
 - Inconsistent



mō “to feel”

扌 hand

莫 mò



mò “sword”

钅 metal

莫 mò

Why compare these?

English

- Computation of phonology from written form
- Does phonology mediate meaning?

Chinese

- Mixed evidence for sub-character phonology
- Directly encodes meaning information

Both languages

- Relative activation of semantics by written form and phonological form
- “Division of labor”

Research



Current research

- Different methods find different processing times
 - Debate over the use of the retrieved phonological information
- Examine the specific role of phonology in word recognition
- Functionally isolate routes
 - Agreement about phonological activation...
 - ...but disagreement about its role

+

The

boy eat

house ran

The

boy

house

very

far

jump

school

eat

ran

car

blue

after

read

SHE

###

RAN

DOG

BED

FOR

SEVEN JUMP

MILES POUNDS

Critical manipulation

- Orthographically similar / Phonological dissimilar

水 永
shuǐ yǒng

English illustration:

beach bench

Critical manipulation

- Orthographically similar / Phonological dissimilar

水 永
shuǐ yǒng

English illustration:

beach

The

###

man

sit

red

dug

in

of

the

out

sail

soil

Motivation



Quantifying Visual Similarity

- A few different methods used in past
- Most common:
 - Same radical
 - Same number of strokes

An imperfect method

Radical: 立

Characters containing this radical (w/ 9 strokes):

亲 竝 𠂔
竝 飒 竖

An imperfect method

Radical: 立

Good pairs:

竝 - 颯

竝 - 竝

Bad pairs:

亲 - 竝

竝 - 竖

The image shows the Japanese character 食 (kashu), which means 'eat' or 'food'. It is a white character on a dark blue background.

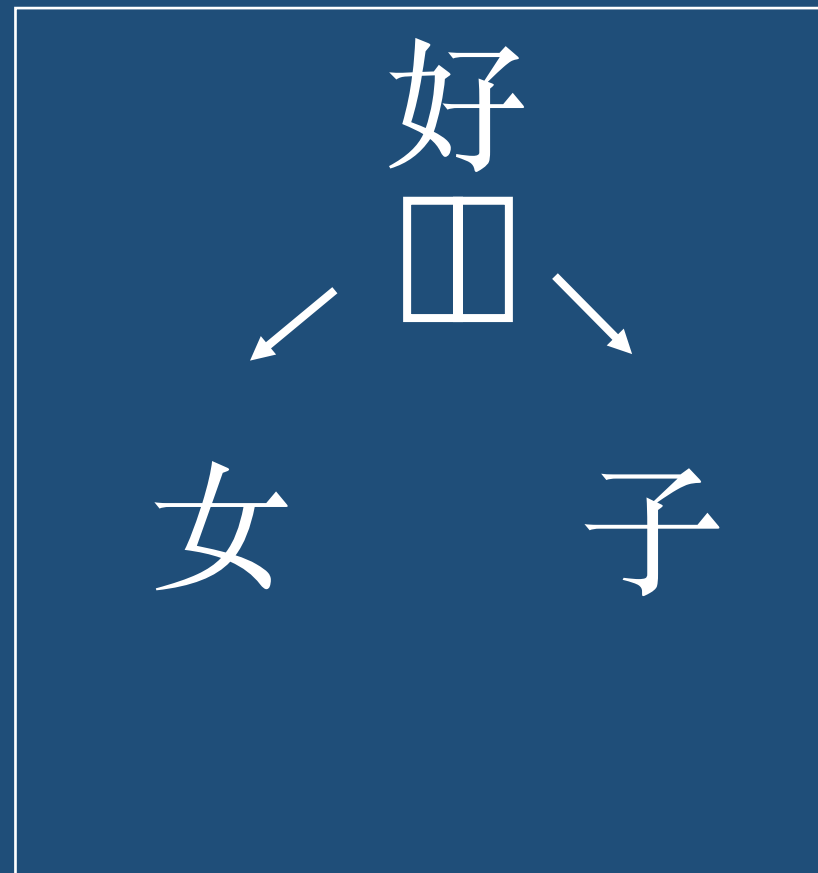
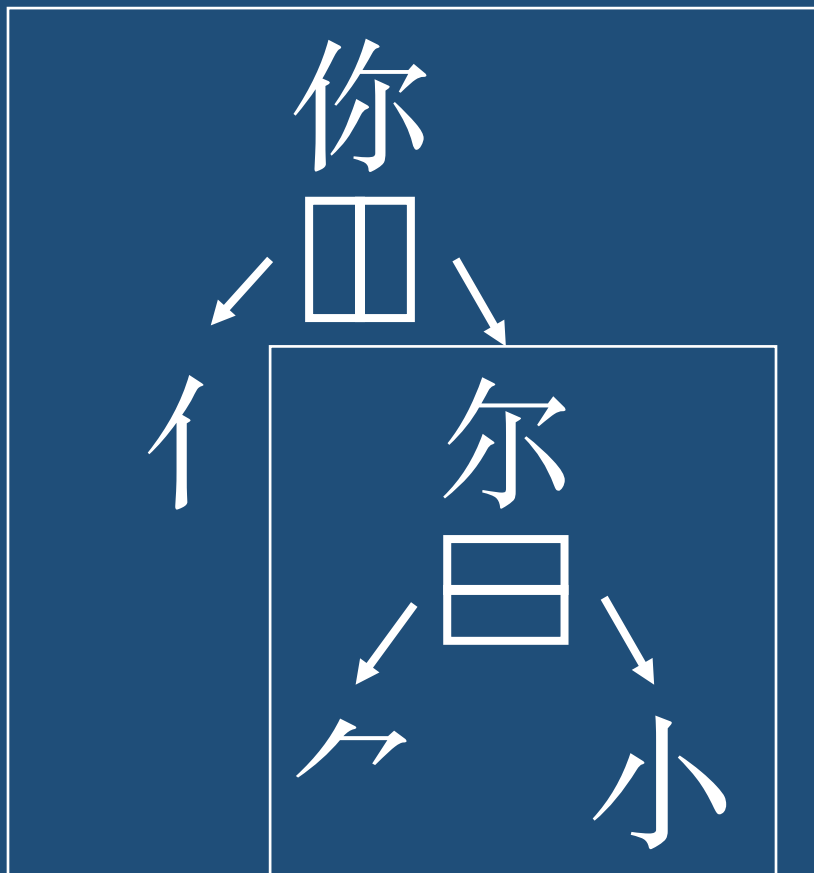
9 strokes

The image shows the English word 'EAT' in a white, sans-serif font on a dark blue background.

9 strokes???

Hierarchical structure

你好!

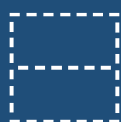


Ideographic Description Characters

- Unicode range U+2FF0 to U+2FFB
- Intended to act as a rough description of characters



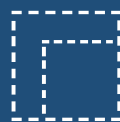
left to right

surround from
below

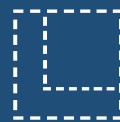
above to below



surround from left

left to middle and
rightsurround from
upper leftabove to middle
and belowsurround from
upper right

full surround

surround from
lower leftsurround from
above

overlaid

Ideographic Description Language

你 →  イ 尔

尔 →  ㄣ 小

兹 →  一  么 么

Algorithm



Parsing IDS BNF

```
char      :: all valid CJK
binary    :: 2FF0 | 2FF1 | 2FF4..2FFB
ternary   :: 2FF2 | 2FF3
sub       :: tergroup|bigroup|char
bigroup   :: binary sub sub
tergroup  :: ternary sub sub sub
expr      :: [ sub ]+
```

```
>>> import idsparser as ids
>>> ids.parse('丽', '㊦—㊦㊦㊦㊦㊦㊦㊦㊦')
['丽', ['㊦', ['—', ['㊦', [['㊦', ['㊦', ['\']],
['㊦', ['㊦', ['\']]]]]]]]]]
>>> ids.pretty_parse('丽', '㊦—㊦㊦㊦㊦㊦㊦㊦㊦')
```

丽

```
L---㊦
    L---—
        L---㊦
            L---㊦
                L---㊦
                    L---㊦
                        L---\
                            L---㊦
                                L---㊦
                                    L---\
```

Similarity algorithm

- Tree generation
- Structural weighting
- Root comparison
- Comparison reuse

Tree generation

```
class IDSTree():

def __init__(self, parse_tree, dictionary=None):
    self.head = parse_tree[0]
    self.ids = idsparser.unparse(parse_tree[1])
    self.tree = IDSNode(parse_tree[1], self._dict)

    self._dict = dictionary
    if self._dict:
        self._dict._nodes[self.head] = self.tree
        self._dict._nodes[self.ids] = self.tree
```

Tree generation

```
class IDSNode():  
  
def __init__(self, parse_tree, dictionary=None):  
    self.head = None  
    self.ids = None  
    self.children = None  
  
    self._dict = dictionary
```


Tree generation

- If head of parse tree is IDS functor:
 - Assign functor to `self.head`
 - Assign flat sequence to `self.ids`
 - Assign tree children to `self.children` list as recursive *IDSNodes*
- If head of parse tree is not functor:
 - Leaf has been reached
 - Assign character to `self.head`
 - Assign character to `self.ids`

Structural weighting

```
self.weights = {  
    '□': (1/3, 2/3),  
    '▣': (1/3, 2/3),  
    '▤': (1/3, 2/3),  
    '▥': (1/3, 2/3),  
    '▦': (1/3, 2/3),  
    '▧': (1/3, 2/3),  
    '▨': (1/3, 2/3),  
    '▩': (1/2, 1/2)  
}  
  
for x in functors.binary:  
    self.weights[x] = (1/2, 1/2)  
for x in functors.ternary:  
    self.weights[x] = (1/3, 1/3, 1/3)
```

Root comparison

```
def compare_nodes(a, b):  
    If a is b... #same tree  
    If a.head == b.head...  
        If a.head not in functors  
            and b.head not in functors  
            Else recursive comparison  
    If a.head != b.head...
```

Comparison reuse

```
class IDSTree():  
  
    def __init__(self, parse_tree, dictionary=None):  
        self.head = parse_tree[0]  
        self.ids = idsparser.unparse(parse_tree[1])  
        self.tree = IDSNode(parse_tree[1], self._dict)  
  
        self._dict = dictionary  
        if self._dict:  
            self._dict._nodes[self.head] = self.tree  
            self._dict._nodes[self.ids] = self.tree
```

Tweaking the algorithm

- Pure structural weighting resulting in many leaves with 0 similarity
- Radical similarity account for 1/3 of similarity
- Stroke proportion accounts for 1/3 of similarity

Examples

Before change:

水 永

0%

After change:

水 永

75%

Distributed processing

- Many individual cloud servers
- Elasticsearch database storing IDS comparisons
- 2000 chars = 1999000 comparisons
- 5000 chars = 12497500 comparisons
- At 10ms per comparison, 34 hours

Creating critical sets

- Pandas dataframe
- Sort by ortho sim, phono sim

Critical:	治	zhi4	Aggregate: 10.402752739414652
Both sim:	沼	zhao3	Score: 78.4634
Ortho sim:	沿	yan2	Score: 78.4203
Phono sim:	置	zhi4	Score: 96.898
Both diff:	淹	yan1	Score: 78.0193

Thoughts



Python in Academia

- Why Python?
 - Easy to experiment
 - Lots of flexibility
 - Plan of attack isn't always clear when exploring a new idea

Python in Academia

 `aggregateseqs.py` `allids.csv` `computesyllables.py` `context_judgement.py` `generate_quartets.py` `global_quartets.py` `handata.py` `maketrees.py` `pretty_quartets.py` `pretty_stim_set.csv` `distances.csv` `downloadhsk.py` `processpatphon.py` `processradicals.py` `read_sqlitetemp.py` `reformatsubtlex.py` `write_db.py` `write_sqlitetemp.py` `write_wordortho.py`

Why do things programmatically?

- Reproducibility!
- A digital trail exists for all tools in the experiment design pipeline
- Flexibility!
- A small change in the pipeline doesn't require starting from scratch

Take away message

- Python can be used to explore new solutions to old problems
- Easily enables experimentation
- Tools exist for every sub-problem!