

# Accessing and Analyzing Linguistic Field Data

Steven Bird

University of Melbourne, AUSTRALIA

July 9, 2006

# Introduction

- variety of data types, e.g.
- **lexicon**: database of words, minimally containing part of speech information and glosses
- **paradigm**: any kind of rational tabulation of words or phrases to illustrate contrasts and systematic variation
- **text**: any larger unit such as a narrative or a conversation
- descriptions: field notes, grammars and analytical papers
- complex inter-relations (new word; lexicon update, paradigm; add to field notes; extend grammar, ...)
- tasks: sorting, tabulating, garnering statistics, searching, verifying transcriptions
- principle challenge is computational

# Introduction

- variety of data types, e.g.
- **lexicon**: database of words, minimally containing part of speech information and glosses
- **paradigm**: any kind of rational tabulation of words or phrases to illustrate contrasts and systematic variation
- **text**: any larger unit such as a narrative or a conversation
- descriptions: field notes, grammars and analytical papers
- complex inter-relations (new word; lexicon update, paradigm; add to field notes; extend grammar, ...)
- tasks: sorting, tabulating, garnering statistics, searching, verifying transcriptions
- principle challenge is computational

# Introduction

- variety of data types, e.g.
- **lexicon:** database of words, minimally containing part of speech information and glosses
- **paradigm:** any kind of rational tabulation of words or phrases to illustrate contrasts and systematic variation
- **text:** any larger unit such as a narrative or a conversation
- descriptions: field notes, grammars and analytical papers
- complex inter-relations (new word; lexicon update, paradigm; add to field notes; extend grammar, ...)
- tasks: sorting, tabulating, garnering statistics, searching, verifying transcriptions
- principle challenge is computational

# Introduction

- variety of data types, e.g.
- **lexicon:** database of words, minimally containing part of speech information and glosses
- **paradigm:** any kind of rational tabulation of words or phrases to illustrate contrasts and systematic variation
- **text:** any larger unit such as a narrative or a conversation
- descriptions: field notes, grammars and analytical papers
- complex inter-relations (new word; lexicon update, paradigm; add to field notes; extend grammar, ...)
- tasks: sorting, tabulating, garnering statistics, searching, verifying transcriptions
- principle challenge is computational

# Introduction

- variety of data types, e.g.
- **lexicon:** database of words, minimally containing part of speech information and glosses
- **paradigm:** any kind of rational tabulation of words or phrases to illustrate contrasts and systematic variation
- **text:** any larger unit such as a narrative or a conversation
- descriptions: field notes, grammars and analytical papers
- complex inter-relations (new word; lexicon update, paradigm; add to field notes; extend grammar, ...)
- tasks: sorting, tabulating, garnering statistics, searching, verifying transcriptions
- principle challenge is computational

# Introduction

- variety of data types, e.g.
- **lexicon:** database of words, minimally containing part of speech information and glosses
- **paradigm:** any kind of rational tabulation of words or phrases to illustrate contrasts and systematic variation
- **text:** any larger unit such as a narrative or a conversation
- descriptions: field notes, grammars and analytical papers
- complex inter-relations (new word; lexicon update, paradigm; add to field notes; extend grammar, ...)
- tasks: sorting, tabulating, garnering statistics, searching, verifying transcriptions
- principle challenge is computational

# Introduction

- variety of data types, e.g.
- **lexicon:** database of words, minimally containing part of speech information and glosses
- **paradigm:** any kind of rational tabulation of words or phrases to illustrate contrasts and systematic variation
- **text:** any larger unit such as a narrative or a conversation
- descriptions: field notes, grammars and analytical papers
- complex inter-relations (new word; lexicon update, paradigm; add to field notes; extend grammar, ...)
- tasks: sorting, tabulating, garnering statistics, searching, verifying transcriptions
- principle challenge is computational



# Introduction

- variety of data types, e.g.
- **lexicon:** database of words, minimally containing part of speech information and glosses
- **paradigm:** any kind of rational tabulation of words or phrases to illustrate contrasts and systematic variation
- **text:** any larger unit such as a narrative or a conversation
- descriptions: field notes, grammars and analytical papers
- complex inter-relations (new word; lexicon update, paradigm; add to field notes; extend grammar, ...)
- tasks: sorting, tabulating, garnering statistics, searching, verifying transcriptions
- principle challenge is computational

# Tools and Technologies

- General-purpose tools
  - word processors
  - spreadsheets
  - databases
- Special purpose tools
  - Shoebox
- broader question: data management in support of best practices
- Bird, Steven and Gary Simons (2003). Seven Dimensions of Portability for Language Documentation and Description, *Language* 79: 557-582.

# Tools and Technologies

- General-purpose tools
  - word processors
  - spreadsheets
  - databases
- Special purpose tools
  - Shoebox
- broader question: data management in support of best practices
- Bird, Steven and Gary Simons (2003). Seven Dimensions of Portability for Language Documentation and Description, *Language* 79: 557-582.

# Tools and Technologies

- General-purpose tools
  - word processors
  - spreadsheets
  - databases
- Special purpose tools
  - Shoebox
- broader question: data management in support of best practices
- Bird, Steven and Gary Simons (2003). Seven Dimensions of Portability for Language Documentation and Description, *Language* 79: 557-582.

# Tools and Technologies

- General-purpose tools
  - word processors
  - spreadsheets
  - databases
- Special purpose tools
  - Shoebox
- broader question: data management in support of best practices
- Bird, Steven and Gary Simons (2003). Seven Dimensions of Portability for Language Documentation and Description, *Language* 79: 557-582.

# Tools and Technologies

- General-purpose tools
  - word processors
  - spreadsheets
  - databases
- Special purpose tools
  - Shoebox
- broader question: data management in support of best practices
- Bird, Steven and Gary Simons (2003). Seven Dimensions of Portability for Language Documentation and Description, *Language* 79: 557-582.

# Tools and Technologies

- General-purpose tools
  - word processors
  - spreadsheets
  - databases
- Special purpose tools
  - Shoebox
- broader question: data management in support of best practices
- Bird, Steven and Gary Simons (2003). Seven Dimensions of Portability for Language Documentation and Description, *Language* 79: 557-582.

# Tools and Technologies

- General-purpose tools
  - word processors
  - spreadsheets
  - databases
- Special purpose tools
  - Shoebox
- broader question: data management in support of best practices
- Bird, Steven and Gary Simons (2003). Seven Dimensions of Portability for Language Documentation and Description, *Language* 79: 557-582.



# Tools and Technologies

- General-purpose tools
  - word processors
  - spreadsheets
  - databases
- Special purpose tools
  - Shoebox
- broader question: data management in support of best practices
- Bird, Steven and Gary Simons (2003). Seven Dimensions of Portability for Language Documentation and Description, *Language* 79: 557-582.

# General Purpose Tools: Word Processors

- office software widely used in computer-based language documentation work
- familiar, ready availability (though may be too expensive?)
- Word processors for managing dictionaries:
  - consistency of content and format — time consuming
  - how would we check the following constraint:  
*each entry must have a part-of-speech field, drawn from a set of 20 possibilities, displayed after the pronunciation field, and rendered in 11-point bold*
  - generally a bad idea!
- save to a non-proprietary format (e.g. RTF, HTML, or XML)
- write programs to do this checking automatically

# General Purpose Tools: Word Processors

- office software widely used in computer-based language documentation work
- familiar, ready availability (though may be too expensive?)
- Word processors for managing dictionaries:
  - consistency of content and format — time consuming
  - how would we check the following constraint:  
*each entry must have a part-of-speech field, drawn from a set of 20 possibilities, displayed after the pronunciation field, and rendered in 11-point bold*
  - generally a bad idea!
- save to a non-proprietary format (e.g. RTF, HTML, or XML)
- write programs to do this checking automatically

# General Purpose Tools: Word Processors

- office software widely used in computer-based language documentation work
- familiar, ready availability (though may be too expensive?)
- Word processors for managing dictionaries:
  - consistency of content and format — time consuming
  - how would we check the following constraint:  
*each entry must have a part-of-speech field, drawn from a set of 20 possibilities, displayed after the pronunciation field, and rendered in 11-point bold*
  - generally a bad idea!
- save to a non-proprietary format (e.g. RTF, HTML, or XML)
- write programs to do this checking automatically

# General Purpose Tools: Word Processors

- office software widely used in computer-based language documentation work
- familiar, ready availability (though may be too expensive?)
- Word processors for managing dictionaries:
  - consistency of content and format — time consuming
  - how would we check the following constraint:  
*each entry must have a part-of-speech field, drawn from a set of 20 possibilities, displayed after the pronunciation field, and rendered in 11-point bold*
  - generally a bad idea!
- save to a non-proprietary format (e.g. RTF, HTML, or XML)
- write programs to do this checking automatically

# General Purpose Tools: Word Processors

- office software widely used in computer-based language documentation work
- familiar, ready availability (though may be too expensive?)
- Word processors for managing dictionaries:
  - consistency of content and format — time consuming
  - how would we check the following constraint:  
*each entry must have a part-of-speech field, drawn from a set of 20 possibilities, displayed after the pronunciation field, and rendered in 11-point bold*
  - generally a bad idea!
- save to a non-proprietary format (e.g. RTF, HTML, or XML)
- write programs to do this checking automatically

# General Purpose Tools: Word Processors

- office software widely used in computer-based language documentation work
- familiar, ready availability (though may be too expensive?)
- Word processors for managing dictionaries:
  - consistency of content and format — time consuming
  - how would we check the following constraint:  
*each entry must have a part-of-speech field, drawn from a set of 20 possibilities, displayed after the pronunciation field, and rendered in 11-point bold*
  - generally a bad idea!
- save to a non-proprietary format (e.g. RTF, HTML, or XML)
- write programs to do this checking automatically

# General Purpose Tools: Word Processors

- office software widely used in computer-based language documentation work
- familiar, ready availability (though may be too expensive?)
- Word processors for managing dictionaries:
  - consistency of content and format — time consuming
  - how would we check the following constraint:  
*each entry must have a part-of-speech field, drawn from a set of 20 possibilities, displayed after the pronunciation field, and rendered in 11-point bold*
  - generally a bad idea!
- save to a non-proprietary format (e.g. RTF, HTML, or XML)
- write programs to do this checking automatically



# General Purpose Tools: Word Processors

- office software widely used in computer-based language documentation work
- familiar, ready availability (though may be too expensive?)
- Word processors for managing dictionaries:
  - consistency of content and format — time consuming
  - how would we check the following constraint:  
*each entry must have a part-of-speech field, drawn from a set of 20 possibilities, displayed after the pronunciation field, and rendered in 11-point bold*
  - generally a bad idea!
- save to a non-proprietary format (e.g. RTF, HTML, or XML)
- write programs to do this checking automatically

## Example: Microsoft Word

- sleep [sli:p] **vi** *condition of body and mind...*
- We can enter this in MSWord, then "Save as Web Page"
- resulting HTML:

```
<p class=MsoNormal>sleep <span style='mso-spacerun:yes'>  
class=SpellE>sli:p</span>]<span style='mso-spacerun:yes'>  
style='font-size:11.0pt'>vi</span></b><span style='font-size:11.0pt'>  
condition of body and mind ...o:p></o:p></i></p>
```

- Observations:
  - entry: HTML paragraph using the `<p>` element
  - pos: inside a `<span style='font-size:11.0pt'>` element

## Example: Validating Dictionaries stored in MSWord

```
>>> import re
>>> legal_pos = set(['n', 'v.t.', 'v.i.', 'adj', 'det'])
>>> pattern = re.compile(r"'font-size:11.0pt'>([a-z.]+)<")
>>> document = open("dict.htm").read()
>>> used_pos = set(re.findall(pattern, document))
>>> illegal_pos = used_pos.difference(legal_pos)
>>> print list(illegal_pos)
['v.intr', 'v.i', 'intrans']
```

# Example: Converting MSWord dictionary to Excel

```
>>> import re
>>> document = open("dict.htm").read()
>>> document = re.sub("[\r\n]", "", document)
>>> word_pattern = re.compile(r">([\w]+)")
>>> pron_pattern = re.compile(r"\[.*>([a-z:]+)<.*\]")
>>> for entry in document.split("<p"):
...     word_match = word_pattern.search(entry)
...     pron_match = pron_pattern.search(entry)
...     if word_match and pron_match:
...         lex = word_match.group(1)
...         pos = pron_match.group(1)
...         print '%s', '%s' % (lex, pos)
"sleep", "sli:p"
"walk", "wo:k"
"wake", "weik"
```

# Spreadsheets

- wordlists, paradigms
- e.g. comparative wordlist  
(<http://www.rosettaproject.org/>):
  - row for each cognate set
  - column for each language
- export spreadsheets in *CSV format*
- read them using Python's `csv` module
- e.g. find cognates having an edit-distance of at least three from each other

# Spreadsheets

- wordlists, paradigms
- e.g. comparative wordlist  
(<http://www.rosettaproject.org/>):
  - row for each cognate set
  - column for each language
- export spreadsheets in *CSV format*
- read them using Python's `csv` module
- e.g. find cognates having an edit-distance of at least three from each other

# Spreadsheets

- wordlists, paradigms
- e.g. comparative wordlist  
(<http://www.rosettaproject.org/>):
  - row for each cognate set
  - column for each language
- export spreadsheets in *CSV format*
- read them using Python's `csv` module
- e.g. find cognates having an edit-distance of at least three from each other

# Spreadsheets

- wordlists, paradigms
- e.g. comparative wordlist  
(<http://www.rosettaproject.org/>):
  - row for each cognate set
  - column for each language
- export spreadsheets in *CSV format*
- read them using Python's `csv` module
- e.g. find cognates having an edit-distance of at least three from each other



# Spreadsheets

- wordlists, paradigms
- e.g. comparative wordlist  
(<http://www.rosettaproject.org/>):
  - row for each cognate set
  - column for each language
- export spreadsheets in *CSV format*
- read them using Python's `csv` module
- e.g. find cognates having an edit-distance of at least three from each other

# Spreadsheets

- wordlists, paradigms
- e.g. comparative wordlist  
(<http://www.rosettaproject.org/>):
  - row for each cognate set
  - column for each language
- export spreadsheets in *CSV format*
- read them using Python's `CSV` module
- e.g. find cognates having an edit-distance of at least three from each other

# Spreadsheets

- wordlists, paradigms
- e.g. comparative wordlist  
(<http://www.rosettaproject.org/>):
  - row for each cognate set
  - column for each language
- export spreadsheets in *CSV format*
- read them using Python's `CSV` module
- e.g. find cognates having an edit-distance of at least three from each other

# Databases

- lexicons are stored in a full-fledged relational databases
- databases can implement many well-formedness constraints  
e.g. ensure that all parts-of-speech come from an *enumerated type*
- However, the relational model is often too restrictive
- optional, repeatable fields (e.g. dictionary sense definitions and example sentences).
- SQL cannot express many linguistically-motivated queries  
*Find all words that appear in example sentences for which no dictionary entry is provided*

# Databases

- lexicons are stored in a full-fledged relational databases
- databases can implement many well-formedness constraints  
e.g. ensure that all parts-of-speech come from an *enumerated type*
- However, the relational model is often too restrictive
- optional, repeatable fields (e.g. dictionary sense definitions and example sentences).
- SQL cannot express many linguistically-motivated queries  
*Find all words that appear in example sentences for which no dictionary entry is provided*

# Databases

- lexicons are stored in a full-fledged relational databases
- databases can implement many well-formedness constraints  
e.g. ensure that all parts-of-speech come from an *enumerated type*
- However, the relational model is often too restrictive
- optional, repeatable fields (e.g. dictionary sense definitions and example sentences).
- SQL cannot express many linguistically-motivated queries  
*Find all words that appear in example sentences for which no dictionary entry is provided*

# Databases

- lexicons are stored in a full-fledged relational databases
- databases can implement many well-formedness constraints  
e.g. ensure that all parts-of-speech come from an *enumerated type*
- However, the relational model is often too restrictive
- optional, repeatable fields (e.g. dictionary sense definitions and example sentences).
- SQL cannot express many linguistically-motivated queries  
*Find all words that appear in example sentences for which no dictionary entry is provided*

# Databases

- lexicons are stored in a full-fledged relational databases
- databases can implement many well-formedness constraints  
e.g. ensure that all parts-of-speech come from an *enumerated type*
- However, the relational model is often too restrictive
- optional, repeatable fields (e.g. dictionary sense definitions and example sentences).
- SQL cannot express many linguistically-motivated queries  
*Find all words that appear in example sentences for which no dictionary entry is provided*



# Database Example

```
>>> import csv
>>> lexemes = []
>>> defn_words = []
>>> for row in csv.reader(open("dict.csv")):
...     lexeme, pron, pos, defn = row
...     lexemes.append(lexeme)
...     defn_words += defn.split()
>>> undefined = list(set(defn_words).difference(set(lexemes)))
>>> undefined.sort()
>>> print undefined
['...', 'a', 'and', 'body', 'by', 'cease', 'condition', 'down',
'each', 'foot', 'lifting', 'mind', 'of', 'progress', 'setting', 'to']
```

# Processing Shoebox Data

- single most popular tool for managing linguistic field data
- many kinds of validation and formatting not supported by Shoebox software
- simple file format
- each file is a collection of *entries* (aka *records*)
- each entry is made up of one or more *fields*
- challenges:
  - different entry types, implications for which fields are present
  - field order only sometimes significant

# Processing Shoebox Data

- single most popular tool for managing linguistic field data
- many kinds of validation and formatting not supported by Shoebox software
- simple file format
- each file is a collection of *entries* (aka *records*)
- each entry is made up of one or more *fields*
- challenges:
  - different entry types, implications for which fields are present
  - field order only sometimes significant

# Processing Shoebox Data

- single most popular tool for managing linguistic field data
- many kinds of validation and formatting not supported by Shoebox software
- simple file format
- each file is a collection of *entries* (aka *records*)
- each entry is made up of one or more *fields*
- challenges:
  - different entry types, implications for which fields are present
  - field order only sometimes significant

# Processing Shoebox Data

- single most popular tool for managing linguistic field data
- many kinds of validation and formatting not supported by Shoebox software
- simple file format
- each file is a collection of *entries* (aka *records*)
- each entry is made up of one or more *fields*
- challenges:
  - different entry types, implications for which fields are present
  - field order only sometimes significant

# Processing Shoebox Data

- single most popular tool for managing linguistic field data
- many kinds of validation and formatting not supported by Shoebox software
- simple file format
- each file is a collection of *entries* (aka *records*)
- each entry is made up of one or more *fields*
- challenges:
  - different entry types, implications for which fields are present
  - field order only sometimes significant

# Processing Shoebox Data

- single most popular tool for managing linguistic field data
- many kinds of validation and formatting not supported by Shoebox software
- simple file format
- each file is a collection of *entries* (aka *records*)
- each entry is made up of one or more *fields*
- challenges:
  - 1 different entry types, implications for which fields are present
  - 2 field order only sometimes significant

# Processing Shoebox Data

- single most popular tool for managing linguistic field data
- many kinds of validation and formatting not supported by Shoebox software
- simple file format
- each file is a collection of *entries* (aka *records*)
- each entry is made up of one or more *fields*
- challenges:
  - 1 different entry types, implications for which fields are present
  - 2 field order only sometimes significant



# Processing Shoebox Data

- single most popular tool for managing linguistic field data
- many kinds of validation and formatting not supported by Shoebox software
- simple file format
- each file is a collection of *entries* (aka *records*)
- each entry is made up of one or more *fields*
- challenges:
  - 1 different entry types, implications for which fields are present
  - 2 field order only sometimes significant

# Shoebox Example

```
\lx kaa  
\ps N.M  
\cl isi  
\ge cooking banana  
\gp banana bilong kukim  
\sf FLORA  
\dt 12/Feb/2005  
\ex Taeavi iria kaa isi kovopaueva kaparapasia.  
\xp Taeavi i bin planim gaden banana bilong kukim tasol  
\xe Taeavi planted banana in order to cook it.
```

## Accessing Shoebox Data: raw() method

- scan the file, processing entries one at a time
- preserves order of fields

```
>>> from nltk_lite.corpora import shoebox
>>> sum_size = num_entries = 0
>>> for entry in shoebox.raw('rotokas'):
...     num_entries += 1
...     sum_size += len(entry)
>>> print sum_size/num_entries
10
```

## Accessing Shoebox Data: dictionary() method

- read each entry into a Python dictionary
- assumes field names are unique

```
>>> entries = list(shoebox.dictionary('rotokas'))
>>> from pprint import pprint
>>> pprint(entries[3])
{'cl': 'isi',
 'dt': '12/Feb/2005',
 'ex': 'Taeavi iria kaa isi kovopaueva kaparapasiasia.',
 'ge': 'cooking banana',
 'gp': 'banana bilong kukim',
 'lx': 'kaa',
 'ps': 'N.M',
 'sf': 'FLORA',
 'xe': 'Taeavi planted banana in order to cook it.',
 'xp': 'Taeavi i bin planim gaden banana bilong kukim tasol long paia.
```

# Simple Entry Processing

- 1 scan through each entry
- 2 for each entry, scan through each field

```
>>> lexemes = []  
>>> for entry in shoebox.raw('rotokas'):  
...     for field in entry:  
...         if field[0] == 'lx':  
...             normalised_lexeme = field[1].lower()  
...             lexemes.append(normalised_lexeme)
```

## Alternative:

```
>>> lexemes = [field[1].lower()  
...             for entry in shoebox.raw('rotokas')  
...             for field in entry if field[0] == 'lx']
```

# Adding New Fields

- Example: add CV field
- Aside: utility function to do CV template

```
>>> import re
>>> def cv(s):
...     s = s.lower()
...     s = re.sub(r'[^a-z]',      r' _', s)
...     s = re.sub(r'[aeiou]',    r' V', s)
...     s = re.sub(r'[^V_]',      r' C', s)
...     return (s)
```

## Adding New Fields (cont)

```
>>> raw_entries = list(shoebox.raw('rotokas'))
>>> for field in raw_entries[50]:
...     print "\\%s %s" % field
...     if field[0] == "lx":
...         print "\\cv %s" % cv(field[1])
\lx kaeviro
\cv CVVCVCV
\ps V.A
\ge lift off
\ge take off
\gp go antap
\nt used to describe action of plane
\dt 12/Feb/2005
\ex Pita kaeviroroe kepa kekesia oa vuripierevo kiuvu.
\xp Pita i go antap na lukim haus win i bagarapim.
\xe Peter went to look at the house that the wind destroyed.
```

# Removing Fields

- sanitise our data before giving it to others

```
>>> retain = ('lx', 'ps')
>>> raw_entries = list(shoebox.raw('rotokas'))
>>> for entry in raw_entries[50:55]:
...     for field in entry:
...         if field[0] in retain:
...             print "\\%s %s" % field
...     print
\\lx kaeviro
\\ps V.A

\\lx kagave
\\ps N.F

\\lx kaie
\\ps V.A
```



# Formatting Entries: dictionary() method

- request specific fields without concern for ordering

```
>>> entries = list(shoebox.dictionary('rotokas'))
>>> for entry in entries[70:80]:
...     lex = entry['lx']
...     pos = entry['ps']
...     dfn = entry['ge']
...     if 'eng' in entry:
...         dfn = entry['eng']
...     print "%s (%s) '%s'" % (lex, pos, dfn)
kakapikoto (N.N2) 'newborn baby'
kakapu (V.B) 'place in sling for purpose of carrying'
kakapua (N.N) 'sling for lifting'
kakara (N.N) 'bracelet'
Kakarapaia (N.PN) 'village name'
kakarau (N.F) 'stingray'
Kakarera (N.PN) 'name'
Kakareraia (N.???) 'name'
kakata (N.F) 'cockatoo'
kakate (N.F) 'bamboo tube for water'
```

# Generating HTML Tables

## ● Publish shoebox lexicon on the web

```
>>> html = "<table>\n"
>>> for entry in entries[70:80]:
...     lex = entry['lx']
...     pos = entry['ps']
...     dfn = entry['ge']
...     if 'eng' in entry:
...         dfn = entry['eng']
...     html += " <tr><td>%s</td><td>%s</td><td>%s</td></tr>\n" % (lex, pos, dfn)
>>> html += "</table>"
>>> print html
```

```
<table>
<tr><td>kakapikoto</td><td>N.N2</td><td>newborn baby</td></tr>
<tr><td>kakapu</td><td>V.B</td><td>place in sling for purpose of carrying</td></tr>
<tr><td>kakapua</td><td>N.N</td><td>sling for lifting</td></tr>
<tr><td>kakara</td><td>N.N</td><td>bracelet</td></tr>
<tr><td>Kakarapaia</td><td>N.PN</td><td>village name</td></tr>
<tr><td>kakarau</td><td>N.F</td><td>stingray</td></tr>
<tr><td>Kakarera</td><td>N.PN</td><td>name</td></tr>
<tr><td>Kakareraia</td><td>N.???</td><td>name</td></tr>
<tr><td>kakata</td><td>N.F</td><td>cockatoo</td></tr>
<tr><td>kakate</td><td>N.F</td><td>bamboo tube for water</td></tr>
```

# Reduplication

- create a table of lexemes and their glosses

```
>>> levgloss = {}  
>>> for entry in shoebox.dictionary('rotokas'):  
...     if 'lx' in entry and entry['ps'][0] == 'V':  
...         levgloss[entry['lx']] = entry['ge']
```

- For each lexeme, check if the lexicon contains the reduplicated form:

```
>>> for lex in levgloss:  
...     if lex+lex in levgloss:  
...         print "%s (%s); %s (%s)" % (lex, levgloss[lex], lex+lex, levgloss[lex+lex])  
kuvu (fill.up); kuvukuvu (stamp the ground)  
kitu (save); kitukitu (scrub clothes)  
kopa (ingest); kopakopa (gulp.down)  
kasi (burn); kasikasi (angry)  
koi (high pitched sound); koikoi (groan with pain)  
kee (chip); keekkee (shattered)  
kauo (jump); kauokauo (jump up and down)  
kea (deceived); keakea (lie)  
kove (drop); kovekove (drip repeatedly)  
kape (unable to meet); kapekape (grip with arms not meeting)  
kapo (fasten.cover.strip); kapokapo (fasten.cover.strips)  
koa (skin); koakoa (remove the skin)  
kipu (paint); kipukipu (rub.on)  
koe (spoon out a solid); koekoe (spoon out)  
kovo (work); kovokovo (surround)
```

# Complex Search Criteria

- for phonological description, identify segments, alternations, syllable canon...
- what syllable types occur in lexemes (MSC, conspiracies)?

```
>>> from nltk_lite.tokenize import regexp
>>> from nltk_lite.probability import FreqDist
>>> fd = FreqDist()
>>> for lex in lexemes:
...     for syl in regexp(lex, pattern=r'^[aeiou][aeiou]'):
...         fd.inc(syl)
```

## Complex Search Criteria (cont)

- Tabulate the results:

```
>>> for vowel in 'aeiou':  
...     for cons in 'ptkvsr':  
...         print '%s%s:%4d ' % (cons, vowel, fd.count(cons+vowel))  
...     print  
pa: 84  ta: 43  ka: 414  va: 87  sa: 0  ra: 185  
pe: 32  te: 8  ke: 139  ve: 25  se: 1  re: 62  
pi: 97  ti: 0  ki: 88  vi: 96  si: 95  ri: 83  
po: 31  to: 140  ko: 403  vo: 42  so: 3  ro: 86  
pu: 49  tu: 35  ku: 169  vu: 44  su: 1  ru: 72
```

- NB *t* and *s* columns
- *ti* not attested, while *si* is frequent: palatalization?
- which lexeme contains *su*? *kasuari*

# Prosodically-motivated search

- segmental and prosodic constraints
- well-formed morphemes and lexemes
- highly-specific query (e.g. find things predicted not to exist by our theoretical model?)
- *Find all trisyllabic lexemes ending in a long vowel*
- bottom-up design: need to count syllables, and test for vowels

```
>>> def num_syls(word):  
...     template = cv(word)  
...     num_cons = template.count('C')  
...     return num_cons  
>>> def is_vowel(char):  
...     return (char in 'aeiou')
```

# Prosodically-motivated search (cont)

```
>>> for entry in shoebox.dictionary('rotokas'):  
...     if 'lx' in entry:  
...         lex = entry['lx']  
...         pos = entry['ps']  
...         if num_syls(lex) == 3 and is_vowel(lex[-1]) and is_vowel(lex[-2]) and pos[0] == '  
...             dfn = entry['ge']  
...             print "%s (%s) '%s'" % (lex, pos, dfn)  
kaetupie (V.B) 'tighten'  
kakupie (V.B) 'yodel'  
kapatau (V.B) 'add to'  
kapuapie (V.B) 'wound'  
kapupie (V.B) 'close tight'  
kapuupie (V.B) 'close'  
karepie (V.B) 'return'  
karivai (V.A) 'have an appetite'  
kasipie (V.B) 'care for'  
kaukaupie (V.B) 'intense sunlight'  
kavorou (V.A) 'intercept'  
kavupie (V.B) 'leave behind'  
kekepie (V.B) 'show'  
keruria (V.A) 'determined'  
ketoopie (V.B) 'make sprout from seed'  
koatapie (V.B) 'accept'  
koetapie (V.B) 'satisfy curiosity'  
kokovae (V.A) 'sing'  
kokovua (V.B) 'shave the hair line'  
kopiipie (V.B) 'kill'  
korupie (V.B) 'take outside'
```

# Finding Minimal Sets

- E.g. mace vs maze, face vs faze
- minimal set parameters: context, target, display

Minimal Set	Context	Target	Display
<i>bib, bid, big</i>	first two letters	third letter	word
<i>deal (N), deal (V)</i>	whole word	pos	word (pos)



# Finding Minimal Sets

- E.g. mace vs maze, face vs faze
- minimal set parameters: context, target, display

Minimal Set	Context	Target	Display
<i>bib, bid, big</i>	first two letters	third letter	word
<i>deal (N), deal (V)</i>	whole word	pos	word (pos)

# Finding Minimal Sets: Example 1

```
>>> lexemes = [entry['lx'] for entry in shoebox.dictionary('rotokas')
...             if 'lx' in entry]
>>> position = 1
>>> parameters = [(lex[:position] + '_' + lex[position+1:],
...                 lex[position],
...                 lex)
...                 for lex in lexemes if len(lex) == 4]

>>> from nltk_lite.utilities import MinimalSet
>>> def build_min_set(parameters):
...     min_set = MinimalSet()
...     for context, target, display in parameters:
...         min_set.add(context, target, display)
...     return min_set
```

# Finding Minimal Sets: Example 1 (cont)

```
>>> ms = build_min_set(parameters)
>>> for context in ms.contexts(3):
...     print context + ':',
...     for target in ms.targets():
...         print "%-4s" % ms.display(context, target, "-"),
...     print
k_si: kasi -    kesi -    kosi
k_ru: karu kiru keru kuru koru
k_pu: kapu kipu -    -    kopu
k_ro: karo kiro -    -    koro
k_ri: kari kiri keri kuri kori
k_pa: kapa -    kepa -    kopa
k_ra: kara kira kera -    kora
k_ku: kaku -    -    kuku koku
k_ki: kaki kiki -    -    koki
```

## Finding Minimal Sets: Example 2

```
>>> parameters = [(entry['lx'],  
                    entry['ps'][0],  
                    "%s (%s)" % (entry['ps'][0], entry['ge'])  
                    )  
...               for entry in shoebox.dictionary('rotokas')  
...               if 'lx' in entry]  
>>> ms = build_min_set(parameters)  
>>> for context in ms.contexts()[:10]:  
...     print "%10s:" % context, "; ".join(ms.display_all(context))  
kokovara: N (unripe coconut); V (unripe)  
kapua: N (sore); V (have sores)  
koie: N (pig); V (get pig to eat)  
kovo: C (garden); N (work); V (work)  
kavori: N (lobster); V (collect crayfish or lobster)  
korita: N (cutlet?); V (cut up meat)  
keru: N (bone); V (harden like bone)  
kirokiro: N (bush used for sorcery); V (write)  
kaapie: N (fishhook); V (capture)  
kou: C (heap); V (defecate)
```

# Example Application: Improving Access to Lexical Resources

- lexicon as a *language resource*
- existence of standardized writing system? literacy?
- lookup by approximate spelling
- soundex
- edit distance

# Example Application: Improving Access to Lexical Resources

- lexicon as a *language resource*
- existence of standardized writing system? literacy?
- lookup by approximate spelling
- soundex
- edit distance

# Example Application: Improving Access to Lexical Resources

- lexicon as a *language resource*
- existence of standardized writing system? literacy?
- lookup by approximate spelling
- soundex
- edit distance

# Example Application: Improving Access to Lexical Resources

- lexicon as a *language resource*
- existence of standardized writing system? literacy?
- lookup by approximate spelling
- soundex
- edit distance



# Example Application: Improving Access to Lexical Resources

- lexicon as a *language resource*
- existence of standardized writing system? literacy?
- lookup by approximate spelling
- soundex
- edit distance

# Computing Soundex

```
>>> group = {  
...     ' ':0,                # blank (for short words)  
...     'p':1, 'b':1, 'v':1, # labials  
...     't':2, 'd':2, 's':2, # alveolars  
...     'l':3, 'r':3,        # sonorant consonants  
...     'i':4, 'e':4,        # high front vowels  
...     'u':5, 'o':5,        # high back vowels  
...     'a':6                # low vowels  
... }
```

# Computing Soundex (cont)

```
>>> def soundex(word):  
...     if len(word) == 0: return word    # sanity check  
...     word += ' '                      # ensure word long enough  
...     c0 = word[0].upper()  
...     c1 = group[word[1]]  
...     cons = filter(lambda x: x in 'pbvtdslr ', word[2:])  
...     c2 = group[cons[0]]  
...     c3 = group[cons[1]]  
...     return "%s%d%d%d" % (c0, c1, c2, c3)  
>>> print soundex('kalosavi')  
K632  
>>> print soundex('ti')  
T400
```

# Soundex Index

```
>>> soundex_idx = {}  
>>> for lex in lexemes:  
...     code = soundex(lex)  
...     if code not in soundex_idx:  
...         soundex_idx[code] = set()  
...     soundex_idx[code].add(lex)
```

# Sorting and Lookup

```
>>> from nltk_lite.utilities import edit_dist
>>> def fuzzy_spell(target):
...     scored_candidates = []
...     code = soundex(target)
...     for word in soundex_idx[code]:
...         dist = edit_dist(word, target)
...         scored_candidates.append((dist, word))
...     scored_candidates.sort()
...     return [w for (d,w) in scored_candidates[:10]]

>>> fuzzy_spell('kokopouto')
['kokopecto', 'kokopuoto', 'kokepato', 'koovoto', 'koepato',
 'kooupato', 'kopato', 'kopiito', 'kovuto', 'koavaato']
>>> fuzzy_spell('kogou')
['kogo', 'koou', 'kokeu', 'koko', 'kokoa', 'kokoi', 'kokoo',
 'koku', 'koee', 'kooku']
```

# Generating Summary Reports

- overall picture of the quality and organisation of the data
- E.g. entry patterns

```
>>> fd = FreqDist()
>>> for entry in shoebox.raw('rotokas'):
...     for field in entry:
...         fd.inc(field[0])
>>> fd.sorted_samples()[:10]
['ge', 'ex', 'xe', 'xp', 'gp', 'lx', 'ps',
'dt', 'rt', 'eng']
```

## Generating Summary Reports (cont)

```
>>> fd = FreqDist()
>>> for entry in shoebox.raw('rotokas'):
...     marker_list = [field[0] for field in entry]
...     markers = ':'.join(marker_list)
...     fd.inc(markers)
>>> top_ten = fd.sorted_samples()[:10]
>>> print '\n'.join(top_ten)
lx:rt:ps:ge:gp:dt:ex:xp:xe
lx:ps:ge:gp:dt:ex:xp:xe
lx:ps:ge:gp:dt:ex:xp:xe:ex:xp:xe
lx:rt:ps:ge:gp:dt:ex:xp:xe:ex:xp:xe
lx:ps:ge:gp:nt:dt:ex:xp:xe
lx:ps:ge:gp:dt
lx:ps:ge:ge:gp:dt:ex:xp:xe:ex:xp:xe
lx:rt:ps:ge:ge:gp:dt:ex:xp:xe:ex:xp:xe
lx:ps:ge:ge:gp:dt:ex:xp:xe
```

# Looking at Timestamps

```
>>> fd = FreqDist()
>>> from string import split
>>> for entry in shoebox.dictionary('rotokas'):
...     if 'dt' in entry:
...         (day, month, year) = split(entry['dt'], '/')
...         fd.inc((month, year))
>>> for time in fd.sorted_samples():
...     print time[0], '/', time[1], ': ', fd.count(time)
Feb / 2005 : 307
Dec / 2004 : 151
Jan / 2005 : 123
Feb / 2004 : 64
Sep / 2004 : 49
May / 2005 : 46
Mar / 2005 : 37
Apr / 2005 : 29
Jul / 2004 : 14
Nov / 2004 : 5
Oct / 2004 : 5
Aug / 2004 : 4
May / 2003 : 2
```