



Annotated Corpora in Linguistic Research

Steven Bird & Stephanie Strassel

Linguistic Data Consortium
University of Pennsylvania
Philadelphia, PA 19104
www.ldc.upenn.edu

The Linguistic Data Consortium

- ◆ a not-for-profit organization
- ◆ serving:
 - researchers, educators, technology developers
- ◆ in language-related fields by:
 - creating/collecting, archiving, distributing language resources

Research, development and teaching

- ◆ vast amounts of data, many languages
- ◆ ‘ecologically valid’ data
- ◆ range of specialized skills
 - ◆ processing of newswire, closed caption video, VOA, ...
 - ◆ intellectual property rights
 - ◆ data transformation and publication: CD, Web
- ◆ collaborative research
- ◆ replication of results



LDC Principles

Publish the data that researchers need

- ◆ sponsored programs (TDT, Hub-4, OLEADA,...)
- ◆ community initiatives (ACL/DCI, Unipen...)
- ◆ non-LDC projects (CSAE, CELEX, Trains...)
- ◆ LDC-funded (Treebank, COMLEX, WordNet...)

Make data available to everyone

- ◆ consortium membership is open to all
- ◆ most databases available to non-members

Promote the idea of shared resources

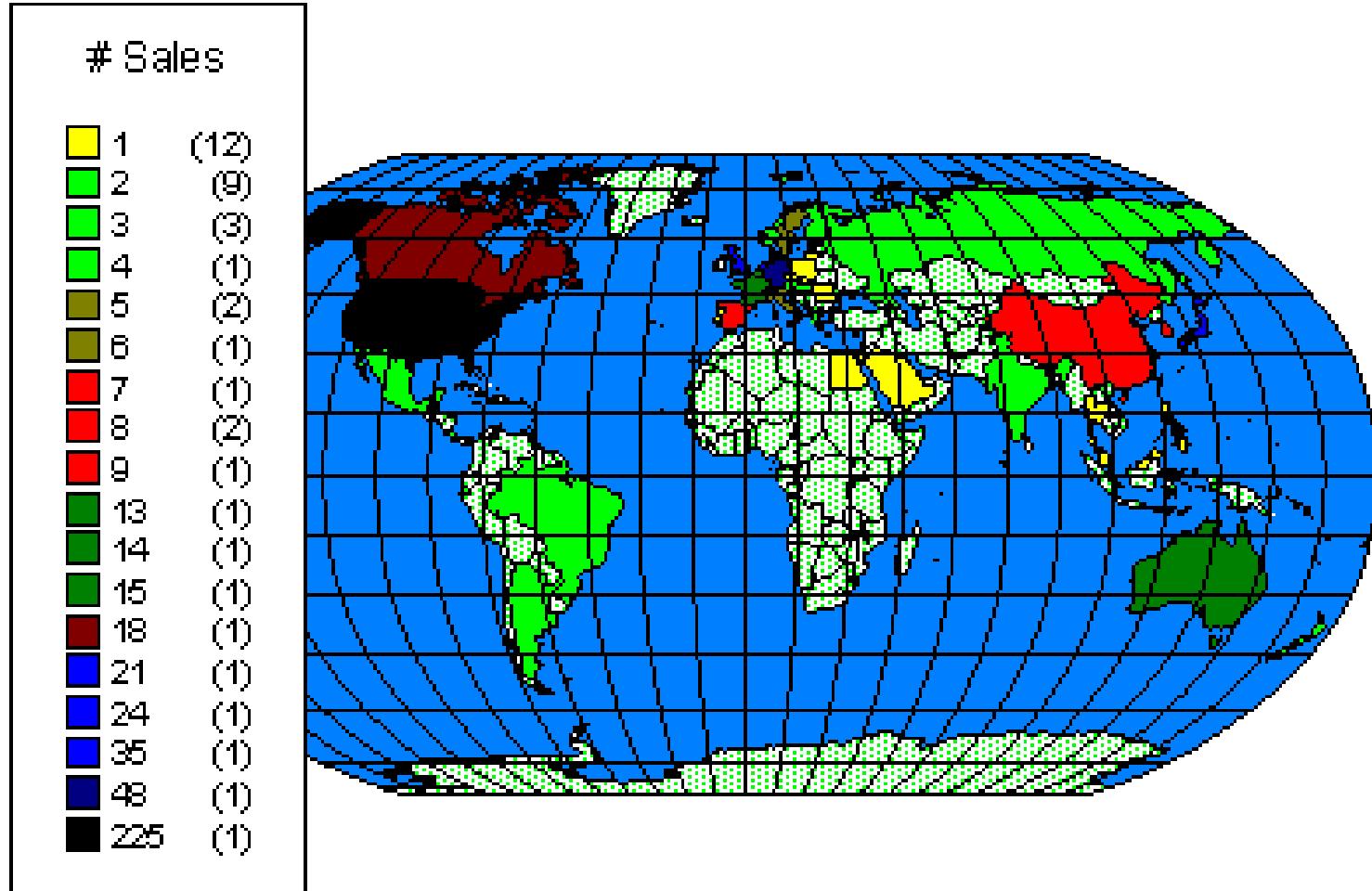
- ◆ IPR intermediary
- ◆ advice on collection, publication and IPR issues
- ◆ develop standards and tools

140 Publications + 12 in press

- ◆ 93 speech corpora + 7 in press
- ◆ 33 text corpora + 5 in press
- ◆ 14 lexicons
- ◆ 250 Gb

Distribution

- ◆ 267 members
- ◆ 520 nonmembers



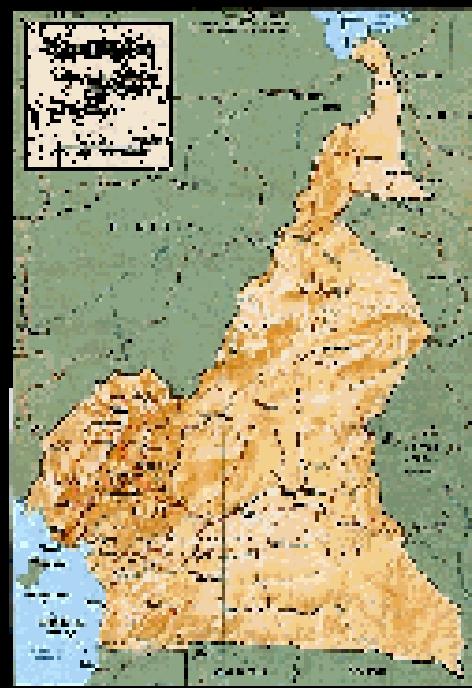
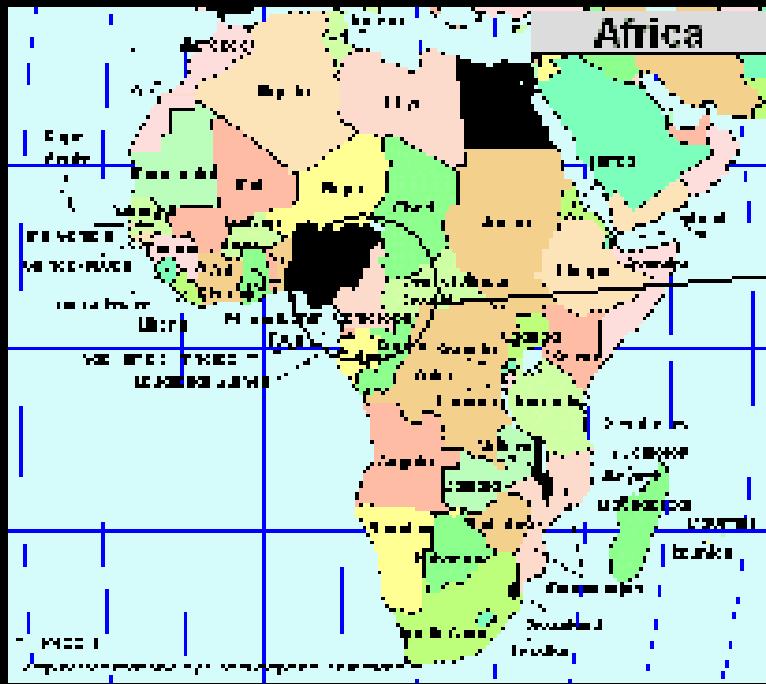
	Speech / Transcripts						
Language	Broadcast	Telephone	Wideband	Parallel Text	Newswire/ Other Text	Lexicon	Traditional Dictionary
Arabic (Egyptian)							
Czech	■						■
Dutch					■	■	
English	■	■	■	■	■	■	
French		■		■	■		■
German		■		■	■	■	
Hindi		■					
Japanese			■	■	■	■	
Korean		■					
Mandarin	■	■		■	■	■	
Persian		■			■		
Portuguese					■		
Russian					■		■
Serbo-Croatian					■		■
Spanish	■	■	■	■	■	■	
Taiwanese			■				
Tamil		■			■		
Tamil					■		
Turkish					■		
Vietnamese		■					

► Afrikaans, Bamileke, Basque, Estonian, Hungarian, Italian, Kazakh, Kurdish, Latvian, Manding, Polish, Slovene, Ukrainian, Uzbek, Xhosa, Yoruba



Linguistic Data
Consortium
University of Pennsylvania

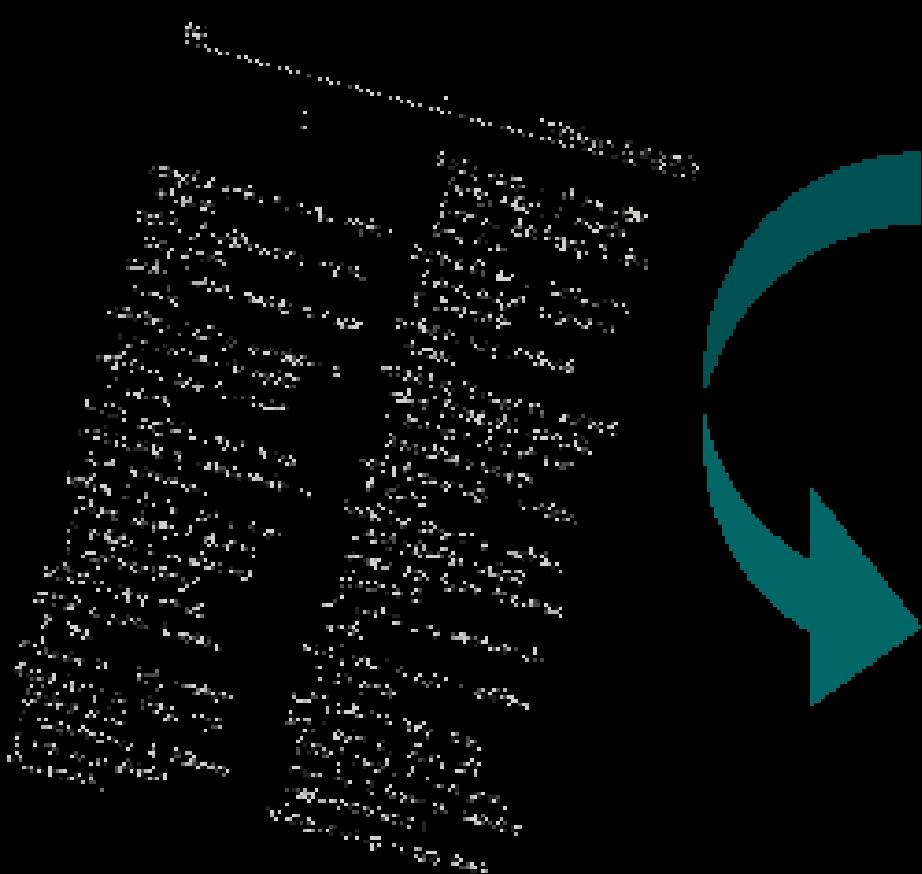
Cameroon





Bamileke
Dictionary

Bamileke Dictionary



Annotated Corpora in Linguistic Research

1500



Linguistic Data Consortium
Switchboard Corpus
Switchboard Corpus
Switchboard Corpus

Switchboard Keyword Search

Switchboard Corpus ToolBox

Define a query or command to search

Box: **Example** | Expansion level: **1** Rows of results: **1000**

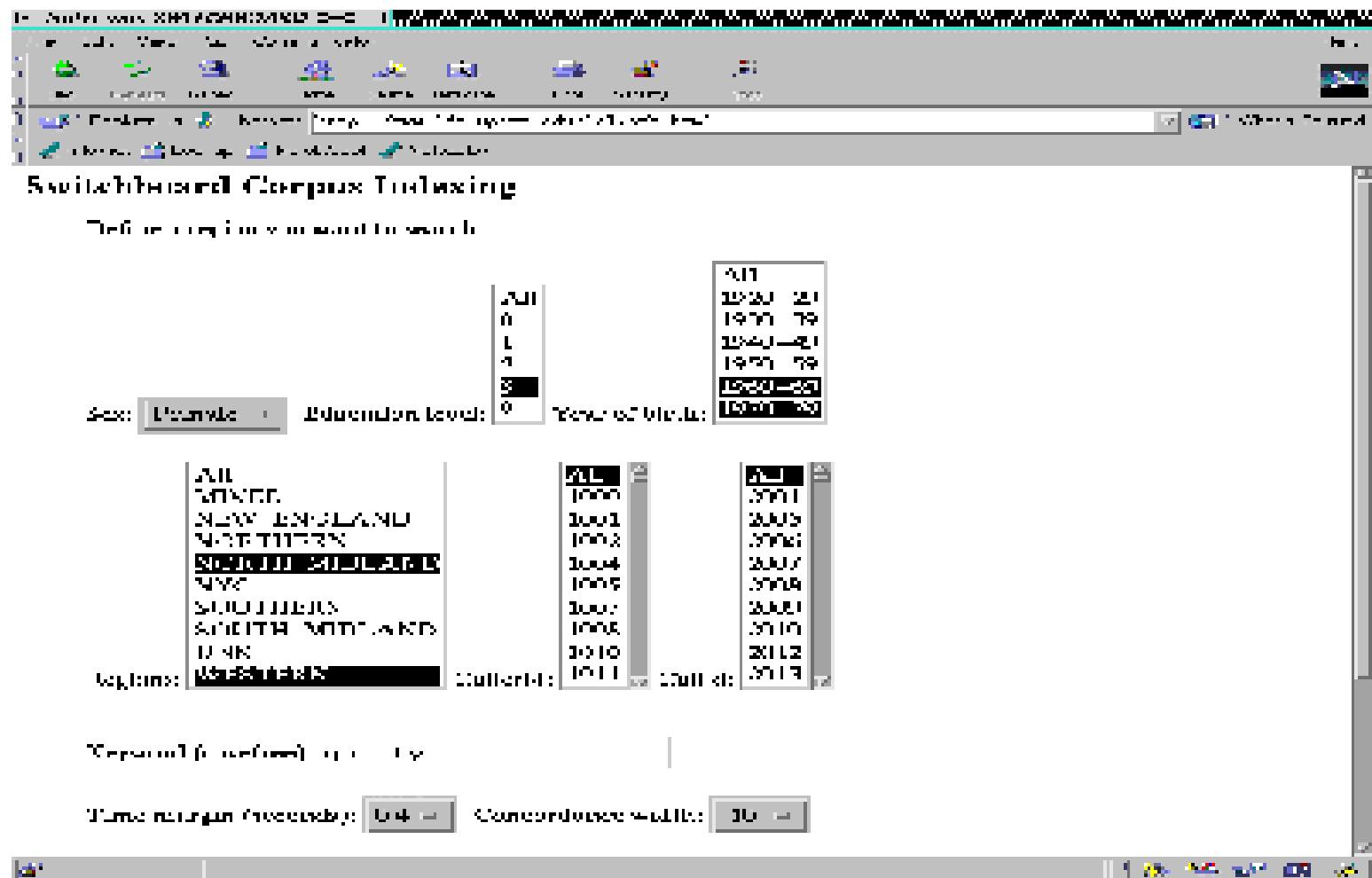
All
0
1
2
3
4
5
6
7
8
9

All
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014

Keywords: **Switchboard** | Context ID: **1000** Count of: **1000**

Version(s) selected: **1**

Time margin (seconds): **0.4** | Concordance width: **30**





Switchboard Keyword Search Hot List

You can check the boxes next to the words to view several words at the same time. By selecting the filter, you can get the whole transcript file. If you find a wrong speech segment, Clicking 'Delete' will automatically remove the problem to the result. (so it's very convenient)

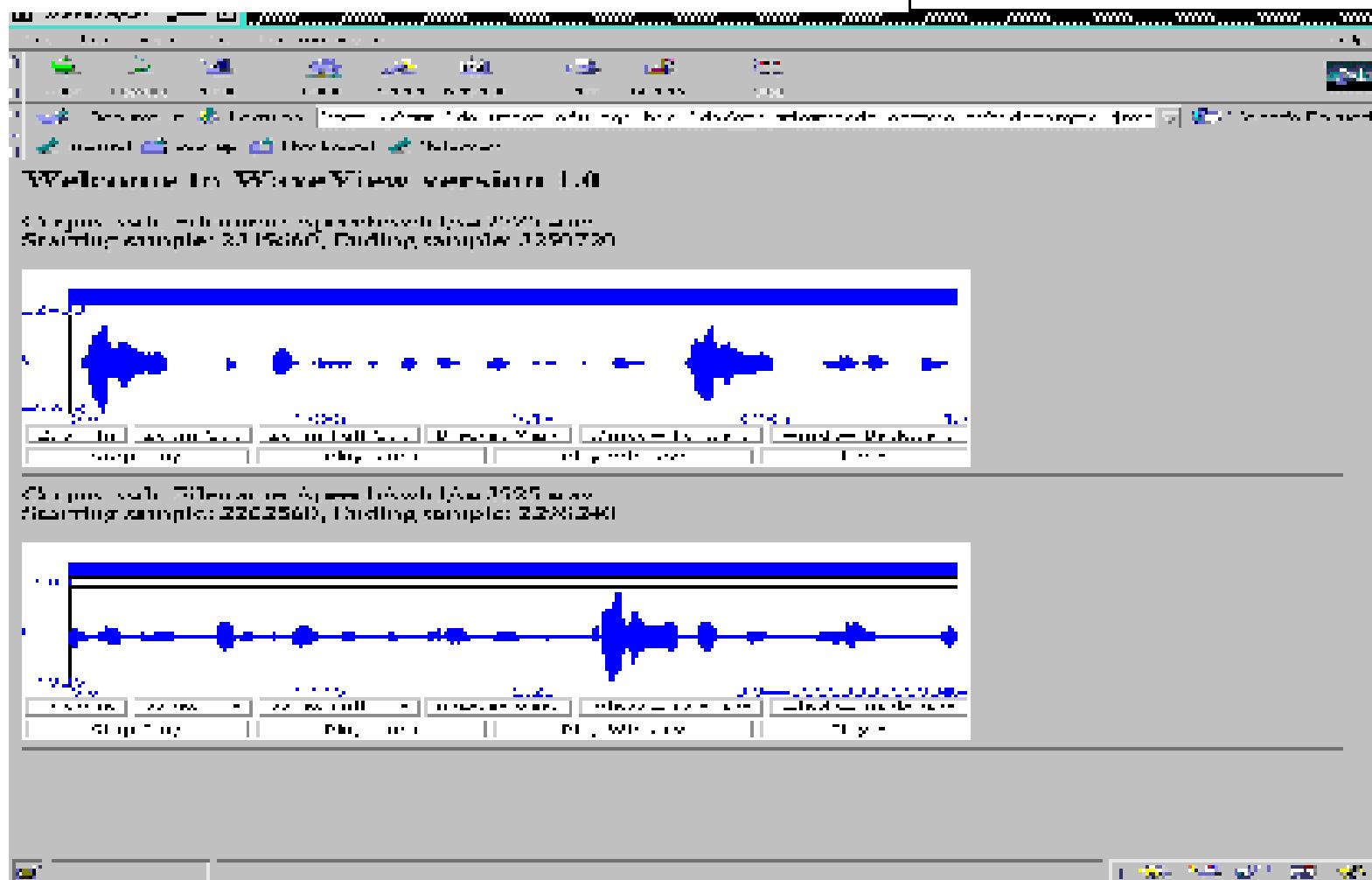
|| Advanced Search || Delete |

Index	Speaker	Text
1	B	There's a lot of bad quality of telephone lines.
2	B	Dependable telephone lines. They're there to sit out quality. Oh well. That's life.
3	B	Dependable telephone lines. I mean, who cares quality or anything like that.
4	B	Dependable telephone lines. I mean, why payed quality with your telephone?
5	B	Dependable telephone lines. I mean, the local phone company quality was always the worst.
6	B	Dependable telephone lines. I mean, they're not good quality, you know what I mean?
7	B	Dependable telephone lines. I mean, nothing, I mean, quality of telephone lines.

|| Advanced Search || Delete |



Switchboard Keyword Search Results





Switchboard Tagged Corpus Search

Pattern Matching Search

Type in the pattern you want to search for:

Page length: Page Number:

Output word types:

Screen columns width: Time margin (sec min's):

Preferred file data type you want to receive:

Select View/Select apply's width: and height:

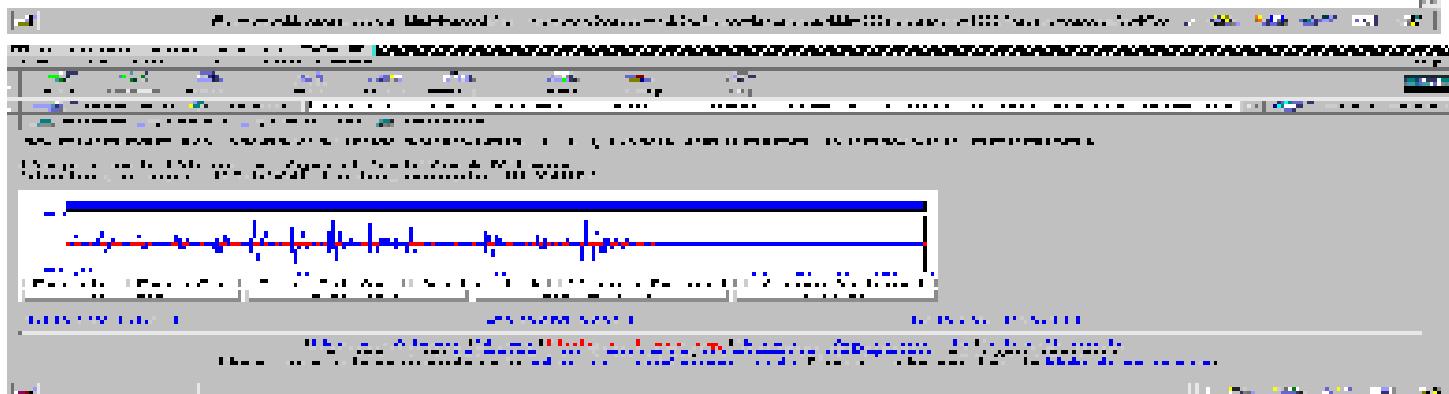
Result transmitted to your browser: or to your help application (client):



Switchboard Tagged Corpus Search: Concordance



Switchboard Tagged Corpus Search Results

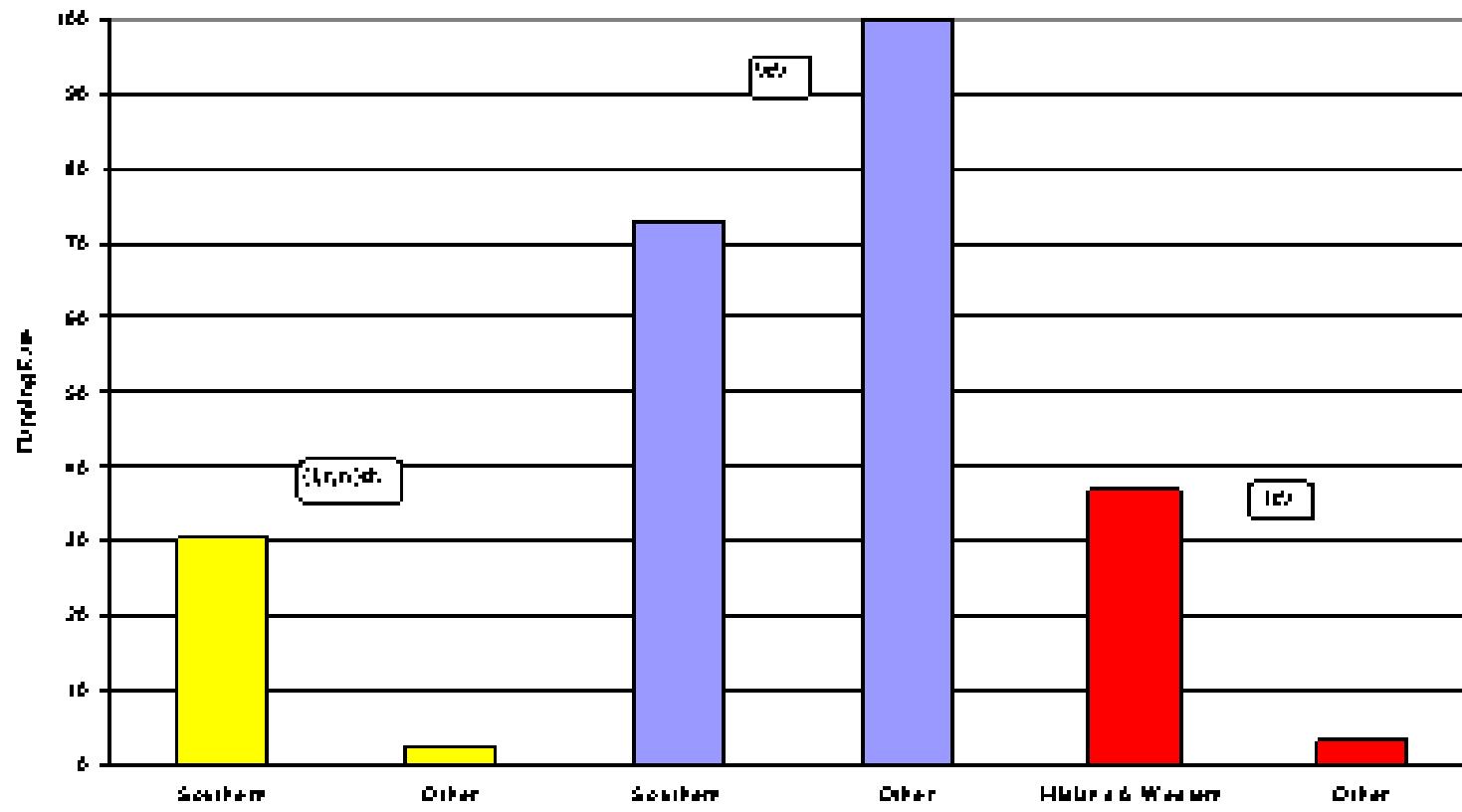




ANSWER **ANSWER** **ANSWER**

*Switchboard Tagged
Corpus Search:
Demographic Information*

Regional Variation in Flapping



(From Sundd 1998)