# The advantages of using relational databases for large historical corpora

Mark Davies
Brigham Young University
http://davies-linguistics.byu.edu/
mark_davies@byu.edu

# Outline

- Introductory comments

- VIEW/BNC (100m words)

- Corpus of Historical English (OED, 37m)

- Corpus do Português (July 2006, 45m)

- Examples of annotation via relational databases

- Proposed 200m word corpus of historical English

# Size

- Function of intended use
- Product of biases of corpus creators
- CED, CEEC, etc. perfectly fine for historical pragmatics
- Bottom line: not make size a limitation

# Accessibility

➢ Via Web

➢ Cost (free?? to end-users)

➢ User interface (may be a much larger and diverse audience than we had first intended)

# Searchability

➢ Spelling, morphology: wildcards

➢ Syntax: POS tagging

➢ Semantics: Collocates, comparing collocates, synonyms, user-defined lists

➢ [ Pragmatics ]

➢ Q: Is what we study a function of the corpus, or what we really want to study?

# Size: Corpus del Español (www.corpusdelespanol.org) ; 100m words 1200s-1900s

Features: **VIEW**: view.byu.edu (British National Corpus; 100m words)

# BNC / VIEW features

Word document

# Corpus of Historical English: view.byu.edu/che



Corpus of Historical English [DAVIES] - Microsoft Internet Explorer

File   Edit   View   Favorites   Tools   Help

Back   |   Search   Favorites

Address http://view.byu.edu/che/x.asp   Go

Google   |   Search   |   3999 blocked   |   Check   |   AutoFill   |   Options

Based on the Oxford English Dictionary (37 million words, Old English – Modern English)

## CORPUS OF HISTORICAL ENGLISH
### Created by Mark Davies / BYU

Click on the number in a column to see [KEYWORD IN CONTEXT] display          [HELP]

**SEARCH**

WORD(S)/PHRASE(S)   HELP

*light*

[CUSTOMIZED LISTS] [+OE]

- ● FREQUENCY (TABLE)   HELP
- ○ FREQUENCY (CHART)
- ○ RELEVANCE
- ○ SURROUNDING WORDS

| | CENTURY | 1000s | 1100s | 1200s | 1300s | 1400s | 1500s | 1600s | 1700s | 1800s | 1900s |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 + | LIGHT | | | | 137 | 149 | 812 | 1988 | 1328 | 4473 | 4078 |
| 2 + | DELIGHT | | | | | 1 | 174 | 418 | 241 | 435 | 218 |
| 3 + | FLIGHT | | | | 13 | 19 | 97 | 273 | 182 | 399 | 904 |
| 4 + | LIGHTS | | | | | 6 | 43 | 218 | 138 | 432 | 658 |
| 5 + | SLIGHT | | | | 20 | 5 | 39 | 189 | 155 | 761 | 544 |
| 6 + | DELIGHTS | | | | | 37 | 144 | 77 | 97 | 64 |
| 7 + | LIGHTNING | | | | 1 | 2 | 51 | 132 | 112 | 242 | 180 |
| 8 + | LIGHTLY | | | | 25 | 42 | 83 | 108 | 51 | 233 | 221 |
| | DELIGHTED | | | | | 28 | 98 | 66 | 164 | 138 |

**LIMITS**

SURROUNDING WORDS   HELP

0 ▼ LEFT/RIGHT 0 ▼

FREQUENCY IN CENTURY   HELP

**DISPLAY**

SORT BY CENTURY   HELP

1600s ▼

FREQ DISPLAY   HELP

RAW # ▼

# HITS   HELP

100

SUBMIT   RESET

**[KEYWORD IN CONTEXT]**          More information... ▼

| 1 | 1509 | BARCLAY: Shyp of Folys | Because the **lightning** or thunder violent .. suffreth thee and thy house to be vnbrent. |
| 2 | 1555 | T. PHAER: neid | And from the skyes the **lightning** fyers do flashe wt grisly steauen. |
| 3 | 1556 | W. TOWERSON: Voy. | The 4.day we had terrible thunder and **lightning**, with exceeding great gusts of raine, called Ternados. |
| 4 | 1558 | PHAR: neid | Me the father of Gods .. Beblasted with his **lightning** wynd. |
| 5 | 1561 | Burn. Paules Ch. | On Wednesday .. was seene a marueilous great fyrie **lightning**. |
| 6 | 1561 | T. NORTON: Calvin's Inst. | To set the aire on fier with **lightning** flames. |
| 7 | 1563 | W. FULKE: Meteors | If **lightning** or any other inflamation be in the upper part of these clouds. |
| 8 | 1563 | W. FULKE: Meteors | The most dangerous, violent and hurtfull kind of **lightning** is called Fulmen. |
| 9 | 1563 | W. FULKE: Meteors | Fulgur is that kinde of **lightning** which followeth thunder. |
| 10 | 1563 | W. FULKE: Meteors | Although the **lightning** appeare unto us, a good pretty while before the thunderclap be heard. |
| 11 | 1563 | FULKE: Meteors | The Sea-Calfe is never hurt with **lightning**: wherefore the Emperours tents were woont to be covered with their skinnes. |

Done          Internet

# Corpus of Historical English: Size and distribution

| Century | # words |
|---------|---------|
| 1000s | 207,594 |
| 1100s | 53,359 |
| 1200s | 313,717 |
| 1300s | 1,077,330 |
| 1400s | 1,388,448 |
| 1500s | 3,182,053 |
| 1600s | 5,127,445 |
| 1700s | 3,688,076 |
| 1800s | 10,422,785 |
| 1900s | 11,063,571 |
| TOTAL | 36,524,378 |

# Simple frequency of word or phrase over time: **turn on**

Corpus of Historical English [DAVIES] - Microsoft Internet Explorer

File   Edit   View   Favorites   Tools   Help

Back   Search   Favorites   Media

Address  http://view.byu.edu/che/x.asp   Go   Links

Based on 37 million words
of text from Old English
to Present-Day English

# CORPUS OF HISTORICAL ENGLISH
## Created by Mark Davies / BYU

**SEARCH**

WORD(S)/PHRASE(S)   HELP
turn on
[CUSTOMIZED LISTS] [+OE]

○ FREQUENCY (TABLE)   HELP
● FREQUENCY (CHART)
○ RELEVANCE
○ SURROUNDING WORDS

**LIMITS**

SURROUNDING WORDS   HELP
0 ▼ LEFT/RIGHT 0 ▼

FREQUENCY IN CENTURY   HELP

**DISPLAY**

SORT BY CENTURY   HELP
1900s ▼

FREQ DISPLAY   HELP
RAW # ▼

# HITS   HELP
100

SUBMIT   RESET

Click on a bar to see [KEYWORD IN CONTEXT] display
Click on a number below the bar chart for table display   [HELP]

| CENTURY | 1000s | 1100s | 1200s | 1300s | 1400s | 1500s | 1600s | 1700s | 1800s | 1900s |
|---|---|---|---|---|---|---|---|---|---|---|
| PER MILLION | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.2 | 1.4 | 2.3 | 5.0 |
| FREQUENCY | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 5 | 24 | 55 |

[KEYWORD IN CONTEXT]   More information...

| 1 | 1900 | The Newlands ground is the most difficult to make runs on in the whole of South Africa, the bowlers always being able to get considerable **turn on** the ball. |
| 2 | 1901 | Please, Nanna, don't **turn on** the dark. |
| 3 | 1902a | While she was gone the man from the gas company called to **turn on** the meter. |
| 4 | 1904 | Such objections to the meaningfulness of the mechanical image as **turn on** the difficulties of defining mass, length, time, and their combinations, cannot be discussed in this connection. |
| 5 | 1909 | I don't know the moment he delivers the ball which way it will **turn on** pitching. |
| 6 | 1917 | If someone will .. kindly re-carbon the arc for us and **turn on** the electric current we shall be ready to start showing. |
| 7 | 1919 | She'll screen well, and she's one of the few that can **turn on** the tears when she wants to. |
| 8 | 1926 | Most fat cows and heavy heifers lost around 25c and in instances the down-**turn on** better grades was even greater. |
| 9 | 1934 | Cover One with a Gun (v. phr.): to **turn on** the heat. |
| 10 | 1942 | Immediately you are grounded, **turn on** your back and unlock the quick release mechanism. |
| 11 | 1946 | The immediate action is to **turn on** the reciprocal course. |
| 12 | 1949 | The majority .. were Zionists .. who at the eve of the war had been waiting for their **turn on** the immigration quota. |
| 13 | 1953 | The trimmer with only one **turn on** his winch trims hand over hand as the sail swings over. |
| 14 | 1953 | Philosopher's arguments which **turn on** these terms are apt, sooner or later, to start to rotate idly. There is nothing |

Internet

# Spelling changes: **vn\*** (chart)



**Corpus of Historical English [DAVIES] - Microsoft Internet Explorer**

File   Edit   View   Favorites   Tools   Help

Back   |   Search   Favorites   Media   |

Address  http://view.byu.edu/che/x.asp                          Go    Links »

**CORPUS OF HISTORICAL ENGLISH**
**Created by Mark Davies / BYU**

Based on 37 million words
of text from Old English
to Present-Day English

**SEARCH**

WORD(S)/PHRASE(S)   HELP
vn*
[CUSTOMIZED LISTS] [+OE]

○ FREQUENCY (TABLE)   HELP
● FREQUENCY (CHART)
○ RELEVANCE
○ SURROUNDING WORDS

**LIMITS**

SURROUNDING WORDS   HELP
0 ▾ LEFT/RIGHT 0 ▾

FREQUENCY IN CENTURY   HELP

**DISPLAY**

SORT BY CENTURY   HELP
1900s ▾

FREQ DISPLAY   HELP
RAW # ▾

# HITS   HELP
100

SUBMIT   RESET

Click on any bar or number below for a list of matching strings        [HELP]

| CENTURY | 1000s | 1100s | 1200s | 1300s | 1400s | 1500s | 1600s | 1700s | 1800s | 1900s |
|---------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| PER MILLION | 38.5 | 187.4 | 1,711.7 | 3,710.1 | 3,917.3 | 4,554.9 | 1,783.1 | 4.9 | 0.8 | 0.6 |
| FREQUENCY | 8 | 10 | 537 | 3997 | 5439 | 14494 | 9143 | 18 | 8 | 7 |

**[KEYWORD IN CONTEXT]**                     More information... ▾

1  1908 Well, I was rared a Carthlick, but I haven't **followed it up** much. To tell ye the truth, I class 'em all alike -- priests, parsons, salvos, and all the lot of 'em.

2  1951 The State Legislature at Albany ordered an inquiry, which met last week. And this week the American Legion has **followed it up** with a narcotic clinic of its own.

3  1955 Beaumont College, who won their juniors at Molesey a week ago, **followed it up** by winning the junior-senior eights.

4  1956 We followed it, up to the spring hole on the edge of the flat land, a no-good bit of bog hole on the edge of arable land.

5  1961 Why she should have **followed it up** beats me. You didn't ask her? .. Of course not! As I told you, she carries too many guns.

6  1965 I .. **followed it up** with a .. right-cross with the school cap.

Internet

# Spelling changes (table): rank-frequency listing by 1500s

# Chart display (by century): **to * [davies:up]** [CUSTOMIZED LISTS]

# Table display (1500s): **to * [davies:up]** [CUSTOMIZED LISTS]



Corpus of Historical English [DAVIES] - Microsoft Internet Explorer

File  Edit  View  Favorites  Tools  Help

Address http://view.byu.edu/che/x.asp

Google | Search | 3999 blocked | Check | AutoFill | Options

**CORPUS OF HISTORICAL ENGLISH**
**Created by Mark Davies / BYU**

Based on the Oxford English Dictionary (37 million words, Old English – Modern English)

**SEARCH**

WORD(S)/PHRASE(S)     HELP
to * [davies:up]
[CUSTOMIZED LISTS] [+OE]

○ FREQUENCY (TABLE)     HELP
● FREQUENCY (CHART)
○ RELEVANCE
○ SURROUNDING WORDS

**LIMITS**

SURROUNDING WORDS     HELP
5 ▼ LEFT/RIGHT 5 ▼

FREQUENCY IN CENTURY     HELP

**DISPLAY**

SORT BY CENTURY     HELP
1900s ▼

FREQ DISPLAY     HELP
RAW # ▼

# HITS     HELP
100

SUBMIT   RESET   ☐

Click on the number in a column to see [KEYWORD IN CONTEXT] display     **[HELP]**

| | CENTURY | 1000s | 1100s | 1200s | 1300s | 1400s | 1500s | 1600s | 1700s | 1800s | 1900s |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | TO TAKE VP | | | | | 3 | 14 | 6 | | | |
| 2 | TO SET VP | | | | | | 13 | 11 | | | |
| 3 | TO MAKE VP | | | | | 6 | 10 | 9 | | | |
| 4 | TO COME VP | | | | | 1 | 7 | 3 | | | |
| 5 | TO LIFT VP | | | | | | 7 | 1 | | | |
| 6 | TO BEARE VP | | | | | | 7 | 4 | | | |
| 7 | TO TAKE UP | | | | | 1 | 7 | 38 | 46 | 70 | 70 |
| 8 | TO SET UP | | | | | | 7 | 43 | 36 | 54 | 117 |
| 9 | TO FILL VP | | | | | | 6 | 7 | | | |
| 10 | TO MAKE UP | | | | | | 6 | 45 | 40 | 93 | 85 |

**CUSTOMIZED LISTS: MODIFY LIST**     More information... ▼

To **modify** a list:

Add or delete words from your list and click on [Submit]

**Return to main menu**
**Overview of wordlists**

LISTS YOU HAVE CREATED

UP  M D    KNOW  M D

**MODIFY LIST**

up     WORDLIST NAME

up
vp
vpp
vppe
upp
uppe

LIST OF WORDS

Submit   Reset

Done                                     Internet

# Morphology: *ly (+1900s -1800s -1700s)

# Morphology/lexicon: word roots: **\*light\*** (+1900s -1600s)

# Lexicon: New words in 1900s: * (+1900s -1800s)

# Relevancy / Z-score like listing: * **seat**

# Lexical bundles: * * (+1900s -1800s)

# Semantics (collocates): **hard** * (+1900s -1500s)



Corpus of Historical English [DAVIES] - Microsoft Internet Explorer

Address: http://view.byu.edu/che/x.asp

Based on 37 million words
of text from Old English
to Present-Day English

## CORPUS OF HISTORICAL ENGLISH
### Created by Mark Davies / BYU

**SEARCH**

WORD(S)/PHRASE(S)   HELP

hard *|

[CUSTOMIZED LISTS] [+OE]

- FREQUENCY (TABLE)   HELP
- FREQUENCY (CHART)
- RELEVANCE
- SURROUNDING WORDS

**LIMITS**

SURROUNDING WORDS   HELP
0 ▾ LEFT/RIGHT 0 ▾

FREQUENCY IN CENTURY   HELP
+1900s -1500s

**DISPLAY**

SORT BY CENTURY   HELP
1900s ▾

FREQ DISPLAY   HELP
RAW # ▾

# HITS   HELP
100

SUBMIT   RESET

Click on the number in a column to see [KEYWORD IN CONTEXT] display          [HELP]

| | CENTURY | 1000s | 1100s | 1200s | 1300s | 1400s | 1500s | 1600s | 1700s | 1800s | 1900s |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | HARD DISK | | | | | | | | | | 29 |
| 2 | HARD TIMES | | | | | | | 2 | 1 | 11 | 18 |
| 3 | HARD CORE | | | | | | | | | 1 | 17 |
| 4 | HARD WAY | | | | | | | | | | 16 |
| 5 | HARD DRIVE | | | | | | | | | | 13 |
| 6 | HARD TIME | | | | | | | | | 2 | 13 |
| 7 | HARD COPY | | | | | | | | | | 11 |
| 8 | HARD STUFF | | | | | | | | 2 | 2 | 11 |
| 9 | HARD KNOCKS | | | | | | | | | 3 | 10 |
| 10 | HARD UP | | | | | | | 1 | | 19 | 10 |

**[KEYWORD IN CONTEXT]**          More information...

1   1917 This Division had already had a very **hard time** .. From the moment of de-busing its life was made very difficult for it.

2   1936 The fellows who have .. small double-dip stands .. are having a **hard time**.

3   1956 Old-time sterno drinkers .. made the flophouse forenoon murky with their hard-time breath.

4   1961 Fresh from the land of sis-boom-bah .. , the Americans had a **hard time** at first learning to applaud good play by either team.

5   1967 It has taken twenty years to get rid of the amateurs, and the professionals .. are having a **hard time** breaking down the sneery reputation gained.

6   1970 Ann-Margret is giving him a **hard time** on the home front, too, tossing out little zingers about his advancing age like Flab is reality.

7   1970 Smith told Petty Officer David Lewis, We are going to have a sit-in and give the Jimmy a **hard time**.

8   1974 Even people who've been out awhile with hepatitis or mononucleosis have a **hard time** making their way back. They lack pep.

9   1978 I was having a **hard time** just breathing. But they've got me padded up pretty good right now and I think I'll be all right.

10   1990 Critics have always had a **hard time** pegging Zush -- some regard him as an outsider artist .. Others prefer to link him to .. art-historical antecedents.

11   1991 I've wondered whether I could ever set this quad before you in some fairness to it and its oddlings, things I've had a

Done                                                                              Internet

# Semantics (collocates): * **meat** (+1900s -1500s)

Corpus of Historical English [DAVIES] - Microsoft Internet Explorer

File   Edit   View   Favorites   Tools   Help

Back | Search | Favorites | Media

Address http://view.byu.edu/che/x.asp

Based on 37 million words of text from Old English to Present-Day English

# CORPUS OF HISTORICAL ENGLISH
## Created by Mark Davies / BYU

**SEARCH**

WORD(S)/PHRASE(S)   HELP
`* meat`
[CUSTOMIZED LISTS] [+OE]

- ⦿ FREQUENCY (TABLE)   HELP
- ○ FREQUENCY (CHART)
- ○ RELEVANCE
- ○ SURROUNDING WORDS

**LIMITS**

SURROUNDING WORDS   HELP
`0 ▾` LEFT/RIGHT `0 ▾`

FREQUENCY IN CENTURY   HELP
`+1900s -1500s`

**DISPLAY**

SORT BY CENTURY   HELP
`1900s ▾`

FREQ DISPLAY   HELP
`RAW # ▾`

# HITS   HELP
`100`

[SUBMIT]  [RESET]

Click on the number in a column to see [KEYWORD IN CONTEXT] display       [HELP]

| | CENTURY | 1000s | 1100s | 1200s | 1300s | 1400s | 1500s | 1600s | 1700s | 1800s | 1900s |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | MINCED MEAT | | | | | | | 7 | 4 | 5 | 16 |
| 2 | RED MEAT | | | | | | | | | 2 | 11 |
| 3 | FRESH MEAT | | | | | | | 3 | 3 | 12 | 10 |
| 4 | IN MEAT | | | | | | | 1 | 4 | 5 | 9 |
| 5 | COLD MEAT | | | | | | | 5 | 4 | 19 | 8 |
| 6 | LEAN MEAT | | | | | | | | 1 | 6 | 7 |
| 7 | COOKED MEAT | | | | | | | | | 1 | 7 |
| 8 | LUNCHEON MEAT | | | | | | | | | | 6 |
| 9 | AS MEAT | | | | | | | 6 | 1 | 4 | 6 |
| 10 | DEAD MEAT | | | | | | | | | 4 | 6 |

**THREE MINUTE TOUR**                          More information...

The Corpus of Historical English will allow you to see:

1. The frequency of **words** (wane, cleped, or fed^er/fether/feather), **phrases** (thou hast or turn on), or **parts of words** (e.g. heart* (sorted by frequency in the 1500s or the 1800s), *ship (by 1400s or 1700s), *friend* (by 1700s), or the substring gl??m* (by 1800s). (Note: charts vs. lists)
2. Which words or phrases have **entered into or left the language** from one century to another, (e.g. *ment 1400s/1500s to 1600s [entering / leaving] or *light* 1600s to 1900s [entering / leaving])
3. **Which words occur most often with others** (often related to the meaning of the word), e.g.:
   -- words immediately after hard (1500s, 1900s) [sorted by uniqueness: 1500s, 1900s]
   -- words within five words of breaks (1700s, 1900s) or face (1600s, 1900s)
4. **Changes in which words occur most often with others** (sometimes a clue to changes in meaning), e.g.
   -- words after hard [+1500s/-1900s, -1500s/+1900s] or before meat [+1500s/-1900s, -1500s/+1900s]
   -- words near market [+1600s/-1900s, -1600s/+1900s] or match [+1600s/-1900s, -1600s/+1900s]

Note that in each case, the link runs an actual query against the corpus and returns the results (i.e. these are not "canned" searches). Feel free to see how the search form has been filled out, and then modify this to create your own queries.

Internet

# Semantics (wide-range collocates): **market** [5L/5R] (+1900s -1600s)

# More information….

# Corpus do Português: 45m words, 1200s-1900s: www.corpusdoportuguês.org

File   Edit   View   Favorites   Tools   Help

Back   Search   Favorites

Address  http://view.byu.edu/cdp/   Go

Google   G Search   3999 blocked   ABC Check   AutoFill   Options

**SEARCH STRING** (HELP)
WORD/PHRASE  *temp*
(INSERT TAG)  -select-   ----
SEARCH   CUSTOMIZED LISTS

**DISPLAY** (HELP)
- ● TABLE
- ○ CHART
- ○ SURROUNDING (HELP)
  -select-   0   0

**SORT BY** (HELP)
- ● FREQUENCY
- ○ PERCENT

**SECTION**   SORT   (HELP)

| 1 | -- ALL -- |
| 2 | -- ALL -- |
1900s-ALL / 1900s-BRAZ / 1900s-PORT / ORAL

MIN FREQ   MIN FREQ
0   0

**OPTIONS** (HELP)
# HITS   100
DISPLAY   RAW FREQ
GROUP BY   WORDS
SAVE LISTS   NO

SEARCH   RESET

## O CORPUS DO PORTUGUÊS (45 MILLION WORDS, 1200s-1900s)
### Mark Davies (Brigham Young Univ) and Michael J. Ferreira (Georgetown Univ)

SEE CONTEXT: CLICK ON WORD (ALL SECTIONS) OR NUMBER (ONE SECTION)   [MORE HELP...]

|   | WORD | 13 | 14 | 15 | 16 | 17 | 18 | PtAc | PtN | PtFc | PtOr | BrAc | BrN | BrFc | BrOr | TOTAL |
|---|------|----|----|----|----|----|----|------|-----|------|------|------|-----|------|------|-------|
| 1 | TEMPO | 542 | 2525 | 6182 | 4150 | 2908 | 12706 | 1744 | 2615 | 3761 | 1239 | 2492 | 2048 | 3709 | 1381 | 48002 |
| 2 | TEMPOS | 44 | 298 | 497 | 526 | 433 | 1620 | 344 | 590 | 630 | 126 | 347 | 299 | 567 | 103 | 6424 |
| 3 | TEMPLO | 9 | 98 | 295 | 259 | 508 | 612 | 162 | 31 | 87 | 6 | 114 | 31 | 72 | 26 | 2310 |
| 4 | TEMPERATURA |  |  |  | 1 | 1 | 89 | 650 | 44 | 30 | 4 | 660 | 94 | 52 | 14 | 1639 |
| 5 | TEMPESTADE | 9 | 28 | 60 | 106 | 65 | 392 | 51 | 45 | 76 | 6 | 30 | 18 | 66 | 14 | 966 |
| 6 | TEMPORAL | 6 | 65 | 183 | 88 | 59 | 125 | 47 | 34 | 42 | 10 | 122 | 18 | 36 | 3 | 838 |
| 7 | TEMPLOS |  | 14 | 164 | 73 | 89 | 144 | 95 | 12 | 27 | 1 | 119 | 5 | 26 | 7 | 776 |
| 8 | TEMPORADA |  |  |  |  | 1 | 10 | 30 | 139 | 8 | 26 | 24 | 366 | 65 | 24 | 693 |
| 9 | TEMPERAMENTO |  | 7 | 1 | 25 | 37 | 289 | 28 | 8 | 78 | 7 | 23 | 18 | 81 | 16 | 618 |
| 10 | TEMPERATURAS |  |  |  |  |  | 7 | 237 | 25 | 9 | 3 | 227 | 20 | 4 |  | 532 |

**KEYWORDS IN CONTEXT**   More information...

LIMIT BY PART OF SPEECH: NO          LIMIT BY SECTION: 1700S  [SEE ALL]          (HELP)

| 1 | 17:Rosario:Frutas | de todas as almas, que estão debaixo do seu governo, & jurisdição **temporal**, & espiritual: qu |
| 2 | 17:Coutinho:Economica | imperio de mayor valor que o mundo todo; no corpo tem ajurisdiçaõ **temporal**, & na alma a e |
| 3 | 17:NavioConselho | meio das brenhas por oficio e por interêsse, fazendo a sua felicidade **temporal** e eterna, e da |
| 4 | 17:Coutto:Brasil | sobre tantos que os miseráveis tinham padecido e padeciam, se levantou um **temporal** tao fo |
| 5 | 17:Coutto:Brasil | Graça divina, devem as Missoens do Maranhao o seu espiritual, e **temporal** augmento* motiv |
| 6 | 17:Pita:America | a Jurisdição de trez. No cível a Relaçao da Bahia, no **Temporal** ao Governo do Maranháo, e n |
| 7 | 17:Aires:Vaidade | a jurisdiçao de três; no espiritual ao bispado de Pernambuco, no **temporal** ao governo do Mar |
| 8 | 17:Aires:Vaidade | acabou com a fugida, ou com a morte, foi a pena **temporal**, e por consequência pena curta, |
| 9 | 17:Barros:Vieira | perpétua, enquanto eles se perpetuam dentro da sua mesma esfera, mas **temporal**, e exting |
| 10 | 17:Barros:Vieira | as mortes, em que se tinham visto, sobreveio de novo o **temporal** aos 14, e durou em desfei |

Done   Internet

# Frequency tables and chart displays (*cujo* 'whose')

# Frequency tables and chart displays (*difícil.\* de* [vr\*] 'hard to VVl')

**[Davies/Ferreira/NEH] O Corpus do Português (45m words, 1200s-1900s) - Microsoft Internet Explorer**

File  Edit  View  Favorites  Tools  Help

Back  →  Search  Favorites  Media

Address  http://www.corpusdoportugues.org/

Google  Search  1752 blocked  Check  AutoLink

**SEARCH STRING**  HELP
WORD/PHRASE  difícil.\* de [vr\*]
(INSERT TAG)  -SELECT-
SEARCH  CUSTOMIZED LISTS

**DISPLAY**  (HELP)
- TABLE
- CHART
- SURROUNDING WORDS  (HELP)
-SELECT-  5  5

**SORT BY**  (HELP)
- FREQUENCY
- PERCENT

**SECTION**  SORT  (HELP)

1  -- ALL --  2  -- ALL --
1900s-ALL  1900s-ALL
1900s-BR  1900s-BR
1900s-PT  1900s-PT
ORAL  ORAL

MIN FREQ  MIN FREQ
5  5

**OPTIONS**  (HELP)
# HITS  100
DISPLAY  RAW FREQ
GROUP BY  WORDS
SAVE LISTS  NO
SEARCH  RESET

## O CORPUS DO PORTUGUÊS (45 MILLION WORDS, 1200s-1900s)

Mark Davies (Brigham Young U) / Michael J. Ferreira (Georgetown U)

CLICK ON ANY BAR TO SEE THE EXAMPLES FROM THAT CENTURY, DIALECT, OR REGISTER

| DECADE | 1300 | 1400 | 1500 | 1600 | 1700 | 1800 | 1900 | PtAc | PtN | PtFc | PtOr | BrAc | BrN | BrFc | BrOr | PT | BR | AC | NEW | FIC | OR |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| # HITS | 0 | 0 | 5 | 2 | 8 | 95 | 370 | 51 | 95 | 40 | 31 | 25 | 39 | 43 | 46 | 217 | 153 | 76 | 134 | 83 | 77 |
| # (MIL) | 1.4 | 3.1 | 4.9 | 4.0 | 2.6 | 11.9 | 24.0 | 3.4 | 3.7 | 3.6 | 1.3 | 3.4 | 3.8 | 3.6 | 1.3 | 12.0 | 12.0 | 6.8 | 7.4 | 7.2 | 2.6 |
| PER MIL | 0.0 | 0.0 | 1.0 | 0.5 | 3.0 | 8.0 | 15.4 | 14.9 | 26.0 | 11.1 | 23.9 | 7.4 | 10.4 | 11.9 | 36.3 | 18.1 | 12.7 | 11.2 | 18.0 | 11.5 | 30.0 |

**OVERVIEW: INTRODUCTION**  Introduction

NOTE: The corpus and corpus interface are still being finalized. The corpus will be online in late June 2006. Before that time, the corpus will be unavailable at times, and your searches may produce inaccurate results. (**List of fixes and improvements by July 2006**)

This website allows you to quickly and easily search more than 45 million words in more than 50,000 Portuguese texts from the 1200s to the 1900s. The interface allows you to search for **exact words** or **phrases, wildcards, lemmas, part of speech**, or any combinations of these. You can also **search for surrounding words** (collocates) within a ten-word window (e.g. all nouns somewhere near *cadeia*, all adjectives near *mulher*, or all nouns near *girar*).

The corpus also allows you to easily compare the frequency of and distribution of words, phrases, and grammatical constructions across texts, in at least three ways:

javascript:x(12)  Internet

# Frequency tables and chart displays (*difícil.\* de* [vr*] 'hard to VVI') [1900s-Pt]

[Davies/Ferreira/NEH] O Corpus do Português (45m words, 1200s-1900s) - Microsoft Internet Explorer

File  Edit  View  Favorites  Tools  Help

Back  |  Search  Favorites  Media  |  1752 blocked  |  ABC Check  AutoLink

Address  http://www.corpusdoportugues.org/  |  Go

Google  |  Search

**SEARCH STRING**  HELP
WORD/PHRASE  difícil.* de [vr*]
(INSERT TAG)  -SELECT-
SEARCH  CUSTOMIZED LISTS

**DISPLAY**  (HELP)
- ● TABLE
- ○ CHART
- ○ SURROUNDING WORDS  (HELP)
  -SELECT-  5  5

**SORT BY**  (HELP)
- ● FREQUENCY
- ○ PERCENT

**SECTION**  SORT  (HELP)

1  -- ALL --          2  -- ALL --
----------             ----------
1900s-ALL            1900s-ALL
1900s-BR             1900s-BR
**1900s-PT**          1900s-PT
----------             ----------
ORAL                   ORAL

MIN FREQ  5          MIN FREQ  5

**OPTIONS**  (HELP)
# HITS  100
DISPLAY  RAW FREQ
GROUP BY  WORDS
SAVE LISTS  NO
SEARCH  RESET

O CORPUS DO PORTUGUÊS (45 MILLION WORDS, 1200s-1900s)

Mark Davies (Brigham Young U) / Michael J. Ferreira (Georgetown U)

SEE CONTEXT: CLICK ON WORD (ALL SECTIONS) OR NUMBER (ONE SECTION)  [MORE HELP...]

| | WORD | 13 | 14 | 15 | 16 | 17 | 18 | PtAc | PtN | PtFc | PtOr | BrAc | BrN | BrFc | BrOr | SUB | TOTAL |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | DIFÍCIL DE EXPLICAR | | | | | | 5 | | 1 | 5 | 2 | | 2 | 3 | 1 | 8 | 19 |
| 2 | DIFÍCIL DE CONTROLAR | | | | | | | 4 | 2 | | | | | | | 6 | 6 |
| 3 | DIFÍCIL DE ENCONTRAR | | | | | | 1 | | 3 | 2 | 1 | | | | 1 | 6 | 8 |
| 4 | DIFÍCIL DE FAZER | | | | | | | 1 | 2 | 1 | 1 | | 2 | | 2 | 5 | 9 |
| 5 | DIFÍCIL DE ENTENDER | | | | 1 | 1 | | | 2 | 1 | 1 | | 2 | 1 | 1 | 4 | 10 |
| 6 | DIFÍCIL DE ESQUECER | | | | | | | 1 | 3 | | | | | 1 | | 4 | 5 |
| 7 | DIFÍCIL DE DIZER | | | | | | 2 | | 1 | | 3 | | | | | 4 | 6 |
| 8 | DIFÍCIL DE COMPREENDER | | | | | | 2 | | 3 | 1 | | | | 1 | | 4 | 7 |
| 9 | DIFICÍLIMOS DE DESFAZER | | | | | | | | | 4 | | | | | | 4 | 4 |
| 10 | DIFÍCIL DE ATINGIR | | | | | | | | 2 | 1 | | | | 1 | | 3 | 4 |

**OVERVIEW: INTRODUCTION**  Introduction

NOTE: The corpus and corpus interface are still being finalized. The corpus will be online in late June 2006. Before that time, the corpus will be unavailable at times, and your searches may produce inaccurate results. (**List of fixes and improvements by July 2006**)

This website allows you to quickly and easily search more than 45 million words in more than 50,000 Portuguese texts from the 1200s to the 1900s. The interface allows you to search for **exact words** or **phrases, wildcards, lemmas, part of speech**, or any combinations of these. You can also **search for surrounding words** (collocates) within a ten-word window (e.g. all nouns somewhere near *cadeia*, all adjectives near *mulher*, or all nouns near *girar*).

The corpus also allows you to easily compare the frequency of and distribution of words, phrases, and grammatical constructions across texts, in at least three ways:

http://www.corpusdoportugues.org/x3.asp?s=n&w4=difícil&c4=&w5=de&c5=&w6=entender&c6=  |  Internet

Lemmatized: **fazer.*** (fazer = to make, do): +1300s-1600s, -1900s

# Tagged for part of speech: **mulheres [aj*]** (ADJ women) (sorted by 1800s)

[Davies/Ferreira/NEH] O Corpus do Português (45m words, 1200s-1900s) - Microsoft Internet Explorer

File   Edit   View   Favorites   Tools   Help

Back   |   Search   Favorites   |   

Address http://view.byu.edu/cdp/   Go

Google   |   Search   |   3999 blocked   |   Check   |   AutoFill   |   Options

## SEARCH STRING   (HELP)

WORD/PHRASE   mulheres [aj*]

(INSERT TAG)   -select-   |   ---- 

SEARCH   CUSTOMIZED LISTS

### DISPLAY   (HELP)
- ⦿ TABLE
- ○ CHART
- ○ SURROUNDING   (HELP)
  -select-   |   0   |   0

### SORT BY   (HELP)
- ⦿ FREQUENCY
- ○ PERCENT

### SECTION   SORT   (HELP)

1  PORT-FICT / PORT-NEWS / PORT-ACAD / ---------- / **1800s** / 1700s / 1600s

2  -- ALL -- / ---------- / 1900s-ALL / 1900s-BRAZ / 1900s-PORT / ---------- / ORAL

MIN FREQ   5  ☐     MIN FREQ   0  ☐

### OPTIONS   (HELP)

# HITS   100
DISPLAY   RAW FREQ
GROUP BY   WORDS
SAVE LISTS   NO

SEARCH   RESET

## O CORPUS DO PORTUGUÊS (45 MILLION WORDS, 1200s-1900s)

### Mark Davies (Brigham Young Univ) and Michael J. Ferreira (Georgetown Univ)

SEE CONTEXT: CLICK ON WORD (ALL SECTIONS) OR NUMBER (ONE SECTION)   [MORE HELP...]

| | WORD | 13 | 14 | 15 | 16 | 17 | 18 | PtAc | PtN | PtFc | PtOr | BrAc | BrN | BrFc | BrOr | SUB | TOTAL |
|---|------|----|----|----|----|----|----|------|-----|------|------|------|-----|------|------|-----|-------|
| 1 | MULHERES BONITAS | | | | | | 15 | | | 6 | | | 5 | 4 | 1 | 15 | 31 |
| 2 | MULHERES PERDIDAS | | | | | | 12 | | | 1 | | | | 2 | | 12 | 15 |
| 3 | MULHERES HONESTAS | | | | 1 | 1 | 9 | | | 1 | | | | 1 | | 9 | 13 |
| 4 | MULHERES BELAS | | | | | | 7 | | | 5 | | | | | | 7 | 12 |
| 5 | MULHERES CASADAS | | | | 5 | | 6 | 1 | | 4 | 1 | 4 | 3 | 3 | | 6 | 27 |
| 6 | MULHERES LIVRES | | | | | | 4 | | | 1 | | 1 | | 1 | | 4 | 7 |
| 7 | MULHERES FEIAS | | | | | | 4 | | | 1 | | | | | | 4 | 5 |
| 8 | MULHERES AJOELHADAS | | | | | | 4 | | | | | | | | | 4 | 4 |
| 9 | MULHERES DESSE | | | | | | 3 | | | 1 | | | | | | 3 | 4 |
| 10 | MULHERES DESTE | | | | | 1 | 3 | | | | | | | | | 3 | 4 |

### KEYWORDS IN CONTEXT   More information...

LIMIT BY PART OF SPEECH: NO          LIMIT BY SECTION: 1800s  [SEE ALL]          (HELP)

| | | |
|---|---|---|
| 1 | 18:Almeida:Gatos5 | - Vivo lá. Há muitos anos - disse-me ele. - Terra de **mulheres bonitas**, segundo me consta |
| 2 | 18:Dinis:Fidalgos | fantástico, ferido no seu orgulho cerebral em eminência de cretino - adeus **mulheres boni** |
| 3 | 18:Dinis:Pupilas | - Maurício tem essa habilidade de ser visto todos os dias por as **mulheres bonitas** da terra |
| 4 | 18:Caminha:Bom-crioulo | reconheceu que era observada, se é que certo instinto, peculiar das **mulheres bonitas**, lho |
| 5 | 18:Alencar:Lucíola | decerto, queria um bem louco ão pequeno, preferia-o a todas as **mulheres bonitas** do mu |
| 6 | 18:Lopes:Viúva | que devora com uma fome canina, como quando contemplava uma multidão de **mulheres l** |
| 7 | 18:Alencar:Correr | de sons, aquela onda de perfumes, de toilletes, e de **mulheres bonitas** que se alastrava p |
| 8 | 18:Alencar:Correr | faremos, se for necessário, como dizia Afonso Karr a propósito das **mulheres bonitas**; far |
| 9 | 18:Machado:Datas | Verdes, na qual são comendadores do número os namorados que desprezam as **mulheres** |
| 10 | 18:Macedo:Mulheres | de ser admirada; mas desculpou-a dizendo que era um desejo natural às **mulheres bonita** |

Internet

# Collocates: **mulheres // [aj:fs]** (ADJ women) (sorted by 1900s)

[Davies/Ferreira/NEH] O Corpus do Português (45m words, 1200s-1900s) - Microsoft Internet Explorer

File   Edit   View   Favorites   Tools   Help

Back → ⊗ ⟳ ⌂ | Search Favorites Media | ⟳ | 🖶 • 🖨 ⊙ • 📄 🔧

Address 🔗 http://www.corpusdoportugues.org/                                        ⏷ ⟳ Go

Google ⏷ | G Search • ⟳ | 🔁 1752 blocked | ABC Check • AutoLink •

## SEARCH STRING                                                    HELP

WORD/PHRASE   mulher

(INSERT TAG)   -SELECT-

SEARCH     CUSTOMIZED LISTS

### DISPLAY                                                          (HELP)
- ○ TABLE
- ○ CHART
- ◉ SURROUNDING WORDS     (HELP)
  - adj.FSG | 0 | 4

### SORT BY                                                          (HELP)
- ○ FREQUENCY
- ◉ PERCENT

### SECTION     DISPLAY                                              (HELP)

| 1 | -- ALL -- |  | 2 | -- ALL -- |
| --------- | | | --------- |
| 1900s-ALL | | | 1900s-ALL |
| 1900s-BR | | | 1900s-BR |
| 1900s-PT | | | 1900s-PT |
| --------- | | | --------- |
| ORAL | | | ORAL |

MIN FREQ       MIN FREQ
4  ☑           0  ☐

### OPTIONS                                                          (HELP)
# HITS         100
DISPLAY        RAW FREQ
GROUP BY       WORDS
SAVE LISTS     NO

SEARCH     RESET

## O CORPUS DO PORTUGUÊS (45 MILLION WORDS, 1200s-1900s)

Mark Davies (Brigham Young U) / Michael J. Ferreira (Georgetown U)

RESULTS ORDERED BY FREQUENCY (HELP)                        [MORE HELP...]
SEE CONTEXT: CLICK ON WORD (ALL SECTIONS) OR NUMBER (ONE SECTION)

|    | WORD | 13 | 14 | 15 | 16 | 17 | 18 | PtAc | PtN | PtFc | PtOr | BrAc | BrN | BrFc | BrOr | SUB | TOT | % |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 1 | GRÁVIDA |  |  |  |  |  | 3 | 7 | 1 |  | 1 | 4 | 3 | 9 | 2 | 27 | 30 | 10.54 |
| 2 | ADÚLTERA |  |  |  |  |  | 4 |  |  | 4 |  |  |  |  |  | 4 | 8 | 7.69 |
| 3 | NUA |  |  |  |  |  | 2 |  |  | 25 |  | 3 |  | 13 | 2 | 43 | 45 | 5.85 |
| 4 | IDOSA |  |  |  |  |  | 6 |  | 2 | 3 |  |  |  | 1 |  | 6 | 12 | 5.71 |
| 5 | GORDA |  |  |  |  | 1 | 13 |  | 1 | 7 |  |  |  | 17 |  | 25 | 39 | 4.16 |
| 6 | DEITADA |  |  |  |  |  | 1 |  |  | 6 |  |  | 1 | 4 |  | 11 | 12 | 3.61 |
| 7 | MAGRA |  |  |  |  |  | 1 |  | 2 | 2 |  |  |  | 6 | 2 | 12 | 13 | 2.89 |
| 8 | PRENHE |  |  | 4 | 2 | 3 | 1 |  |  | 4 |  |  |  | 1 |  | 5 | 15 | 2.56 |
| 9 | CASADA |  |  |  |  |  | 40 | 2 | 2 | 5 | 1 | 2 | 4 | 11 | 1 | 28 | 68 | 2.55 |
| 10 | BONITA |  |  |  |  | 2 | 72 |  |  | 10 |  | 3 | 4 | 31 | 3 | 51 | 125 | 2.54 |

## KEYWORDS IN CONTEXT                              Introduction

LIMIT BY PART OF SPEECH: NO                      LIMIT BY SECTION: NO

| 1 | 18:Machado:Jucanda | Quando entrou na sala de espera, viu uma **mulher** de pé, **magra**, amarelada, envolvida |
| 2 | 19N:Pt:Expr | durante o ano anterior. # Esta jovem **mulher magra**, delicada e de voz tranquila é cons |
| 3 | 19N:Pt:Público | para a protegerem ver caixa Rosa é uma **mulher magra** e frágil, a quem o cabelo claro |
| 4 | 19:Fic:Pt:Tavares:Insubmissos | . No lugar fronteiro ão dele ia aquela **mulher magra**, nova, mas de idade quase indefiní |
| 5 | 19:Fic:Pt:Amorim:Mascara | . Veio abrir a porta a Rita, **mulher** do Longuinhos, **magra**, vestida de farrapos sujos. |
| 6 | 19:Fic:Br:Carvalho:Bebados | , no caso, da morta, uma **mulher** muito branca e **magra**, com os cabelos escorridos sob |
| 7 | 19:Fic:Br:Holanda:Burro | a alma do distrito. Mafra, a **mulher magra**, levantou-se, postou-se atrás do poeta e con |
| 8 | 19:Fic:Br:Louzeiro:Devotos | olhando o táxi que se distanciava. Uma **mulher** jovem e **magra**, pés descalços, atraves |
| 9 | 19:Fic:Br:Louzeiro:Pixote | Era tudo que sabia. Lembra-se vagamente da **mulher**, **magra** e pálida, a voz sumida. U |
| 10 | 19:Fic:Br:Rego:Pedra | a vida na casa paroquial. Era uma **mulher magra** e alta. De voz seca e autoritária, |
| 11 | 19:Fic:Br:Neto:Turbilhão | , esmagava o bolo de feijão; a **mulher magra**, triste, comia lentamente, com ar |

Internet

# Collocates: **mulheres // [aj:fs]** (ADJ women) (1900s vs. 1800s)



Screenshot of "O Corpus do Português (45 million words, 1200s-1900s)" by Mark Davies (Brigham Young U) / Michael J. Ferreira (Georgetown U).

Search string: WORD/PHRASE: mulher. Display: SURROUNDING WORDS, adj.FSG, 0, 4. Sort by: PERCENT. Section: 1800s / 1900s-ALL.

RESULTS ORDERED BY FREQUENCY (HELP)
SEE CONTEXT: CLICK ON WORD

| | WORD | # REG 1 | # REG 2 | % REG 1 | | WORD | # REG 2 | # REG 1 | % REG 2 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | INFAME | 6 | 0 | 1.00 | 1 | BRASILEIRA | 14 | 0 | 1.00 |
| 2 | INSENSÍVEL | 6 | 0 | 1.00 | 2 | SOZINHA | 14 | 0 | 1.00 |
| 3 | ADORADA | 5 | 0 | 1.00 | 3 | DIFERENTE | 8 | 0 | 1.00 |
| 4 | DESMAIADA | 5 | 0 | 1.00 | 4 | DESEJADA | 8 | 0 | 1.00 |
| 5 | INFERNAL | 5 | 0 | 1.00 | 5 | VERDE | 6 | 0 | 1.00 |
| 6 | TANTA | 5 | 0 | 1.00 | 6 | CLARA | 5 | 0 | 1.00 |
| 7 | SEGUINTE | 4 | 0 | 1.00 | 7 | IMPORTANTE | 5 | 0 | 1.00 |
| 8 | PERVERSA | 4 | 0 | 1.00 | 8 | EXCEPCIONAL | 5 | 0 | 1.00 |

**KEYWORDS IN CONTEXT**   Introduction

LIMIT BY PART OF SPEECH: NO    LIMIT BY SECTION: NO

1  19N:Pt:Beira — há seis anos. Agora, ela era uma **mulher sozinha**, firme, de olhar verde e cabelo apanhado

2  19N:Pt:Jornal — como é que vai ser agora? A **mulher** não fica **sozinha**. Apesar de sermos humildes, estamos

3  19:Fic:Pt:Melo:Autópsia — patrões depois de comer, e mandam a **mulher** trabalhar **sozinha** em cima deles. Anica não

4  19:Fic:Pt:Melo:Autópsia — de libertação. Assim, que podia fazer **mulher sozinha**, com marido ausente, filho pequeno p

5  19:Fic:Pt:Melo:Autópsia — , euê, pois claro. E mesmo **mulher sozinha** tem é o medo à espera na solidão antiga

6  19Ac:Br:Enc — não é raro que o parceiro abandone a **mulher sozinha**, que tem que tomar uma decisão, na

7  19:Fic:Br:Amaral:Amigos — pares cujo perfil era muito semelhante: uma **mulher sozinha** e um homossexual que não to

8  19:Fic:Br:Dantas:Cartilha — condenada ão inferno? É do ditado: **mulher sozinha** é andorinha. Ah, Virgem de Guadalupe!

9  19:Fic:Br:Dantas:Cartilha — porque não tem estofo próprio, fibra de **mulher** que sabe viver **sozinha**? Embora, muitas ve

10  19:Fic:Br:Queirós:Dora — o guarda me segurou o pulso: - **Mulher sozinha** na rua, tarde da noite, a ordem

11  19:Fic:Br:Teixeira:Rua — prosperidade cada vez maior. Mas deixei minha **mulher sozinha** no hotel. Com licença. Nisto

# Collocates: **cadeia // [nn*]** (ADJ string) (FICT vs. ACAD)

# Word comparisons: **[nn*] {agudo/aguçado/afilado}** 'sharp N'

# Synonyms: **[=gritar]**: synonyms of 'to shout'

File   Edit   View   Favorites   Tools   Help

Back | Search | Favorites | Media

Address http://www.corpusdoportugues.org/

Google — | Search | 1752 blocked | Check | AutoLink

**SEARCH STRING** HELP
WORD/PHRASE [=gritar]
(INSERT TAG) -SELECT-
SEARCH   CUSTOMIZED LISTS

**DISPLAY** (HELP)
- TABLE
- CHART
- SURROUNDING WORDS (HELP)
  -SELECT- | 5 | 5

**SORT BY** (HELP)
- FREQUENCY
- PERCENT

**SECTION** SORT (HELP)
1 -- ALL --        2 -- ALL --
----------           ----------
1900s-ALL          1900s-ALL
1900s-BR           1900s-BR
1900s-PT           1900s-PT
----------           ----------
ORAL                 ORAL

MIN FREQ            MIN FREQ
5                       5

**OPTIONS** (HELP)
# HITS            100
DISPLAY          RAW FREQ
GROUP BY       WORDS
SAVE LISTS      NO
SEARCH   RESET

## O CORPUS DO PORTUGUÊS (45 MILLION WORDS, 1200s-1900s)

Mark Davies (Brigham Young U) / Michael J. Ferreira (Georgetown U)

SEE CONTEXT: CLICK ON WORD (ALL SECTIONS) OR NUMBER (ONE SECTION)   [MORE HELP...]

| | WORD | 13 | 14 | 15 | 16 | 17 | 18 | PtAc | PtN | PtFc | PtOr | BrAc | BrN | BrFc | BrOr | TOTAL |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | GRITAR [S] | | | 13 | 38 | 18 | 222 | 2 | 23 | 209 | 23 | | 20 | 163 | 9 | 740 |
| 2 | RECLAMAR [S] | | 1 | 1 | 3 | 3 | 78 | 10 | 33 | 20 | 5 | 26 | 39 | 48 | 17 | 284 |
| 3 | PROTESTAR [S] | | | 2 | 3 | 6 | 84 | 8 | 34 | 30 | 6 | 5 | 37 | 18 | 2 | 235 |
| 4 | BRADAR [S] | | 18 | 48 | 27 | 3 | 73 | | 2 | 22 | | | 1 | 4 | 2 | 200 |
| 5 | QUEIXAR-SE [S] | | 1 | 2 | 7 | 13 | 77 | 1 | 7 | 31 | 1 | | 3 | 10 | | 153 |
| 6 | APREGOAR [S] | 36 | 20 | 34 | 19 | 7 | 17 | | 4 | 7 | | 1 | 3 | 3 | | 151 |
| 7 | BERRAR [S] | | | 9 | 16 | | 46 | | 1 | 34 | 10 | | | 26 | 1 | 143 |
| 8 | LADRAR [S] | 2 | 2 | 9 | 16 | 13 | 25 | 1 | | 43 | | | 1 | 7 | | 119 |
| 9 | RALHAR [S] | | 1 | | 1 | 1 | 81 | | 8 | 3 | | | | 12 | | 107 |
| 10 | CHIAR [S] | | | 4 | 18 | 8 | 27 | | 2 | 16 | 1 | | 1 | 7 | | 84 |

### KEYWORDS IN CONTEXT                                          Introduction

| 31 | 18:Azevedo:Cortiço | de saias. Jerônimo perdeu a paciência e ia **protestar** brutalmente contra semelhante invas... |
| 32 | 18:Azevedo:Cortiço | com vontade de afastar-se, mas sem animo de **protestar**, por acanhamento, tentou reatar |
| 33 | 18:Azevedo:Mulato | vontade até de queixar se ão pai; de **protestar** contra aquelas contrariedades que a faziam |
| 34 | 18:Macedo:Luneta | Há só uma voz que pode e há de **protestar**, é a minha, a voz suspeita, |
| 35 | 18:Pompéia:Ateneu | Esta caranguejola, enorme e pesada, que parecia **protestar**, a cada solavanco, contra o ca... |
| 36 | 18:Taunay:Inocência | algumas pautadas de incitamento, pareceu querer o cargueiro **protestar** contra o tratamen... |
| 37 | 18:Taunay:Inocência | mesmo. Oh! Sr. Garcia! quis **protestar** Pereira. --Nada;.. digo-lhe isto do coração.. |
| 38 | 18:Alencar:Correr | no ator que caíra em falta, nem de **protestar** contra o ato dos diretores por uma semelhan... |
| 39 | 18:Cunha:Artigos | de Lourenço de Médicis. Precisamos, porém, **protestar** também contra a introdução deste p... |
| 40 | 18:Alencar:Garatuja | naquela capitania, que vinham todos unidos em corpo **protestar** contra a violência inaudita |
| 41 | 18:Alencar:Garatuja | e revoltado nele o sentimento do belo, ia **protestar**, quando pareceu-lhe que de novo entre... |
| 42 | 18:Machado:Textos | dessa confissão, vê-se que o poeta queria principalmente **protestar** contra o caminho que l... |

Internet

# Customized lists: **querer.\* [davies@gritar-syn]:** want to "shout"



[Davies/Ferreira/NEH] O Corpus do Português (45m words, 1200s-1900s) - Microsoft Internet Explorer

File   Edit   View   Favorites   Tools   Help

Back ··· · ⊗ ⟳ ⌂ | ⚲Search ⚑Favorites ⚙Media ⚙ | ⎙· ⎙ ⎙ · ⎕ ⎘

Address ⎙ http://www.corpusdoportugues.org/

Google · | G Search · ⚙ | 🗗 1752 blocked | ᴬᴮᶜ Check · AutoLink ·

## SEARCH STRING                        HELP

WORD/PHRASE   querer.* [davies@g

(INSERT TAG)   -SELECT-

SEARCH    CUSTOMIZED LISTS

### DISPLAY                             (HELP)

- ⦿ TABLE
- ○ CHART
- ○ SURROUNDING WORDS   (HELP)
  -SELECT-   5   5

### SORT BY                             (HELP)

- ⦿ FREQUENCY
- ○ PERCENT

### SECTION   SORT                      (HELP)

[1]  -- ALL --     [2]  -- ALL --
     ----------          ----------
     1900s-ALL          1900s-ALL
     1900s-BR           1900s-BR
     1900s-PT           1900s-PT
     ----------          ----------
     ORAL               ORAL

MIN FREQ          MIN FREQ
5   ☐             5   ☐

### OPTIONS                             (HELP)

# HITS        100
DISPLAY       RAW FREQ
GROUP BY      WORDS
SAVE LISTS    NO

SEARCH    RESET

## O CORPUS DO PORTUGUÊS (45 MILLION WORDS, 1200s-1900s)

Mark Davies (Brigham Young U) / Michael J. Ferreira (Georgetown U)

SEE CONTEXT: CLICK ON WORD (ALL SECTIONS) OR NUMBER (ONE SECTION)   [MORE HELP...]

| | WORD | 13 | 14 | 15 | 16 | 17 | 18 | PtAc | PtN | PtFc | PtOr | BrAc | BrN | BrFc | BrOr | TOTAL |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | QUIS GRITAR | | | | | | 8 | | | 1 | | | | 5 | | 14 |
| 2 | QUIS PROTESTAR | | | | | | 3 | | 1 | 2 | | | | 1 | | 7 |
| 3 | QUER GRITAR | | | | | | 4 | | | | | | | | | 4 |
| 4 | QUERIA PROTESTAR | | | | | | 1 | | | 1 | | | | | | 2 |
| 5 | QUERIA GRITAR | | | | | | 1 | | | 1 | | | | | | 2 |
| 6 | QUEREMOS RECLAMAR | | | | | | | | 1 | | | | | | 1 | 2 |
| 7 | QUERO GRITAR | | | | | | | | | 1 | | | | 1 | | 2 |
| 8 | QUERO CLAMAR | | | | | | | | | | | | 1 | | | 1 |
| 9 | QUIZ GRITAR | | | | | | | | | 1 | | | | | | 1 |
| 10 | QUISESSE GRITAR | | | | | | | | | | | | | 1 | | 1 |

### KEYWORDS IN CONTEXT                  Introduction

| 31 | 18:Azevedo:Cortiço | de saias. Jerônimo perdeu a paciência e ia **protestar** brutalmente contra semelhante invasã |
| 32 | 18:Azevedo:Cortiço | com vontade de afastar-se, mas sem animo de **protestar**, por acanhamento, tentou reatar |
| 33 | 18:Azevedo:Mulato | vontade até de queixar se ão pai; de **protestar** contra aquelas contrariedades que a faziam |
| 34 | 18:Macedo:Luneta | Há só uma voz que pode e há de **protestar**, é a minha, a voz suspeita, |
| 35 | 18:Pompéia:Ateneu | Esta caranguejola, enorme e pesada, que parecia **protestar**, a cada solavanco, contra o ca |
| 36 | 18:Taunay:Inocência | algumas pautadas de incitamento, pareceu querer o cargueiro **protestar** contra o tratament |
| 37 | 18:Taunay:Inocência | mesmo. Oh! Sr. Garcia! quis **protestar** Pereira. --Nada;.. digo-lhe isto do coração.. |
| 38 | 18:Alencar:Correr | no ator que caíra em falta, nem de **protestar** contra o ato dos diretores por uma semelhant |
| 39 | 18:Cunha:Artigos | de Lourenço de Médicis. Precisamos, porém, **protestar** também contra a introdução deste p |
| 40 | 18:Alencar:Garatuja | naquela capitania, que vinham todos unidos em corpo **protestar** contra a violência inaudita |
| 41 | 18:Alencar:Garatuja | e revoltado nele o sentimento do belo, ia **protestar**, quando pareceu-lhe que de novo entre |
| 42 | 18:Machado:Textos | dessa confissão, vê-se que o poeta queria principalmente **protestar** contra o caminho que l |

# Customized lists: **estou + emotions** (I am + emotions): 1900s

# Customized lists: User-created lists (can correct our "errors"): **sair 'to leave'**

# Advantages of relational databases

- Size, speed, and annotation
- Unlimited, modular tables (synonyms, customized lists, etc)
- Helps in annotating the corpus

# POS tagging

- Bootstrap from existing tagged corpus
- Probabilistic: *citie*, then *emptie* (update tens of millions of tokens at a time)
- Rule-based: to emptie, the emptie CONJ
- Even words that are not in the existing corpus: DET XX PREP/CONJ

# Lemmatization

➢ English: not much of a problem (*work, flye, returnes*, but *saw, ground, wound, etc*)

➢ Much worse for other languages (Sp trabajo 'work/trabajo, I-work/trabajar)

➢ Sp hubiese = 56/72 possible forms

# Lemmatization (historical)

➢ Apply modern lemmas (finish)

➢ Apply spelling rules:  *tie: emptie / empty

➢ Use frequency information from database to find overly frequent lemmas (*goo / go, had lane / had lain*)

➢ Use frequency information to target highly frequent unknown forms first
(students / see context / provide modern form / words acquire modern features)

➢ *a pease of*:  compare to PDE:
        *a piece / peas / peace of*

➢ Can also use contextual words (*pease*)

# Interface for annotating Old Portuguese (blocks of text)

# Interface for annotating Old Portuguese (by word, POS, lemma)

# Interface for annotating Old Portuguese (by lemma frequency I)

**BY TEXT / PASSAGE**
12_coimbra_text-msc12 ▾
OFFSET 200
SUBMIT    RESET

**BY WORD**
WORD [        ]  LEM CONTAR    POS [      ]
CORPUS OLD ▾ # 200
SUBMIT    RESET

**HAS  LEMMA**
FREQ 50 -10000
OPT > 2    x MPT
GO

**NO LEMMA**
500    POS [    ]
GO

KWIC                                        CHANGE

| 1  | 2345 | dano    | nn:ms | >> |
| 2  | 997  | danos   | nn:mp | >> |
| 3  | 251  | dãno    | nn:ms | >> |
| 4  | 228  | dapno   | nn:ms | >> |
| 5  | 224  | damno   | nn:ms | >> |
| 6  | 217  | dampno  | nn:ms | >> |
| 7  | 106  | danno   | nn:ms | >> |
| 8  | 100  | dapnos  | nn:mp | >> |
| 9  | 87   | damnos  | nn:mp | >> |
| 10 | 62   | dãpno   | nn:ms | >> |
| 11 | 45   | dannos  | nn:mp | >> |

| CONTEXT | [        ] | WORD ▾ | // | [        ] | WORD ▾ |
| WORD    |            |        |    | dãno       |        |
| LEMMA   |            |        |    | dano       |        |
| POS     |            |        |    | nn:ms      |        |

Submit    Reset

MORE…

| 1 | 38153155 | dano | nn:ms | a elle : & assi vendolhe eu este **dãno** me contento, & quero que me arrinquem |
| 2 | 38157207 | dano | nn:ms | que vem de dar esmolla, & o **dãno** que soccede aos ingratos. Trata de hum |
| 3 | 38169988 | dano | nn:ms | que o posera, & em emmenda do **dãno** que lhe fez em sua casa, dando |
| 4 | 38214492 | dano | nn:ms | permittio Deos que pagassë o grãde mal & **dãno** que fizerã à Rainha sua yrmã, & |
| 5 | 39344895 | dano | nn:ms | ate o passarem, para se euitar o **dãno** que com isso se faz, Há o |
| 6 | 39348315 | dano | nn:ms | sem pagarem coimas : comtal que não fação **dãno** nas nouidades : porque fazendoo pagarão os dãnos |
| 7 | 39348350 | dano | nn:ms | trazer os ditos gados. E quando fizerem **dãno** nas nouidades, farão S 2 penhora Quarta |
| 8 | 39348369 | dano | nn:ms | em tanto gado que baste para pagar o **dãno** & coimas. E não prenderão os pastores |
| 9 | 39361752 | dano | nn:ms | dinheiro : de maneira que não soomente recebem **dãno** no preço em que as comprão fiadas, |

# Composition of proposed 200 million word historical corpus

| | SPOKEN/INFORMAL | | FICTION | | NEWS | | ACADEMIC | | OED | TOTAL |
|---|---|---|---|---|---|---|---|---|---|---|
| | SIZE | SOURCES | SIZE | SOURCES | SIZE | SOURCES | SIZE | SOURCES | | |
| **1500s** | 5 | EEBO-Dr CED | 10 | EEBO-Fic | 1 | ?? | 13 | EEBO-NF | 1 | 30 |
| **1600s** | 5 | EEBO-Dr BAILEY CED | 10 | EEBO-Fic | 1 | NEWSBOOKS (NEWDIGATE) | 13 | EEBO-NF | 1 | 30 |
| **1700s** | 5 | LION-Dr BAILEY CED LATE 18C | 10 | LION-Fic | 1 | LAMPETER ZEN | 13 | LION-NF | 1 | 30 |
| **1800s** | 9 | LION-Dr PG-Dr BAILEY | 16 | LION-Fic PG-Fic | 7 | MOA-Jrnl SCAN-News (LON TIMES) (NY TIMES) | 16 | LION-NF PG-NF | 2 | 50 |
| British | 4 | | 7 | | 3 | | 8 | | | 23 |
| American | 4 | | 7 | | 3 | | 8 | | | 23 |
| Misc | 1 | | 2 | | 1 | | 0 | | | 4 |
| **1900s** | 9 | MOVIES RADIO ORALHIST BNC, LLC, SEC, (LDC) PG-Dr SCAN-Dr | 21 | PG-Fic SCAN-Fic BNC-Fic | 14 | SCAN-Per (LON TIMES) (NY TIMES) BNC-News NDNP | 14 | PG-NF SCAN-NF BNC-NF | 2 | 60 |
| British | 4 | | 9 | | 6 | | 7 | | | 27 |
| American | 4 | | 9 | | 6 | | 7 | | | 27 |
| Misc | 1 | | 3 | | 2 | | 0 | | | 6 |
| **TOTAL** | 33 | | 67 | | 24 | | 69 | | 7 | 200 |
| **PERCENT** | 16.5 | | 33.5 | | 12.0 | | 34.5 | | 3.5 | 100% |

# Conclusion

- Relational databases: size, speed, and annotation

- Very wide range of searches

- Information from database can be used as an integral part of the tagging process

- Already working in several online corpora, more in planning stages